











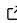


PatientProfiles: An R package to identify patient characteristics in data mapped to the OMOP common data model

Martí Català¹, Mike Du¹, Yuchen Guo¹, Kim Lopez-Guell¹, Núria Mercadé-Besora¹, Xihang Chen¹, Marta Alcalde-Herraz¹, Xintong Li¹, Daniel Prieto-Alhambra^{1,2}, and Edward Burn¹

¹ Health Data Sciences Group, Nuffield Department of Orthopaedics, Rheumatology and Musculoskeletal Sciences, University of Oxford, United Kingdom ² Department of Medical Informatics, Erasmus University Medical Center, Rotterdam, The Netherlands ¶ Corresponding author

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: [Nick Golding](#)  

Reviewers:

- [@javier-gracia-tabuenca-tuni](#)
- [@pacoramon](#)

Submitted: 09 July 2025

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](https://creativecommons.org/licenses/by/4.0/)).

Summary

Real-world data (RWD) mapped to the Observational Medical Outcomes Partnership Common Data Model (OMOP CDM) offers a standardised framework for conducting observational health research across diverse data sources. However, identifying and summarising patient-level characteristics within this model often requires custom code, limiting efficiency and reproducibility. To address this, we developed the open-source PatientProfiles R package. This package streamlines the process of extracting demographic characteristics, computing intersections between cohorts and clinical events, and generating standard summaries of patient populations in OMOP CDM datasets.

Built on the tidyverse and omopgenerics infrastructure, PatientProfiles supports SQL translation for scalable database operations and includes comprehensive test coverage across multiple database systems. It provides a suite of functions grouped into demographics, intersections, summaries, utility, and mock data generation. The package is designed for transparency, modularity, and reusability in epidemiological workflows and is available via CRAN and GitHub, along with documentation and vignettes to support users.

Statement of need

Real-world data (RWD), routinely collected health data such as GP records, hospital data, and insurance claims data are valuable resources for conducting epidemiological research studies. However, with such data typically not collected primarily for research, different RWD sources can vary substantially in format and clinical coding systems. To overcome this difficulty a common data model (CDM) is often used. A CDM helps standardising data structures across various sources, enhancing data consistency, quality, and interoperability. A particularly popular data model is the Observational Medical Outcomes Partnership (OMOP) CDM, with more than 800 million patients' health care data transformed into this format ([Overhage et al., 2011](#)).

The OMOP CDM is a person-centric relational data model. Patients' data is spread across various tables related to different clinical domains with, for example, the *condition occurrence* table containing diagnoses while the *drug exposure* table contains drug prescriptions. These different clinical tables are all linked back to the *person* table which contains a unique identifier for each individual along with some key demographic data such as their date of birth.

40 Meanwhile, records in the *observation period* table define the period of calendar time over
41 which an individual is followed-up.(Blacketer, 2025)

42 One of the principal benefits of mapping data to a CDM is that it allows for the same analytic
43 code to be run across different datasets. Developing well-tested and easy to use software for
44 common analytic tasks can therefore bring significant benefits, both improving the speed in
45 which analyses can be performed and improved quality by reducing the amount of study-specific
46 bespoke code needing to be written.

47 Obtaining the characteristics of individuals is one of the most common first tasks when working
48 with patient-level data. In almost all analyses specific characteristics of individuals will need to
49 be identified, after which groups of individuals who share some specific common condition or
50 characteristic need to be identified and relationships between these groups are described (for
51 example the time between a given diagnosis and a health outcome of interest).

52 We created the PatientProfiles R package to support identifying patient characteristics in data
53 mapped to the OMOP CDM. It provides functionality to obtain demographic information (such
54 as age, sex, prior observation time, future observation time, and so on), describe intersections
55 between different groups of patients, and summarise the results in a standard output format.

56 Design principles

57 PatientProfiles was designed to adhere to the tidyverse tidy design principles. The tidyverse
58 is a collection of R packages designed for data science, offering a cohesive and consistent
59 syntax for data manipulation, and analysis (Wickham et al., 2019). The dplyr package defines
60 multiple methods that can be implemented to many different sources of data. Of particular
61 relevance to working with OMOP CDM data which is typically stored in a database, the dbplyr
62 package provides translations of dplyr methods to SQL.

63 The core dependency of PatientProfiles is the omopgenerics package (Català & Burn, 2024),
64 which provides methods, classes and basic operations for packages working with data in the
65 OMOP CDM format. It defines a central object, a *cdm_reference*, that provides a central
66 reference to all the different OMOP CDM tables, along with various other S3 classes and
67 methods that facilitate working with the data contained in this reference.

68 Development of the PatientProfiles R package

69 PatientProfiles was developed in accordance with best practices for R packages with the
70 devtools and usethis R packages used for common development tasks. The core, general
71 dependencies of the package include dplyr and tidyr for common data manipulations and dbplyr
72 which provides translations to SQL. In addition the core dependency related to OMOP CDM
73 data is the omopgenerics package which provides core classes and methods specific to this
74 data format.

75 The PatientProfiles package includes functionality to create its own mock data in the OMOP
76 CDM format. This mock data is used to test the package using the testthat framework
77 (Wickham, 2011). Every line of the packages is tested multiple times trying to account for
78 various edge cases. Currently, the package is tested iteratively against different database
79 management systems: PostgreSQL, SQL Server, Amazon Redshift, and DuckDB. In addition
80 to unit tests, end-to-end integration tests of the package have been conducted to ensure the
81 face validity of results.

82 The package is open-source and released via CRAN: [https://CRAN.R-project.org/package=](https://CRAN.R-project.org/package=PatientProfiles)
83 [PatientProfiles](https://CRAN.R-project.org/package=PatientProfiles) (Català et al., 2025) (version 1.4.1 as of 9th July 2025) and also available
84 on github: <https://github.com/darwin-eu/PatientProfiles> with its own website with more
85 documentation and vignettes that cover the content of the package more in depth.

Overview of the PatientProfiles R package

PatientProfiles contains three main groups of functions (Figure 1). **Demographics** functions are used to add information contained in person and observation period tables to other tables or objects of interest. **Intersections** are used to intersect a table with an object of interest (it can be another table, a cohort of patients or a paritciar clinical concept). The **summarise** functions are used to create standard objects that summarise the content of a table of interest. Finally, the package also contains some complementary utility functions to for example create mock data.

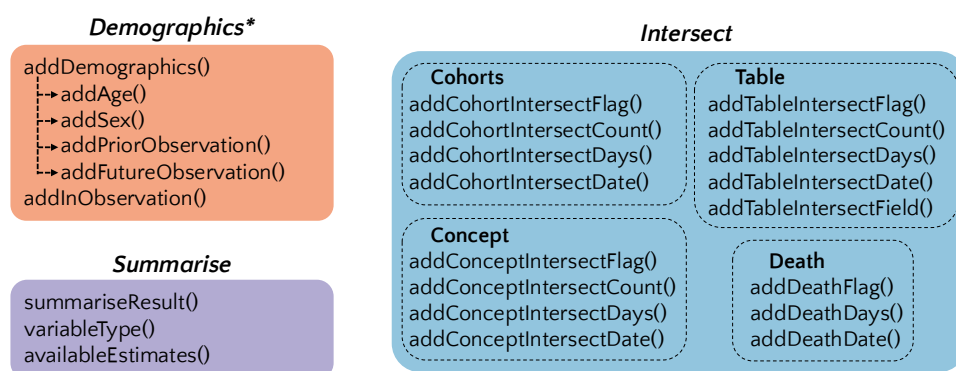


Figure 1: PatientProfiles functions bloks. Note that each demographic function has its own analogous *query* function to only add a query to the data, e.g. `addAge()` -> `addAgeQuery()`.

Mock data

A reference to an OMOP CDM instance is needed to use PatientProfiles. In this simple tutorial we will use mock toy data produced by the same package. By default this toy data is copied into an in-process duckdb database.

```

library(PatientProfiles)
cdm <- mockPatientProfiles(numberIndividuals = 1000)
# to customise cohorts
cdm$my_flu_cohort <- cdm$cohort1 |>
  dplyr::filter(cohort_definition_id == 1L) |>
  omopgenerics::newCohortTable(cohortSetRef = dplyr::tibble(
    cohort_definition_id = 1L, cohort_name = "flu"
  ))
cdm$target <- cdm$cohort2 |>
  omopgenerics::newCohortTable(cohortSetRef = dplyr::tibble(
    cohort_definition_id = c(1L, 2L, 3L),
    cohort_name = c("covid_test", "flu_test", "asthma")
  ))
  
```

Demographics

`addDemographics()` is used to characterise the demographics of a table. The table is needed to be part of a `cdm_reference` object and to contain a person identifier column (either `person_id` or `subject_id`). There are multiple columns that can be added with this function:

- *age*: the age at a certain *indexDate*. You can also add an *age group* column grouping individuals for different age group ranges.

```

117   ■ sex: the sex of the individual.
118   ■ prior observation: the number of days between start of observation and indexDate.
119   ■ future observation: the number of days between indexDate and end of observation.
120   ■ date of birth: the birth date of the individual.

121 An example to add the demographics to a mock cohort table is:

122 cdm$my_flu_cohort |>
123   addDemographics(
124     indexDate = "cohort_start_date",
125     ageGroup = list("children" = c(0, 17), "adult" = c(18, Inf))
126   ) |>
127   dplyr::glimpse()

128 ## Rows: ??
129 ## Columns: 10
130 ## Database: DuckDB v1.0.0 [root@Darwin 23.4.0:R 4.4.1/:memory:]
131 ## $ cohort_definition_id <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1...
132 ## $ subject_id           <int> 117, 818, 634, 729, 886, 245, 761, 385, 597, 53, ...
133 ## $ cohort_start_date    <date> 1959-06-03, 1936-09-03, 1979-08-13, 2010-03-
134 14, ...
135 ## $ cohort_end_date      <date> 1960-11-29, 1981-03-23, 2053-12-02, 2044-04-
136 28, ...
137 ## $ age                  <int> 27, 2, 16, 36, 58, 20, 38, 82, 94, 77, 18, 5, 5, ...
138 ## $ age_group            <chr> "adult", "children", "children", "adult", "adult"...
139 ## $ sex                  <chr> "Female", "Female", "Female", "Male", "Male", "Ma...
140 ## $ prior_observation     <int> 10015, 976, 6068, 13221, 21371, 7397, 13961, 3005...
141 ## $ future_observation    <int> 768, 41462, 32133, 40592, 20967, 10212, 20892, 19...
142 ## $ date_of_birth         <date> 1932-01-01, 1934-01-01, 1963-01-01, 1974-01-
143 01, ...

144 For each one of the functionalities there exist individual functions: addAge(), addSex(),
145 addPriorObservation(), addFutureObservation() and addDateOfBirth().

```

146 Observation period id

147 The *observation_period* contains the period of time that an individual in the database is in
 148 observation. There might be multiple individual periods per person, but they can not overlap
 149 each other. When doing analysis it can be of interest knowing if a certain date is in observation,
 150 whether the individual will be in observation after a certain time, and from which observation
 151 period is an observation. To do so we have two functions:

- 152 ■ addInObservation() to identify if an individual is in observation in a certain *window*
 153 respect an *indexDate*.
- 154 ■ addObservationPeriodId() to identify in which observation period ordinal is that date
 155 from.

```

cdm$gibleed |>
  addInObservation(
    indexDate = "cohort_start_date",
    window = list("obs_index_date" = c(0, 0), "in_1_year" = c(365, 365)),
    nameStyle = "{window_name}"
  ) |>
  addObservationPeriodId() |>
  dplyr::glimpse()

156 ## Rows: ??
157 ## Columns: 7

```

```

158 ## Database: DuckDB v1.0.0 [root@Darwin 23.4.0:R 4.4.1/:memory:]
159 ## $ cohort_definition_id <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ...
160 ## $ subject_id <int> 962, 1158, 4462, 351, 3556, 320, 1965, 2105, 259...
161 ## $ cohort_start_date <date> 1995-07-09, 2016-12-27, 1990-10-23, 2018-06-
162 28,...
163 ## $ cohort_end_date <date> 2019-06-14, 2017-02-15, 2018-04-27, 2018-06-
164 29,...
165 ## $ obs_index_date <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ...
166 ## $ in_1_year <int> 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, ...
167 ## $ observation_period_id <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ...

```

168 Query functions

169 Usually OMOP CDM instances are stored in SQL databases. The functions that we have
 170 seen take the original table add the new columns and save the result into a new table (name
 171 argument). Each function has its own homologous one (same name terminated with 'Query')
 172 that instead of saving the result to a table only returns a query to generate the table. For
 173 local instances both functions provide exactly the same result.

174 Intersections

175 PatientProfiles has 15 functions that are used to compute intersections between tables.
 176 Common functions parameters:

- 177 ■ **indexDate** Name of the column that contains the date that will be the origin time of
 178 our calculations.
- 179 ■ **censorDate** Name of the column that contains the date to censor the observation
 180 window.
- 181 ■ **window** Window of time respect to the index date that we will consider relevant events
 182 on.

183 There exist 4 different function types:

- 184 ■ **Flag**: It creates a new integer column that can have 3 possible values: 1 whether an
 185 event of interest is observed; 0 if the event is not observed; NA if the individual is not in
 186 observation within that window.
- 187 ■ **Count**: It creates a new integer column with the number of observed events, NA is
 188 reported if the individual is not in observation in that window.
- 189 ■ **Date**: It creates a new date column that contains the date of a certain event, NA is
 190 reported if the event is not observed or the individual is not in observation in that
 191 window.
- 192 ■ **Days**: It creates a new integer date with the time difference with a certain event, NA
 193 is reported if the event is not observed or the individual is not in observation in that
 194 window.

195 For the *Flag* and *Count* functions there are 2 extra parameters: - **targetStartDate** Nome of
 196 the column that identifies the start of the event. - **targetEndDate** Name of the end of the
 197 episode, if NULL the event is considered to start and end on the targetStartDate.

198 With the following code you can add the number of visits recorded in the prior year
 199 (**number_visits**) and a flag to see if there is a record of a asthma test any time prior to the
 200 index date.

```

201 cdm$cohort1 |>
202   addTableIntersectCount(
203     tableName = "visit_occurrence",
204     window = c(-365, 0),
205     nameStyle = "number_visits"

```

```

206   ) |>
207   addCohortIntersectFlag(
208     cohortTableName = "cohort2",
209     cohortId = 3,
210     window = c(-Inf, 0),
211     nameStyle = "prior_asthma"
212   ) |>
213   dplyr::glimpse()

214   ## Rows: ??
215   ## Columns: 6
216   ## Database: DuckDB v1.0.0 [root@Darwin 23.4.0:R 4.4.1/:memory:]
217   ## $ cohort_definition_id <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1...
218   ## $ subject_id          <int> 117, 818, 634, 886, 245, 761, 597, 53, 124, 285, ...
219   ## $ cohort_start_date   <date> 1959-06-03, 1936-09-03, 1979-08-13, 1996-07-
220   06, ...
221   ## $ cohort_end_date     <date> 1960-11-29, 1981-03-23, 2053-12-02, 2048-10-
222   18, ...
223   ## $ number_visits       <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...
224   ## $ prior_asthma        <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0...

```

For the *Date* and *Days* functions there are 2 extra parameters: - *targetDate* Name of the column that contains the event of interest. - *order* Whether we are interested with the “first” or “last” event in the window.

With the following code you would add which is the date of the next test (of flu or covid) after the index date:

```

228   cdm$cohort1 |>
229     addCohortIntersectDate(
230       targetCohortTable = "cohort2",
231       targetCohortId = c(1, 2),
232       window = c(1, Inf)
233     ) |>
234     dplyr::glimpse()

230   ## Rows: ??
231   ## Columns: 6
232   ## Database: DuckDB v1.0.0 [root@Darwin 23.4.0:R 4.4.1/:memory:]
233   ## $ cohort_definition_id <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1...
234   ## $ subject_id          <int> 761, 597, 124, 285, 157, 348, 297, 919, 830, 741,...
235   ## $ cohort_start_date   <date> 2037-03-23, 2093-08-24, 1981-12-13, 1993-01-
236   28, ...
237   ## $ cohort_end_date     <date> 2058-09-28, 2141-02-14, 2020-03-08, 2006-09-
238   29, ...
239   ## $ covid_test_1_to_inf <date> 2077-08-29, NA, NA, NA, 2043-06-27, NA, NA, 2045...
240   ## $ flu_test_1_to_inf   <date> NA, 2194-01-24, 2058-08-07, 2093-09-06, NA, 2041...

```

NOTE that each function has some arguments related to the intersecting target (cohort, concept or clinical table).

Summarise data

`summariseResult()` is a function that allow the user to summarise multiple columns into multiple estimates (see `availableEstimates()`) into a standard format output, see the below example:

```

247   cdm$cohort1 |>

```



```

248 addCohortName() |>
249 summariseResult(
250   group = "cohort_name",
251   strata = list("sex", c("sex", "prior_asthma")),
252   variables = list(c("number_visits", "age"), c("covid_test_1_to_inf", "flu_test_1_to_inf")),
253   estimates = list(c("median", "q25", "q75"), c("min", "max"))
254 ) |>
255 dplyr::glimpse()

256 ## Rows: 84
257 ## Columns: 13
258 ## $ result_id      <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,...
259 ## $ cdm_name        <chr> "PP MOCK", "PP MOCK", "PP MOCK", "PP MOCK", "PP MOCK"...
260 ## $ group_name      <chr> "cohort_name", "cohort_name", "cohort_name", "cohort_name"...
261 ## $ group_level     <chr> "flu", "flu", "flu", "flu", "flu", "flu", "flu", "flu"...
262 ## $ strata_name     <chr> "overall", "overall", "overall", "overall", "overall"...
263 ## $ strata_level    <chr> "overall", "overall", "overall", "overall", "overall"...
264 ## $ variable_name   <chr> "number records", "number subjects", "age", "age", "a..."
265 ## $ variable_level  <chr> NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA, NA...
266 ## $ estimate_name   <chr> "count", "count", "median", "q25", "q75", "median", "...
267 ## $ estimate_type    <chr> "integer", "integer", "integer", "integer", "integer"...
268 ## $ estimate_value  <chr> "350", "350", "49", "23", "87", "0", "0", "0", "1945-
269 ...
270 ## $ additional_name <chr> "overall", "overall", "overall", "overall", "overall"...
271 ## $ additional_level <chr> "overall", "overall", "overall", "overall", "overall"...

```

Conclusions

The PatientProfiles R package provides functionality to assist researchers working with data mapped to the OMOP CDM format. By basing the package around this data model which has a known structure the package could be developed with simple interfaces yet deep functionality. The package has already been used in published studies (Català et al., 2024; Mercadé-Besora et al., 2024) and is freely available to be used in future research.

Funding information

Development of the PatientProfiles R package was funded by the European Medicines Agency as part of the Data Analysis and Real World Interrogation Network (DARWIN EU®). This manuscript represents the views of the DARWIN EU® Coordination Centre only and cannot be interpreted as reflecting the views of the European Medicines Agency or the European Medicines Regulatory Network.

References

- Blacketer, C. (2025). *Definition and DDLs for the OMOP common data model (CDM)*.
- Català, M., & Burn, E. (2024). *Omopgenerics: Methods and classes for the OMOP common data model*. <https://darwin-eu.github.io/omopgenerics/>, <https://github.com/darwin-eu/omopgenerics>
- Català, M., Guo, Y., Du, M., Lopez-Guell, K., Burn, E., & Mercade-Besora, N. (2025). *PatientProfiles: Identify characteristics of patients in the OMOP common data model*. <https://darwin-eu.github.io/PatientProfiles/>, <https://github.com/darwin-eu/PatientProfiles>

- 293 Català, M., Mercadé-Besora, N., Kolde, R., Trinh, N. T. H., Roel, E., Burn, E., Rathod-
294 Mistry, T., Kostka, K., Man, W. Y., Delmestri, A., Nordeng, H. M. E., Uusküla, A.,
295 Duarte-Salles, T., Prieto-Alhambra, D., & Jödicke, A. M. (2024). The effectiveness
296 of COVID-19 vaccines to prevent long COVID symptoms: Staggered cohort study of
297 data from the UK, Spain, and Estonia. *The Lancet Respiratory Medicine*, 12, 226–240.
298 [https://doi.org/10.1016/S2213-2600\(23\)00414-9](https://doi.org/10.1016/S2213-2600(23)00414-9)
- 299 Mercadé-Besora, N., Li, X., Kolde, R., Trinh, N. T., Sanchez-Santos, M. T., Man, W. Y., Roel,
300 E., Reyes, C., Delmestri, A., Nordeng, H. M. E., Uusküla, A., Duarte-Salles, T., Prats,
301 C., Prieto-Alhambra, D., Jödicke, A. M., & Català, M. (2024). The role of COVID-19
302 vaccines in preventing post-COVID-19 thromboembolic and cardiovascular complications.
303 *Heart*, 110(9), 635–643. <https://doi.org/10.1136/heartjnl-2023-323483>
- 304 Overhage, J. M., Ryan, P. B., Reich, C. G., Hartzema, A. G., & Stang, P. E. (2011).
305 Validation of a common data model for active safety surveillance research. *Journal of*
306 *the American Medical Informatics Association*, 19(1), 54–60. [https://doi.org/10.1136/](https://doi.org/10.1136/amiajnl-2011-000376)
307 [amiajnl-2011-000376](https://doi.org/10.1136/amiajnl-2011-000376)
- 308 Wickham, H. (2011). Testthat: Get started with testing. *The R Journal*, 3, 5–10. https://journal.r-project.org/archive/2011-1/RJournal_2011-1_Wickham.pdf
- 310 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., Golemund,
311 G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T. L., Miller, E., Bache, S. M.,
312 Müller, K., Ooms, J., Robinson, D., Seidel, D. P., Spinu, V., ... Yutani, H. (2019). Welcome
313 to the tidyverse. *Journal of Open Source Software*, 4(43), 1686. [https://doi.org/10.21105/](https://doi.org/10.21105/joss.01686)
314 [joss.01686](https://doi.org/10.21105/joss.01686)