

Preprocessing Electronic Health Records for Analysis-Ready Data in an Asthma Cohort

Kimberly Lactaoen, MS
Thursday, June 12, 2025

EHR data

- EHR data captures extensive patient data for large-scale studies
- However, it is complex, subject to inaccuracies, and prone to missingness which may not represent a patient's true health status¹
- Demonstrate specific challenges when preparing data for analysis
 - Tidyverse (dplyr, lubridate, stringr, and tidyr)

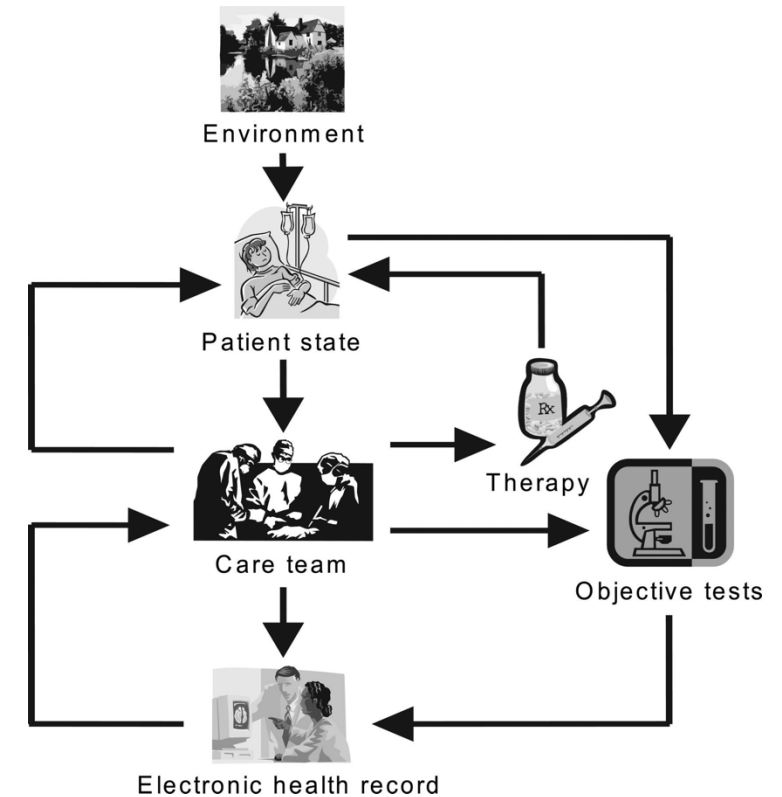


Figure 1. Feedback loops in the electronic health record.

J Am Med Inform Assoc, Volume 20, Issue 1, January 2013, Pages 117–121,
<https://doi.org/10.1136/amiajnl-2012-001145>

¹Hubbard RA, Lou C, Himes BE: The Effective Sample Size of EHR-Derived Cohorts Under Biased Sampling. *Modern Statistical Methods for Health Research*. Zhao Y, Chen D (eds.). Springer International Publishing, 2021

EHR data collection

- Data collected from the Clarity database from PennChart (EPIC)
- Request data relevant to asthma research
 - Demographics
 - Diagnosis codes
 - Medications
 - Laboratory test results



Figure 2. Screenshot of PennChart login page.

Demographics

- **Challenge:** understanding how data was collected
- **Background info:** demographic information is presented by encounter level
- **Resolution:** communication with data analyst
 - *Sex, Race, Ethnicity, and Smoking* status were captured from Patient Table
 - Used the most recent entry
 - *Insurance and BMI* variables change over time

Patient_ID	Enc_ID	Date	Sex	Race	BMI	Ethnicity	Smoking	Insurance
1	kdl1	2023-01-05	F	White	28	non-Hispanic/Latino	Ever	Private
2	kdl2	2025-10-13	M	Black	31	non-Hispanic/Latino	Never	Medicare
1	kdl3	2025-03-22	F	White	29	non-Hispanic/Latino	Ever	Private
1	kdl4	2025-05-30	F	White	29	non-Hispanic/Latino	Ever	Unknown
3	kdl5	2025-06-12	M	White	31	non-Hispanic/Latino	Unknown	Private
4	kdl6	2025-04-27	M	Asian	32	Hispanic/Latino	Unknown	Medicaid

Figure 3. Example of demographic characteristics from EHR data.

Demographics

```
1 `{{r}}`
2 sex <- enc |>
3   group_by(patient_ID) |>
4   summarise(sex_combined = paste(unique(sex), collapse = ", ")) |>
5   ungroup() ##same
6 table(sex$sex_combined, useNA = "ifany")
7
8 race <- enc |>
9   group_by(patient_ID) |>
10  summarise(race_combined = paste(unique(race), collapse = ", ")) |>
11  ungroup() ##same
12 table(race$race_combined, useNA = "ifany")
13
14 ethnicity <- enc |>
15   group_by(patient_ID) |>
16   summarise(ethnicity_combined = paste(unique(ethnicity), collapse = ", ")) |>
17   ungroup() ##same
18 table(ethnicity$ethnicity_combined, useNA = "ifany")
19
20 smoking <- enc |>
21   group_by(patient_ID) |>
22   summarise(smoking_combined = paste(unique(smoking), collapse = ", ")) |>
23   ungroup() ##same
24 table(smoking$smoking_combined, useNA = "ifany")
25
26 bmi <- enc |>
27   group_by(patient_ID) |>
28   summarise(BMI_combined = paste(unique(BMI), collapse = ", ")) |>
29   ungroup() ##different
30 table(bmi$BMI_combined, useNA = "ifany")
31
32 insurance <- enc |>
33   group_by(patient_ID) |>
34   summarise(insurance_combined = paste(unique(insurance), collapse = ", ")) |>
35   ungroup() ##different
36 table(insurance$insurance_combined, useNA = "ifany")
37
38 `{{r}}`
```



```
1 `{{r}}`
2 Table1 <- enc |>
3   group_by(patient_ID) |>
4   mutate(Date = as_date(Date)) |>
5   filter(Date == (max(Date))) |>
6   filter(row_number()==1) |>
7   select(patient_ID, sex, race, ethnicity, smoking) |>
8   distinct()
9 `{{r}}`
```

Figure 4. Example of shortening R script after understanding how the data was collecting, eliminating much of the exploratory analysis.

Demographics

- **Challenge:** conflicting patient information
- **Background info:** The *Race* variable has Hispanic/Latino-Black(HLB) and Hispanic/Latino-White(HLW) values where some patients with these values also have non-Hispanic/Latino as their *Ethnicity*
- **Resolution:** those that were denoted HLB or HLW in *Race* AND Hispanic/Latino in *Ethnicity*, were recoded as Black or White, respectively, in *Race*. Those with conflicting information were recoded as Unknown.

Patient_ID	Enc_ID	Race	Ethnicity
5	kdI7	HLB	non-Hispanic/Latino
6	kdI8	White	non-Hispanic/Latino
7	kdI9	Black	non-Hispanic/Latino
8	kdI10	HLW	Hispanic/Latino

Figure 5. Example of Race and Ethnicity variables from EHR data.

Demographics

```
1  ````{r}
2  enc_update <- enc |>
3    mutate(Race = case_when(patient_ID %in% Race_Ethnicity_conflict$patient_ID ~ "Unknown",
4                                .default = as.character(Race))) |>
5    mutate(Ethnicity = case_when(patient_ID %in% Race_Ethnicity_conflict$patient_ID ~ "Unknown",
6                                  .default = as.character(Ethnicity))) |>
7    mutate(Race = case_when(patient_ID %in%
8                              no_conflict_HLW$patient_ID ~ "White",
9                              .default = as.character(Race))) |>
10   mutate(Race = case_when(patient_ID %in%
11                             no_conflict_HLB$patient_ID ~ "Black",
12                             .default = as.character(Race)))
13  ````
```

Figure 6. Example of using `case_when` function from `dplyr` package.

Diagnostic codes

- **Challenge:** Diagnostic code descriptions are not always consistent with diagnostic code
- **Background info:** EHR data has diagnostic codes and diagnostic code descriptions. As an example, J45.2 is the ICD-10 code for Mild intermittent asthma, uncomplicated
- **Resolution:** link diagnostic code information directly from government-based ICD-10 files

Patient_ID	Enc_ID	ICD-10	Description
5	kdl7	J45.2	Mild intermittent asthma without complication
6	kdl8	J45.2	Mild intermittent asthma
7	kdl9	J45.2	Mild intermittent asthma, uncomplicated

Figure 7. Example of ICD-10 codes and their descriptions in EHR data.

Diagnostic codes

```
1  ```{r}
2  dx <- read_delim("/project/diagnosis.txt") |>
3    mutate(icd_code =
4      str_replace_all(CODE, "\\.", ""))
5
6  lines <- read_lines("/project/icd10cm_order_2025.txt")
7
8  parts_df <- tibble(
9    icd_code = map_chr(str_split(lines, "\\s+"), ~ .x[2]),
10    icd_definition = str_sub(lines, 78)
11  )
12
13  dx_updated <- left_join(dx, parts_df, by = "icd_code") |>
14    select(patient_ID, Date, icd_code, icd_definition)
15  ```
```

Figure 8. Example of using functions from readr, stringr, and dplyr packages.

Diagnostic codes

```
1  ```{r}
2  dx_bypatient <- dx_updated |>
3    select(patient_ID, icd_code) |>
4    distinct() |>
5    mutate(value = 1) |>
6    pivot_wider(
7      names_from = icd_code,
8      values_from = value,
9      values_fill = list(value = 0))
10  ```
```

patient_ID	icd_code
1	J45.2
1	J45.3
1	J45.5
2	J45.2
2	D50.9
...	...



patient_ID	J45.2	J45.3	J45.5	D50.9	M54.2	...
1	1	1	1	0	0	...
2	1	0	0	1	0	...
3	1	1	1	1	1	...

Figure 9. Example of using `pivot_wider` function from `tidyr` packages to create data frame of diagnostic codes by patient.

Medications

- **Challenge:** selecting medications relevant to asthma treatment
- **Background info:** We searched for therapeutic and pharmacological classes like respiratory agents and corticosteroids and narrowed our search to medications like inhaled corticosteroids (ICS) and oral corticosteroids (OCS). However, not all specific medication names under these categories were likely to treat asthma.
- **Resolution:** Individually reviewed each medication name

NAME	MED_ROUTE
AEROBID INHALATION	INHALATION
<u>ALLERGY RELIEF NASAL</u>	<u>NASAL</u>
ALVESCO 160 MCG/ACT INHALATION AERS	INHALATION
ALVESCO 160 MCG/ACT INHALATION AERS	NULL
ALVESCO 80 MCG/ACT INHALATION AERS	INHALATION
ALVESCO 80 MCG/ACT INHALATION AERS	NULL
ALVESCO 80 MCG/ACT INHALATION AERS	ORAL
ALVESCO INHALATION	INHALATION

Figure 10. Example of medication names and medication routes.

Medications

NAME	MED_ROUTE	Keep
BENRALIZUMAB 30 MG/ML SC SOAJ	SUBCUTANEOUS	yes
BENRALIZUMAB 30 MG/ML SC SOSY	INJECTION	yes
BENRALIZUMAB 30 MG/ML SC SOSY	NULL	yes
BENRALIZUMAB 30 MG/ML SC SOSY	SUBCUTANEOUS	yes
BENRALIZUMAB 30 MG/ML SC SOSY	SUBCUTANEOUS INFUSION	yes
BENRALIZUMAB SC	SUBCUTANEOUS	yes
CINQAIR 100 MG/10ML IV SOLN	INTRAVENOUS	yes
CINQAIR 100 MG/10ML IV SOLN	NULL	yes
CINQAIR IV	INTRAVENOUS	yes
CINQAIR IV	NULL	yes

►

LAMA_route

biologic_route

ICS_LABA_route

ICS_LABA

```
1  `` {r}
2  med_route <- c("OCS_route", "ICS_route", "SABA_route", "SAMA_route", "LABA_route",
3                "LAMA_route", "biologic_route", "ICS_LABA_route",
4                "ICS_LABA_LAMA_route", "SABA_SAMA_route", "LABA_LAMA_route")
5
6  meds_to_keep <- map_dfr(med_route, ~
7    read_excel("/project/unique_med_route.xlsx", sheet = .x) |>
8    select(NAME, MED_ROUTE, Keep)
9  )
10 ``
```

Figure 11. Example of manually reviewing medication data then using `map_dfr` function from the `purrr` package to create one data frame of medications that will be used to filter for relevant asthma medications.

Laboratory test results

- **Challenge:** not all facilities report the same units of measurement
- **Background info:** Eosinophil units of measures vary in this codified dataset and some units are not reported
- **Resolution:** Convert units of measures to one unit. For those eosinophil measures without a unit of measure, either remove from data, make best assumptions, or chart review.

```
eosinophils_updated <- eosinophils |>
  mutate(REFERENCE_UNIT = case_when(REFERENCE_UNIT %in% c("TH0/uL",
    "K/uL",
    "Thousand/uL",
    "x10E3/uL",
    "Thou/uL",
    "x10(3)/mcl",
    "x10*3/uL",
    "x10E3/uL",
    "k/uL",
    "K/UL",
    "x10E3/u") ~ "10*3/uL",
    REFERENCE_UNIT == "cells/uL" ~ "cells/uL",
    REFERENCE_UNIT == "K/cu mm" ~ "k/cumm",
    .default = REFERENCE_UNIT)) |>
  mutate(ORD_VALUE = as.numeric(ORD_VALUE)) |>
  mutate(ORD_VALUE = case_when(REFERENCE_UNIT == "cells/uL" ~ ORD_VALUE/1000,
    .default = ORD_VALUE)) |>
  mutate(REFERENCE_UNIT = case_when(REFERENCE_UNIT %in% c("k/cumm",
    "cells/uL") ~ "10*3/uL",
    .default = REFERENCE_UNIT)) |>
  rename(eos_value = ORD_VALUE)
```

Figure 12. Example script for converting multiple units of measure to one unite in eosinophil data.

Conclusion

- Understand how the data was collected
 - Example: Patient Table and encounter-level demographic variables
- Be on the lookout for conflicting information
 - Example: Race and Ethnicity variables
- Be aware of unspecified information
 - Example: ICD-10 codes
- Remember to collaborate with experts
 - Example: medication names
- Simplify data when you can
 - Example: eosinophil units of measure

Thank you!

- Blanca Himes, PhD
- Gary Weissman, MD, MS



Email: kimberly.lactaoen@pennmedicine.upenn.edu