

NIH *All of Us* Allostatic (over)Load Dataset Construction

CHORDS Lab – Washington State University
Authors: Shawna Beese, Trey DeLong, Jason Cross
Last updated: 2025-08-15

Overview

This Jupyter notebook walks through the construction of a clean, analysis-ready allostatic (over)load dataset using biomarker and demographic data. The notebook extracts key measures including:

- C-reactive protein (CRP)
- HDL cholesterol
- Hemoglobin A1c (A1C)
- Heart rate
- Waist/height ratio

It removes outliers and merges these measures with demographic variables, including age, sex at birth, gender, race, and ethnicity.

Binary and ordinal allostatic load scores are derived from these measures, providing ready-to-use outcomes for downstream analyses. The notebook also calculates age at measurement and generates a summary table of demographic characteristics by sex.

🚩 **Note:** This notebook requires access to the NIH *All of Us* Researcher Workbench and is intended **only** for use within that secure environment. See [NIH_access_guide.md](#) for details on becoming a registered NIH *All of Us* researcher.

This notebook relies on the R packages `dplyn`, `tidyr`, `misty`, `eeptools`, and `table1`.

Setup and Biomarker Extraction

Setup and Initial Libraries

Load necessary R packages (`tidyr`, `dplyn`, `misty`) and prepare the dataset for biomarker analysis.

```
In [ ]: install.packages(c("tidyr", "dplyn", "misty"))

library(tidyr)
library(dplyn)
library(misty)
```

Extract and Process CRP Measurements

Filter the dataset for C-reactive protein measurements. For each person, identify the most recent measurement, remove data older than 1 year, and calculate per-person mean values. Check basic statistics using `misty::describe`.

```
In [ ]: CRP = dataset_XXXXXXX_measurement_df[dataset_XXXXXXX_measurement_df$standard_concept_name ==
      "C-reactive protein [mass/volume] in Serum or Plasma",]
CRP = CRP%>group_by(person_id)%>%mutate(recent_datetime = max(measurement_datetime))
CRP$greater_year_ago<-difftime(CRP$recent_datetime, CRP$measurement_datetime)>360
CRP$less_three_apart<-difftime(CRP$recent_datetime, CRP$measurement_datetime)
CRP_tr1m1 = CRP[CRP$greater_year_ago,]
CRP_value_means = CRP_tr1m1%>%summarise(counts = n(),mean_val = mean(value_as_number))
head(CRP_value_means)
misty::describe(CRP_value_means$mean_val)
```

Remove CRP Outliers

Filter out extreme CRP values using ± 4 SD from the mean to create a cleaned dataset. Ensure each person_id is unique and rename columns for clarity.

```
In [ ]: #save CRP dataframe with outliers removed
library(dplyn)

CRP_OR <- CRP_filtered %>%
  group_by(person_id) %>%
  summarise(mean_val = mean(mean_val, na.rm = TRUE), counts = n()) %>%
  ungroup() %>%
  filter(
    mean_val < mean(mean_val, na.rm = TRUE) + 4 * sd(mean_val, na.rm = TRUE) &
    mean_val > mean(mean_val, na.rm = TRUE) - 4 * sd(mean_val, na.rm = TRUE)
  )

# Check duplicate person_id counts
sum(duplicated(CRP_OR$person_id))

#re-label and reduce data object
colnames(CRP_OR)<-c('person_id','CRP','counts')
misty::describe(CRP_OR)
CRP_OR[CRP_OR[,3]>3]
colnames(CRP_OR)
sum(duplicated(CRP_OR$person_id))
```

Extract and Process HDL Measurements

Filter for HDL cholesterol, keep measurements within 1 year, calculate per-person mean, and remove outliers using ± 4 SD.

```
In [ ]: HDL <- dataset_XXXXXXX_measurement_df %>%
  filter(standard_concept_name == "Cholesterol in HDL [Mass/volume] in Serum or Plasma") %>%
  select(person_id, value_as_number, measurement_datetime) %>%
  group_by(person_id) %>%
  mutate(recent_datetime = max(measurement_datetime)) %>% # Find most recent measurement per person
  filter(difftime(recent_datetime, measurement_datetime, units = "days") <= 365) %>% # Keep measurements within 1 year
  summarise(counts = n(), mean_val = mean(value_as_number, na.rm = TRUE)) %>% # Compute mean and count per person
  ungroup()

head(HDL)
misty::describe(HDL)

In [ ]: HDL_OR <- HDL %>%
  group_by(person_id) %>%
  summarise(mean_val = mean(mean_val, na.rm = TRUE)) %>%
  ungroup() %>%
  filter(
    mean_val < mean(mean_val, na.rm = TRUE) + 4 * sd(mean_val, na.rm = TRUE) &
    mean_val > mean(mean_val, na.rm = TRUE) - 4 * sd(mean_val, na.rm = TRUE)
  )

head(HDL_OR)
sum(duplicated(HDL_OR$person_id))
colnames(HDL_OR)<-c('person_id','HDL')
misty::describe(HDL_OR)
```

Extract and Process A1C Measurements

Filter for Hemoglobin A1c, keep recent measurements, compute mean per person, and remove outliers.

```
In [ ]: A1C <- dataset_XXXXXXX_measurement_df %>%
  filter(standard_concept_name == "Hemoglobin A1c/Hemoglobin total in Blood") %>%
  select(person_id, value_as_number, measurement_datetime) %>%
  group_by(person_id) %>%
  mutate(recent_datetime = max(measurement_datetime)) %>%
  filter(difftime(recent_datetime, measurement_datetime, units = "days") <= 365) %>%
  summarise(counts = n(), mean_val = mean(value_as_number, na.rm = TRUE)) %>%
  ungroup()

head(A1C)
misty::describe(A1C)

In [ ]: A1C_OR <- A1C %>%
  group_by(person_id) %>%
  summarise(mean_val = mean(mean_val, na.rm = TRUE)) %>%
  ungroup() %>%
  filter(
    mean_val < mean(mean_val, na.rm = TRUE) + 4 * sd(mean_val, na.rm = TRUE) &
    mean_val > mean(mean_val, na.rm = TRUE) - 4 * sd(mean_val, na.rm = TRUE)
  )

# Check for duplicate person_ids
sum(duplicated(A1C_OR$person_id))
colnames(A1C_OR)<-c('person_id','A1C')
misty::describe(A1C_OR)
```

Waist/Height Ratio and Heart Rate

Extract Waist and Height Measurements

Select waist circumference and height measurements to calculate waist-to-height ratios later. Outliers are not removed at this stage.

```
In [ ]: #extract waist and body height rows - no outliers removed
Waist = dataset_XXXXXXX_measurement_df[dataset_XXXXXXX_measurement_df$standard_concept_name ==
      "Computed waist circumference, mean of closest two measures",
      c('person_id','value_as_number')]
colnames(Waist)=c('person_id','Waist')
misty::describe(Waist)

Height = dataset_XXXXXXX_measurement_df[dataset_XXXXXXX_measurement_df$standard_concept_name ==
      "Body height",
      c('person_id','value_as_number')]
colnames(Height)=c('person_id','Height')
misty::describe(Height)
```

Calculate Waist-to-Height Ratio

Merge waist and height data and compute the waist/height ratio for each person.

```
In [ ]: #calculate waist/height ratio
WaistHeight=merge(Waist,Height,by="person_id")
WaistHeight$WaistHeight=WaistHeight$Waist/WaistHeight$Height
misty::describe(WaistHeight)
```

Remove Waist/Height Outliers

Clean the waist/height ratio data by removing extreme values using ± 4 SD.

```
In [ ]: WaistHeight_OR <- WaistHeight %>%
  group_by(person_id) %>%
  summarise(mean_WaistHeight = mean(WaistHeight, na.rm = TRUE)) %>%
  ungroup() %>%
  filter(
    mean_WaistHeight < mean(mean_WaistHeight, na.rm = TRUE) + 4 * sd(mean_WaistHeight, na.rm = TRUE) &
    mean_WaistHeight > mean(mean_WaistHeight, na.rm = TRUE) - 4 * sd(mean_WaistHeight, na.rm = TRUE)
  )

# Check for duplicates
sum(duplicated(WaistHeight_OR$person_id))
colnames(WaistHeight_OR)=c('person_id','WaistHeight')
misty::describe(WaistHeight_OR)
```

Extract Heart Rate Measurements

Select mean heart rate measurements (2nd and 3rd measures).

```
In [ ]: #extract heart rate rows - no outliers removed
HR = dataset_XXXXXXX_measurement_df[dataset_XXXXXXX_measurement_df$standard_concept_name ==
      "Computed heart rate, mean of 2nd and 3rd measures",
      c('person_id','value_as_number')]
colnames(HR)=c('person_id','HR')
misty::describe(HR)
```

Remove Heart Rate Outliers

Clean heart rate data using ± 4 SD and save for analysis.

```
In [ ]: #remove HR outliers and save

HR_OR <- HR %>%
  group_by(person_id) %>%
  summarise(mean_HR = mean(HR, na.rm = TRUE)) %>%
  ungroup() %>%
  filter(
    mean_HR < mean(mean_HR, na.rm = TRUE) + 4 * sd(mean_HR, na.rm = TRUE) &
    mean_HR > mean(mean_HR, na.rm = TRUE) - 4 * sd(mean_HR, na.rm = TRUE)
  )

# Check duplicate person_id counts
sum(duplicated(HR_OR$person_id))
colnames(HR_OR)=c('person_id','HR')
misty::describe(HR_OR)
```

Prepare Measurement Dates

Create a dataset linking person_id with measurement dates.

```
In [ ]: colnames(dataset_XXXXXXX_measurement_df)

In [ ]: MeasureDate = cbind(dataset_XXXXXXX_measurement_df$person_id,
      dataset_XXXXXXX_measurement_df$measurement_datetime)

colnames(MeasureDate)=c('person_id','MeasureDate')
head(MeasureDate)
```

Merge All Biomarkers into One Dataset (AL5)

Combine cleaned CRP, HDL, A1C, heart rate, and waist/height data into a single allostatic load dataset.

```
In [ ]: #merge allostatic load measures
A = merge(CRP_OR,HR_OR,by="person_id")
B = merge(A,HDL_OR,by="person_id")
C = merge(B,WaistHeight_OR,by="person_id")
AL5 = merge(C,A1C_OR,by="person_id")
misty::describe(AL5)
```

```
In [ ]: sum(duplicated(AL5$person_id))
```

```
In [ ]: head(AL5)
```

Allostatic Load Scores

Create Allostatic Load Scores

Convert biomarkers to binary "high-risk" flags based on quartiles. Sum the flags for an ordinal allostatic load score and create a binary "high AL" indicator.

```
In [ ]: ##create Allostatic ordinal and binary
#top quartile for AL5 items
#calculate thresholds
CRP_quart=quantile(AL5$CRP,probs=.75,na.rm=T)
HR_quart=quantile(AL5$HR,probs=.75,na.rm=T)
HDL_quart=quantile(AL5$HDL,probs=.25,na.rm=T)
WaistHeight_quart=quantile(AL5$WaistHeight,probs=.75,na.rm=T)
A1C_quart=quantile(AL5$A1C,probs=.75,na.rm=T)
#encode items to binary
AL5$CRP_load=ifelse(AL5$CRP>CRP_quart,1,0)
AL5$HR_load=ifelse(AL5$HR>HR_quart,1,0)
AL5$HDL_load=ifelse(AL5$HDL>HDL_quart,1,0)
AL5$WaistHeight_load=ifelse(AL5$WaistHeight>WaistHeight_quart,1,0)
AL5$A1C_load=ifelse(AL5$A1C>A1C_quart,1,0)
#sum for allostatic ordinal
AL5$AL5_sum=AL5$CRP_load+AL5$HR_load+AL5$HDL_load
AL5$AL5_high=ifelse(AL5$AL5_sum>2,1,0)
misty::describe(AL5[,c('CRP_load','HR_load','HDL_load','WaistHeight_load','A1C_load','AL5_sum',
      'AL5_high')])
colnames(AL5)
```

Demographics and Final Preparation

Merge Demographics with Measurements

Combine measurement dates and AL5 scores with demographic information for each person.

```
In [ ]: DemoMeasureDate <-merge(MeasureDate, dataset_XXXXXXX_person_df, by = "person_id")
DemoAL <- merge(AL5, dataset_XXXXXXX_person_df, by = "person_id")
head(DemoMeasureDate)
```

Remove Duplicate Records

Ensure each person_id appears only once to avoid duplication in downstream analysis.

```
In [ ]: clean_DemoMeasureDate <- DemoMeasureDate %>% distinct(person_id ,.keep_all = TRUE)
nrow(clean_DemoMeasureDate)

In [ ]: nrow(AL5)

In [ ]: AL5_DemoReady <- merge(clean_DemoMeasureDate, AL5, by = "person_id")
nrow(AL5_DemoReady)
head(AL5_DemoReady)
colnames(AL5_DemoReady)
```

Calculate Age at Measurement

Use `eeptools::age_calc` to compute age at the time of measurement and add to the dataset.

```
In [ ]: install.packages("eeptools")

In [ ]: AL5_Age<- eeptools::age_calc(as.Date(AL5_DemoReady$date_of_birth),
      as.Date(AL5_DemoReady$MeasureDate),
      units="years")

boxplot(AL5_Age)
```

```
In [ ]: AL5_DemoReadyAge <-cbind(AL5_DemoReady, AL5_Age)
```

Create Summary Table

Use `table1` to summarize demographics (race, ethnicity, age) stratified by sex. Includes a helper function to display HTML tables in Jupyter notebooks.

```
In [ ]: #Ready<-na.omit(AL5_DemoReadyAge)
output <- table1::table1(~ race + ethnicity + AL5_Age | sex_at_birth, data=AL5_DemoReadyAge)

require(htmltools)
require(IRdisplay)
require(table1)

display_jupyter <- function(x){
  css <- system.file("table1_defaults.1.0/table1_defaults.css", package="table1")
  css <- paste(readlines(css), collapse="\n")
  x <- htmltools::tagList(htmltools::tags$style(css), htmltools::tags$div(class="Rtable1", x))
  IRdisplay::display_html(as.character(x))
}

display_jupyter(output)
```