

process_cleaning_merging

Approach:

My analytical approach addresses the distinct characteristics of the data, which is organised into three main temporal categories: daily, hourly and minute-level datasets. I will process each of these categories independently. A critical initial step for each dataset will be to thoroughly examine its integrity and completeness to ensure data quality, before proceeding with further analysis. The same approach will be implemented for the other two data sets representing heart rate and user weight data.

Packages used:

```
library(janitor)
library(skimr)
library(tidyverse)
library(lubridate)
library(here)
library(readr)
```

Daily data sets:

Importing datasets:

```
activity <- read.csv(here("mturkfitbit_export_4.12.16-5.12.16",
                          "Fitabase Data 4.12.16-5.12.16",
                          "dailyActivity_merged.csv"))

steps <- read.csv(here("mturkfitbit_export_4.12.16-5.12.16",
                      "Fitabase Data 4.12.16-5.12.16",
                      "dailySteps_merged.csv"))

intensities <- read_csv(here("mturkfitbit_export_4.12.16-5.12.16",
                             "Fitabase Data 4.12.16-5.12.16",
                             "dailyIntensities_merged.csv"))
```

```
## Rows: 940 Columns: 10
## -- Column specification -----
## Delimiter: ","
## chr (1): ActivityDay
## dbl (9): Id, SedentaryMinutes, LightlyActiveMinutes, FairlyActiveMinutes, Ve...
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
```

```
sleep <- read_csv(here("mturkfitbit_export_4.12.16-5.12.16",
                      "Fitabase Data 4.12.16-5.12.16",
                      "sleepDay_merged.csv"))

## Rows: 413 Columns: 5
## -- Column specification -----
## Delimiter: ","
## chr (1): SleepDay
## dbl (4): Id, TotalSleepRecords, TotalMinutesAsleep, TotalTimeInBed
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.
```

Upon initial inspection, it is evident that the 'activity' dataset already contains information on calories, intensity and steps, thus eliminating the need to load a separate dataset for these metrics. The 'sleep' dataset will be kept separate due to its distinct nature.

Data cleaning:

It's easier to me to work with snake_case so I 'clear_names()' in all datasets.

```
daily_all_data_frames <- list(
  activity = activity,
  steps = steps,
  intensities = intensities,
  sleep = sleep
)

daily_cleaned_data_frames <- daily_all_data_frames %>%
  map(~ .x %>% janitor::clean_names())

list2env(daily_cleaned_data_frames, env = .GlobalEnv)
```

```
## <environment: R_GlobalEnv>
```

Now let's have a quick overview of the data using *skim*.

```
activity %>% skimr::skim_without_charts()
```

Table 1: Data summary

Name	Piped data
Number of rows	940
Number of columns	15
Column type frequency:	
character	1
numeric	14
Group variables	None

Variable type: character

skim_variable	n_missing	complete_rate	min	max	empty	n_unique	whitespace
activity_date	0	1	8	9	0	31	0

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100
id	0	1	4.855407e-20	2.4805e-16	0.39603	0.6320127e-40	0.945115e-60	0.962181e-80	0.977689e+09
total_steps	0	1	7.637910e-50	8.7150e+03	0	3.789750e-70	0.55500e-10	0.72700e-30	0.01900e+04
total_distance	0	1	5.490000e-30	2.0000e+00	0	2.620000e-50	0.40000e-70	0.10000e-20	0.03000e+01
tracker_distance	0	1	5.480000e-30	1.0000e+00	0	2.620000e-50	0.40000e-70	0.10000e-20	0.03000e+01
logged_activities_distance	0	1	1.100000e-6	2.00000e-01	0	0.000000e-9	0.000000e-9	0.000000e-10	0.040000e+00
very_active_distance	0	1	1.500000e-20	6.0000e+00	0	0.000000e-20	0.00000e-2	0.050000e-20	0.092000e+01
moderately_active_distance	0	1	5.700000e-01	8.80000e-01	0	0.000000e-20	0.00000e-8	0.000000e-6	0.480000e+00
light_active_distance	0	1	3.340000e-20	4.0000e+00	0	1.950000e-30	0.60000e-40	0.80000e-10	0.071000e+01
sedentary_active_distance	0	1	0.000000e-10	0.00000e-02	0	0.000000e-9	0.000000e-9	0.000000e-10	0.000000e-01
very_active_minutes	0	1	2.116000e-30	2.84000e+01	0	0.000000e-40	0.00000e-30	0.200000e-20	0.100000e+02
fairly_active_minutes	0	1	1.356000e-10	9.9000e+01	0	0.000000e-60	0.00000e-10	0.000000e-10	0.130000e+02
lightly_active_minutes	0	1	1.928100e-10	2.91700e+02	0	1.270000e-10	0.90000e-20	0.240000e-50	0.280000e+02
sedentary_minutes	0	1	9.912100e-30	2.12700e+02	0	7.297500e-10	0.257500e-10	0.29500e-10	0.040000e+03
calories	0	1	2.303610e-70	8.1700e+02	0	1.828500e-20	0.34000e-20	0.33250e-40	0.000000e+03

```
sleep %>% skimr::skim_without_charts()
```

Table 4: Data summary

Name	Piped data
Number of rows	413
Number of columns	5
Column type frequency:	
character	1
numeric	4
Group variables	None

Variable type: character

skim_variable	n_missing	complete_rate	min	max	empty	n_unique	whitespace
sleep_day	0	1	20	21	0	31	0

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100
id	0	1	5.000979e+00	0.6036e+01	503960363	977333714	702921684	962181068	792009665
total_sleep_records	0	1	1.120000e+01	0.50000e-01	1	1	1	1	3
total_minutes_asleep	0	1	4.194700e+02	2.8340e+02	58	361	433	490	796
total_time_in_bed	0	1	4.586400e+02	2.7100e+02	61	403	463	526	961

There are no missing values in each data set. Let's check if the number of people represented by the 'id' column is the same in every dataset.

```
n_distinct(activity)
```

```
## [1] 940
```

```
n_distinct(sleep)
```

```
## [1] 410
```

```
n_distinct(activity$id)
```

```
## [1] 33
```

```
n_distinct(sleep$id)
```

```
## [1] 24
```

It's now evident that the sleep dataset contains a smaller number of users. This needs to be taken into consideration during the analysis.

Daily data has been gathered from 33 people; however, only 24 users are represented in the sleep data.

Fixing the data format:

```
activity$activity_date=as.POSIXct(activity$activity_date, format="%m/%d/%Y", tz=Sys.timezone())
```

Hourly Data Sets:

Importing datasets:

```
hourly_calories <- read.csv(here("mturkfitbit_export_4.12.16-5.12.16",  
                                "Fitabase Data 4.12.16-5.12.16",  
                                "hourlyCalories_merged.csv"))
```

```
hourly_intensities <- read.csv(here("mturkfitbit_export_4.12.16-5.12.16",
                                   "Fitabase Data 4.12.16-5.12.16",
                                   "hourlyIntensities_merged.csv"))

hourly_steps <- read_csv(here("mturkfitbit_export_4.12.16-5.12.16",
                              "Fitabase Data 4.12.16-5.12.16",
                              "hourlySteps_merged.csv"))
```

Upon initial review the datasets should be merge together.

Data cleaning:

```
hourly_all_data_frames <- list(
  hourly_calories = hourly_calories,
  hourly_intensities = hourly_intensities,
  hourly_steps = hourly_steps
)

hourly_cleaned_data_frames <- hourly_all_data_frames %>%
  map(~ .x %>% janitor::clean_names())

list2env(hourly_cleaned_data_frames, envir = .GlobalEnv)
```

```
## <environment: R_GlobalEnv>
```

Let's again have a quick overview plus check the distinct 'id' value.

```
hourly_calories %>% skimr::skim_without_charts()
```

Table 7: Data summary

Name	Piped data
Number of rows	22099
Number of columns	3
Column type frequency:	
character	1
numeric	2
Group variables	None

Variable type: character

skim_variable	n_missing	complete_rate	min	max	empty	n_unique	whitespace
activity_hour	0	1	19	21	0	736	0

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100
id	0	1	4.848235e+00	4.225e+00	503960362	320127002	445114986	696218106	8877689391
calories	0	1	9.739000e+01	610700e+01	42	63	83	108	948

```
hourly_intensities %>% skimr::skim_without_charts()
```

Table 10: Data summary

Name	Piped data
Number of rows	22099
Number of columns	4
Column type frequency:	
character	1
numeric	3
Group variables	None

Variable type: character

skim_variable	n_missing	complete_rate	min	max	empty	n_unique	whitespace
activity_hour	0	1	19	21	0	736	0

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100
id	0	1	4.848235e+00	4.225e+00	503960362	320127002	445114986	696218106	8877689391
total_intensity	0	1	1.204000e+01	1.130e+01	0	0	3.000000e+00	6.000000e+01	180
average_intensity	0	1	2.000000e-01	3.5000e-01	0	0	5.000000e-02	2.700000e-01	3

```
hourly_steps %>% skimr::skim_without_charts()
```

Table 13: Data summary

Name	Piped data
Number of rows	22099
Number of columns	3
Column type frequency:	
character	1
numeric	2
Group variables	None

Variable type: character

skim_variable	n_missing	complete_rate	min	max	empty	n_unique	whitespace
activity_hour	0	1	19	21	0	736	0

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100
id	0	1	4.848235e+01	4.225e+01	503960360	232012700	244511498	696218106	8877689391
step_total	0	1	3.201700e+02	29038e+02	0	0	40	357	10554

```
n_distinct(hourly_calories)
```

```
## [1] 22099
```

```
n_distinct(hourly_intensities)
```

```
## [1] 22099
```

```
n_distinct(hourly_steps)
```

```
## [1] 22099
```

```
n_distinct(hourly_calories$id)
```

```
## [1] 33
```

```
n_distinct(hourly_intensities$id)
```

```
## [1] 33
```

```
n_distinct(hourly_steps$id)
```

```
## [1] 33
```

Datasets have the same length, let's merge it.

Datasets merging:

```
hourly_cis_temp <- merge(hourly_calories, hourly_intensities, by = c('id', 'activity_hour'))
hourly_cis <- merge(hourly_cis_temp, hourly_steps, by = c('id', 'activity_hour'))
head(hourly_cis)
```

	id	activity_hour	calories	total_intensity	average_intensity
## 1	1503960366	4/12/2016 1:00:00 AM	61	8	0.133333
## 2	1503960366	4/12/2016 1:00:00 PM	66	6	0.100000
## 3	1503960366	4/12/2016 10:00:00 AM	99	29	0.483333
## 4	1503960366	4/12/2016 10:00:00 PM	65	9	0.150000
## 5	1503960366	4/12/2016 11:00:00 AM	76	12	0.200000
## 6	1503960366	4/12/2016 11:00:00 PM	81	21	0.350000
##	step_total				
## 1	160				
## 2	221				
## 3	676				
## 4	89				
## 5	360				
## 6	338				

I ended up with one dataset reflecting the hourly data, called 'hourly_cis'.

Fixing the timestamp.

I would like to separate day and time from 'activity_hour'.

```
hourly_cis$activity_hour=as.POSIXct(hourly_cis$activity_hour, format="%m/%d/%Y %I:%M:%S %p", tz=Sys.time())
hourly_cis$time <- format(hourly_cis$activity_hour, format = "%H:%M:%S")
hourly_cis$date <- format(hourly_cis$activity_hour, format = "%m/%d/%y")
```

Minute Data Sets:

Importing Data:

```
minute_calories <- read.csv(here("mturkfitbit_export_4.12.16-5.12.16",
                                "Fitabase Data 4.12.16-5.12.16",
                                "minuteCaloriesNarrow_merged.csv"))

minute_intensities <- read.csv(here("mturkfitbit_export_4.12.16-5.12.16",
                                    "Fitabase Data 4.12.16-5.12.16",
                                    "minuteIntensitiesNarrow_merged.csv"))

minute_METs <- read.csv(here("mturkfitbit_export_4.12.16-5.12.16",
                              "Fitabase Data 4.12.16-5.12.16",
                              "minuteMETsNarrow_merged.csv"))

minute_sleep <- read.csv(here("mturkfitbit_export_4.12.16-5.12.16",
                               "Fitabase Data 4.12.16-5.12.16",
                               "minuteSleep_merged.csv"))

minute_steps <- read.csv(here("mturkfitbit_export_4.12.16-5.12.16",
                               "Fitabase Data 4.12.16-5.12.16",
                               "minuteStepsNarrow_merged.csv"))
```


Data cleaning:

```
minute_all_data_frames <- list(  
  minute_calories = minute_calories,  
  minute_intensities = minute_intensities,  
  minute_METs = minute_METs,  
  minute_sleep = minute_sleep,  
  minute_steps = minute_steps  
)  
  
minute_cleaned_data_frames <- minute_all_data_frames %>%  
  map(~ .x %>% janitor::clean_names())  
  
list2env(minute_cleaned_data_frames, envir = .GlobalEnv)
```

```
## <environment: R_GlobalEnv>
```

Overview of data.

```
minute_calories %>% skimr::skim_without_charts()
```

Table 16: Data summary

Name	Piped data
Number of rows	1325580
Number of columns	3
Column type frequency:	
character	1
numeric	2
Group variables	None

Variable type: character

skim_variable	n_missing	complete_rate	min	max	empty	n_unique	whitespace
activity_minute	0	1	19	21	0	44160	0

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100
id	0	1	4.847898e+29	2.22313e+09	0.396036e+32	0.320127e+32	0.445115e+32	0.962181e+32	8.877689e+09
calories	0	1	1.620000e+01	1.100000e+00	0	9.400000e-01	1.220000e+01	0.430000e+01	0.975000e+01

```
minute_intensities %>% skimr::skim_without_charts()
```

Table 19: Data summary

Name	Piped data
Number of rows	1325580
Number of columns	3
Column type frequency:	
character	1
numeric	2
Group variables	None

Variable type: character

skim_variable	n_missing	complete_rate	min	max	empty	n_unique	whitespace
activity_minute	0	1	19	21	0	44160	0

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100
id	0	1	4847897691.2	4.222313e+05	03960360320	1270024451	1498696218	1067877689	391
intensity	0	1	0.2	5.200000e-01	0	0	0	0	3

```
minute_METs %>% skimr::skim_without_charts()
```

Table 22: Data summary

Name	Piped data
Number of rows	1325580
Number of columns	3
Column type frequency:	
character	1
numeric	2
Group variables	None

Variable type: character

skim_variable	n_missing	complete_rate	min	max	empty	n_unique	whitespace
activity_minute	0	1	19	21	0	44160	0

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100
id	0	1	4.847898e+09	2422313e+05	03960366	32012700	24451149	86962181	0678877689391
me_ts	0	1	1.469000e+01	206000e+01	0	10	10	11	157

```
minute_sleep %>% skimr::skim_without_charts()
```

Table 25: Data summary

Name	Piped data
Number of rows	188521
Number of columns	4
Column type frequency:	
character	1
numeric	3
Group variables	None

Variable type: character

skim_variable	n_missing	complete_rate	min	max	empty	n_unique	whitespace
date	0	1	19	21	0	49773	0

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100
id	0	1	4.996595e+09	66950e+09	03960366	3977333714	702921684	962181067	8792009665
value	0	1	1.100000e+00	300000e-01	1	1	1	1	3
log_id	0	1	1.149611e+06	822863e+07	372227280	1439308639	1501142214	1552534115	1616251768

```
minute_steps %>% skimr::skim_without_charts()
```

Table 28: Data summary

Name	Piped data
Number of rows	1325580
Number of columns	3
Column type frequency:	
character	1
numeric	2
Group variables	None

Variable type: character

skim_variable	n_missing	complete_rate	min	max	empty	n_unique	whitespace
activity_minute	0	1	19	21	0	44160	0

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100
id	0	1	4.847898e+09	4.22313e+09	0.03960360	0.32012700	2.44511498	6.96218106	7.87768939
steps	0	1	5.340000e+00	1.3000e+01	0	0	0	0	220

Check unique 'id'.

```
n_distinct(minute_calories)
```

```
## [1] 1325580
```

```
n_distinct(minute_intensities)
```

```
## [1] 1325580
```

```
n_distinct(minute_METs)
```

```
## [1] 1325580
```

```
n_distinct(minute_sleep)
```

```
## [1] 187978
```

```
n_distinct(minute_steps)
```

```
## [1] 1325580
```

```
n_distinct(minute_calories$id)
```

```
## [1] 33
```

```
n_distinct(minute_intensities$id)
```

```
## [1] 33
```

```
n_distinct(minute_METs$id)
```

```
## [1] 33
```

```
n_distinct(minute_sleep$id)
```

```
## [1] 24
```

```
n_distinct(minute_steps$id)
```

```
## [1] 33
```

As expected, the ‘minute_sleep’ data frame contains data from a smaller number of distinct people.

Datasets merging:

```
minute_ci_temp <- merge(minute_calories, minute_intensities,  
                        by = c("id", "activity_minute"))
```

```
minute_cim_temp <- merge(minute_ci_temp, minute_METs,  
                        by = c("id", "activity_minute"))
```

```
minute_ciMs <- merge(minute_cim_temp, minute_steps,  
                    by = c("id", "activity_minute"))
```

```
head(minute_ciMs)
```

```
##           id      activity_minute calories intensity me_ts steps  
## 1 1503960366 4/12/2016 1:00:00 AM    0.9438          0    12     0  
## 2 1503960366 4/12/2016 1:00:00 PM    0.9438          0    12     0  
## 3 1503960366 4/12/2016 1:01:00 AM    2.6741          1    34    36  
## 4 1503960366 4/12/2016 1:01:00 PM    0.9438          0    12     0  
## 5 1503960366 4/12/2016 1:02:00 AM    2.0449          1    26     9  
## 6 1503960366 4/12/2016 1:02:00 PM    0.9438          0    12     0
```

I ended up with one dataset reflecting the minutes data, called ‘minute_ciMs’.

Fixing the timestep.

```
minute_ciMs$activity_minute=as.POSIXct(minute_ciMs$activity_minute, format="%m/%d/%Y %I:%M:%S %p", tz=Sys.timezone())  
minute_ciMs$time <- format(minute_ciMs$activity_minute, format = "%H:%M:%S")  
minute_ciMs$date <- format(minute_ciMs$activity_minute, format = "%m/%d/%y")
```

```
minute_sleep$date=as.POSIXct(minute_sleep$date, format="%m/%d/%Y %I:%M:%S %p", tz=Sys.timezone())  
minute_sleep$time <- format(minute_sleep$date, format = "%H:%M:%S")  
minute_sleep$date <- format(minute_sleep$date, format = "%m/%d/%y")
```

Weight and Heartrate Data Sets:

Importing Data:

```
weight <- read_csv(here("mturkfitbit_export_4.12.16-5.12.16",
                        "Fitabase Data 4.12.16-5.12.16",
                        "weightLogInfo_merged.csv"))

heartrate <- read_csv(here("mturkfitbit_export_4.12.16-5.12.16",
                          "Fitabase Data 4.12.16-5.12.16",
                          "heartrate_seconds_merged.csv"))
```

Data cleaning:

```
weight <- clean_names(weight)
heartrate <- clean_names(heartrate)

weight %>% skimr::skim_without_charts()
```

Table 31: Data summary

Name	Piped data
Number of rows	67
Number of columns	8
Column type frequency:	
character	1
logical	1
numeric	6
Group variables	None

Variable type: character

skim_variable	n_missing	complete_rate	min	max	empty	n_unique	whitespace
date	0	1	19	21	0	56	0

Variable type: logical

skim_variable	n_missing	complete_rate	mean	count
is_manual_report	0	1	0.61	TRU: 41, FAL: 26

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100
id	0	1.00	7.009282e+09	50322e+09	503960e+09	62181e+09	62181e+09	877689e+09	877689e+09
weight_kg	0	1.00	7.204000e+01	92000e+01	126000e+01	140000e+01	150000e+01	1505000e+01	1835000e+02
weight_pounds	0	1.00	1.588100e+01	270000e+01	159600e+01	253600e+01	277900e+01	2875000e+01	2943200e+02
fat	65	0.03	2.350000e+01	120000e+01	200000e+01	2275000e+01	2350000e+01	2425000e+01	2500000e+01
bmi	0	1.00	2.519000e+01	1070000e+01	2045000e+01	2396000e+01	2439000e+01	2556000e+01	2754000e+01
log_id	0	1.00	1.461772e+12	229948e+12	460444e+12	1261079e+12	1261802e+12	1262375e+12	1263098e+12

We can see that for ‘fat’ column most values are not complete (97.015%).

```
heartrate %>% skimr::skim_without_charts()
```

Table 35: Data summary

Name	Piped data
Number of rows	2483658
Number of columns	3
Column type frequency:	
character	1
numeric	2
Group variables	None

Variable type: character

skim_variable	n_missing	complete_rate	min	max	empty	n_unique	whitespace
time	0	1	19	21	0	961274	0

Variable type: numeric

skim_variable	n_missing	complete_rate	mean	sd	p0	p25	p50	p75	p100
id	0	1	5.513765e+09	50223761e+09	2022484408e+09	3388161845e+09	5553957448e+09	8962181068e+09	8877689391e+09
value	0	1	7.733000e+01	19.4	36	63	73	88	203

```
n_distinct(weight$id)
```

```
## [1] 8
```

```
n_distinct(heartrate$id)
```

```
## [1] 14
```

Unfortunately, the group of users from which the ‘weight’ data is taken is too small to draw any conclusions. Even the ‘heart rate’ data is based on a small sample.

Fixsing taamesteps:

```
weight$date=as.POSIXct(weight$date, format="%m/%d/%Y %I:%M:%S %p", tz=Sys.timezone())
weight$time <- format(weight$date, format = "%H:%M:%S")
weight$date <- format(weight$date, format = "%m/%d/%y")
```

```
heartrate$time=as.POSIXct(heartrate$time, format="%m/%d/%Y %I:%M:%S %p", tz=Sys.timezone())
heartrate$date <- format(heartrate$time, format = "%m/%d/%y")
heartrate$time_sec <- format(heartrate$time, format = "%H:%M:%S")
```

Summary:

Following data cleaning and merging operations, I've consolidated user activity information into three primary datasets, encompassing 33 unique users: 'daily_ais', 'hourly_cis' and 'minute_ciMs'. Sleep data requires independent analysis due to its user base, comprising 24 users with both daily and minute-level granularity. Heart rate data shows promise for limited integration into the overall analysis, despite its restricted user base. Contrarily, the weight dataset is currently deemed to have insufficient user representation for meaningful inclusion.

List of dataframes that will be used for futher analysis:

- daily_ais
- hourly_cis
- minute_ciMs
- sleep
- minute_sleep
- heartrate