

Be-Healthy

Pooja Rathee

2022-04-08

First we install required packages and go through it

```
library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.1 --

## v ggplot2 3.3.5    v purrr  0.3.4
## v tibble  3.1.6    v dplyr  1.0.8
## v tidyr   1.2.0    v stringr 1.4.0
## v readr   2.1.2    v forcats 0.5.1

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

Read csv file

```
daily_activity <- read_csv("dailyActivity_merged.csv")

## Rows: 940 Columns: 15
## -- Column specification -----
## Delimiter: ","
## chr  (1): ActivityDate
## dbl (14): Id, TotalSteps, TotalDistance, TrackerDistance, LoggedActivitiesDi...
##
## 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_day <- read.csv("sleepDay_merged.csv")
```

Take a look at the daily activity data

```
head(daily_activity)
```

```
## # A tibble: 6 x 15
##       Id ActivityDate TotalSteps TotalDistance TrackerDistance LoggedActivitie~
##       <dbl> <chr>         <dbl>         <dbl>         <dbl>         <dbl>
## 1  1.50e9 4/12/2016         13162          8.5           8.5           0
## 2  1.50e9 4/13/2016         10735          6.97          6.97          0
## 3  1.50e9 4/14/2016         10460          6.74          6.74          0
## 4  1.50e9 4/15/2016          9762          6.28          6.28          0
## 5  1.50e9 4/16/2016         12669          8.16          8.16          0
## 6  1.50e9 4/17/2016          9705          6.48          6.48          0
## # ... with 9 more variables: VeryActiveDistance <dbl>,
## #   ModeratelyActiveDistance <dbl>, LightActiveDistance <dbl>,
## #   SedentaryActiveDistance <dbl>, VeryActiveMinutes <dbl>,
## #   FairlyActiveMinutes <dbl>, LightlyActiveMinutes <dbl>,
## #   SedentaryMinutes <dbl>, Calories <dbl>
```

Identify all the columns in daily__activity

```
colnames(daily_activity)
```

```
## [1] "Id" "ActivityDate"
## [3] "TotalSteps" "TotalDistance"
## [5] "TrackerDistance" "LoggedActivitiesDistance"
## [7] "VeryActiveDistance" "ModeratelyActiveDistance"
## [9] "LightActiveDistance" "SedentaryActiveDistance"
## [11] "VeryActiveMinutes" "FairlyActiveMinutes"
## [13] "LightlyActiveMinutes" "SedentaryMinutes"
## [15] "Calories"
```

Take a look at the sleep data

```
head(sleep_day)
```

```
##       Id SleepDay TotalSleepRecords TotalMinutesAsleep
## 1 1503960366 4/12/2016 12:00:00 AM 1 327
## 2 1503960366 4/13/2016 12:00:00 AM 2 384
## 3 1503960366 4/15/2016 12:00:00 AM 1 412
## 4 1503960366 4/16/2016 12:00:00 AM 2 340
## 5 1503960366 4/17/2016 12:00:00 AM 1 700
## 6 1503960366 4/19/2016 12:00:00 AM 1 304
## TotalTimeInBed
## 1 346
## 2 407
## 3 442
## 4 367
## 5 712
## 6 320
```

Identify all the columns in sleep day

```
colnames(sleep_day)
```

```
## [1] "Id"          "SleepDay"      "TotalSleepRecords"  
## [4] "TotalMinutesAsleep" "TotalTimeInBed"
```

Understanding some summary statistics

```
n_distinct(daily_activity$Id)
```

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

```
n_distinct(sleep_day$Id)
```

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

Calculate the observations

```
nrow(daily_activity)
```

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

```
nrow(sleep_day)
```

```
## [1] 413
```

summary statistics of daily activity

```
daily_activity %>%  
  select(TotalSteps,  
         TotalDistance,  
         SedentaryMinutes) %>%  
  summary()
```

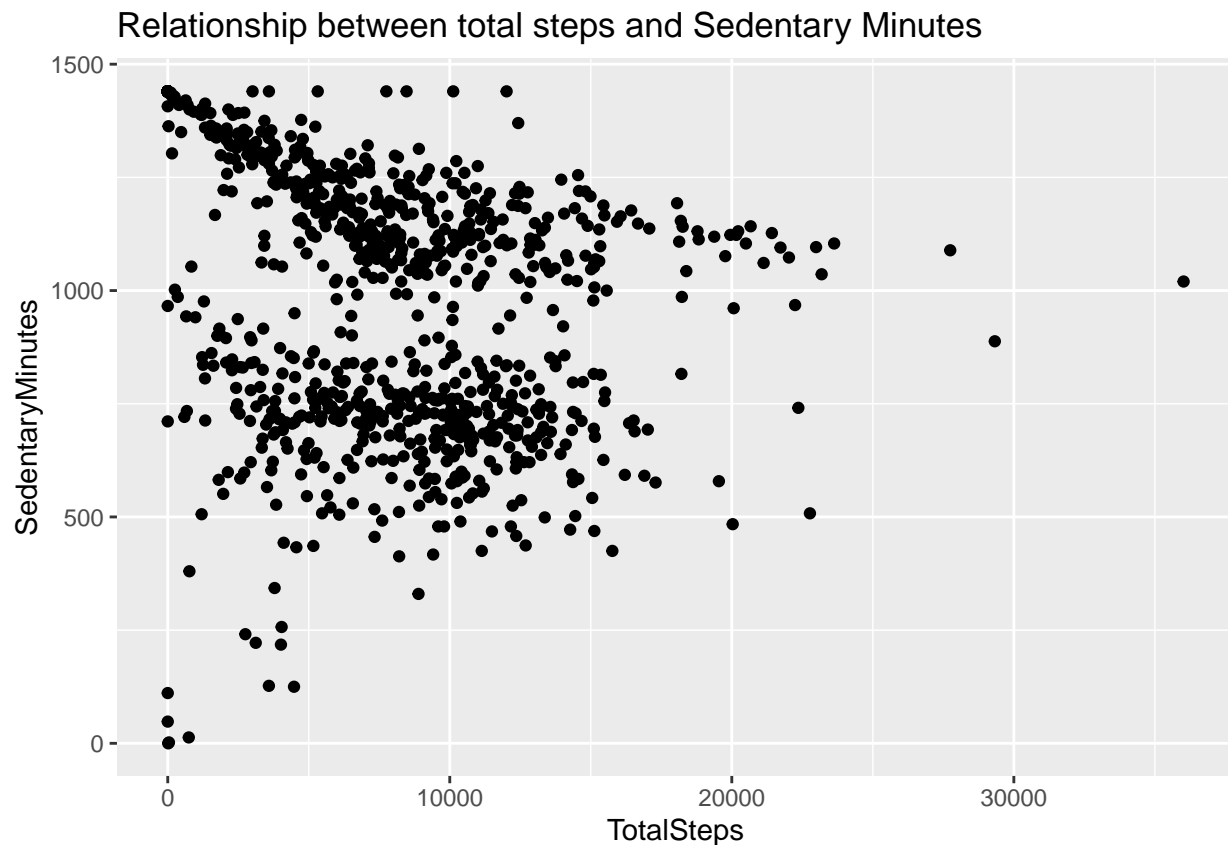
```
##   TotalSteps   TotalDistance   SedentaryMinutes  
##   Min.      :    0   Min.      : 0.000   Min.      :    0.0  
##   1st Qu.: 3790   1st Qu.: 2.620   1st Qu.: 729.8  
##   Median : 7406   Median : 5.245   Median :1057.5  
##   Mean    : 7638   Mean    : 5.490   Mean     : 991.2  
##   3rd Qu.:10727   3rd Qu.: 7.713   3rd Qu.:1229.5  
##   Max.    :36019   Max.    :28.030   Max.     :1440.0
```

For the sleep dataframe

```
sleep_day %>%  
  select(TotalSleepRecords,  
         TotalMinutesAsleep,  
         TotalTimeInBed) %>%  
  summary()
```

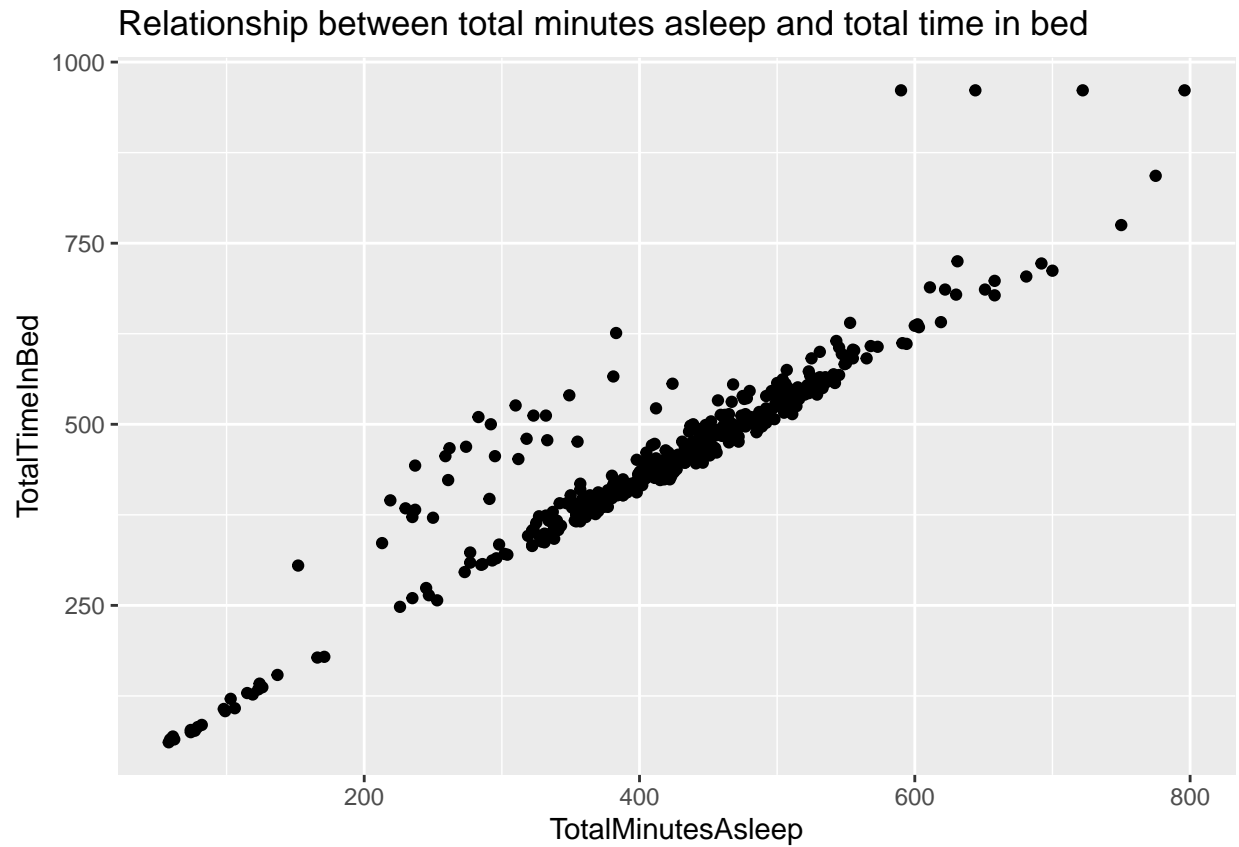
```
## TotalSleepRecords TotalMinutesAsleep TotalTimeInBed  
## Min. :1.000      Min. : 58.0      Min. : 61.0  
## 1st Qu.:1.000    1st Qu.:361.0    1st Qu.:403.0  
## Median :1.000    Median :433.0    Median :463.0  
## Mean :1.119      Mean :419.5      Mean :458.6  
## 3rd Qu.:1.000    3rd Qu.:490.0    3rd Qu.:526.0  
## Max. :3.000      Max. :796.0      Max. :961.0
```

```
ggplot(data=daily_activity, aes(x=TotalSteps, y=SedentaryMinutes)) + geom_point()+  
  labs(title = "Relationship between total steps and Sedentary Minutes")
```



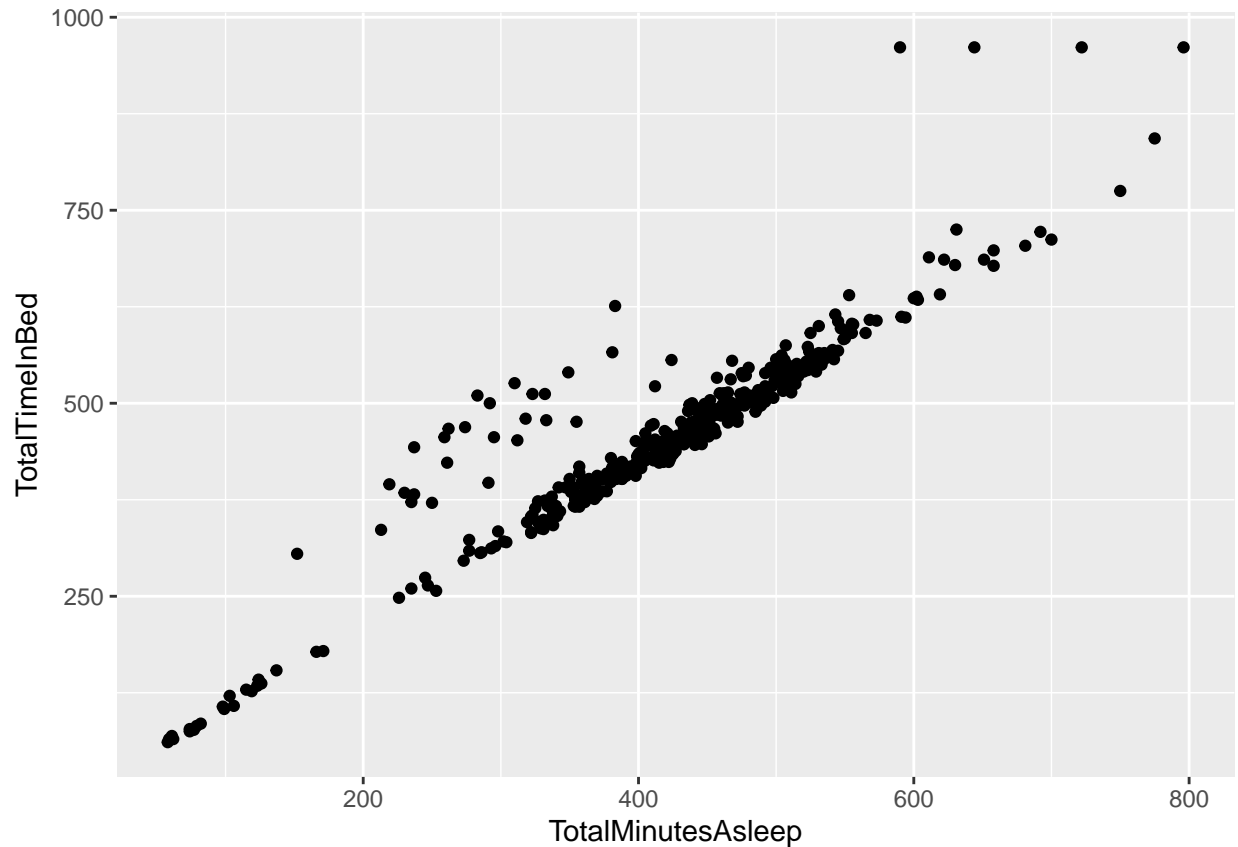
```
# plotting the graph for sleep_day
```

```
ggplot(data=sleep_day, aes(x=TotalMinutesAsleep, y=TotalTimeInBed)) + geom_point() +  
  labs(title = "Relationship between total minutes asleep and total time in bed")
```



What's the relationship between minutes asleep and time in bed? You might expect it to be almost completely linear - are there any unexpected trends?

```
ggplot(data=sleep_day, aes(x=TotalMinutesAsleep, y=TotalTimeInBed)) + geom_point()
```



What could these trends tell you about how to help market this product? Or areas where you might want to explore further?

Merging these two datasets together

```
combined_data <- merge(sleep_day, daily_activity, by="Id")
```

How many participants are there in data

```
n_distinct(combined_data$Id)
```

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

There were more participant Ids in the daily activity dataset that have been filtered out using merge.