

google capstone 2

Amal

2024-04-02

Install and load the tidyverse

```
options(repos = "https://cran.rstudio.com/")
install.packages('tidyverse')
```

```
## Installing package into 'C:/Users/AMD/AppData/Local/R/win-library/4.3'
## (as 'lib' is unspecified)
```

```
## package 'tidyverse' successfully unpacked and MD5 sums checked
##
```

```
## The downloaded binary packages are in
## C:\Users\AMD\AppData\Local\Temp\RtmpcTT8JG\downloaded_packages
```

```
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
```

```
## v dplyr      1.1.4      v readr      2.1.5
## v forcats    1.0.0      v stringr   1.5.1
## v ggplot2    3.5.0      v tibble    3.2.1
## v lubridate  1.9.3      v tidyr     1.3.1
## v purrr      1.0.2
```

```
## -- Conflicts ----- tidyverse_conflicts() --
```

```
## x dplyr::filter() masks stats::filter()
```

```
## x dplyr::lag()     masks stats::lag()
```

```
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

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

```
sleep_day <- read.csv("sleepDay_merged.csv")
```

```
head(daily_activity)
```

```
##           Id ActivityDate TotalSteps TotalDistance TrackerDistance
## 1 1503960366   4/12/2016     13162           8.50           8.50
## 2 1503960366   4/13/2016     10735           6.97           6.97
## 3 1503960366   4/14/2016     10460           6.74           6.74
## 4 1503960366   4/15/2016      9762           6.28           6.28
## 5 1503960366   4/16/2016     12669           8.16           8.16
## 6 1503960366   4/17/2016      9705           6.48           6.48
##   LoggedActivitiesDistance VeryActiveDistance ModeratelyActiveDistance
## 1                        0                1.88                   0.55
## 2                        0                1.57                   0.69
## 3                        0                2.44                   0.40
## 4                        0                2.14                   1.26
## 5                        0                2.71                   0.41
```

```
## 6          0          3.19          0.78
##   LightActiveDistance SedentaryActiveDistance VeryActiveMinutes
## 1          6.06          0          25
## 2          4.71          0          21
## 3          3.91          0          30
## 4          2.83          0          29
## 5          5.04          0          36
## 6          2.51          0          38
##   FairlyActiveMinutes LightlyActiveMinutes SedentaryMinutes Calories
## 1          13          328          728    1985
## 2          19          217          776    1797
## 3          11          181         1218    1776
## 4          34          209          726    1745
## 5          10          221          773    1863
## 6          20          164          539    1728
```

identify columns

```
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"
```

```
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
```

```
colnames(sleep_day)
```

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

How many unique participants are there in each dataframe? It looks like there may be more participants in the daily activity dataset than the sleep dataset.

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

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

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

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

How many observations are there in each dataframe?

```
nrow(daily_activity)
```

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

```
nrow(sleep_day)
```

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

What are some quick summary statistics we'd want to know about each data frame? For the daily activity dataframe:

```
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
```

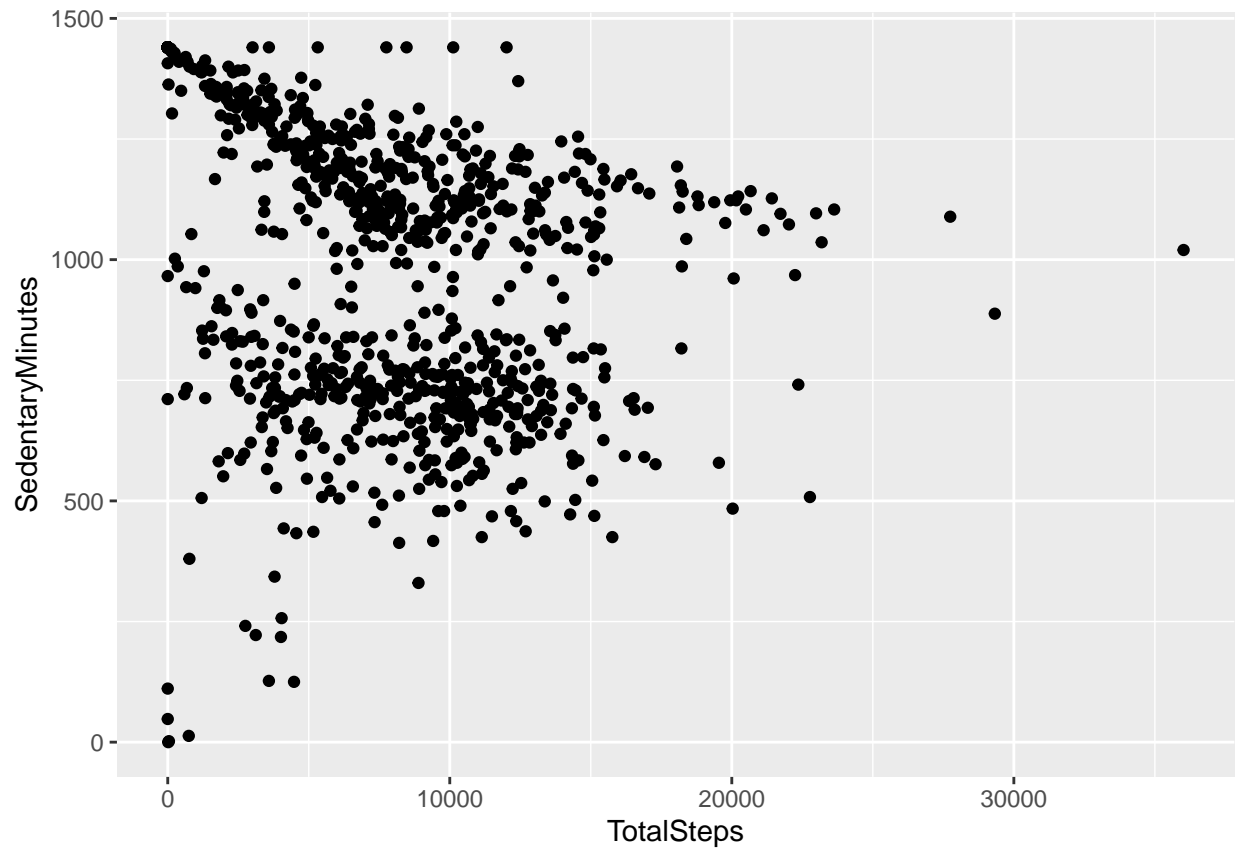
For the sleep dataframe:

```
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
```

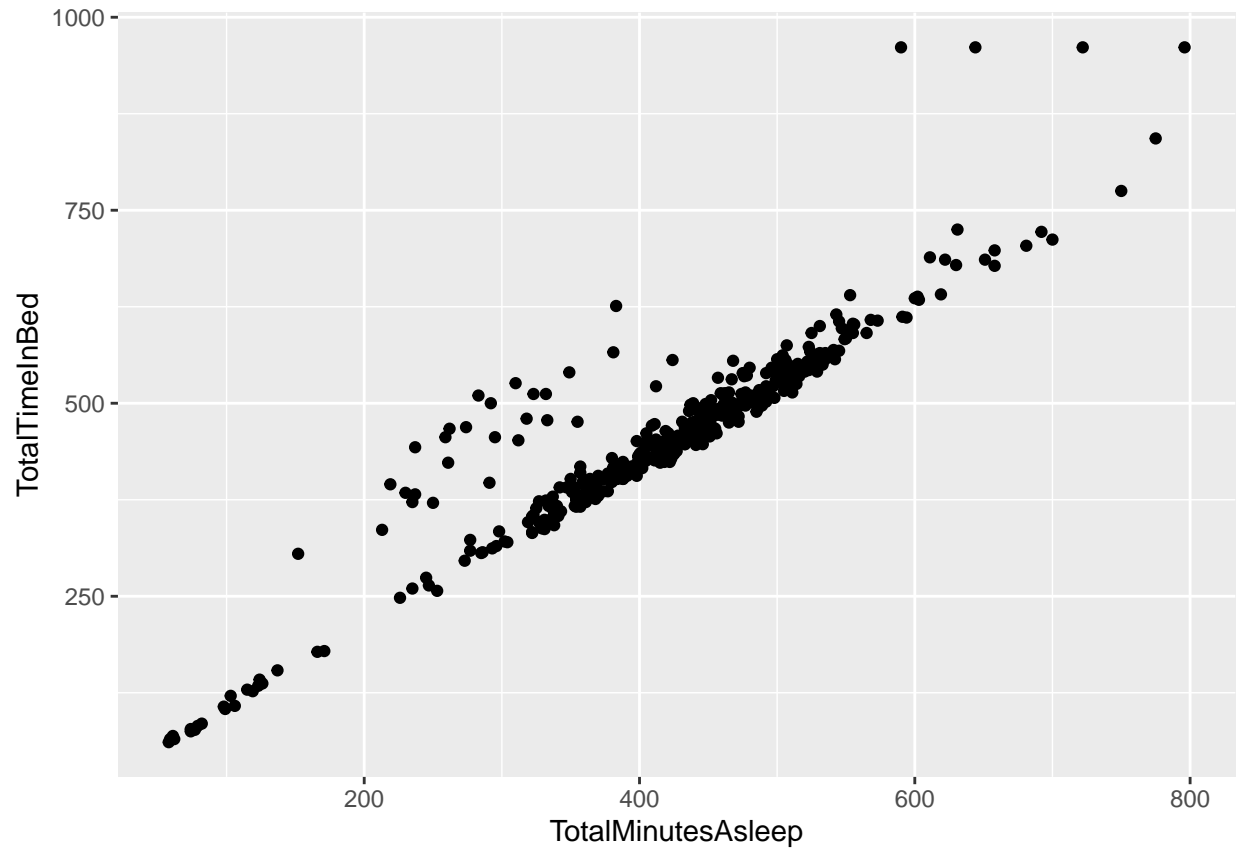
What's the relationship between steps taken in a day and sedentary minutes? How could this help inform the customer segments that we can market to? E.g. position this more as a way to get started in walking more? Or to measure steps that you're already taking?

```
ggplot(data=daily_activity, aes(x=TotalSteps, y=SedentaryMinutes)) + geom_point()
```



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()
```



merging two dataset together

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

Take a look at how many participants are in this data set.

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

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