

# Bellabeat Casestudy

## Introduction

### How Can a Wellness Technology Company Play It Smart?

#### Step 1: Ask

**Background** Bellabeat is a high tech company that manufactures health focused smart products. They offer different smart devices that collect data on activity, sleep, stress, and reproductive health to empower women with knowledge about their own health and habits.

The main focus of this case is to analyze smart devices fitness data and determine how it could help unlock new growth opportunities for Bellabeat. We will focus on one of Bellabeat's products: Bellabeat app.

The Bellabeat app provides users with health data related to their activity, sleep, stress, menstrual cycle, and mindfulness habits. This data can help users better understand their current habits and make healthy decisions. The Bellabeat app connects to their line of smart wellness products.

#### Key Stakeholders

- Urška Sršen Bellabeat cofounder and Chief Creative Officer
- Sando Mur Bellabeat cofounder and key member of Bellabeat executive team
- Bellabeat Marketing Analytics team

#### Business Task

Given the previous facts, the business task is defined as searching for user patterns of usage of their smart devices in order to gain insights that would later better orientate marketing decisions. So, in one phrase it would be:

How do our users use our smart devices?. Identify trends in how consumers use non Bellabeat smart devices to apply insights into Bellabeat's marketing strategy

#### Step 2: Prepare

**Dataset used** The data source used for this case study is **FitBit Fitness Tracker Data**. This dataset is stored in Kaggle and was made available through Mobius and generated by respondents to a distributed survey via Amazon Mechanical Turk between 03.12.2016 05.12.2016.

#### Accessibility and privacy of data

The data is licensed under CC0: Public Domain, waiving all of his or her rights to the work worldwide under copyright law, including all related and neighboring rights, to the extent by law. The work can be copied, modified, distributed and perform the work, even for commercial purposes, all without asking permission

#### Data organization and verification

The dataset is a collection of 18 .csv files. 15 in long format, 3 in wide format. The datasets consists of wide ranging information from activity metrics, calories, sleep records, metabolic equivalent of tasks (METs), heart rate and steps; in timeframes of seconds, minutes, hours and days

## Data limitations

The data has some limitations which could Undermine the results of the analysis Such limitations to take into consideration are:

- Missing demographics
- Small sample size
- Short time period of Data collection

## Step 3: Process

```
install.packages("tidyverse")
```

### Loading Libraries

```
## Installing package into '/cloud/lib/x86_64-pc-linux-gnu-library/4.3'  
## (as 'lib' is unspecified)
```

```
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ----- tidyverse 2.0.0 --  
## v dplyr      1.1.2      v readr      2.1.4  
## v forcats    1.0.0      v stringr   1.5.0  
## v ggplot2    3.4.3      v tibble    3.2.1  
## v lubridate  1.9.2      v tidyr     1.3.0  
## 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
```

```
library(lubridate)  
library(dplyr)  
library(ggplot2)  
library(tidyr)
```

### Importing datasets

For this project, I will use FitBit Fitness Tracker Data

```
install.packages("here")
```

```
## Installing package into '/cloud/lib/x86_64-pc-linux-gnu-library/4.3'  
## (as 'lib' is unspecified)
```

```
library(readr)
```

```
activity <- read_csv("Fitabase Data 4.12.16-5.12.16/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.
```

```

calories <- read_csv("Fitabase Data 4.12.16-5.12.16/hourlyCalories_merged.csv")

## Rows: 22099 Columns: 3
## -- Column specification -----
## Delimiter: ","
## chr (1): ActivityHour
## dbl (2): Id, Calories
##
## 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.
intensities <- read_csv("Fitabase Data 4.12.16-5.12.16/hourlyIntensities_merged.csv")

## Rows: 22099 Columns: 4
## -- Column specification -----
## Delimiter: ","
## chr (1): ActivityHour
## dbl (3): Id, TotalIntensity, AverageIntensity
##
## 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("Fitabase Data 4.12.16-5.12.16/sleepDay_merged.csv")
weight <- read_csv("Fitabase Data 4.12.16-5.12.16/weightLogInfo_merged.csv")

```

I already checked the data in Google Sheets. I just need to make sure that everything were imported correctly by using View() and head() functions.

```

head(activity)

## # A tibble: 6 x 15
##       Id ActivityDate TotalSteps TotalDistance TrackerDistance
##       <dbl> <chr>         <dbl>         <dbl>         <dbl>
## 1 1503960366 4/12/2016         13162           8.5           8.5
## 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
## # i 10 more variables: LoggedActivitiesDistance <dbl>,
## #   VeryActiveDistance <dbl>, ModeratelyActiveDistance <dbl>,
## #   LightActiveDistance <dbl>, SedentaryActiveDistance <dbl>,
## #   VeryActiveMinutes <dbl>, FairlyActiveMinutes <dbl>,
## #   LightlyActiveMinutes <dbl>, SedentaryMinutes <dbl>, Calories <dbl>

```

I spotted some problems with the timestamp data. So before analysis, I need to convert it to date time format and split to date and time.

```

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

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

# activity

```

```

activity$ActivityDate=as.POSIXct(activity$ActivityDate, format="%m/%d/%Y", tz=Sys.timezone())
activity$date <- format(activity$ActivityDate, format = "%m/%d/%y")
# sleep
sleep$SleepDay=as.POSIXct(sleep$SleepDay, format="%m/%d/%Y %I:%M:%S %p", tz=Sys.timezone())
sleep$date <- format(sleep$SleepDay, format = "%m/%d/%y")

```

Warning message in system("timedatectl", intern = TRUE): "running command 'timedatectl' had status 1"

Now that everything is ready, I can start exploring data sets.

## Exploring and summarizing data

```
n_distinct(activity$Id)
```

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

```
n_distinct(calories$Id)
```

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

```
n_distinct(intensities$Id)
```

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

```
n_distinct(sleep$Id)
```

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

```
n_distinct(weight$Id)
```

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

This information tells us about number participants in each data sets.

There is 33 participants in the activity, calories and intensities data sets, 24 in the sleep and only 8 in the weight data set. 8 participants is not significant to make any recommendations and conclusions based on this data.

Let's have a look at summary statistics of the data sets:

```

# activity
activity %>%
  select(TotalSteps,
         TotalDistance,
         SedentaryMinutes, Calories) %>%
  summary()

```

```

##      TotalSteps      TotalDistance      SedentaryMinutes      Calories
##  Min.       :    0      Min.       : 0.000      Min.       :  0.0      Min.       :    0
##  1st Qu.: 3790      1st Qu.:  2.620      1st Qu.: 729.8      1st Qu.: 1828
##  Median : 7406      Median :  5.245      Median :1057.5      Median : 2134
##  Mean   : 7638      Mean   :  5.490      Mean   : 991.2      Mean   : 2304
##  3rd Qu.:10727      3rd Qu.:  7.713      3rd Qu.:1229.5      3rd Qu.: 2793
##  Max.   :36019      Max.   :28.030      Max.   :1440.0      Max.   : 4900

```

```

# explore num of active minutes per category
activity %>%
  select(VeryActiveMinutes, FairlyActiveMinutes, LightlyActiveMinutes) %>%
  summary()

```

```
## VeryActiveMinutes FairlyActiveMinutes LightlyActiveMinutes
## Min. : 0.00 Min. : 0.00 Min. : 0.0
## 1st Qu.: 0.00 1st Qu.: 0.00 1st Qu.:127.0
## Median : 4.00 Median : 6.00 Median :199.0
## Mean : 21.16 Mean : 13.56 Mean :192.8
## 3rd Qu.: 32.00 3rd Qu.: 19.00 3rd Qu.:264.0
## Max. :210.00 Max. :143.00 Max. :518.0
```

```
# calories
calories %>%
  select(Calories) %>%
  summary()
```

```
##      Calories
## Min. : 42.00
## 1st Qu.: 63.00
## Median : 83.00
## Mean : 97.39
## 3rd Qu.:108.00
## Max. :948.00
```

```
# sleep
sleep %>%
  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
```

```
# weight
weight %>%
  select(WeightKg, BMI) %>%
  summary()
```

```
##      WeightKg      BMI
## Min. : 52.60 Min. :21.45
## 1st Qu.: 61.40 1st Qu.:23.96
## Median : 62.50 Median :24.39
## Mean : 72.04 Mean :25.19
## 3rd Qu.: 85.05 3rd Qu.:25.56
## Max. :133.50 Max. :47.54
```

### Some interesting discoveries from this summary:

- Average sedentary time is 991 minutes or 16 hours. Definately needs to be reduced!
- The majority of the participants are lightly active.
- On the average, participants sleep 1 time for 7 hours.
- Average total steps per day are 7638 which a little bit less for having health benefits for according to the CDC research. They found that taking 8,000 steps per day was associated with a 51% lower risk

for all-cause mortality (or death from all causes). Taking 12,000 steps per day was associated with a 65% lower risk compared with taking 4,000 steps.

## Merging data

Before beginning to visualize the data, I need to merge two data sets. I'm going to merge (inner join) activity and sleep on columns Id and date (that I previously created after converting data to date time format).

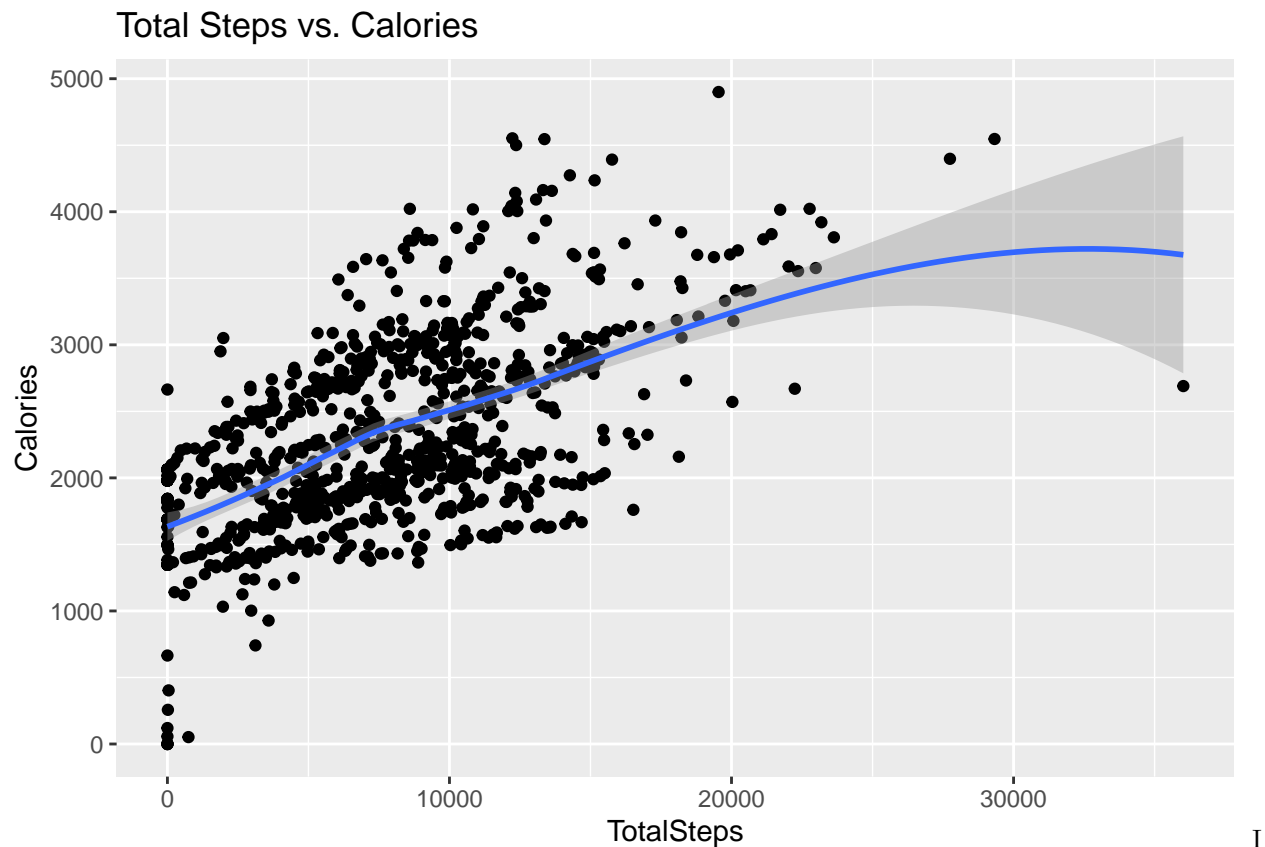
```
merged_data <- merge(sleep, activity, by=c('Id', 'date'))
head(merged_data)
```

```
##           Id      date SleepDay TotalSleepRecords TotalMinutesAsleep
## 1 1503960366 04/12/16 2016-04-12                1                327
## 2 1503960366 04/13/16 2016-04-13                2                384
## 3 1503960366 04/15/16 2016-04-15                1                412
## 4 1503960366 04/16/16 2016-04-16                2                340
## 5 1503960366 04/17/16 2016-04-17                1                700
## 6 1503960366 04/19/16 2016-04-19                1                304
##   TotalTimeInBed ActivityDate TotalSteps TotalDistance TrackerDistance
## 1             346   2016-04-12     13162           8.50           8.50
## 2             407   2016-04-13     10735           6.97           6.97
## 3             442   2016-04-15      9762           6.28           6.28
## 4             367   2016-04-16     12669           8.16           8.16
## 5             712   2016-04-17      9705           6.48           6.48
## 6             320   2016-04-19     15506           9.88           9.88
##   LoggedActivitiesDistance VeryActiveDistance ModeratelyActiveDistance
## 1                      0                1.88                0.55
## 2                      0                1.57                0.69
## 3                      0                2.14                1.26
## 4                      0                2.71                0.41
## 5                      0                3.19                0.78
## 6                      0                3.53                1.32
##   LightActiveDistance SedentaryActiveDistance VeryActiveMinutes
## 1                6.06                0                25
## 2                4.71                0                21
## 3                2.83                0                29
## 4                5.04                0                36
## 5                2.51                0                38
## 6                5.03                0                50
##   FairlyActiveMinutes LightlyActiveMinutes SedentaryMinutes Calories
## 1                13                328                728     1985
## 2                19                217                776     1797
## 3                34                209                726     1745
## 4                10                221                773     1863
## 5                20                164                539     1728
## 6                31                264                775     2035
```

## Visualization

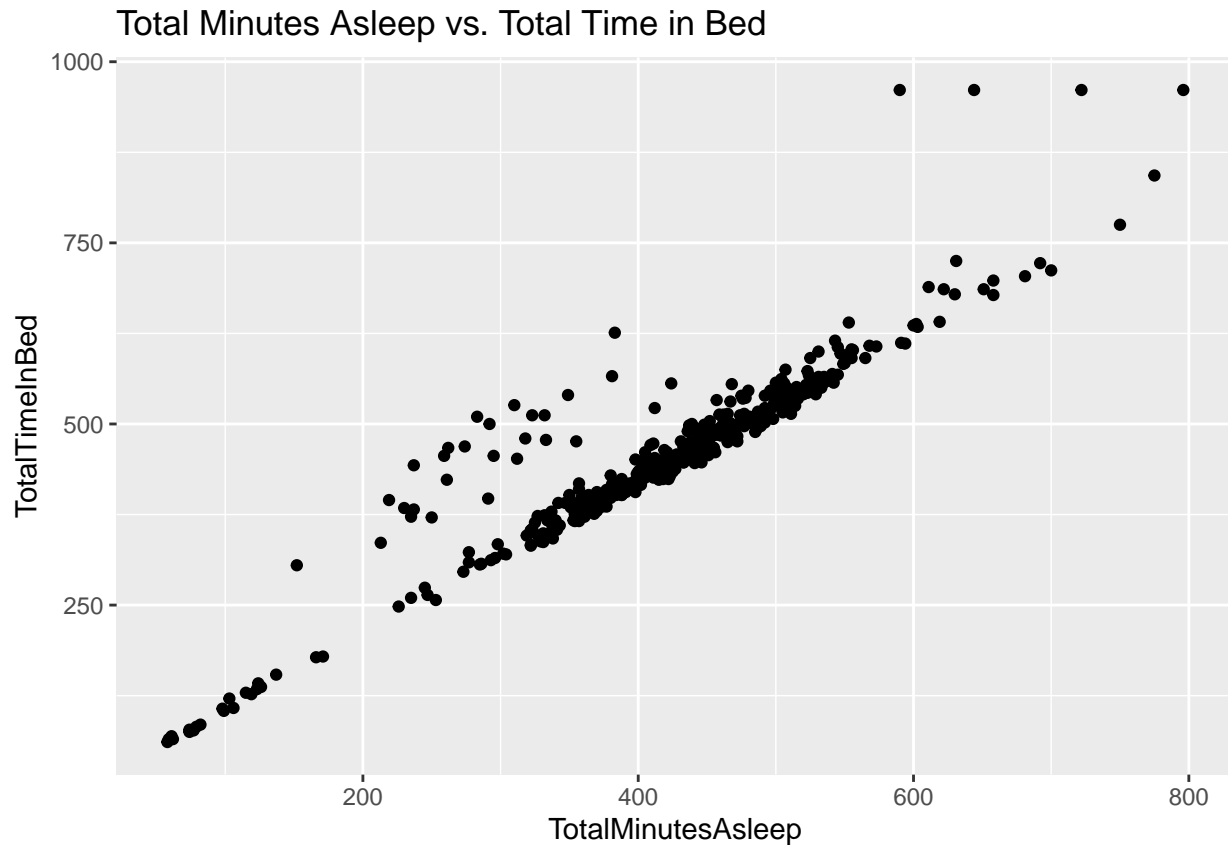
```
ggplot(data=activity, aes(x=TotalSteps, y=Calories)) +
  geom_point() + geom_smooth() + labs(title="Total Steps vs. Calories")
```

```
## `geom_smooth()` using method = 'loess' and formula = 'y ~ x'
```



see positive correlation here between Total Steps and Calories, which is obvious - the more active we are, the more calories we burn.

```
ggplot(data=sleep, aes(x=TotalMinutesAsleep, y=TotalTimeInBed)) +  
  geom_point() + labs(title="Total Minutes Asleep vs. Total Time in Bed")
```



The relationship between Total Minutes Asleep and Total Time in Bed looks linear. So if the Bellabeat users want to improve their sleep, we should consider using notification to go to sleep.

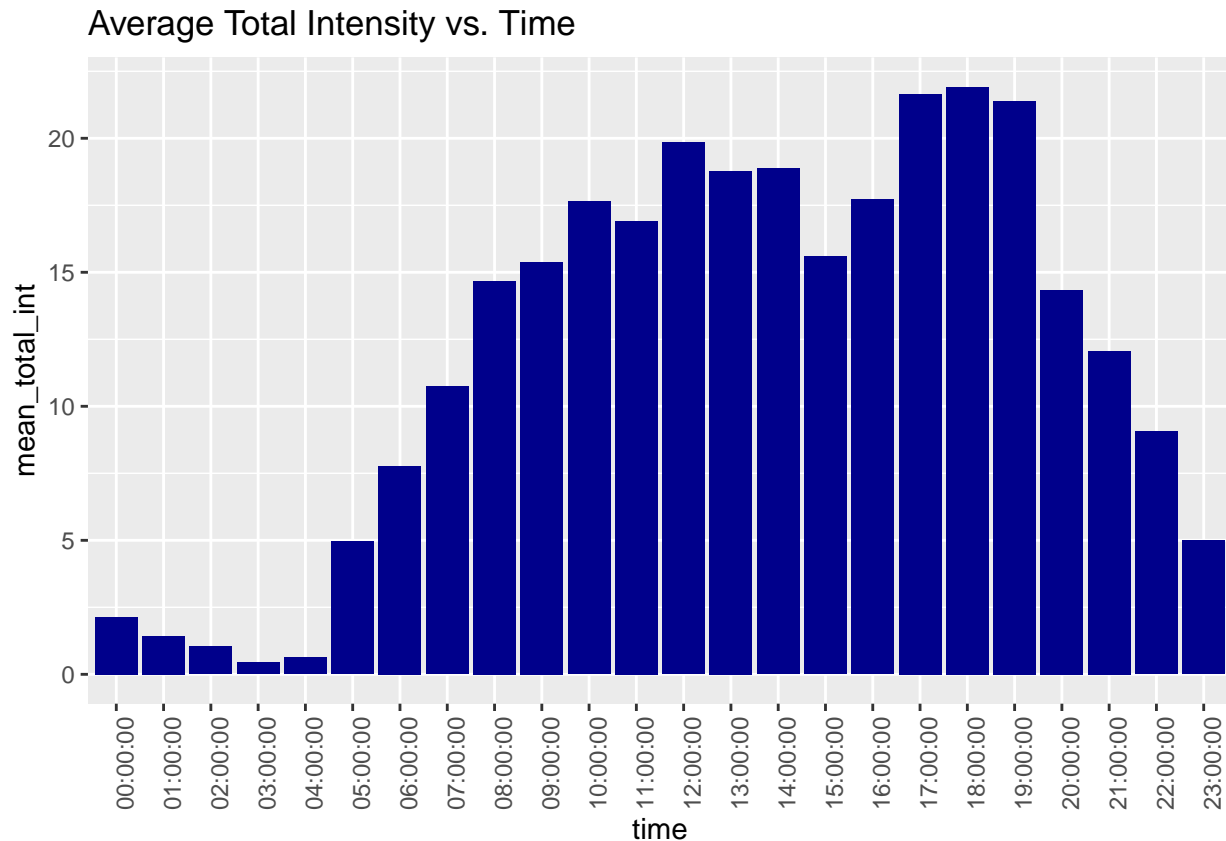
Let's look at intensities data over time (hourly).

```
int_new <- intensities %>%
  group_by(time) %>%
  drop_na() %>%
  summarise(mean_total_int = mean(TotalIntensity))

ggplot(data=int_new, aes(x=time, y=mean_total_int)) + geom_histogram(stat = "identity", fill='darkblue') +
  theme(axis.text.x = element_text(angle = 90)) +
  labs(title="Average Total Intensity vs. Time")
```

```
## Warning in geom_histogram(stat = "identity", fill = "darkblue"): Ignoring
## unknown parameters: `binwidth`, `bins`, and `pad`
```





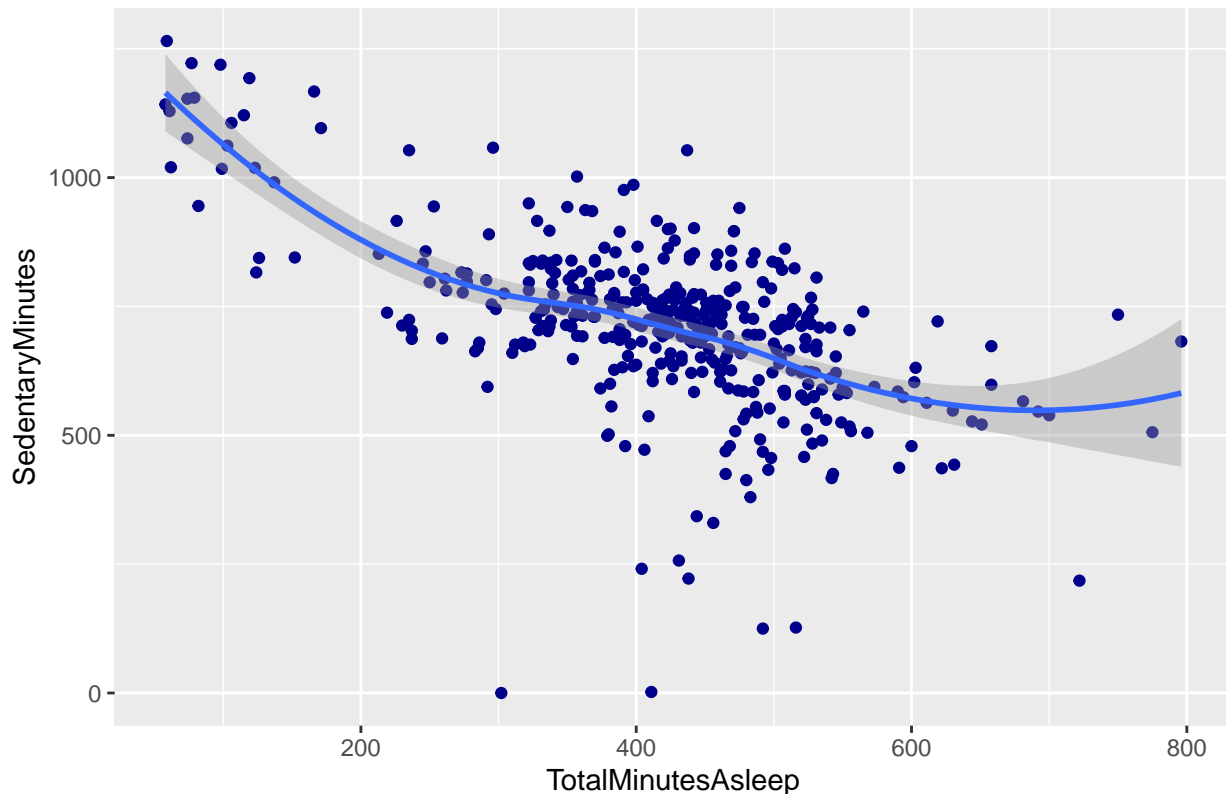
- After visualizing Total Intensity hourly, I found out that people are more active between 5 am and 10pm.
- Most activity happens between 5 pm and 7 pm - I suppose, that people go to a gym or for a walk after finishing work. We can use this time in the Bellabeat app to remind and motivate users to go for a run or walk.

Let's look at the relationship between Total Minutes Asleep and Sedentry Minutes.

```
ggplot(data=merged_data, aes(x=TotalMinutesAsleep, y=SedentryMinutes)) +
  geom_point(color='darkblue') + geom_smooth() +
  labs(title="Minutes Asleep vs. Sedentry Minutes")
```

```
## `geom_smooth()` using method = 'loess' and formula = 'y ~ x'
```

## Minutes Asleep vs. Sedentary Minutes



- Here we can clearly see the negative relationship between Sedentary Minutes and Sleep time.
- As an idea: if Bellabeat users want to improve their sleep, Bellabeat app can recommend reducing sedentary time.
- Keep in mind that we need to support this insights with more data, because correlation between some data doesn't mean causation.

### Summarizing recommendations for the business

As we already know, collecting data on activity, sleep, stress, and reproductive health has allowed Bellabeat to empower women with knowledge about their own health and habits. Since it was founded in 2013, Bellabeat has grown rapidly and quickly positioned itself as a tech-driven wellness company for women.

After analyzing FitBit Fitness Tracker Data, I found some insights that would help influence Bellabeat marketing strategy.

#### Target audience

Women who work full-time jobs (according to the hourly intensity data) and spend a lot of time at the computer/in a meeting/ focused on work they are doing (according to the sedentary time data).

These women do some light activity to stay healthy (according to the activity type analysis). Even though they need to improve their everyday activity to have health benefits. They might need some knowledge about developing healthy habits or motivation to keep going.

- As there is no gender information about the participants, I assumed that all genders were presented and balanced in this data set.

### The key message for the Bellabeat online campaign

The Bellabeat app is not just another fitness activity app. It's a guide (a friend) who empowers women to balance full personal and professional life and healthy habits and routines by educating and motivating them through daily app recommendations.

### **Ideas for the Bellabeat app**

1. Average total steps per day are 7638 which is a little bit less for having health benefits for according to the CDC research. They found that taking 8,000 steps per day was associated with a 51% lower risk for all-cause mortality (or death from all causes). Taking 12,000 steps per day was associated with a 65% lower risk compared with taking 4,000 steps. Bellabeat can encourage people to take at least 8 000 explaining the benefits for their health.
2. If users want to lose weight, it's probably a good idea to control daily calorie consumption. Bellabeat can suggest some ideas for low-calorie lunch and dinner.
3. If users want to improve their sleep, Bellabeat should consider using app notifications to go to bed.
4. Most activity happens between 5 pm and 7 pm - I suppose, that people go to a gym or for a walk after finishing work. Bellabeat can use this time to remind and motivate users to go for a run or walk.
5. As an idea: if users want to improve their sleep, the Bellabeat app can recommend reducing sedentary time.