

BellaBeat Case Study

Zubair Matani

11/5/2021



Figure 1: bellabeat logo

Introduction

This is a Google Data Analytics Capstone project. The purpose of this project is to learn, understand and apply the concepts we learnt in the Google Analytics Course.

Goals and Business Statement

This case study tasks us with assisting a wearable fitness technology company, BellaBeat, improve their marketing strategies for their products by investigating customer activity with other fitness trackers like FitBit.

We will look at datasets to find the following:

- How customers use fitness trackers in their everyday life?
- What features are most popularly used?
- Which of these features does BellBeat already have, and how can we improve our marketing skills for those?
- What additional features can BellBeat introduce to add more customers?

Data Usage

What Data are We Using?

The data provided to us by BellBeats is <https://www.kaggle.com/arashnic/fitbit>. This data website ranges from their daily activities, to their steps to their heart rate, calories intake and much more. All of the data is stored in different .csv files which we will be importing to analyze and support our statements.

Loading CSV Files

The following data sets will be used:

- Daily Activity
- Daily Calories
- Daily Sleep
- Weight Log Info
- Daily intensities

```

dailyActivity <- read.csv("dailyActivity_merged.csv")
dailyCalories <- read.csv("dailyCalories_merged.csv")
sleepDay <- read.csv("sleepDay_merged.csv")
dailyIntensities <- read.csv("dailyIntensities_merged.csv")
weightLog <- read.csv("weightLogInfo_merged.csv")

```

Exploring the Tables

For each of the tables we have decided to work our analysis with, we will take a closer look at them using the `head()`, `glimpse()` and `colnames()` function. This would allow us to look at the first six values of each table and see each table with it's distributed columns respectively. `### Daily Activities`

```
head(dailyActivity)
```

```

##           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

```

```
glimpse(dailyActivity)
```

```

## Rows: 940
## Columns: 15
## $ Id <dbl> 1503960366, 1503960366, 1503960366, 150396036~
## $ ActivityDate <chr> "4/12/2016", "4/13/2016", "4/14/2016", "4/15/~
## $ TotalSteps <int> 13162, 10735, 10460, 9762, 12669, 9705, 13019~
## $ TotalDistance <dbl> 8.50, 6.97, 6.74, 6.28, 8.16, 6.48, 8.59, 9.8~
## $ TrackerDistance <dbl> 8.50, 6.97, 6.74, 6.28, 8.16, 6.48, 8.59, 9.8~
## $ LoggedActivitiesDistance <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ VeryActiveDistance <dbl> 1.88, 1.57, 2.44, 2.14, 2.71, 3.19, 3.25, 3.5~
## $ ModeratelyActiveDistance <dbl> 0.55, 0.69, 0.40, 1.26, 0.41, 0.78, 0.64, 1.3~

```

```
## $ LightActiveDistance      <dbl> 6.06, 4.71, 3.91, 2.83, 5.04, 2.51, 4.71, 5.0~
## $ SedentaryActiveDistance  <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ VeryActiveMinutes       <int> 25, 21, 30, 29, 36, 38, 42, 50, 28, 19, 66, 4~
## $ FairlyActiveMinutes     <int> 13, 19, 11, 34, 10, 20, 16, 31, 12, 8, 27, 21~
## $ LightlyActiveMinutes    <int> 328, 217, 181, 209, 221, 164, 233, 264, 205, ~
## $ SedentaryMinutes        <int> 728, 776, 1218, 726, 773, 539, 1149, 775, 818~
## $ Calories                <int> 1985, 1797, 1776, 1745, 1863, 1728, 1921, 203~
```

```
colnames(dailyActivity)
```

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

Daily Carlories

```
head(dailyCalories)
```

```
##           Id ActivityDay Calories
## 1 1503960366  4/12/2016    1985
## 2 1503960366  4/13/2016    1797
## 3 1503960366  4/14/2016    1776
## 4 1503960366  4/15/2016    1745
## 5 1503960366  4/16/2016    1863
## 6 1503960366  4/17/2016    1728
```

```
glimpse(dailyCalories)
```

```
## Rows: 940
## Columns: 3
## $ Id      <dbl> 1503960366, 1503960366, 1503960366, 1503960366, 1503960366~
## $ ActivityDay <chr> "4/12/2016", "4/13/2016", "4/14/2016", "4/15/2016", "4/16/~
## $ Calories  <int> 1985, 1797, 1776, 1745, 1863, 1728, 1921, 2035, 1786, 1775~
```

```
colnames(dailyCalories)
```

```
## [1] "Id"          "ActivityDay" "Calories"
```

Daily Intensities

```
head(dailyIntensities)
```

```
##           Id ActivityDay SedentaryMinutes LightlyActiveMinutes
## 1 1503960366  4/12/2016           728           328
## 2 1503960366  4/13/2016           776           217
## 3 1503960366  4/14/2016          1218           181
## 4 1503960366  4/15/2016           726           209
## 5 1503960366  4/16/2016           773           221
## 6 1503960366  4/17/2016           539           164
## FairlyActiveMinutes VeryActiveMinutes SedentaryActiveDistance
## 1              13              25              0
```

```
## 2          19          21          0
## 3          11          30          0
## 4          34          29          0
## 5          10          36          0
## 6          20          38          0
##   LightActiveDistance ModeratelyActiveDistance VeryActiveDistance
## 1          6.06          0.55          1.88
## 2          4.71          0.69          1.57
## 3          3.91          0.40          2.44
## 4          2.83          1.26          2.14
## 5          5.04          0.41          2.71
## 6          2.51          0.78          3.19
```

```
glimpse(dailyIntensities)
```

```
## Rows: 940
## Columns: 10
## $ Id <dbl> 1503960366, 1503960366, 1503960366, 150396036~
## $ ActivityDay <chr> "4/12/2016", "4/13/2016", "4/14/2016", "4/15/~
## $ SedentaryMinutes <int> 728, 776, 1218, 726, 773, 539, 1149, 775, 818~
## $ LightlyActiveMinutes <int> 328, 217, 181, 209, 221, 164, 233, 264, 205, ~
## $ FairlyActiveMinutes <int> 13, 19, 11, 34, 10, 20, 16, 31, 12, 8, 27, 21~
## $ VeryActiveMinutes <int> 25, 21, 30, 29, 36, 38, 42, 50, 28, 19, 66, 4~
## $ SedentaryActiveDistance <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ LightActiveDistance <dbl> 6.06, 4.71, 3.91, 2.83, 5.04, 2.51, 4.71, 5.0~
## $ ModeratelyActiveDistance <dbl> 0.55, 0.69, 0.40, 1.26, 0.41, 0.78, 0.64, 1.3~
## $ VeryActiveDistance <dbl> 1.88, 1.57, 2.44, 2.14, 2.71, 3.19, 3.25, 3.5~
```

```
colnames(dailyIntensities)
```

```
## [1] "Id" "ActivityDay"
## [3] "SedentaryMinutes" "LightlyActiveMinutes"
## [5] "FairlyActiveMinutes" "VeryActiveMinutes"
## [7] "SedentaryActiveDistance" "LightActiveDistance"
## [9] "ModeratelyActiveDistance" "VeryActiveDistance"
```

Sleep Day

```
head(sleepDay)
```

```
##           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
```

```
glimpse(sleepDay)
```

```
## Rows: 413
## Columns: 5
## $ Id          <dbl> 1503960366, 1503960366, 1503960366, 1503960366, 150~
## $ SleepDay    <chr> "4/12/2016 12:00:00 AM", "4/13/2016 12:00:00 AM", "~
## $ TotalSleepRecords <int> 1, 2, 1, 2, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, ~
## $ TotalMinutesAsleep <int> 327, 384, 412, 340, 700, 304, 360, 325, 361, 430, 2~
## $ TotalTimeInBed   <int> 346, 407, 442, 367, 712, 320, 377, 364, 384, 449, 3~
```

```
colnames(sleepDay)
```

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

Weight Log

```
head(weightLog)
```

```
##           Id           Date WeightKg WeightPounds Fat  BMI
## 1 1503960366 5/2/2016 11:59:59 PM    52.6    115.9631 22 22.65
## 2 1503960366 5/3/2016 11:59:59 PM    52.6    115.9631 NA 22.65
## 3 1927972279 4/13/2016 1:08:52 AM   133.5    294.3171 NA 47.54
## 4 2873212765 4/21/2016 11:59:59 PM    56.7    125.0021 NA 21.45
## 5 2873212765 5/12/2016 11:59:59 PM    57.3    126.3249 NA 21.69
## 6 4319703577 4/17/2016 11:59:59 PM    72.4    159.6147 25 27.45
##   IsManualReport   LogId
## 1             True 1.462234e+12
## 2             True 1.462320e+12
## 3             False 1.460510e+12
## 4             True 1.461283e+12
## 5             True 1.463098e+12
## 6             True 1.460938e+12
```

```
glimpse(weightLog)
```

```
## Rows: 67
## Columns: 8
## $ Id          <dbl> 1503960366, 1503960366, 1927972279, 2873212765, 2873212~
## $ Date        <chr> "5/2/2016 11:59:59 PM", "5/3/2016 11:59:59 PM", "4/13/2~
## $ WeightKg    <dbl> 52.6, 52.6, 133.5, 56.7, 57.3, 72.4, 72.3, 69.7, 70.3, ~
## $ WeightPounds <dbl> 115.9631, 115.9631, 294.3171, 125.0021, 126.3249, 159.6~
## $ Fat         <int> 22, NA, NA, NA, NA, 25, NA, NA, NA, NA, NA, NA, NA, ~
## $ BMI         <dbl> 22.65, 22.65, 47.54, 21.45, 21.69, 27.45, 27.38, 27.25, ~
## $ IsManualReport <chr> "True", "True", "False", "True", "True", "True", "True"~
## $ LogId       <dbl> 1.462234e+12, 1.462320e+12, 1.460510e+12, 1.461283e+12, ~
```

```
colnames(weightLog)
```

```
## [1] "Id"          "Date"          "WeightKg"      "WeightPounds"
## [5] "Fat"         "BMI"           "IsManualReport" "LogId"
```

Short Summary

Using the `glimpse` and the `colnames` functions, it is easily noticeable that the ID column is common all 5 data sets in this analysis.

The daily activity table gives us a hint that it contains values for calories and intensities as well which would allow us to extract data only using the ID column.

In order to shorten our summary table we will create a new data frame selecting only ID, ActivityDate and Calories.

```
dailyActivityNewFrame <- dailyActivity %>%  
  select(Id, ActivityDate, Calories)  
  
head(dailyActivityNewFrame)
```

```
##           Id ActivityDate Calories  
## 1 1503960366    4/12/2016    1985  
## 2 1503960366    4/13/2016    1797  
## 3 1503960366    4/14/2016    1776  
## 4 1503960366    4/15/2016    1745  
## 5 1503960366    4/16/2016    1863  
## 6 1503960366    4/17/2016    1728
```

To make sure the new data frame we have just created has the correct number of rows, we will use an SQL query to check the number of rows.

```
sqlCheck <- sqldf('SELECT * FROM dailyActivityNewFrame INTERSECT SELECT * FROM dailyCalories')  
head(sqlCheck)
```

```
##           Id ActivityDate Calories  
## 1 1503960366    4/12/2016    1985  
## 2 1503960366    4/13/2016    1797  
## 3 1503960366    4/14/2016    1776  
## 4 1503960366    4/15/2016    1745  
## 5 1503960366    4/16/2016    1863  
## 6 1503960366    4/17/2016    1728
```

```
nrow(sqlCheck)
```

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

The nrow() function shows 940 shows which is the same as we noticed earlier hence our data check is verified and we are good to go ahead. ## Analysis For the analysis stage, we will consider distinct data from the tables which would allow us to analyze data for each ID and not the same repititative ones.

Repitative vs Distinct Rows

```
n_distinct(dailyActivity$Id)
```

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

```
nrow(dailyActivity)
```

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

```
n_distinct(sleepDay$Id)
```

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

```
nrow(sleepDay)
```

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

```
n_distinct(weightLog$Id)
```

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

```
nrow(weightLog)
```

```
## [1] 67
```

Quick Statistics

Daily Activity

```
dailyActivity %>%  
  select(TotalSteps, TotalDistance, SedentaryMinutes, VeryActiveMinutes) %>%  
  summary()
```

```
##      TotalSteps      TotalDistance      SedentaryMinutes      VeryActiveMinutes  
##  Min.       :    0      Min.       : 0.000      Min.       :    0.0      Min.       :    0.00  
## 1st Qu.: 3790      1st Qu.: 2.620      1st Qu.: 729.8      1st Qu.:    0.00  
## Median : 7406      Median : 5.245      Median :1057.5      Median :    4.00  
## Mean   : 7638      Mean   : 5.490      Mean    : 991.2      Mean    :   21.16  
## 3rd Qu.:10727      3rd Qu.: 7.713      3rd Qu.:1229.5      3rd Qu.:   32.00  
## Max.    :36019      Max.    :28.030      Max.     :1440.0      Max.     :  210.00
```

Sleep

```
sleepDay %>%  
  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 Log

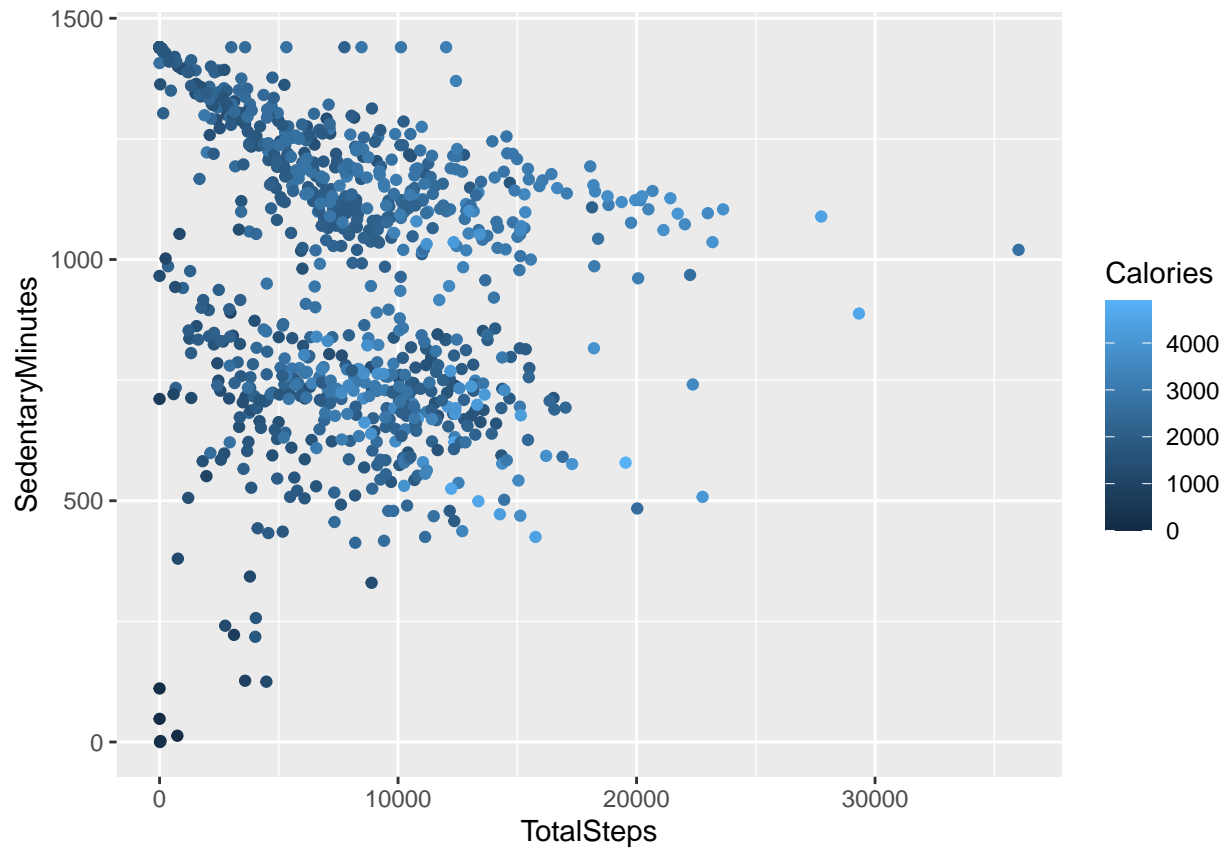
```
weightLog %>%  
  select(WeightPounds, BMI) %>%  
  summary()
```

```
##      WeightPounds      BMI  
##  Min.       :116.0      Min.       :21.45  
## 1st Qu.:135.4      1st Qu.:23.96  
## Median :137.8      Median :24.39  
## Mean   :158.8      Mean   :25.19  
## 3rd Qu.:187.5      3rd Qu.:25.56  
## Max.    :294.3      Max.    :47.54
```

Analysis: Plot

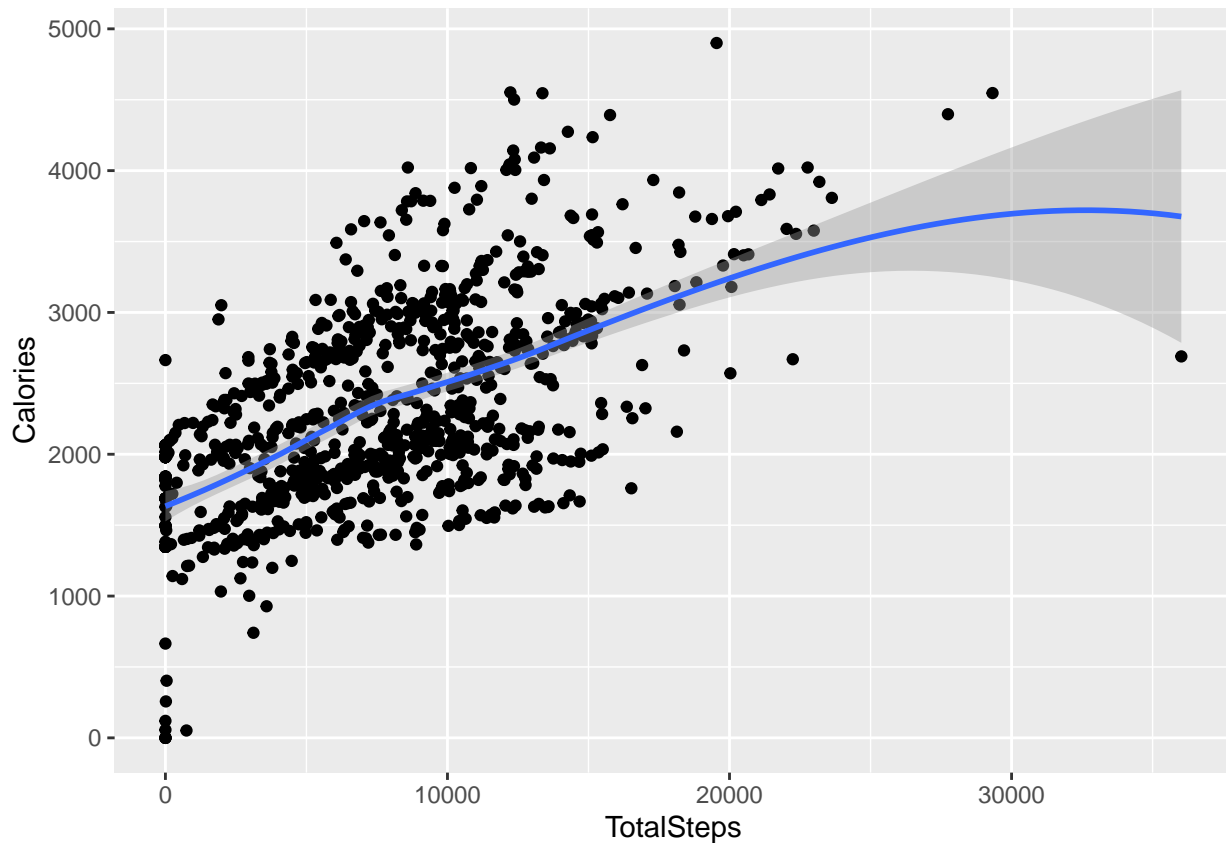
The plot for Total Steps vs Sedentary Minutes is as follows:

```
ggplot(data = dailyActivity, aes(x = TotalSteps, y = SedentaryMinutes, color = Calories)) + geom_point()
```



```
ggplot(data = dailyActivity, aes(x = TotalSteps, y = Calories)) + geom_point() + stat_smooth()
```

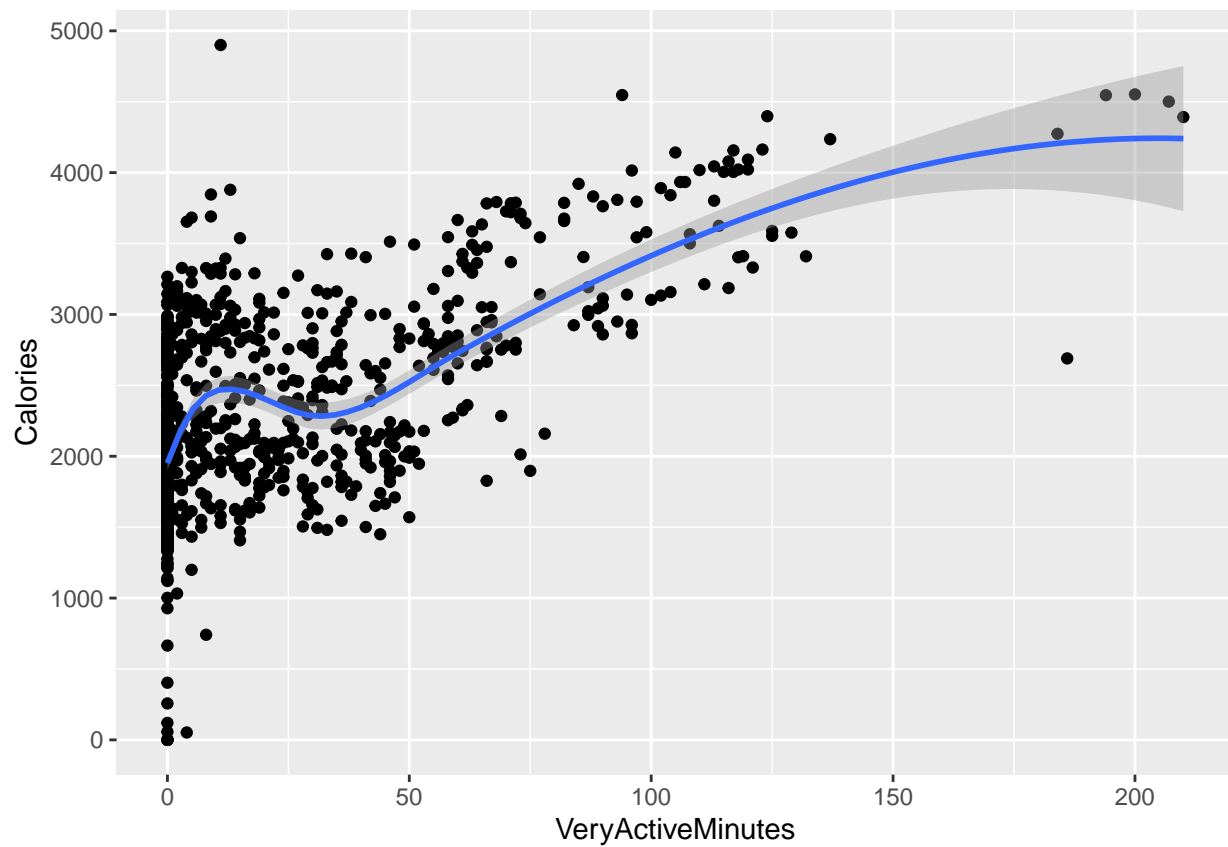
```
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```

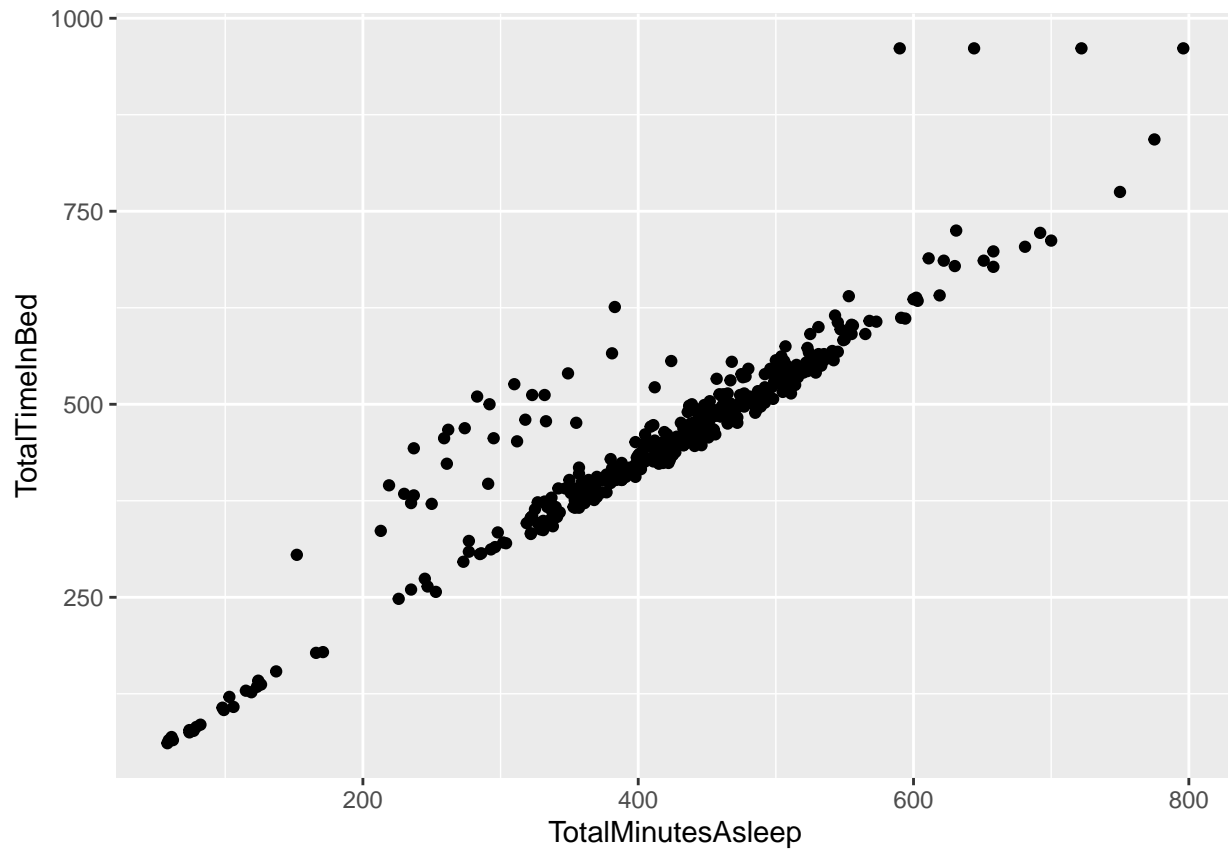
The graph above clearly demonstrates that the relation between the number of steps walked and the amount of calories burned is linear. The more the number of steps, the higher the calories burned.

```
ggplot(data = dailyActivity, aes(x = VeryActiveMinutes, y = Calories)) + geom_point() + stat_smooth()

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



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



This plot shows that people have not been logging their sleep hours in the proper manner and therefore this data is not too useful to make use of. ## Conclusion

Some of the things I noticed when analyzing this data:

- Fitbit is not collecting any hydration data, Bellabeat does which makes it more user friendly and something that regular users would enjoy working with.
- People need to log their sleep times properly
- Another feature I would like Bellabeat to add is the different modes of exercise or to create a program that would allow the sensors to monitor and recognize what sort of exercise the user is doing and then log it in the watch rather than entering or starting the watch feature to collect data.