

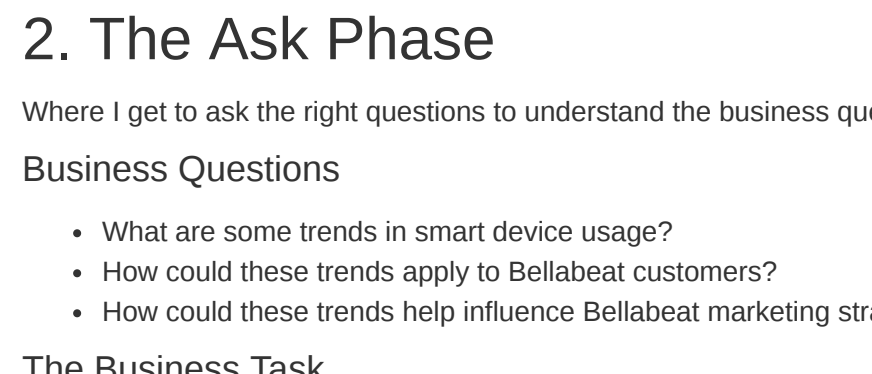
# BellaBeat Case Study Notes

OlaJuwon Aina

2022-09-12

## 1. About the Company

BellaBeat is a high-tech manufacturer of women's health products. BellaBeat is a successful little business with the potential to grow into a major player in the global smart device industry. Uka Sren, cofounder and Chief Creative Officer of BellaBeat, believes that examining smart device fitness data could help the company discover new development prospects.



## 2. The Ask Phase

Where I get to ask the right questions to understand the business questions and also identify key stakeholders on the project.

### Business Questions

- What are some trends in smart device usage?
- How could these trends apply to BellaBeat customers?
- How could these trends help influence BellaBeat marketing strategy?

### The Business Task

How consumers use non-BellaBeat smart devices to gain insights

### Stakeholders Involved

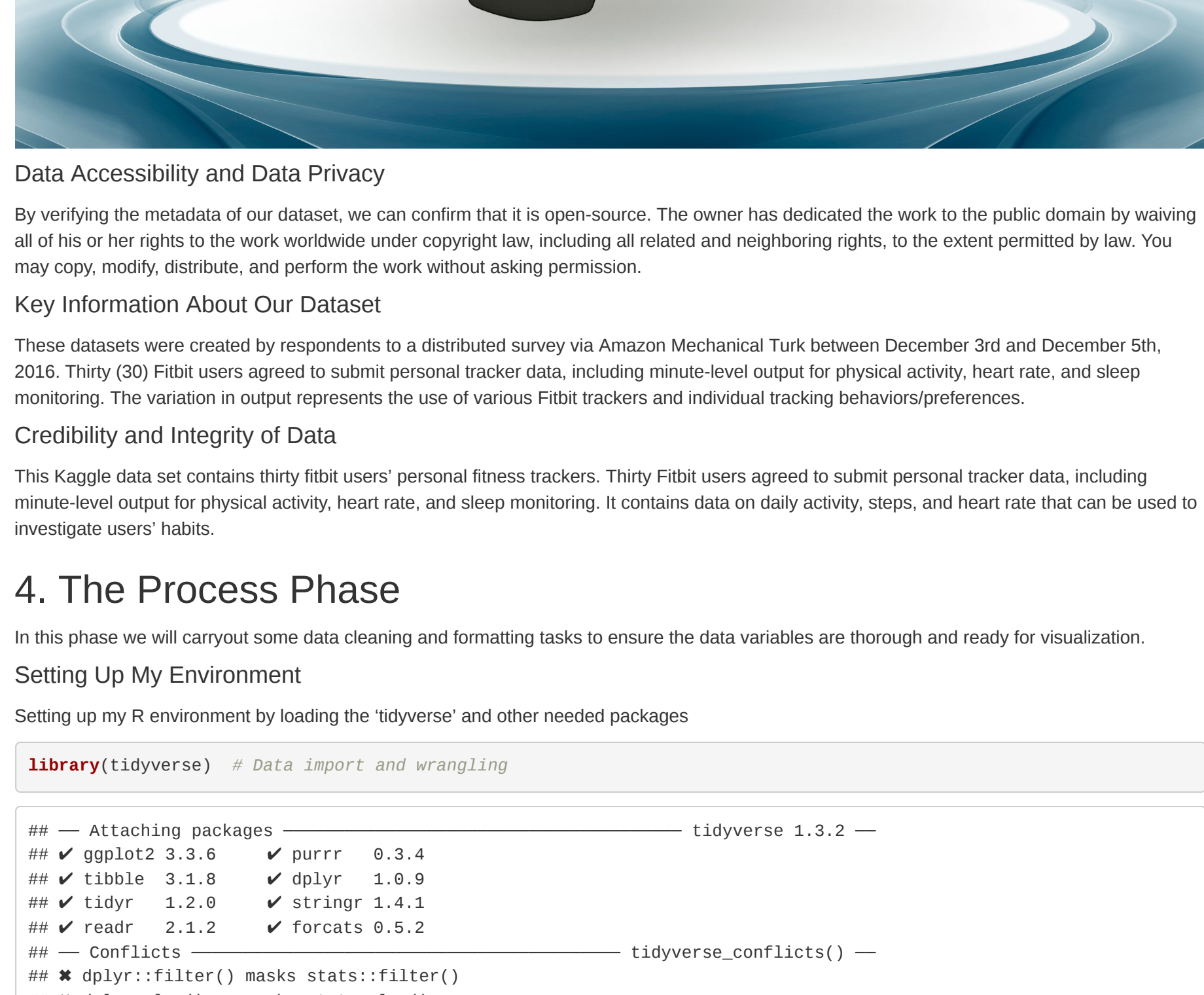
- Uka Sren - The cofounder and Chief Creative Officer of BellaBeat.
- Satish Reddy - BellaBeat cofounder and key member of BellaBeat executive team
- The Marketing Analytics team at BellaBeat

## 3. The Prepare Phase

Here's where I get to gather the dataset to use, identify the source, the security, credibility and integrity.

### Dataset used

The fibit fitness tracker public data will be used for this analysis. [Here](#)



### Data Accessibility and Data Privacy

By verifying the metadata of our dataset, we can confirm that it is open-source. The owner has dedicated the work to the public domain by waiving all of his/her rights to the work worldwide under copyright law, including all related and neighboring rights, to the extent permitted by law. You may copy, modify, distribute, and perform the work without asking permission.

### Key Information About Our Dataset

These datasets were created by respondents to a distributed survey via Amazon Mechanical Turk between December 3rd and December 5th, 2016. Thirty (30) Fitbit users agreed to submit personal fitness data, including minute-level output for physical activity, heart rate, and sleep monitoring. The variation in output represents the use of various Fitbit trackers and individual tracking behaviors/preferences.

### Credibility and Integrity of Data

This Kaggle data set contains thirty fitbit users' personal fitness trackers. Thirty Fitbit users agreed to submit personal tracker data, including minute-level output for physical activity, heart rate, and sleep monitoring. It contains data on daily activity, steps, and heart rate that can be used to investigate users' habits.

## 4. The Process Phase

In this phase we will carryout some data cleaning and formatting tasks to ensure the data variables are thorough and ready for visualization.

### Setting up My R Environment

Setting up R Environment by loading the 'tidyverse' and other needed packages

```
library(tidyverse) # Data Import and wrangling

## --- Attaching packages --- tidyverse 1.3.2 ---
## ggplot2 3.3.6      ✓ purrr   0.3.4
## lubridate 3.1.8    ✓ dplyr   1.0.9
## tidyr  1.2.0       ✓ stringr 1.4.1
## readr  2.1.2       ✓ forcats  0.5.1
## --- Conflicts --- tidyverse_conflicts() ---
## dplyr::filter() masks stats::filter()
## dplyr::lag()    masks stats::lag()

library(ggplot2) # For Data Visualization
library(tidy)
library(scales) # For transforming numbers in percentage

##
## Attaching package: 'scales'
##
## The following object is masked from 'package:purrr':
##
##   discard
##
## The following object is masked from 'package:readr':
##
##   col_factor

## To know factor
```

### Get To Working Directory

```
getwd() # Displays the working directory
```

```
## [1] "C:/Users/Ola/Documents"
```

### Importing The Datasets

There are 18 csv files in the dataset. Each of them displays data related to the device's various functions: calories, activity level, daily steps, and so on.

To simplify the analysis, we will concentrate on daily data in this study.

#### Daily Activity

```
daily_activity <- read_csv("Fitabase Data 4.12.16-5.12.16/dailyActivity_merged.csv")
```

```
## Rows: 948 Columns: 15
## --- Column specification ---
## Delimiter: ","
## chr (2): 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.
```

```
View(daily_activity)
```

#### Daily Calories

```
daily_calories <- read_csv("Fitabase Data 4.12.16-5.12.16/dailyCalories_merged.csv")
```

```
## Rows: 948 Columns: 5
## --- Column specification ---
## Delimiter: ","
## chr (1): ActivityDay
## dbl (4): 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.
```

#### Daily Intensities

```
daily_intensities <- read_csv("Fitabase Data 4.12.16-5.12.16/dailyIntensities_merged.csv")
```

```
## Rows: 948 Columns: 18
## --- Column specification ---
## Delimiter: ","
## chr (1): ActivityDay
## dbl (17): 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.
```

#### Daily Steps

```
daily_steps <- read_csv("Fitabase Data 4.12.16-5.12.16/dailySteps_merged.csv")
```

```
## Rows: 948 Columns: 3
## --- Column specification ---
## Delimiter: ","
## chr (1): ActivityDay
## dbl (2): Id, StepTotal
##
## 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.
```

#### Daily Sleep

```
daily_sleep <- read_csv("Fitabase Data 4.12.16-5.12.16/dailySleep_merged.csv")
```

```
## Rows: 433 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.
```

#### Weight

```
weight_info <- read_csv("Fitabase Data 4.12.16-5.12.16/weightLogInfo_merged.csv")
```

```
## Rows: 67 Columns: 8
## --- Column specification ---
## Delimiter: ","
## chr (1): Date
## dbl (6): Id, WeightKg, WeightPounds, Fat, BMI, LogId
## lgl (1): IsManualReport
##
## 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.
```

### Preview Datasets

Let's have a look at the various datasets and have a clear understanding of how they look, similarities and cohesion between the various datasets.

#### Daily Activity

```
View(daily_activity)
```

#### Daily Calories

```
View(daily_calories)
```

#### Daily Intensities

```
View(daily_intensities)
```

#### Daily Steps

```
View(daily_steps)
```

#### Daily Sleep

```
View(daily_sleep)
```

#### Weight

```
View(weight_info)
```

### Cleaning and Formatting Our Dataset

After examining the various data sets, it is possible to conclude that table 1 (Daily activity) already contains information from table 2 (Daily calories), table 3 (Daily steps), and table 4 (Daily intensities). Another observation is that each dataset has the same number of observations. As a result, those dataframes will be removed.

```
rm(daily_calories, daily_intensities, daily_steps) #removing tables
```

### Transforming the data to be homogeneous

Before merging the datasets, let's clean the date columns to make them homogeneous and transform them to right date type.

```
# Cleaning the variables
daily_activity <- daily_activity %>%
  rename(Date = ActivityDate) %>%
  mutate(Date = as.Date(Date, format = "%m/%d/%y"))

daily_sleep <- daily_sleep %>%
  rename(Date = SleepDay) %>%
  mutate(Date = as.Date(Date, format = "%m/%d/%y"))

weight_info <- weight_info %>%
  select(-LogId) %>%
  mutate(Date = as.Date(Date, format = "%m/%d/%y")) %>%
  mutate(IsManualReport = as.factor(IsManualReport))
```

### Merging the Datasets

```
final_data <- merge(merge(daily_activity, daily_sleep, by=c('Id', 'Date'), all = TRUE), weight_info, by = c('Id', 'Date'), all = TRUE)
```

### Viewing the Merged dataframe (final\_data)

```
View(final_data)
```

### Removing extra/irrelevant variables

```
final_data <- final_data %>%
  select(-c(TrackerDistance, LoggedActivitiesDistance, TotalSleepRecords, WeightPounds, Fat, BMI, IsManualReport))
```

### Reviewing the Merged dataframe (final\_data) again after removing unwanted variables

```
View(final_data)
```

### Checking the variables & data types

```
str(final_data)

## 'data.frame':   943 obs. of 16 variables:
## $ Id: num  1.5e+09 1.5e+09 1.5e+09 1.5e+09 1.5e+09 1.5e+09 ...
## $ Date: Date, format: "2020-04-12" "2020-04-13" ...
## $ TotalDistance: num 13162 10735 10480 9762 2660 ...
## $ TotalSteps: num 8.5 6.97 6.74 6.28 3.16 ...
## $ VerActiveDistance: num 1.88 1.57 2.44 2.14 2.73 ...
## $ ModerateActiveDistance: num 0.55 0.69 0.4 1.08 0.42 ...
## $ LightActiveDistance: num 6.86 4.71 3.93 2.83 5.04 ...
## $ SedentaryActiveDistance: num 0 0 0 0 0 0 0 ...
## $ VerActiveMinutes: num 25 21 30 29 36 36 42 50 28 19 ...
## $ FairlyActiveMinutes: num 33 19 11 34 10 20 16 31 12 0 ...
## $ LightlyActiveMinutes: num 328 217 181 209 221 164 233 264 295 211 ...
## $ SedentaryMinutes: num 728 776 1218 726 773 ...
## $ Calories: num 1085 1797 1776 1745 1863 ...
## $ TotalMinutesAsleep: num 327 384 NA 412 340 708 NA 304 360 325 ...
## $ TotalTimeInBed: num 346 407 NA 442 267 712 NA 320 377 364 ...
## $ WeightKg: num NA NA NA NA NA NA NA NA NA NA ...
```

We can see that majority of the variables are numerical.

```
summary(final_data)

##      Id      Date      TotalSteps      TotalDistance
## Min.   :1.594e+09   Min.   :2020-04-12   Min.   :0   Min.   :0.000
## 1st Qu.:1.320e+09   1st Qu.:2020-04-10   1st Qu.:3795   1st Qu.:1.628
## Median :1.445e+09   Median :2020-04-26   Median :7439   Median :5.268
## Mean   :1.488e+09   Mean   :2020-04-26   Mean   :7652   Mean   :5.593
## 3rd Qu.:1.862e+09   3rd Qu.:2020-05-04   3rd Qu.:10734   3rd Qu.:7.728
## Max.   :1.878e+09   Max.   :2020-05-12   Max.   :98819   Max.   :128.638
##
## VerActiveDistance ModerateActiveDistance LightActiveDistance
## Min.   :0.000   Min.   :0.0000   Min.   :0.000
## 1st Qu.:0.000   1st Qu.:0.0000   1st Qu.:1.950
## Median :0.228   Median :0.2400   Median :3.380
## Mean   :1.504   Mean   :0.5700   Mean   :3.349
## 3rd Qu.:2.865   3rd Qu.:0.8058   3rd Qu.:4.798
## Max.   :121.928   Max.   :6.4808   Max.   :10.718
##
## SedentaryActiveDistance VerActiveMinutes FairlyActiveMinutes
## Min.   :0.0000000   Min.   :0.00   Min.   :0.00
## 1st Qu.:0.0000000   1st Qu.:0.00   1st Qu.:0.00
## Median :0.0000000   Median :4.00   Median :7.80
## Mean   :0.091801    Mean   :21.24   Mean   :13.63
## 3rd Qu.:0.0000000   3rd Qu.:32.80   3rd Qu.:19.88
## Max.   :0.110060    Max.   :216.00   Max.   :143.80
##
## LightlyActiveMinutes SedentaryMinutes Calories TotalMinutesAsleep
## Min.   :0   Min.   :0.0   Min.   :0   Min.   :58.0
## 1st Qu.:127   1st Qu.:729.0   1st Qu.:1830   1st Qu.:361.0
## Median :159   Median :1057.0   Median :2140   Median :433.0
## Mean   :193   Mean   :998.4   Mean   :2388   Mean :419.5
## 3rd Qu.:264   3rd Qu.:1229.0   3rd Qu.:2796   3rd Qu.:499.0
## Max.   :518   Max.   :1444.0   Max.   :4480   Max.   :796.0
##
## TotalTimeInBed WeightKg
## Min.   :61.0   Min.   :52.68
## 1st Qu.:493.0   1st Qu.:61.49
## Median :463.0   Median :62.59
## Mean   :458.6   Mean   :72.04
## 3rd Qu.:526.0   3rd Qu.:85.05
## Max.   :961.0   Max.   :133.58
## NA's :539   NA's :17
```

## 5. The Analyze and Share Phase

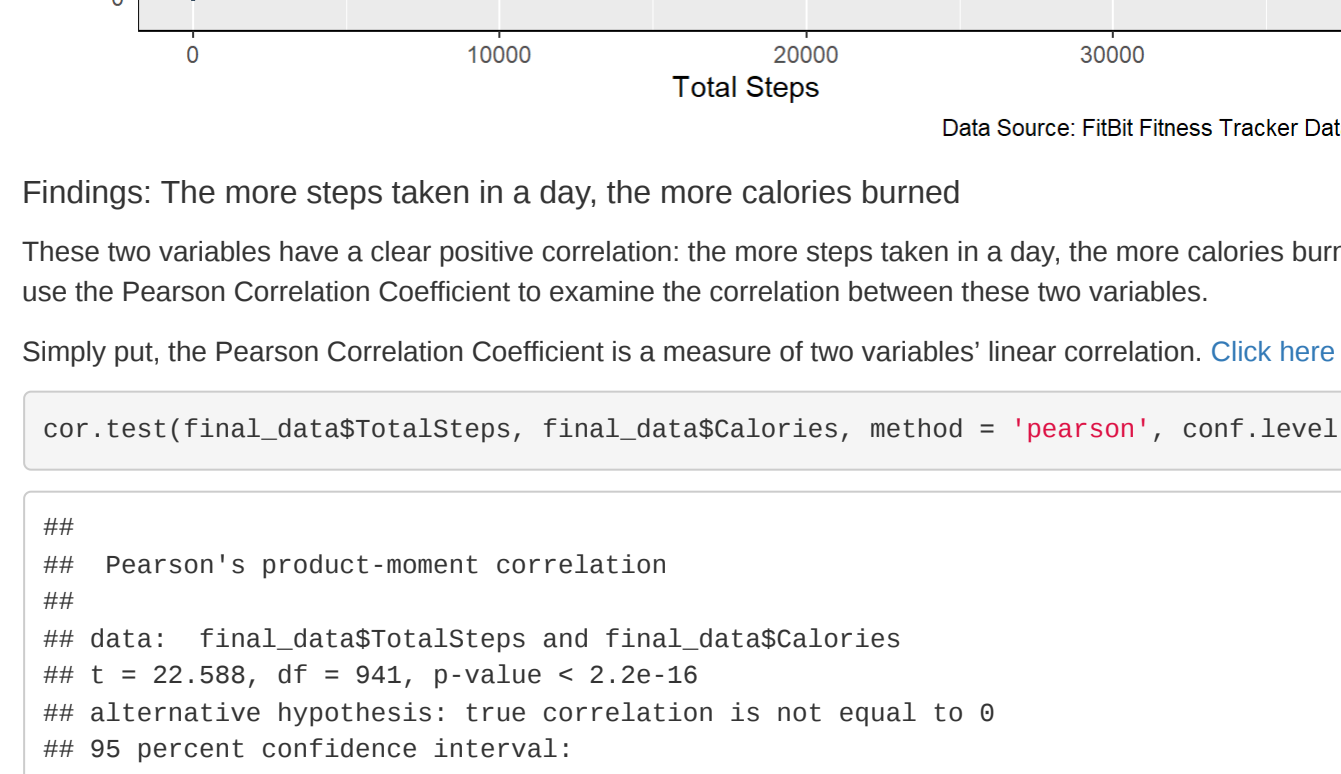
In this phase we will be plotting various graphs to analyze our dataset for possible findings.

### Users Daily Activity

Now with data merged, we can check for Users daily activities in a simple box plot

```
final_data %>%
  mutate(weekdays = weekdays(Date)) %>%
  select(weekdays, TotalSteps) %>%
  mutate(weekdays = factor(weekdays, levels = c('Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday', 'Saturday', 'Sunday'))) %>%
  drop_na() %>%
  ggplot(aes(weekdays, TotalSteps, fill = weekdays)) +
  geom_boxplot() +
  scale_fill_manual(values=c("Set1") +
  theme(legend.position="none") +
  labs(title = "Users' activity by day", x = "Day of the week", y = "Steps",
  caption = "Data Source: Fitbit Fitness Tracker Data")

## Users' activity by day
```



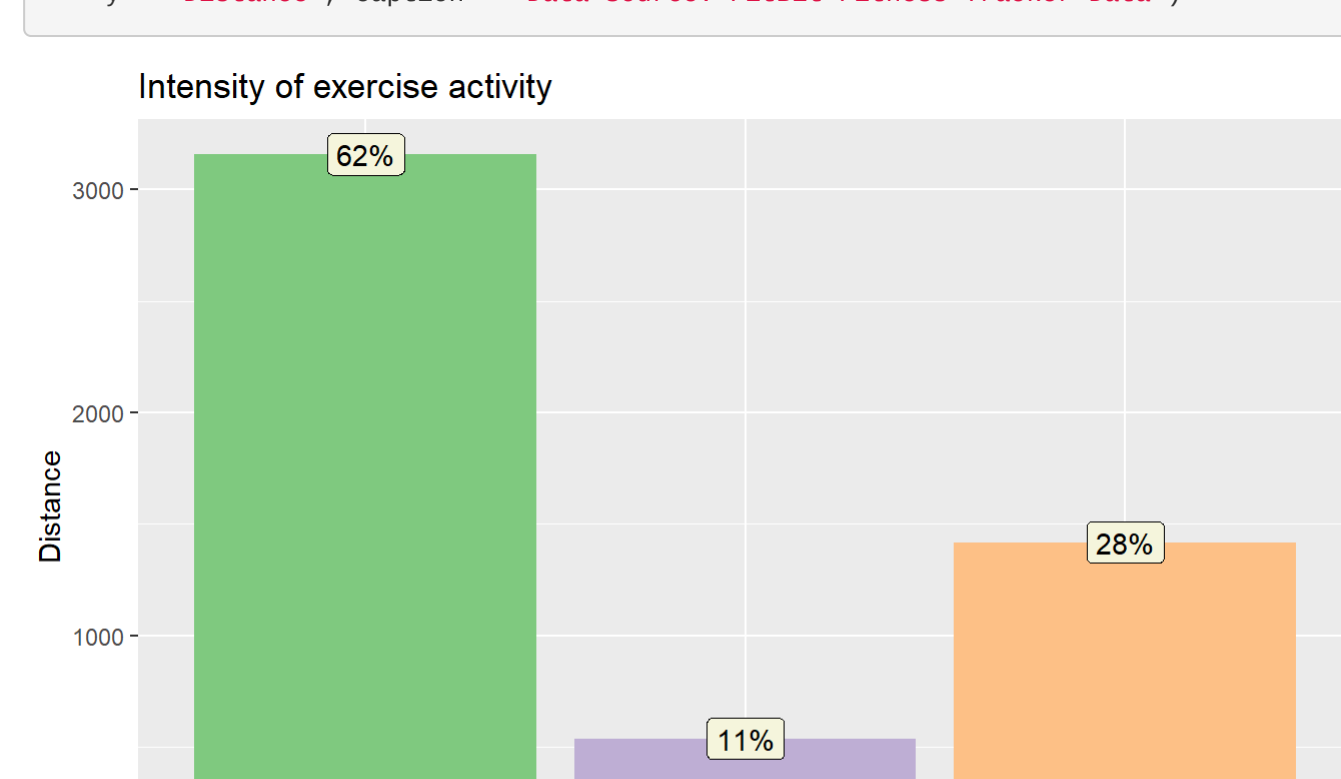
### Next, Check for Calories burned by Steps Taken

Check for calories calories burned by steps (i.e Calories vs Total Steps)

```
final_data %>%
  group_by(TotalSteps, Calories) %>%
  ggplot(aes(x = TotalSteps, y = Calories, color = Calories)) +
  geom_point() +
  theme(legend.position = c(18, 3),
  legend.spacing.y = unit(1, "mm"),
  panel.border = element_rect(colour = "black", fill=NA),
  legend.background = element_blank(),
  legend.box.background = element_rect(colour = "black")) +
  labs(title = "Calories burned by total steps taken", y = "Calories",
  x = "Total Steps", caption = "Data Source: Fitbit Fitness Tracker Data")

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

## Calories burned by total steps taken
```



### Findings: The more steps taken in a day, the more calories burned

These two variables have a clear positive correlation: the more steps taken in a day, the more calories burned. To verify this assumption, we can use the Pearson Correlation Coefficient to examine the correlation between these two variables.

Simply put, the Pearson Correlation Coefficient is a measure of two variables' linear correlation. [Click here](#) for more information.

```
cor.test(final_data$TotalSteps, final_data$Calories, method = "pearson", conf.level = 0.95)
```

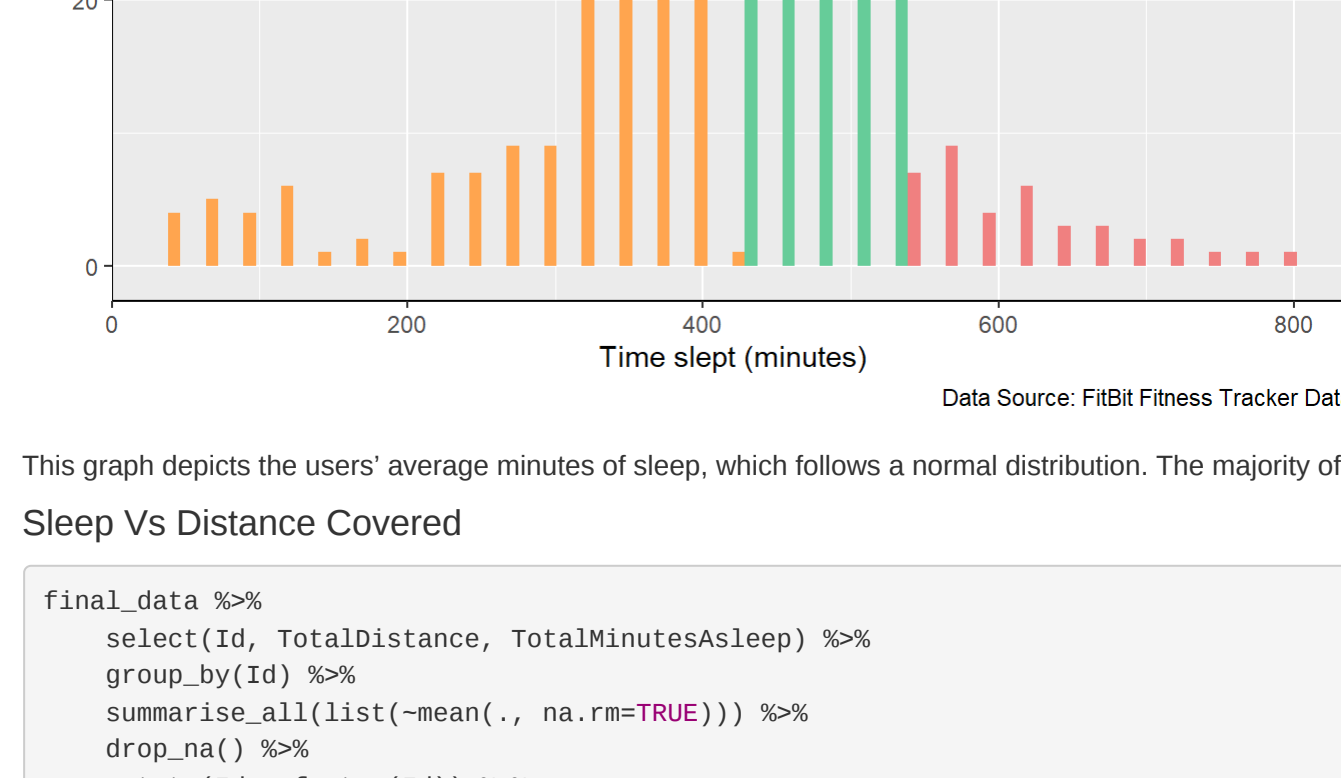
```
##
## Pearson's product-moment correlation
##
## data: final_data$TotalSteps and final_data$Calories
## t = 22.588, df = 841, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.5490551 0.6283441
## sample estimates:
##      cor
## 0.5929493
```

With a confidence level of 95%, the correlation between the variables is almost 0.6. This means that there is a strong relationship between the variables.

### Next, Check for Intensity of Exercise Activity

```
final_data %>%
  select(VerActiveDistance,
  LightActiveDistance) %>%
  summarise(across(everything(), list(sum))) %>%
  gather(activities, value) %>%
  mutate(ratio = value / sum(value),
  label = percent(ratio %>% round(1))) %>%
  mutate(activities = factor(activities, labels = c('Light Activity', 'Moderate Activity', 'Heavy Activity'))) %>%
  ggplot(aes(x = activities), y = value, label = label, fill = activities)) +
  geom_bar(stat="identity") +
  theme(legend.position="none") +
  scale_fill_manual(values=c("beige", colour = "black", yjust = 0.5) +
  scale_fill_brewer(palette="Accent") +
  theme(legend.position="none") +
  labs(title = "Intensity of exercise activity", x = "Activity Level",
  y = "Distance", caption = "Data Source: Fitbit Fitness Tracker Data")

## Intensity of exercise activity
```

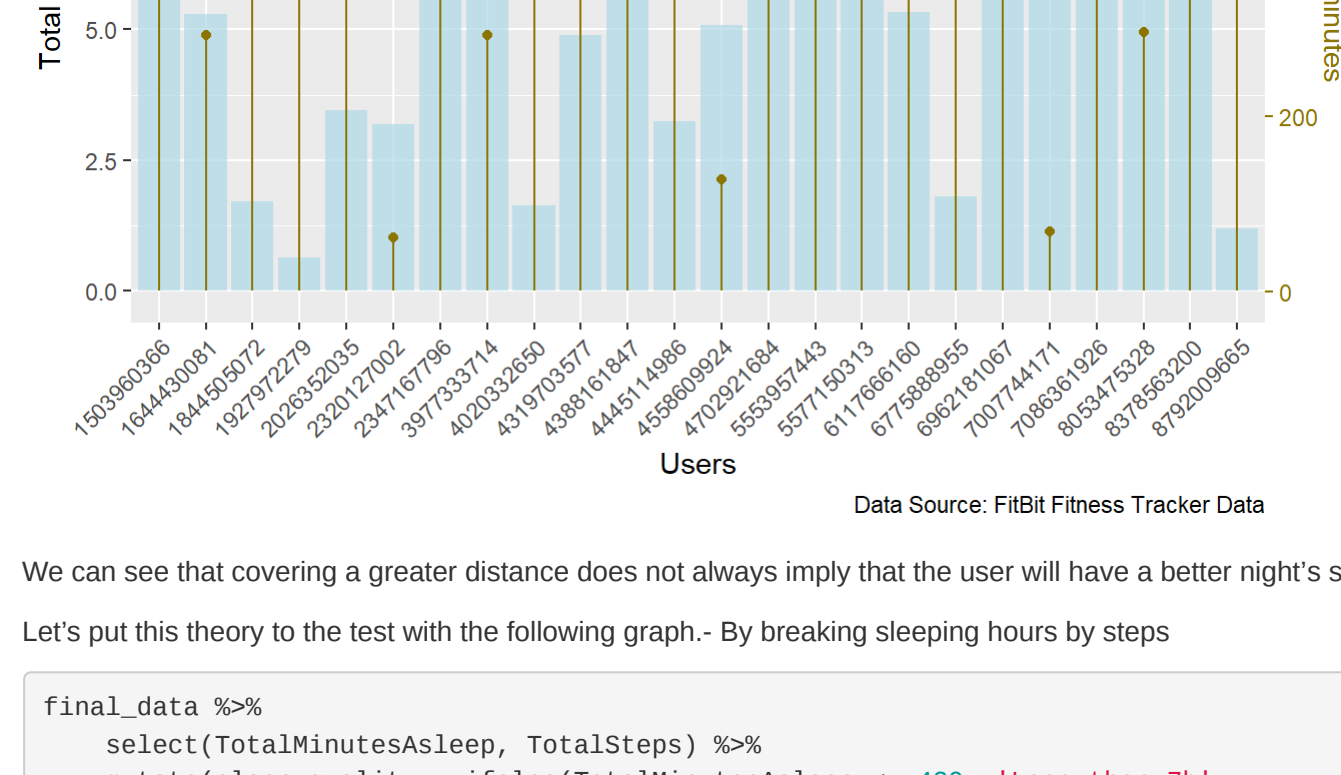


From the analysis above, the most common level of activity during exercise is light.

### Next, Sleep Distribution

```
final_data %>%
  select(TotalMinutesAsleep) %>%
  drop_na() %>%
  mutate(sleep_quality = ifelse(TotalMinutesAsleep <= 420, 'Less than 7h',
  ifelse(TotalMinutesAsleep <= 540, '7h to 9h',
  'More than 9h'))) %>%
  mutate(sleep_quality = factor(sleep_quality,
  levels = c('Less than 7h', '7h to 9h',
  'More than 9h'))) %>%
  ggplot(aes(x = TotalMinutesAsleep, fill = sleep_quality)) +
  geom_histogram(position = "dodge", bins = 30) +
  scale_fill_manual(values=c("tan1", "#66c099", "lightcoral1")) +
  theme(legend.position = c(18, 30), legend.title = element_blank(), legend.spacing.y = unit(1, "mm"),
  panel.border = element_blank(), legend.background = element_rect(colour = "black", fill=NA),
  legend.box.background = element_rect(colour = "black"), legend.box.background = element_rect(colour = "black")) +
  labs(title = "Sleep distribution", x = "Time slept (minutes)", y = "Count",
  caption = "Data Source: Fitbit Fitness Tracker Data")

## Sleep distribution
```



This graph depicts the users' average minutes of sleep, which follows a normal distribution. The majority of users sleep for 320 to 530 minutes.

### Sleep Vs Distance Covered

```
final_data %>%
  select(Id, WeightKg, TotalDistance) %>%
  group_by(Id) %>%
  summarise_all(list(-mean(.), na.rm=TRUE))) %>%
  drop_na() %>%
  mutate(Id = factor(Id)) %>%
  ggplot() +
  geom_bar(aes(x = Id, y = TotalDistance), stat = "identity", fill = 'lightblue', alpha = 0.7) +
  geom_point(aes(x = Id, y = TotalMinutesAsleep/60), color = 'gold4', group = 1) +
  scale_y_continuous(limits=c(0, 12), name = "Total Distance",
  sec.axis = sec.axis(-40, name = "Sleep in minutes")) +
  theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
  theme(axis.title.y.right = element_text(color = "gold4"), axis.ticks.y.right = element_line(color = "gold4"),
  axis.text.y.right = element_text(color = "gold4")) +
  labs(
  title = "Average distance vs average sleep by user", x = "Users",
  caption = "Data Source: Fitbit Fitness Tracker Data")

## Average distance vs average sleep by user
```

