

Peer-graded Assignment: Course Project 1

Introduction

It is now possible to collect a large amount of data about personal movement using activity monitoring devices such as a Fitbit (<http://www.fitbit.com/>), Nike Fuelband (http://www.nike.com/us/en_us/c/nikeplus-fuelband), or Jawbone Up (<https://jawbone.com/up>). These type of devices are part of the quantified self movement a group of enthusiasts who take measurements about themselves regularly to improve their health, to find patterns in their behavior, or because they are tech geeks. But these data remain under-utilized both because the raw data are hard to obtain and there is a lack of statistical methods and software for processing and interpreting the data.

This assignment makes use of data from a personal activity monitoring device. This device collects data at 5 minute intervals through out the day. The data consists of two months of data from an anonymous individual collected during the months of October and November, 2012 and include the number of steps taken in 5 minute intervals each day.

The data for this assignment can be downloaded from the course web site:

- Dataset: [Activity monitoring data]
(<https://d396qusza40orc.cloudfront.net/repdata%2Fdata%2Factivity.zip>
(<https://d396qusza40orc.cloudfront.net/repdata%2Fdata%2Factivity.zip>))

The variables included in this dataset are:

- steps: Number of steps taking in a 5-minute interval (missing values are coded as NA)
- date: The date on which the measurement was taken in YYYY-MM-DD format
- interval: Identifier for the 5-minute interval in which measurement was taken

The dataset is stored in a comma-separated-value (CSV) file and there are a total of 17,568 observations in this dataset.

Loading and preprocessing the data

- Setting the Working Directory

```
path <- setwd("C:/Users/KHAWLA/Downloads/")
```

- Unzipping the zip file

```
unzip("repdata_data_activity.zip", exdir = path)
```

- Loading the data

```
activity <- read.csv("activity.csv")  
head(activity)
```

##	steps	date	interval
## 1	NA	2012-10-01	0
## 2	NA	2012-10-01	5
## 3	NA	2012-10-01	10
## 4	NA	2012-10-01	15
## 5	NA	2012-10-01	20
## 6	NA	2012-10-01	25

- Setting data variables into date format

```
activity$date <- as.POSIXct(activity$date, "%Y%m%d")
```

```
## Warning in strptime(xx, f, tz = tz): unknown timezone '%Y%m%d'
```

```
## Warning in as.POSIXct.POSIXlt(x): unknown timezone '%Y%m%d'
```

```
## Warning in strptime(xx, f, tz = tz): unknown timezone '%Y%m%d'
```

```
## Warning in as.POSIXct.POSIXlt(x): unknown timezone '%Y%m%d'
```

```
## Warning in strptime(xx, f, tz = tz): unknown timezone '%Y%m%d'
```

```
## Warning in as.POSIXct.POSIXlt(x): unknown timezone '%Y%m%d'
```

```
## Warning in strptime(xx, f, tz = tz): unknown timezone '%Y%m%d'
```

```
## Warning in as.POSIXct.POSIXlt(x): unknown timezone '%Y%m%d'
```

```
## Warning in strptime(xx, f, tz = tz): unknown timezone '%Y%m%d'
```

```
## Warning in as.POSIXct.POSIXlt(x): unknown timezone '%Y%m%d'
```

```
## Warning in strptime(x, f, tz = tz): unknown timezone '%Y%m%d'
```

```
## Warning in as.POSIXct.POSIXlt(as.POSIXlt(x, tz, ...), tz, ...): unknown timezone
## '%Y%m%d'
```

- Create a new column 'day'

```
activity$day <- weekdays(activity$date)
```

```
## Warning in as.POSIXlt.POSIXct(x, tz): unknown timezone '%Y%m%d'
```

- Summarise and structure of the dataset

```
summary(activity)
```

```
## Warning in as.POSIXlt.POSIXct(x, tz): unknown timezone '%Y%m%d'
```

```
##      steps      date      interval      day
## Min.   : 0.00   Min.   :2012-10-01   Min.   : 0.0   Length:17568
## 1st Qu.: 0.00   1st Qu.:2012-10-16   1st Qu.: 588.8   Class :character
## Median : 0.00   Median :2012-10-31   Median :1177.5   Mode  :character
## Mean   : 37.38   Mean   :2012-10-31   Mean    :1177.5
## 3rd Qu.: 12.00   3rd Qu.:2012-11-15   3rd Qu.:1766.2
## Max.   :806.00   Max.   :2012-11-30   Max.    :2355.0
## NA's   :2304
```

```
str(activity)
```

```
## 'data.frame': 17568 obs. of 4 variables:
## $ steps : int NA NA NA NA NA NA NA NA NA NA ...
## $ date : POSIXct, format:
```

```
## Warning in as.POSIXlt.POSIXct(x, tz): unknown timezone '%Y%m%d'
```

```
## "2012-10-01" "2012-10-01" ...
## $ interval: int 0 5 10 15 20 25 30 35 40 45 ...
## $ day : chr "lundi" "lundi" "lundi" "lundi" ...
```

What is the mean total number of steps taken per day?

1. Calculate the total number of steps taken per day

```
NumberOfSteps <- with(data = activity, aggregate(steps, by=list(date), sum, na.rm=TRUE))
```

```
## Warning in as.POSIXlt.POSIXct(x, tz): unknown timezone '%Y%m%d'
```

```
## Warning in as.POSIXlt.POSIXct(x, tz): unknown timezone '%Y%m%d'
```

```
names(NumberOfSteps) <- c("Date","Steps")
TotalSteps <- data.frame(NumberOfSteps)
head(TotalSteps)
```

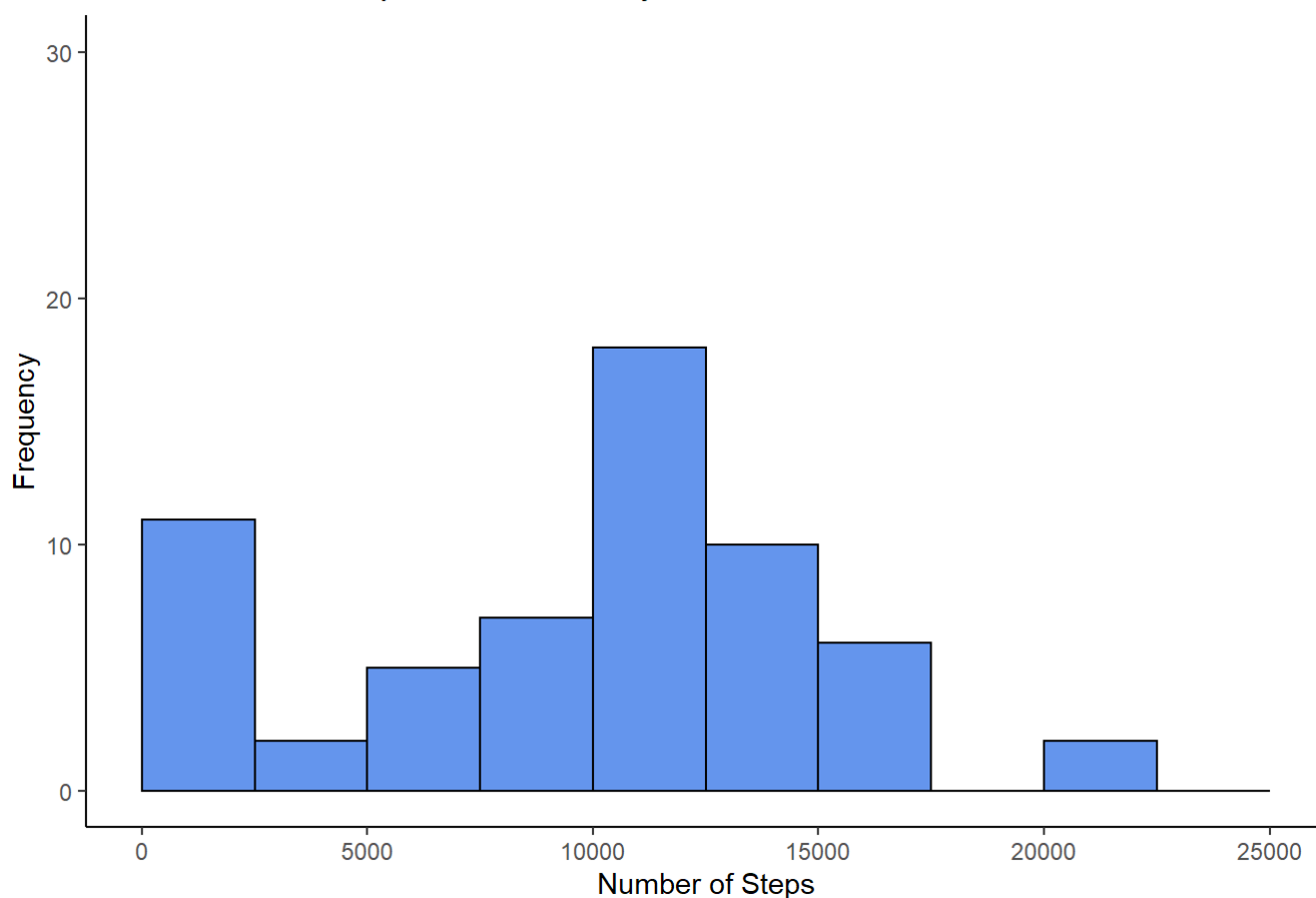
```
## Warning in as.POSIXlt.POSIXct(x, tz): unknown timezone '%Y%m%d'
```

```
##      Date Steps
## 1 2012-10-01    0
## 2 2012-10-02  126
## 3 2012-10-03 11352
## 4 2012-10-04 12116
## 5 2012-10-05 13294
## 6 2012-10-06 15420
```

2. Make a histogram of the total number of steps taken each day

```
library(ggplot2)
ggplot(TotalSteps, aes(x = Steps))+
  geom_histogram(breaks=seq(0,25000, by=2500), fill="cornflowerblue", col="black")+theme_classic()+ ylim(c(0,30))+
  ggtitle("Total number of steps taken each day")+xlab("Number of Steps")+ylab("Frequency")
```

Total number of steps taken each day



3. Calculate and report the mean and median of the total number of steps taken per day

```
mean(TotalSteps$Steps)
```

```
## [1] 9354.23
```

```
median(TotalSteps$Steps)
```

```
## [1] 10395
```

What is the average daily activity pattern?

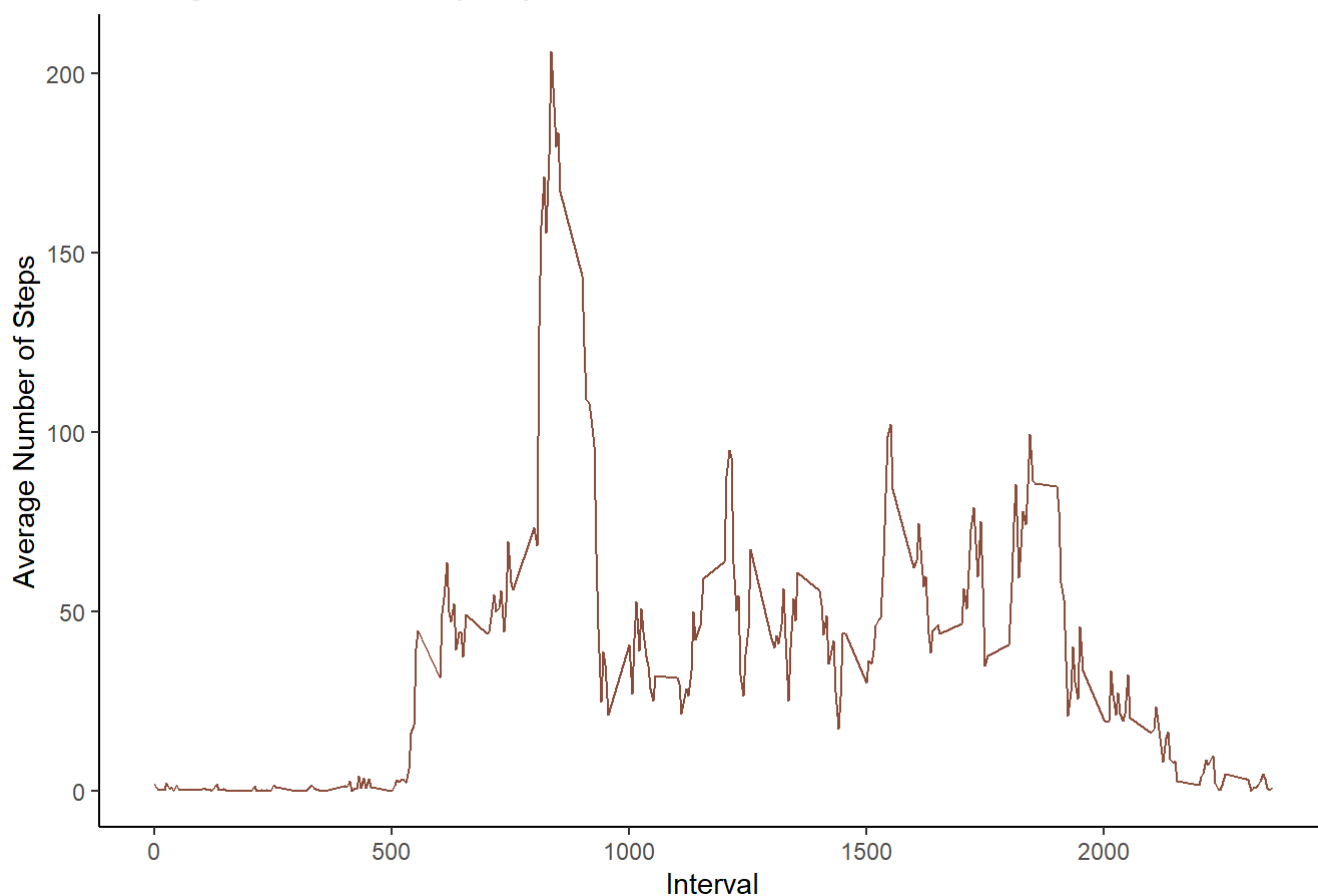
1. Make a time series plot (i.e.type = "l") of the 5-minute interval (x-axis) and the average number of steps taken, averaged across all days (y-axis)

```
AvgDailyActivity <- with(data= activity, aggregate(steps, by= list(interval), mean, na.rm= TRUE))
names(AvgDailyActivity) <- c("Interval", "Average")
AvgActivity <- data.frame(AvgDailyActivity)
head(AvgActivity)
```

```
##   Interval   Average
## 1         0 1.7169811
## 2         5 0.3396226
## 3        10 0.1320755
## 4        15 0.1509434
## 5        20 0.0754717
## 6        25 2.0943396
```

```
ggplot(AvgActivity, aes(x = Interval, y = Average))+
  geom_line(col="salmon4")+ xlab("Interval")+
  ylab("Average Number of Steps")+ ggtitle("Average Number of Steps by Interval")+
  theme_classic()
```

Average Number of Steps by Interval



2. Which 5-minute interval, on average across all the days in the dataset, contains the maximum number of steps?

```
AvgDailyActivity[which.max(AvgDailyActivity$Average),]$Interval
```

```
## [1] 835
```

Imputing missing values

1. Calculate and report the total number of missing values in the dataset (i.e. the total number of rows with NAs)

```
sum(is.na(activity$steps))
```

```
## [1] 2304
```

2. Devise a strategy for filling in all of the missing values in the dataset. The strategy does not need to be sophisticated. For example, you could use the mean/median for that day, or the mean for that 5-minute interval, etc.

```
imputedSteps <- AvgDailyActivity$Average[match(activity$interval, AvgDailyActivity$Interval)]  
  
activityImputed <- transform(activity, steps = ifelse(is.na(activity$steps), yes = imputedSteps, no = activity$steps))
```

3. Create a new dataset that is equal to the original dataset but with the missing data filled in.

```
TotalActImputed <- aggregate(steps~date,activityImputed, sum)
```

```
## Warning in as.POSIXlt.POSIXct(x, tz): unknown timezone '%Y%m%d'
```

```
## Warning in as.POSIXlt.POSIXct(x, tz): unknown timezone '%Y%m%d'
```

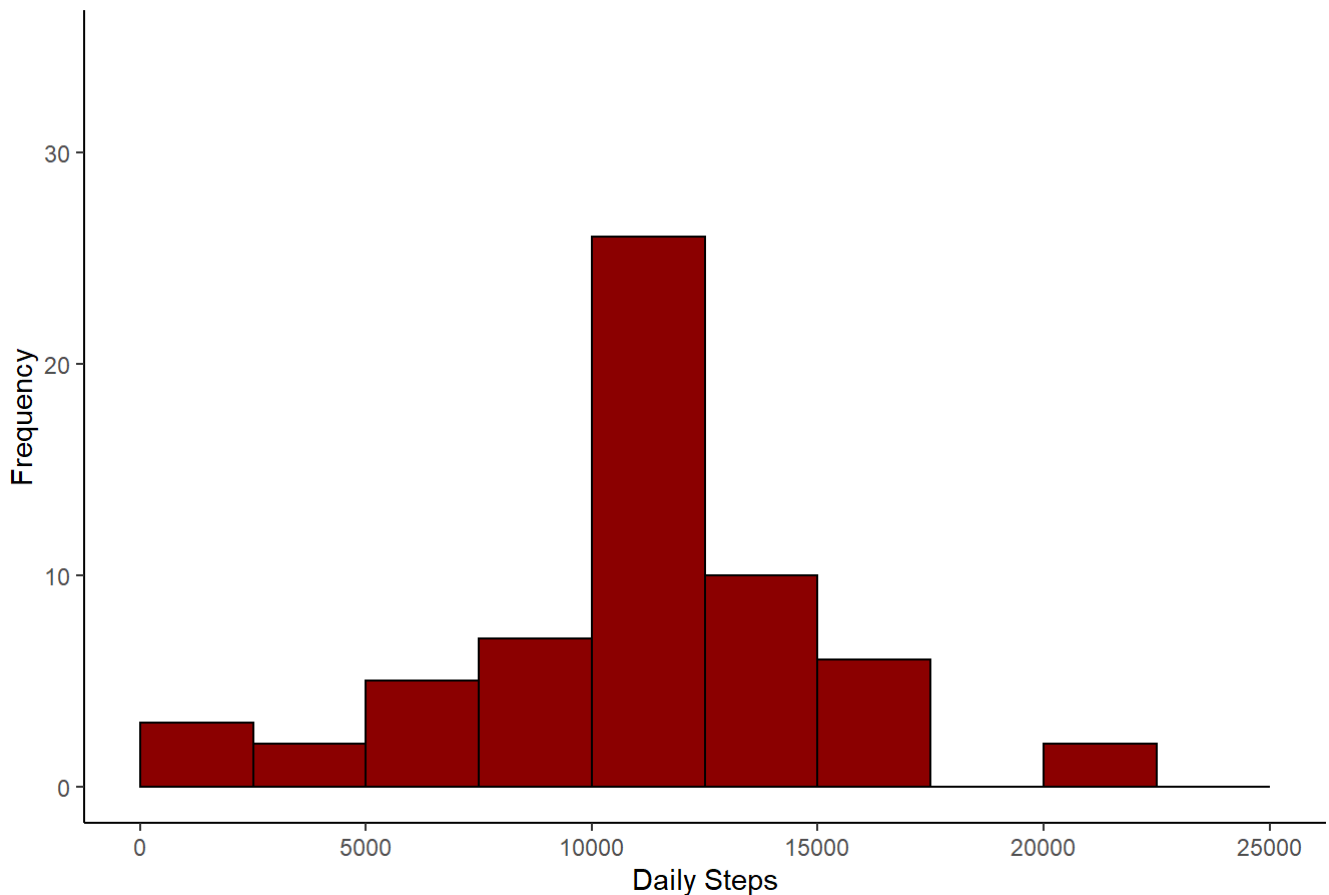
```
names(TotalActImputed) <- c("Date", "Daily_Steps")  
sum(is.na(TotalActImputed$Daily_Steps))
```

```
## [1] 0
```

4. Make a histogram of the total number of steps taken each day and Calculate and report the mean and median total number of steps taken per day. Do these values differ from the estimates from the first part of the assignment? What is the impact of imputing missing data on the estimates of the total daily number of steps?

```
ImputedSteps <- data.frame(TotalActImputed)  
ggplot(ImputedSteps, aes(x = Daily_Steps))+geom_histogram(breaks = seq(0,25000, by = 2500), fill = "darkred", col = "black")+ylim(c(0,35))+theme_classic()+ xlab("Daily Steps")+ylab("Frequency")+ggtitle("Total number of Steps taken by day")
```

Total number of Steps taken by day



```
mean(ImputedSteps$Daily_Steps)
```

```
## [1] 10766.19
```

```
median(ImputedSteps$Daily_Steps)
```

```
## [1] 10766.19
```

Are there differences in activity patterns between weekdays and weekends?

1. Create a new factor variable in the dataset with two levels weekdays and weekends indicating whether a given date is a weekday or weekend day.

```
activity$dayType <- sapply(activity$date, function(x) {  
  if(weekdays(x) == "samedi" | weekdays(x) == "dimanche")  
  {y <- "Weekend"}  
  else {y <- "Weekday"}  
  y  
})
```



```
activityByDay <- aggregate(steps ~ interval + dayType, activity, mean, na.rm = TRUE)

ggplot(activityByDay, aes(x = interval , y = steps, color = dayType)) + geom_line() +
  ggtitle("Average Daily Steps by Day Type") + xlab("Interval") + ylab("Average Number of S
  teps")+
  facet_wrap(~dayType, ncol = 1, nrow = 2)+theme_classic()
```

