

Reproducible Research: Programming Assignment 1

Felix P. Muga II

March 2, 2015

Calling the Libraries

We shall be using the following libraries: **dplyr**, **ggplot2**, **lattice** and **stringr** in this report.

```
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
##
## The following object is masked from 'package:stats':
##
##   filter
##
## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union
```

```
library(ggplot2)
library(lattice)
library(stringr)
```

Loading and Preprocessing the Data

The dataset were downloaded from the **COURSE WEBSITE** and named as **activity.csv** which is a comma-separated value file.

The dataset has 17,568 observations with 3 variables. These variables are:

- **steps** : Number of steps taken 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.
There are 61 unique **dates** in the dataset which starts in 2012-10-01 and ends in 2012-11-30.
- **interval** : **Identifier** for the 5-minute interval in which measurement was taken.
Each **identifier** shall be expressed in HHMM format where HH is one of 0, 1, ..., 23 and MM is one of 0, 5, ..., 55.
The dataset has 288 unique 5-minute interval **identifiers** in a day.
The interval identifier 0000 is the 1st 5-minute interval of the first hour,
The interval identifier 0005 is the 2nd 5-minute interval of the first hour,
:
The interval identifier 0075 is the 12th or the last 5-minute interval of the first hour,
:
The interval identifier 2300 is the 1st 5-minute interval of the 24th hour,
The interval identifier 2305 is the 2nd 5-minute interval of the 24th hour,
:
The interval identifier 2375 is the 12th or the last 5-minute interval of the 24th hour.

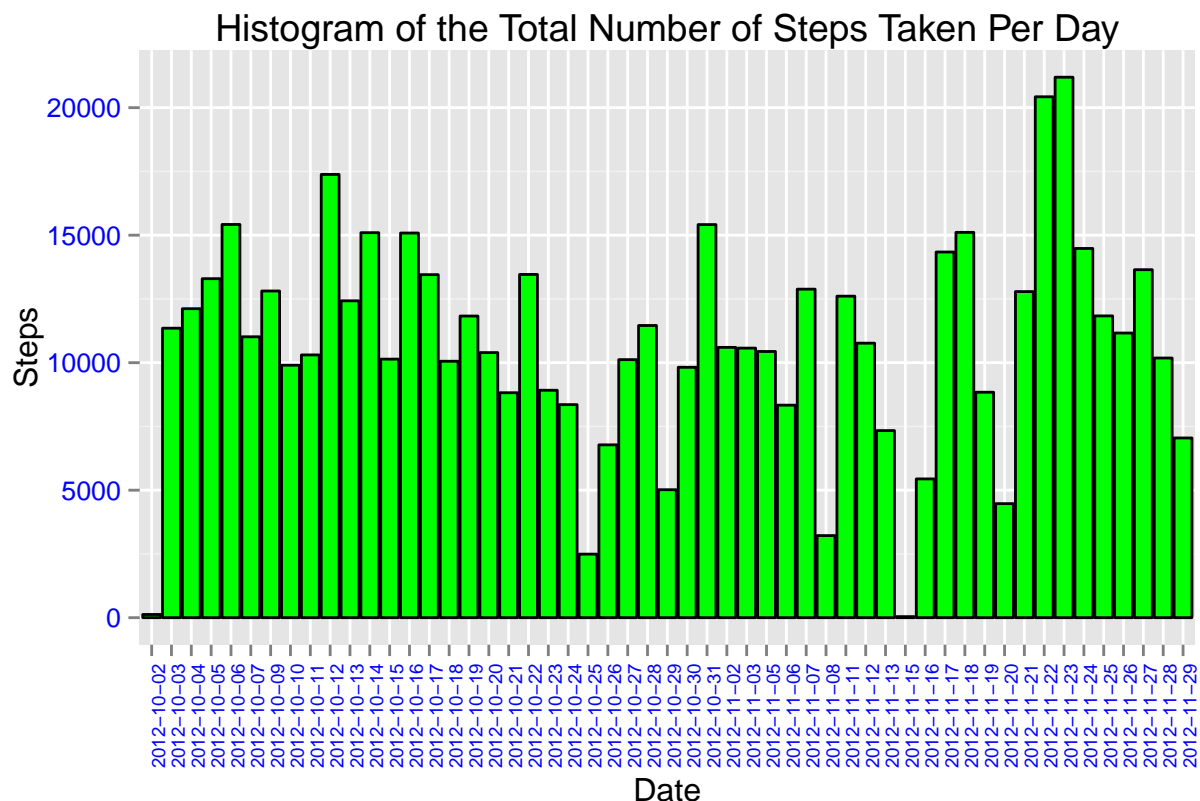
Thus, the total number of observations in the dataset is $288 \times 61 = 17,568$.

```
# Loading the CSV file as a data.format
rawdf <- read.csv("activity.csv")
# Transforming the interval column in HHMM format
rawdf <- transform(rawdf, interval = sprintf("%0004d", interval))
# The dataset is processed where observations with missing values are removed.
df <- rawdf[complete.cases(rawdf),]
```

Problem 1 - What is mean total number of steps taken per day?

1.1. Make a histogram of the total number of steps taken each day. (Days with missing values are not included.)

```
## R Script of the histogram
qplot(date, data = df, weight = steps,
      xlab = "Date", ylab = "Steps",
      main = "Histogram of the Total Number of Steps Taken Per Day") +
  geom_histogram(colour = "black", fill = "green") +
  theme(axis.text.x = element_text(angle = 90, colour="blue",size=7),
        axis.text.y = element_text(colour = "blue",size=10))
```



```
## R Script of the mean and median of the total number of steps per day
numSteps <- summarise(group_by(df, date), total=sum(steps))
```

```
mean <- format(round(mean(numSteps$total)), big.mark = ",")
median <- format(round(median(numSteps$total)), big.mark = ",")
```

1.2. Calculate and report the mean and median total number of steps taken per day. The mean and the median of the total number of steps taken per day are respectively :

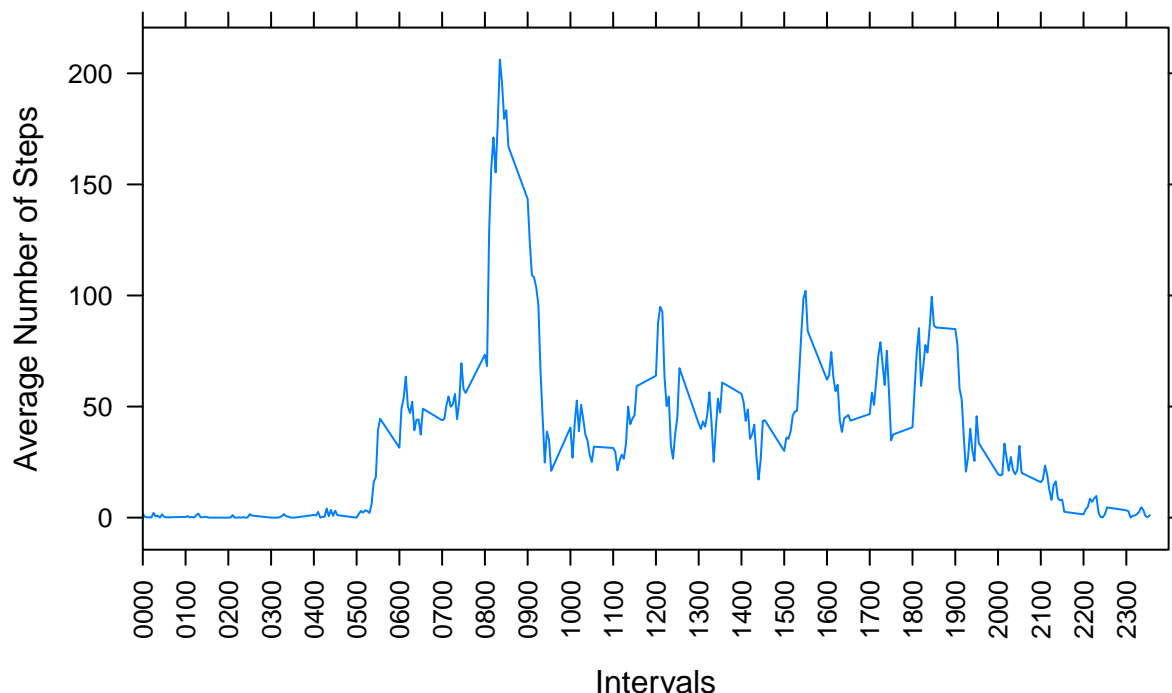
mean = 10,766 and median = 10,765.

Problem 2 - What is the average daily activity pattern?

2.1. 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)

```
#      R Script for the time series plot of the average daily activity pattern
dailyPattern <- summarise(group_by(df, interval), averageSteps = mean(steps))
labels <- sprintf("%0004d",seq(0,2400,by=100))
xyplot(
  averageSteps ~ interval,
  data = dailyPattern,
  type = "l",
  main = "A Time-Series Plot of the Average Activity Pattern
  (There are 12 five-minute intervals between x-ticks)",
  xlab = "Intervals",
  ylab = "Average Number of Steps",
  scales = list(x=list(tick.number=25, labels = labels, rot = 90)),
  xlim = c(0,2399)
)
```

**A Time-Series Plot of the Average Activity Pattern
(There are 12 five-minute intervals between x-ticks)**



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

```
index <- dailyPattern[which.max(dailyPattern$averageSteps),]  
index <- as.character(index[1,1])  
hr <- as.numeric(str_sub(index, start=1, end=2))  
min <- as.numeric(str_sub(index, start=3, end=4))
```

Hence, the 5-minute interval that contains the maximum number of steps is equal to 0835 which is the 8th 5-minute interval of the 9th hour of the day.

Problem 3 - Imputing missing values

```
missingValues <- rawdf[is.na(rawdf$step),]  
missingDays <- summarise(group_by(missingValues, date), missing = n())
```

```
missingNumber = format(nrow(missingValues), big.mark=",")  
missingNumber
```

3.1. Calculate and report the total number of missing values in the dataset.

```
## [1] "2,304"
```

```
missingDays
```

```
## Source: local data frame [8 x 2]  
##  
##      date missing  
## 1 2012-10-01    288  
## 2 2012-10-08    288  
## 3 2012-11-01    288  
## 4 2012-11-04    288  
## 5 2012-11-09    288  
## 6 2012-11-10    288  
## 7 2012-11-14    288  
## 8 2012-11-30    288
```

There are 2,304 missing values in the dataset which are uniformly distributed among the 8 days, since $288 \times 8 = 2,304$.

3.2. Devise a strategy for filling in all of the missing values in the dataset.

- The number of 5-minute intervals in a day is 288.
- Each of the 8 days mentioned above has 288 missing values.
- We computed the average number of steps per interval for each of the 5-minute intervals of a day from the data.
- Hence, we shall fill up the missing value of each of the 5-minute intervals of each of the 8 days with the average number of steps per interval.

```

#      R code to create a new dataset equal to the original with missing data supplied.
#      subset with complete data
lower <- rawdf[!is.na(rawdf$steps),]
#      subset with incomplete data due missing steps
upper <- rawdf[is.na(rawdf$steps),]
#      incomplete data is replaced by mean of number of steps
upper <- data.frame(dailyPattern[match(upper$interval, dailyPattern$interval),2],select(upper, -steps))
#      change the name of the first column to 'steps'
names(upper)[1] <- 'steps'
#      rowbinding the two data frames 'upper' and 'lower'
newdf <- arrange(rbind(upper,lower), date, interval)

```

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

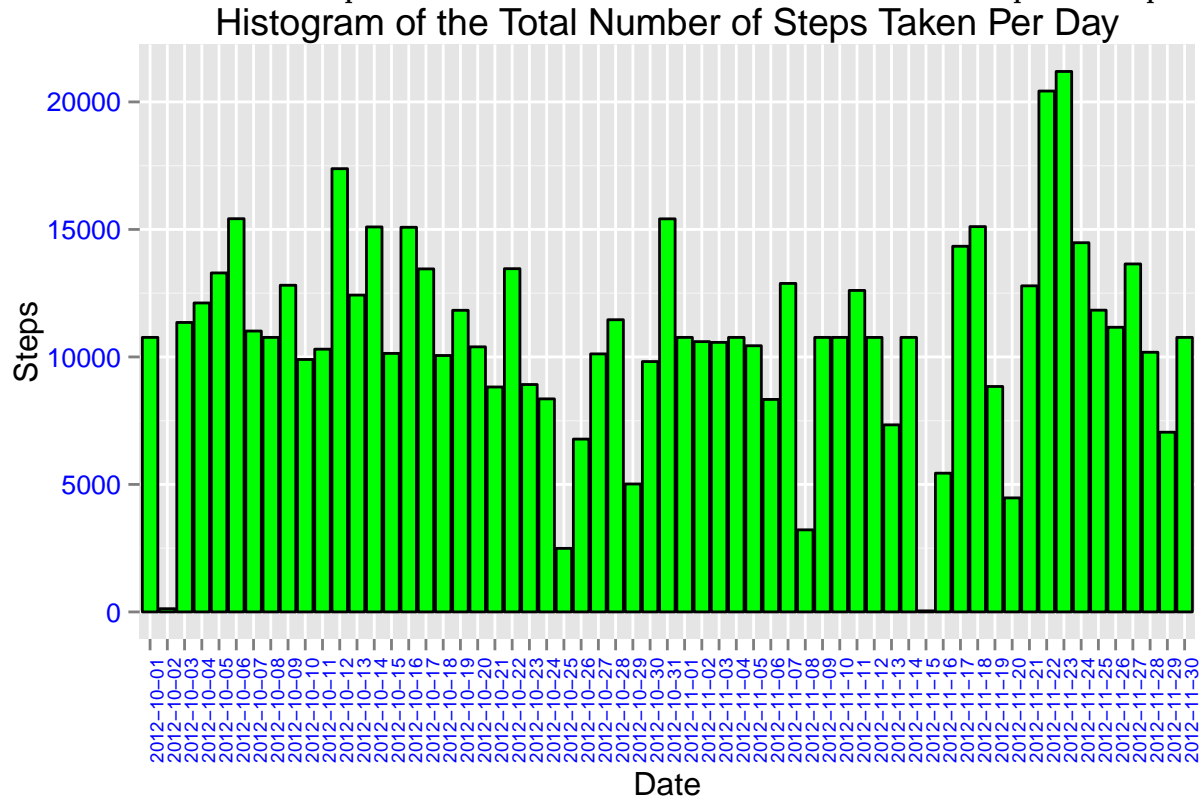
3.4. Make a histogram of the total number of steps taken each day.

```

##      R Script of the histogram
qplot(date, data = newdf, weight = steps,
      xlab = "Date", ylab = "Steps",
      main = "Histogram of the Total Number of Steps Taken Per Day") +
  geom_histogram(colour = "black", fill = "green") +
  theme(axis.text.x = element_text(angle = 90, colour="blue",size=7),
        axis.text.y = element_text(colour = "blue",size=10))

```

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



```
##      R Script of the mean and median of the total number of steps per day
numSteps2 <- summarise(group_by(newdf, date), total=sum(steps))
mean2 <- format(round(mean(numSteps2$total)), big.mark = ",")
median2 <- format(round(median(numSteps2$total)), big.mark = ",")
```

Comparison	original dataset	modified dataset
mean	10,766	10,766
median	10,765	10,766

3.4.2. Do these values differ from the estimates from the first part of the assignment? The error in the computation of the mean between the original and modified datasets is zero.

However, the error with respect to the median between the original and the modified datasets is 0.0001104207.

If we use the median instead of the mean in filling up the missing values, the error could be zero.

3.4.3. What is the impact of imputing missing data on the estimates of the total daily number of steps? If the values that are used to replace the missing data are meaningless or invalid from the “true” values which are missing, then these “illegal” values may cause problems in the statistical analyses.

Problem 4 - Are there differences in activity patterns between weekdays and weekends?

4.1 Create a new factor variable in the dataset with two levels – “weekday” and “weekend” indicating whether a given date is a weekday or weekend day.

4.2. Make a panel plot containing 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 weekday days or weekend days (y-axis).