

Reproducible Research - Activity monitoring devices

Jatin chawda

2/26/2020

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>]

The variables included in this dataset are:

1. steps: Number of steps taking in a 5-minute interval (missing values are coded as NA)
2. date: The date on which the measurement was taken in YYYY-MM-DD format
3. 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.

What is mean total number of steps taken per day? For this part of the assignment, you can ignore the missing values in the dataset.

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

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

What is the average daily activity pattern?

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)

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

Imputing missing values

Note that there are a number of days/intervals where there are missing values (coded as NA). The presence of missing days may introduce bias into some calculations or summaries of the data.

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

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.

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

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?

Are there differences in activity patterns between weekdays and weekends? For this part the `weekdays()` function may be of some help here. Use the dataset with the filled-in missing values for this part.

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.

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).

Loading and preprocessing the data

```
download.file("https://d396qusza40orc.cloudfront.net/repdata%2Fdata%2Factivity.zip", destfile = "activity.zip")
unzip("activity.zip")
```

```
stepdata <- read.csv("activity.csv", header = TRUE)
```

```
head(stepdata)
```

```
##   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
```

Calculate total number of steps taken each day

```
library(magrittr)
```

```
## Warning: package 'magrittr' was built under R version 3.6.2
```

```
library(dplyr)
```

```
## Warning: package 'dplyr' was built under R version 3.6.2
```

```
##
```

```
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
##   filter, lag
```

```
## The following objects are masked from 'package:base':
```

```
##
```

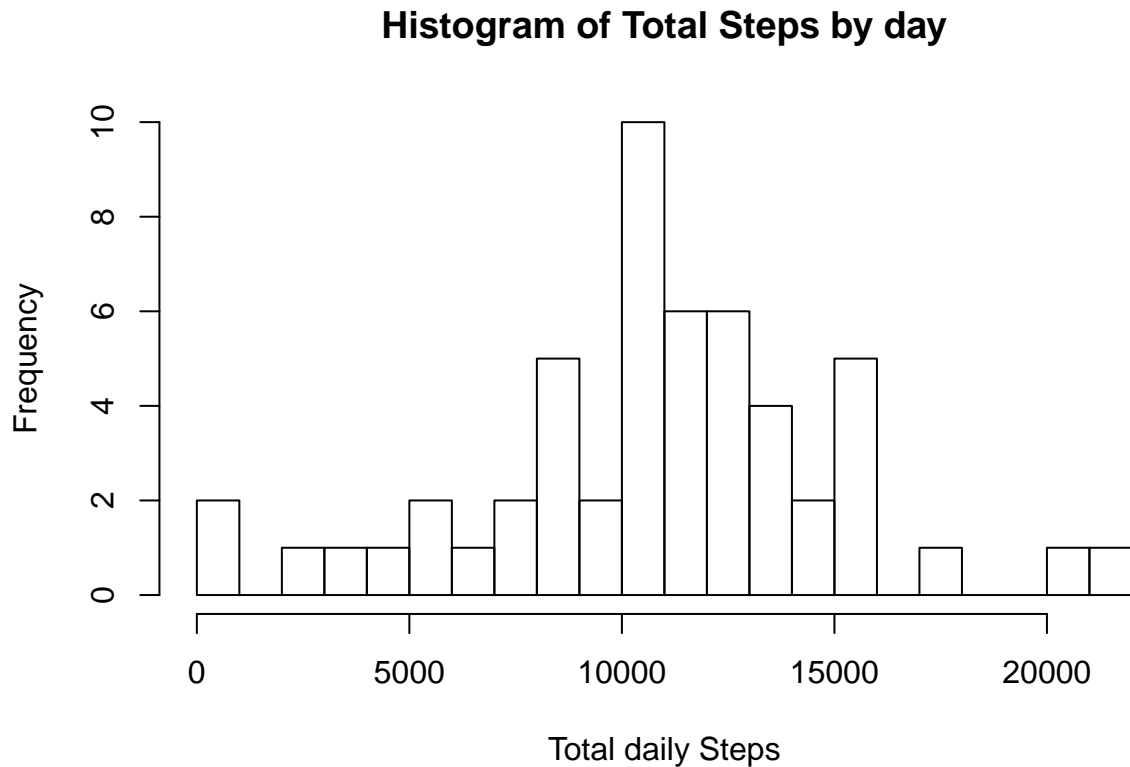
```
##   intersect, setdiff, setequal, union
```

```

databydate <- stepdata %>%
  select(date, steps) %>%
  group_by(date) %>%
  summarize(tsteps= sum(steps)) %>%
  na.omit()

hist(databydate$tsteps, xlab = "Total daily Steps",main="Histogram of Total Steps by day", breaks = 20)

```



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

```
mean(databydate$tsteps)
```

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

```
median(databydate$tsteps)
```

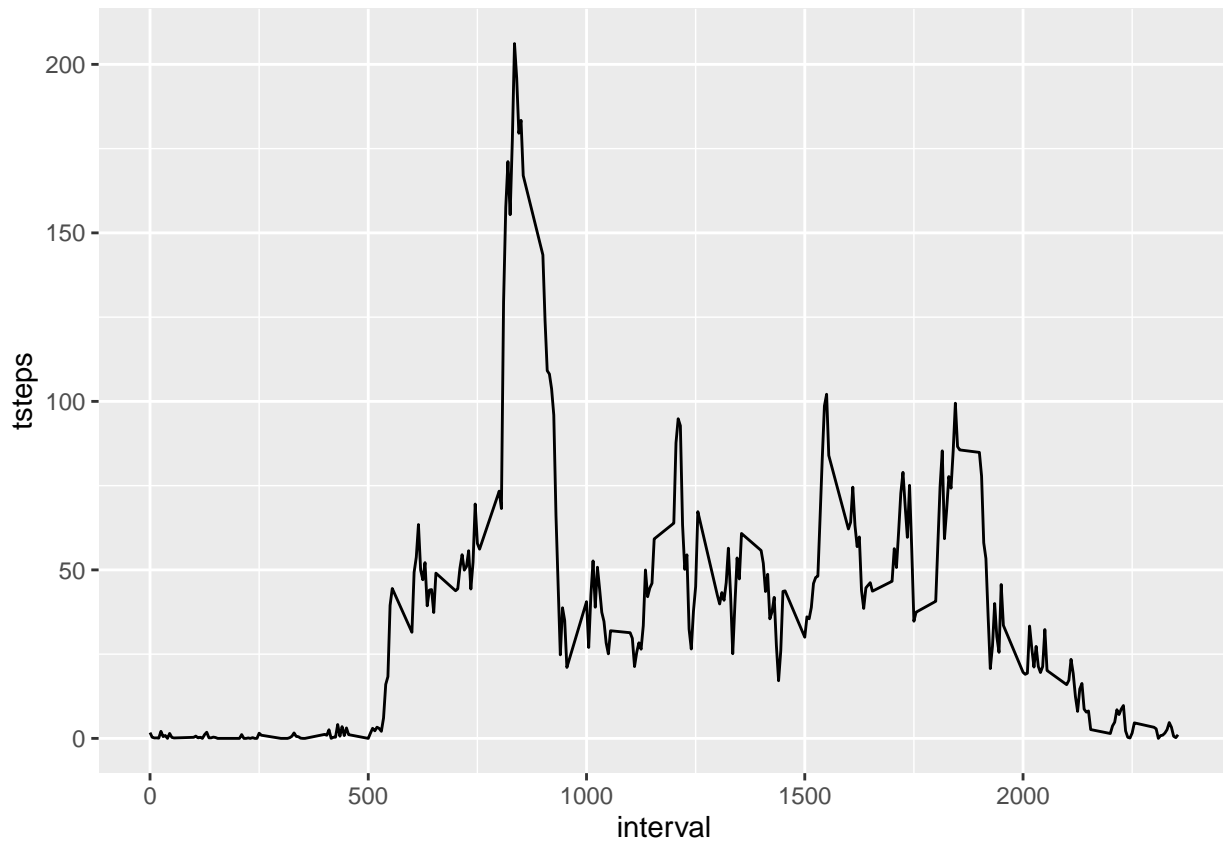
```
## [1] 10765
```

What is the average daily activity pattern?

```
library(ggplot2)
```

```
## Warning: package 'ggplot2' was built under R version 3.6.2
```

```
databyinterval <- stepdata%>%  
  select(interval, steps) %>%  
  na.omit() %>%  
  group_by(interval) %>%  
  summarize(tsteps= mean(steps))  
  
ggplot(databyinterval, aes(x=interval, y=tsteps))+ geom_line()
```



The 5-minute interval that, on average, contains the maximum number of steps

```
databyinterval[which(databyinterval$tsteps== max(databyinterval$tsteps)),]
```

```
## # A tibble: 1 x 2  
##   interval tsteps  
##   <int>   <dbl>  
## 1     835    206.
```

Looking for missing values

```
missingVals <- sum(is.na(data))
```

```
## Warning in is.na(data): is.na() applied to non-(list or vector) of type
## 'closure'
```

```
table(missingVals)
```

```
## missingVals
## 0
## 1
```

```
#Checking and Replacing all the missing Values
```

```
library(magrittr)
library(dplyr)

replacewithmean <- function(x) replace(x, is.na(x), mean(x, na.rm = TRUE))
meandata <- stepdata%>%
  group_by(interval) %>%
  mutate(steps= replacewithmean(steps))
```

```
head(meandata)
```

```
## # A tibble: 6 x 3
## # Groups:   interval [6]
##   steps date      interval
##   <dbl> <fct>      <int>
## 1 1.72  2012-10-01         0
## 2 0.340 2012-10-01         5
## 3 0.132 2012-10-01        10
## 4 0.151 2012-10-01        15
## 5 0.0755 2012-10-01       20
## 6 2.09  2012-10-01       25
```

```
#Make a histogram of the total number of steps taken each day
```

```
FullSummedDataByDay <- aggregate(meandata$steps, by=list(meandata$date), sum)
```

```
names(FullSummedDataByDay)[1] = "date"
names(FullSummedDataByDay)[2] = "totalsteps"
head(FullSummedDataByDay, 15)
```

```
##           date totalsteps
## 1 2012-10-01  10766.19
## 2 2012-10-02   126.00
## 3 2012-10-03  11352.00
## 4 2012-10-04  12116.00
## 5 2012-10-05  13294.00
## 6 2012-10-06  15420.00
## 7 2012-10-07  11015.00
## 8 2012-10-08  10766.19
## 9 2012-10-09  12811.00
## 10 2012-10-10   9900.00
```

```
## 11 2012-10-11 10304.00
## 12 2012-10-12 17382.00
## 13 2012-10-13 12426.00
## 14 2012-10-14 15098.00
## 15 2012-10-15 10139.00
```

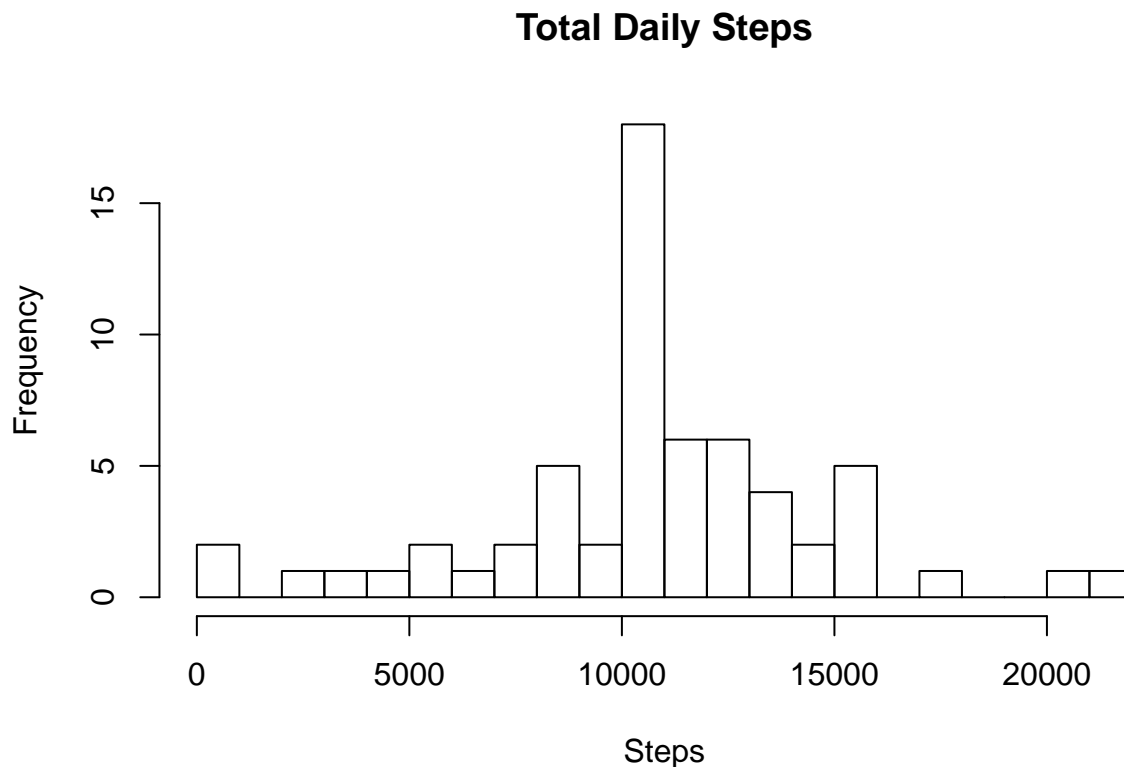
Summary of new data : mean & median

```
summary(FullSummedDataByDay)
```

```
##          date      totalsteps
## 2012-10-01: 1  Min.   :   41
## 2012-10-02: 1  1st Qu.: 9819
## 2012-10-03: 1  Median :10766
## 2012-10-04: 1  Mean   :10766
## 2012-10-05: 1  3rd Qu.:12811
## 2012-10-06: 1  Max.   :21194
## (Other)      :55
```

Making a histogram

```
hist(FullSummedDataByDay$totalsteps, xlab = "Steps", ylab = "Frequency", main = "Total Daily Steps", br
```



#Compare the mean and median of Old and New data

```
oldmean <- mean(databydate$steps, na.rm = TRUE)
newmean <- mean(FullSummedDataByDay$totalsteps)
```

```
oldmean
```

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

```
newmean
```

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

```
oldmedian <- median(databydate$steps, na.rm = TRUE)
newmedian <- median(FullSummedDataByDay$totalsteps)
```

```
oldmedian
```

```
## [1] 10765
```

```
newmedian
```

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

Mean and median values are higher after imputing missing data. The reason is that in the original data, there are some days with steps values NA for any interval. The total number of steps taken in such days are set to 0s by default. However, after replacing missing steps values with the mean steps of associated interval value, these 0 values are removed from the histogram of total number of steps taken each day.

Are there differences in activity patterns between weekdays and weekends?

```
meandata$date <- as.Date(meandata$date)
meandata$weekday <- weekdays(meandata$date)
meandata$weekend <- ifelse(meandata$weekday=="Saturday" | meandata$weekday=="Sunday", "Weekend", "Weekday")

library(ggplot2)
meandataweekendweekday <- aggregate(meandata$steps , by= list(meandata$weekend, meandata$interval), na.rm=T)
names(meandataweekendweekday) <- c("weekend", "interval", "steps")

ggplot(meandataweekendweekday, aes(x=interval, y=steps, color=weekend)) + geom_line() +
  facet_grid(weekend ~.) + xlab("Interval") + ylab("Mean of Steps") +
  ggtitle("Comparison of Average Number of Steps in Each Interval")
```

