

# PA1\_template.Rmd

Saydaliev

4/1/2020

#1 Code for reading in the dataset and/or processing the data

```
library("data.table")  
library(ggplot2)
```

```
activityDT <- data.table::fread(input = "activity.csv")
```

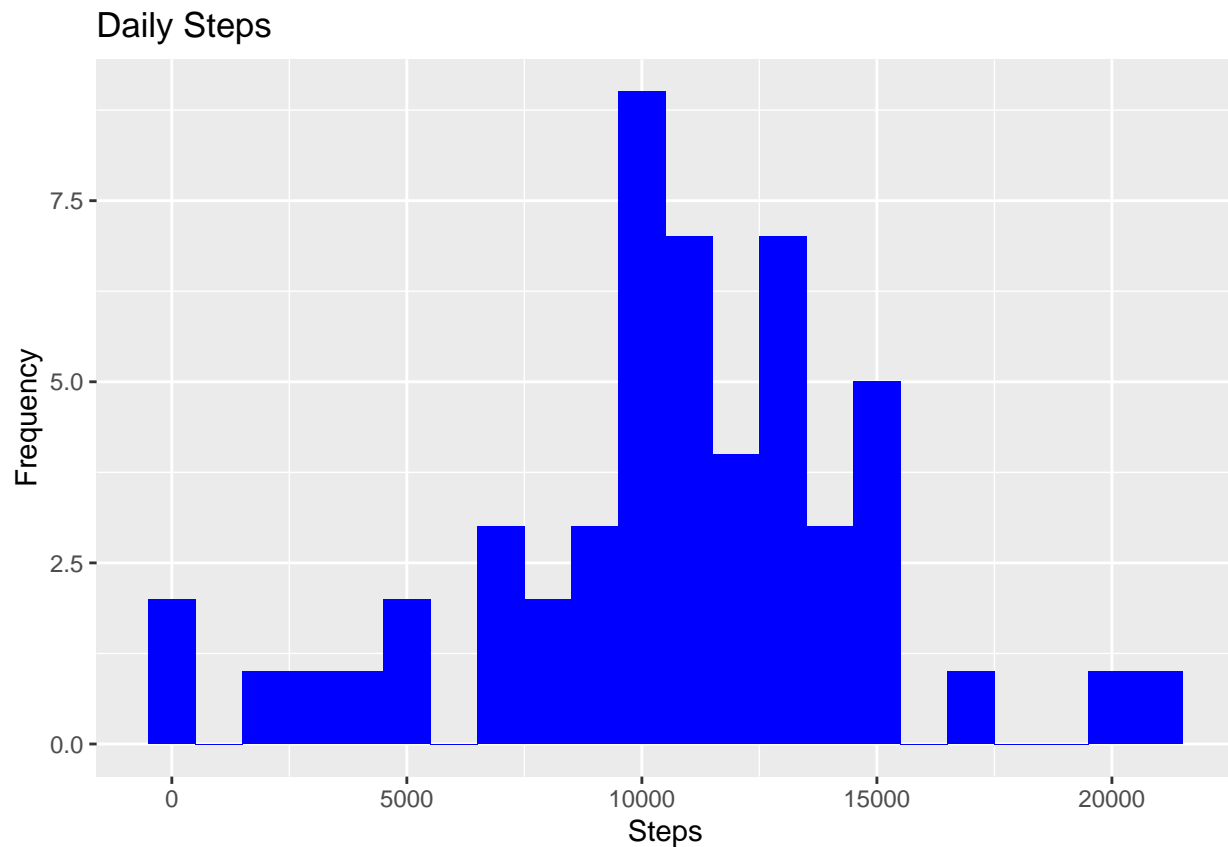
#2 Histogram of the total number of steps taken each day

```
Total_Steps <- activityDT[, c(lapply(.SD, sum, na.rm = FALSE)), .SDcols = c("steps"), by = .(date)]  
head(Total_Steps, 10)
```

```
##           date steps  
##  1: 2012-10-01    NA  
##  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  
##  7: 2012-10-07 11015  
##  8: 2012-10-08    NA  
##  9: 2012-10-09 12811  
## 10: 2012-10-10  9900
```

```
ggplot(Total_Steps, aes(x = steps)) +  
  geom_histogram(fill = "blue", binwidth = 1000) +  
  labs(title = "Daily Steps", x = "Steps", y = "Frequency")
```

```
## Warning: Removed 8 rows containing non-finite values (stat_bin).
```



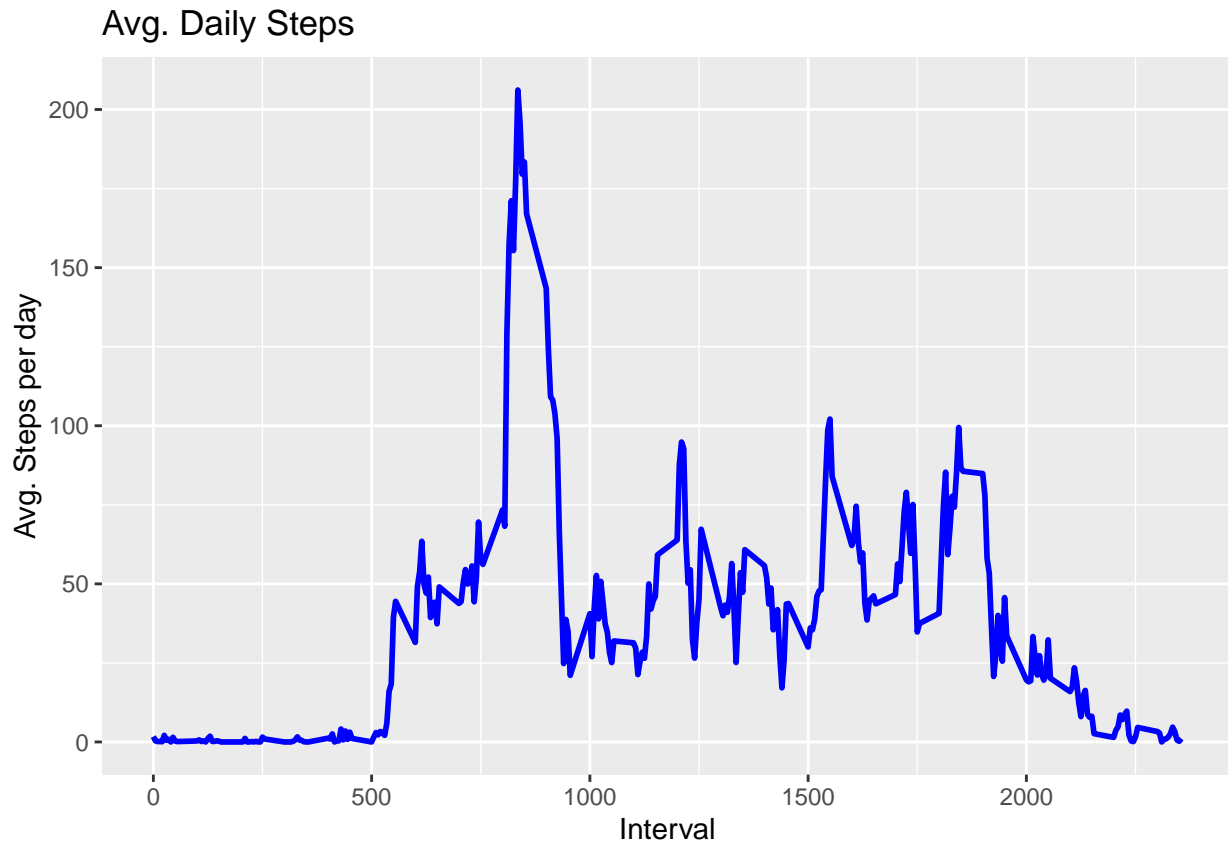
#3 Mean and median number of steps taken each day

```
Total_Steps[, .(Mean_Steps = mean(steps, na.rm = TRUE), Median_Steps = median(steps, na.rm = TRUE))]
```

```
##      Mean_Steps Median_Steps
## 1:    10766.19      10765
```

#4 Time series plot of the average number of steps taken

```
IntervalDT <- activityDT[, c(lapply(.SD, mean, na.rm = TRUE)), .SDcols = c("steps"), by = .(interval)]
ggplot(IntervalDT, aes(x = interval , y = steps)) + geom_line(color="blue", size=1) + labs(title = "Avg
```



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

```
IntervalDT[steps == max(steps), .(max_interval = interval)]
```

```
##      max_interval
## 1:              835
```

#6 Code to describe and show a strategy for imputing missing data

```
activityDT[is.na(steps), .N ]
```

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

#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.

```
activityDT[is.na(steps), "steps"] <- activityDT[, c(lapply(.SD, median, na.rm = TRUE)), .SDcols = c("steps")]
```

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

```
data.table::fwrite(x = activityDT, file = "tidyData.csv", quote = FALSE)
```

#7 Histogram of the total number of steps taken each day after missing values are imputed

## Total number of steps taken per day

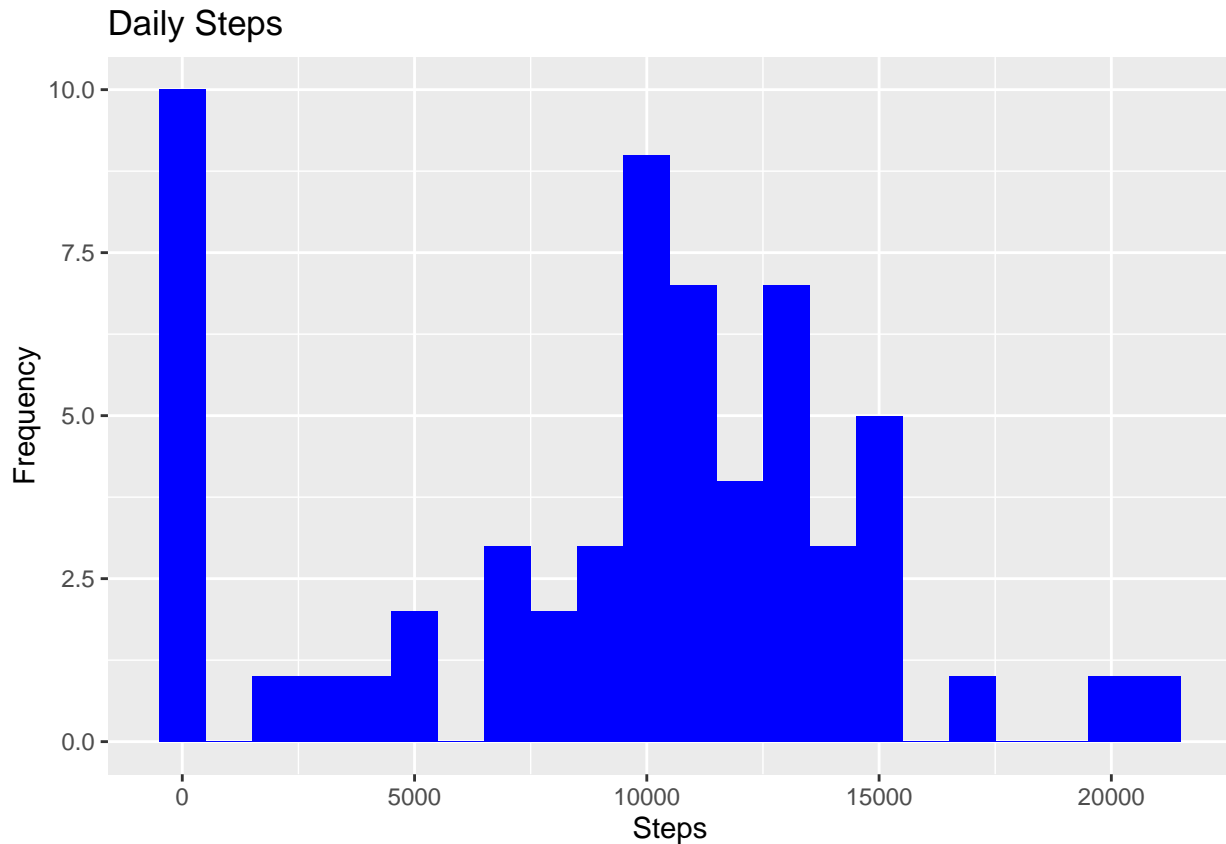
```
Total_Steps <- activityDT[, c(lapply(.SD, sum)), .SDcols = c("steps"), by = .(date)]
```

## mean and median total number of steps taken per day

```
Total_Steps[, .(Mean_Steps = mean(steps), Median_Steps = median(steps))]
```

```
##      Mean_Steps Median_Steps
## 1:      9354.23      10395
```

```
ggplot(Total_Steps, aes(x = steps)) + geom_histogram(fill = "blue", binwidth = 1000) + labs(title = "Daily Steps")
```



## 8 Panel plot comparing the average number of steps taken per 5-minute interval across weekdays and weekends

```
activityDT <- data.table::fread(input = "activity.csv")
activityDT[, date := as.POSIXct(date, format = "%Y-%m-%d")]
activityDT[, `Day of Week` := weekdays(x = date)]
activityDT[grepl(pattern = "Monday|Tuesday|Wednesday|Thursday|Friday", x = `Day of Week`), "weekday or weekend"] <- "weekday"
activityDT[grepl(pattern = "Saturday|Sunday", x = `Day of Week`), "weekday or weekend"] <- "weekend"
activityDT[, `weekday or weekend` := as.factor(`weekday or weekend`)]
head(activityDT, 10)
```

```
##      steps      date interval Day of Week weekday or weekend
## 1:      NA 2012-10-01         0    Monday      weekday
## 2:      NA 2012-10-01         5    Monday      weekday
## 3:      NA 2012-10-01        10    Monday      weekday
## 4:      NA 2012-10-01        15    Monday      weekday
## 5:      NA 2012-10-01        20    Monday      weekday
```

```
## 6:    NA 2012-10-01    25    Monday    weekday
## 7:    NA 2012-10-01    30    Monday    weekday
## 8:    NA 2012-10-01    35    Monday    weekday
## 9:    NA 2012-10-01    40    Monday    weekday
## 10:   NA 2012-10-01    45    Monday    weekday
```

```
#9
```

```
activityDT[is.na(steps), "steps"] <- activityDT[, c(lapply(.SD, median, na.rm = TRUE)), .SDcols = c("steps")]
IntervalDT <- activityDT[, c(lapply(.SD, mean, na.rm = TRUE)), .SDcols = c("steps"), by = .(interval, `weekday or weekend`)]
ggplot(IntervalDT, aes(x = interval, y = steps, color = `weekday or weekend`)) + geom_line() + labs(title = "Avg. Daily Steps by Weektype")
```

