

# UCSL: R Challenge 1\*

Christopher Prince

07/09/2015

## 1 Summary statistics

### 1.1 Complete January 2015 data set

For the first part of the challenge, we are asked to provide the basic descriptive statistics for the Citibike data set from January 2015. We begin by importing the .csv data into a data frame named `cbdata`:

```
#import Jan 2015 CitiBike data from csv file, include a header line
fname <- '~/UCSL/R/ch1/201501-citibike-tripdata.csv'
cbdata <- read.csv(fname, header = TRUE)
```

The statistics for the total data set are calculated below.

```
#calculate descriptive stats and save in variables
#tdur is shorthand for tripduration
cbdata.count <- dim(cbdata)[1] #1st element of dim() vector is the number of observations

tdur.mean <- mean(cbdata$tripduration)
tdur.var <- var(cbdata$tripduration)
tdur.sd <- sd(cbdata$tripduration)
tdur.min <- min(cbdata$tripduration)
tdur.max <- max(cbdata$tripduration)
tdur.range <- tdur.max - tdur.min
tdur.median <- median(cbdata$tripduration)
tdur.qtile <- quantile(cbdata$tripduration)
```

`tdur.qtile` is the vector of values associated with the quantiles. After reporting them, we use the fourth and second elements (corresponding to the first and third quartiles) to calculate the *interquartile range*:

```
tdur.qtile
```

```
##      0%    25%    50%    75%   100%
##      60    334    504    772  43023
```

```
tdur.iqr <- tdur.qtile[4]-tdur.qtile[2]
```

symbol	description	value
<code>cbdata.count</code>	number of observations in <code>cbdata</code>	285552
<code>tdur.mean</code>	mean trip duration	654.33
<code>tdur.var</code>	variance of trip duration	811397
<code>tdur.sd</code>	standard deviation of trip duration	900.78
<code>tdur.median</code>	median trip duration	504

\*Source for this document is at <https://github.com/cmprince/UCSL/blob/master/R/ch1/ch1.Rmd>

symbol	description	value
tdur.min	minimum trip duration	60
tdur.max	maximum trip duration	43023
tdur.range	total range of trip duration values	42963
tdur.iqr	interquartile range of trip duration values	438

## 1.2 Removing outliers

Outliers are then removed by *z*-score. Observations with a *z*-score greater than 3 (more than three standard deviations from the mean) are removed from the set using the `subset` function:

```
#define z-score
z_sc <- (cbdata$tripduration-tdur.mean)/tdur.sd
#creating subset of data for z_sc<3
cbdata.z3 <- subset(cbdata, z_sc<3)
```

The same statistics are calculated for the subset:

```
# Computing same calculation for data without outliers

cbdata.z3.count <- dim(cbdata.z3)[1]
tdur.z3.mean <- mean(cbdata.z3$tripduration)
tdur.z3.var <- var(cbdata.z3$tripduration)
tdur.z3.sd <- sd(cbdata.z3$tripduration)
tdur.z3.min <- min(cbdata.z3$tripduration)
tdur.z3.max <- max(cbdata.z3$tripduration)
tdur.z3.range <- tdur.z3.max - tdur.z3.min
tdur.z3.median <- median(cbdata.z3$tripduration)
tdur.z3.qtile <- quantile(cbdata.z3$tripduration)

tdur.z3.qtile
```

```
##  0%  25%  50%  75% 100%
##  60  333  502  766 3355
```

```
tdur.z3.iqr <- tdur.z3.qtile[4]-tdur.z3.qtile[2]
```

symbol	description	value
cbdata.z3.count	number of observations in <code>cbdata.z3</code>	284255
tdur.z3.mean	mean trip duration	616.47
tdur.z3.var	variance of trip duration	177353
tdur.z3.sd	standard deviation of trip duration	421.13
tdur.z3.median	median trip duration	502
tdur.z3.min	minimum trip duration	60
tdur.z3.max	maximum trip duration	3355
tdur.z3.range	total range of trip duration values	3295
tdur.z3.iqr	interquartile range of trip duration values	433

## 1.3 Discussion

The **central tendency** for the data in January, after removing the  $z > 3$  outliers, is that the average (mean) trip duration is slightly more than 10 minutes (616.47 sec). Half of the trips took less than the median time of 502 sec, about 8 1/2 minutes.

The **dispersion** of the data in January indicate that the middle half of all rides (again, after removing the  $z > 3$  outliers), given by the IQR spanned a range of 433 sec. The total range for all of the data is a little under an hour, 3295 sec.

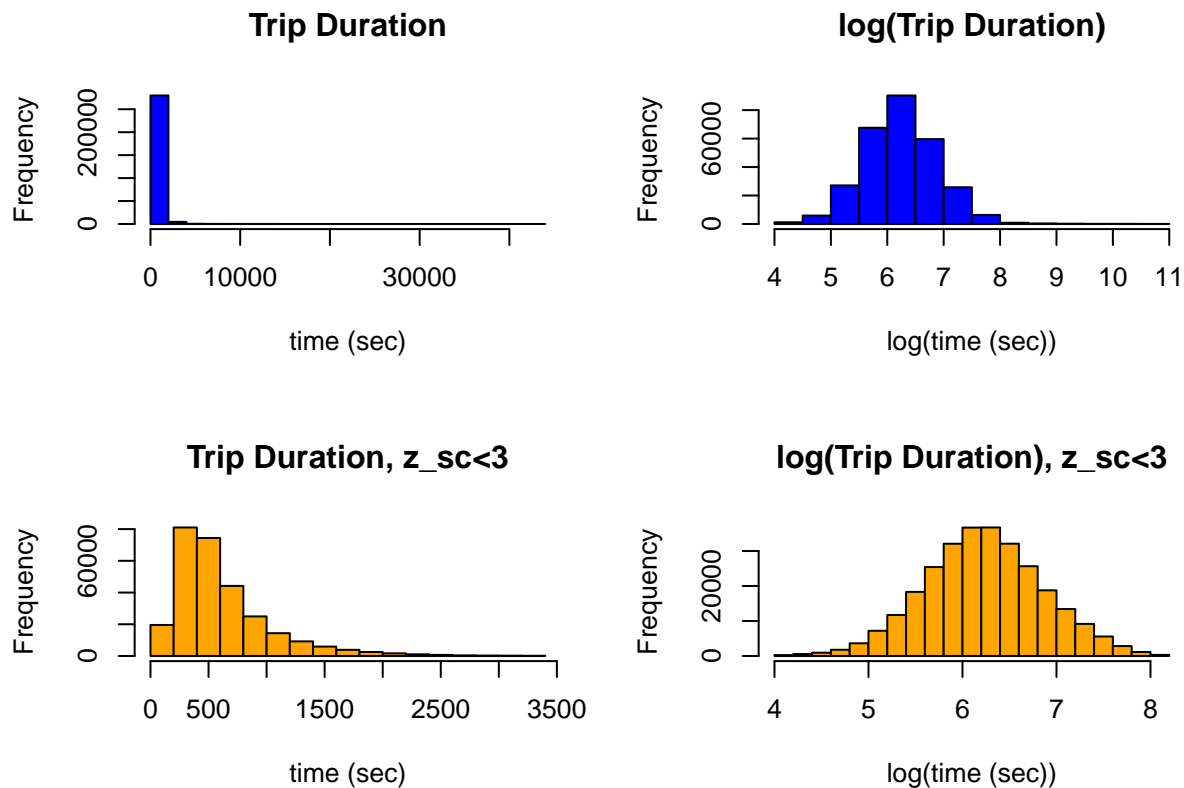
**Removing outliers** had little effect on the quartile calculations: the median shifted just `tdur.median - tdur.z3.median = 2` sec, and the IQR was reduced by only `tdur.iqr - tdur.z3.iqr = 5` sec. However, the mean shifted by a significant amount, `tdur.mean - tdur.z3.mean = 37.85` sec. This is due to removing `cbdata.count - cbdata.z3.count = 1297` data points skewing the mean. The overall range reduced from 42963 sec to 3295 sec.

## 2 Visualization

### 2.1 Histograms for `cbdata` and `cbdata.z3`

Here we plot histograms for both the full set and  $z$ -score reduced set. The log-transformed data is also plotted, which is particularly useful for the full data set.

```
# histograms of tdur and log(tdur) for all data points and the z_sc<3 subset
par(mfrow = c(2, 2))
hist(cbdata$tripduration, main = "Trip Duration", xlab = "time (sec)", col = "blue")
hist(log(cbdata$tripduration), main = "log(Trip Duration)", xlab = "log(time (sec))",
     col = "blue")
hist(cbdata.z3$tripduration, main = "Trip Duration, z_sc<3", xlab = "time (sec)",
     col = "orange")
hist(log(cbdata.z3$tripduration), main = "log(Trip Duration), z_sc<3", xlab = "log(time (sec))",
     col = "orange")
```



Visually, the log-transformed data appears to fit a normal distribution, though other tests (which we'll no doubt learn and are out of this assignment's scope) can assess this.

## 2.2 Visualization by usertype

First we use the `subset` function to parse out new data frames by `usertype`. There are two `usertypes`, `Subscriber` and `Customer`. After creating the subsets for both the full data set and the  $z$ -score reduced set, we count the observations and verify that we haven't missed any blanks or mislabeled points.

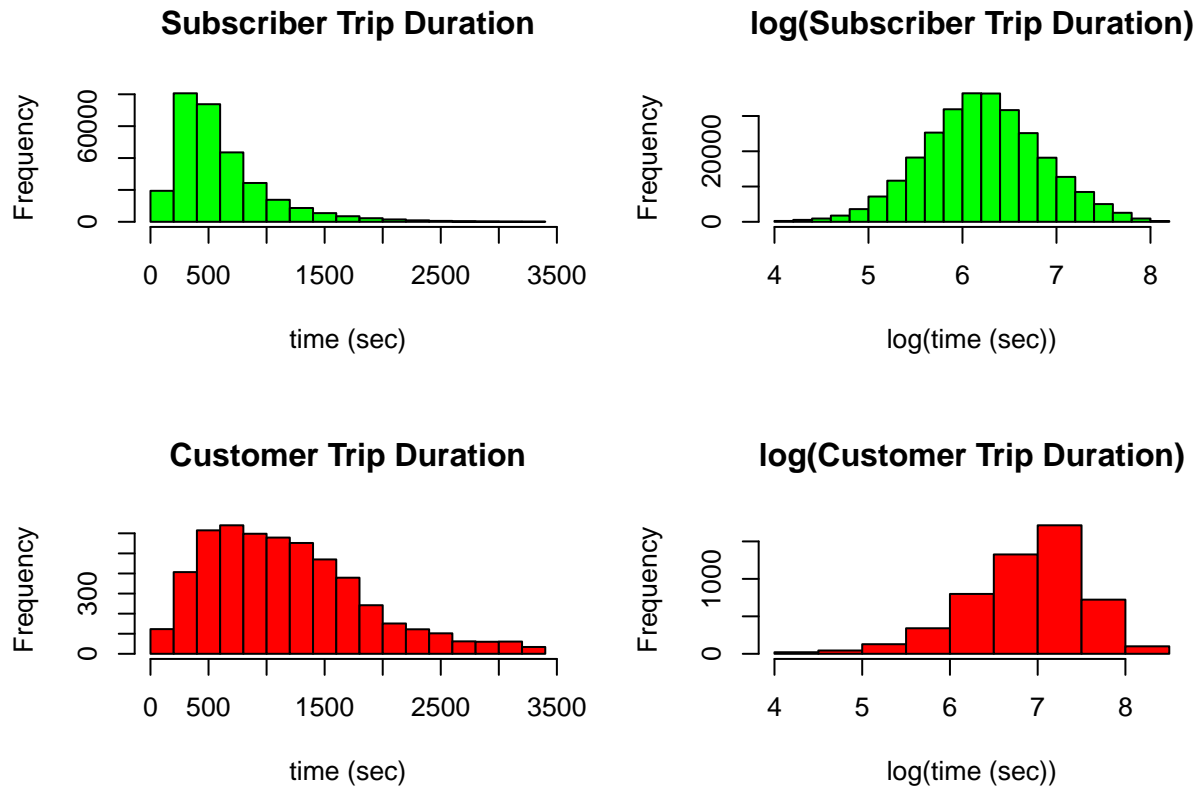
```
#create subsets by usertype
cbdata.subs <- subset(cbdata, cbdata$usertype=='Subscriber')
cbdata.cust <- subset(cbdata, cbdata$usertype=='Customer')
cbdata.z3.subs <- subset(cbdata.z3, cbdata.z3$usertype=='Subscriber')
cbdata.z3.cust <- subset(cbdata.z3, cbdata.z3$usertype=='Customer')

#count the elements of the subsets
cbdata.cust.count <- dim(cbdata.cust)[1]
cbdata.subs.count <- dim(cbdata.subs)[1]
cbdata.z3.cust.count <- dim(cbdata.z3.cust)[1]
cbdata.z3.subs.count <- dim(cbdata.z3.subs)[1]
```

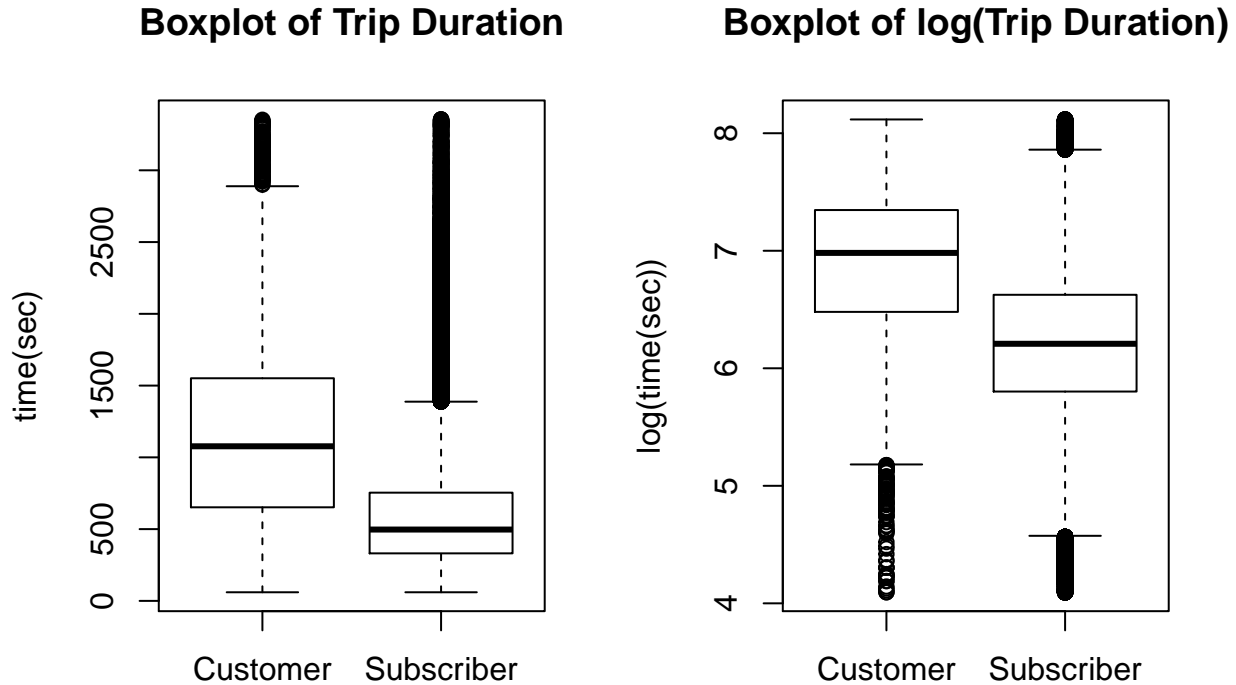
	cbdata	cbdata.z3
subset sums	285552	284255
total counts	285552	284255

Now we visualize the subsets by producing histograms and boxplots.

```
# histograms of tdur by usertype for z_sc<3 subsets
par(mfrow = c(2, 2))
hist(cbddata.z3.subs$tripduration, main = "Subscriber Trip Duration", xlab = "time (sec)",
     col = "green")
hist(log(cbddata.z3.subs$tripduration), main = "log(Subscriber Trip Duration)",
     xlab = "log(time (sec))", col = "green")
hist(cbddata.z3.cust$tripduration, main = "Customer Trip Duration", xlab = "time (sec)",
     col = "red")
hist(log(cbddata.z3.cust$tripduration), main = "log(Customer Trip Duration)",
     xlab = "log(time (sec))", col = "red")
```



```
# and boxplots of same
par(mfrow = c(1, 2))
boxplot(tripduration ~ usertype, data = cbddata.z3, main = "Boxplot of Trip Duration",
        ylab = "time(sec)")
boxplot(log(tripduration) ~ usertype, data = cbddata.z3, main = "Boxplot of log(Trip Duration)",
        ylab = "log(time(sec))")
```



## 2.3 Discussion

Before answering the posed questions, we note that the shapes of the histograms for the two different usertypes are markedly different. Indeed, the log-transform of the **Subscriber** subset has the same normal-looking shape as the outlier-removed data at large. This is not surprising since **Subscribers** account for 98.03% of the data. The histogram of the **Customer** subset, however, has a much different shape, and the log-transform does *not* appear normally-shaped. This suggests that the trip patterns of **Customers** and **Subscribers** are significantly different.

Outliers in data sets will create longer tails in histograms and more data points beyond the boxplot whiskers. To visualize the complete data set, there will be a loss of resolution in both types of graphs.

In a histogram, there will be many sparsely populated bins in the tails, with most of the observations piled into just a few bins around the median.

In a boxplot, the outliers will dominate the axis along which the values are plotted due to their range. Effectively this compresses the IQR into a smaller space on the graph, making the visualization less effective.

From the discussion of the descriptive statistics above, we can hypothesize that the histograms and boxplots for the data sets including outliers will suffer the effects identified above.