

# Data PreProcessing with R

## Normalising Data

Let us examine how we can generate summary statistics for a variable. You will also see how to carry out Min-Max normalisation and Z Score Standardisation using R.

**#Input dataset numeric weather into a dataframe (from last week)**

```
weathernum <- read.table(file = "c:/weathernumeric.txt", stringsAsFactors=FALSE, sep = ",", header=TRUE)
```

**#Show the first 10 records and all the columns**

```
weathernum[1:10,]
```

**#Examine the 5 number summary statistics with mean for temperature and humidity. From this calculate the interquartile range**

```
summary(weathernum$temperature)
```

```
summary(weathernum$humidity)
```

**For you to do!**

- By hand, calculate the IQR for these dimensions. Are there any outliers for these variables?
- Use the IQR method to detect outliers.

**#Min-Max normalisation for the weathernum\$humidity attribute**

```
mmnorm.humidity <- (weathernum$humidity - min(weathernum$humidity))/(max(weathernum$humidity) - min(weathernum$humidity))
```

```
summary(mmnorm.humidity)
```

**#zScore Standardisation**

```
zscore.humidity <- (weathernum$humidity - mean(weathernum$humidity))/sd(weathernum$humidity)
```

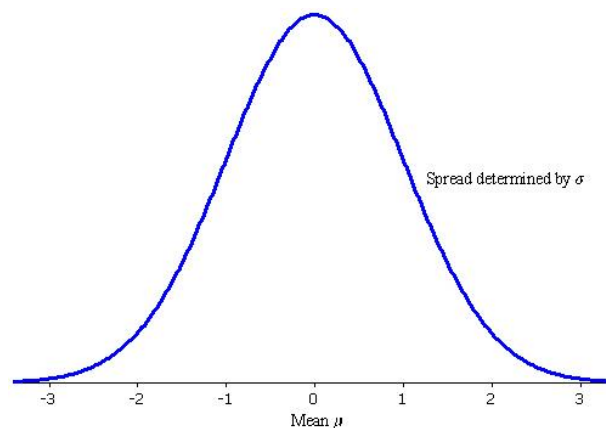
```
summary(zscore.humidity)
```

**For you to do!**

- What are the ranges for this variable for both normalisation methods? Do you think there are any outliers?
- Write the R formula to decimal scale weathernum\$humidity attribute. Hint use nchar function. Examine the summary statistics afterwards. Note any observations

## Transforming the data in an attempt to achieve normality#

Some data mining and statistical methods require that the variables be normally distributed. The normal distribution is a continuous probability distribution commonly known as the bell curve, which is symmetric. It is centred at the mean and spread is determined by the standard deviation. Standard normal Z distribution has a mean of 0 and a standard deviation of 1.



Standard normal Z-distribution  
with  $\mu=0$  and  $\sigma=1$

However, variables are often skewed to the right (+) or left (-). We can measure use the following statistic to measure the skewness of a distribution

$$Skewness = \frac{3(\text{mean} - \text{median})}{\text{standard deviation}}$$

For perfectly symmetric data, the mean, median and mode = 0

We can transform the distribution in different ways in an attempt to reduce skewness and achieve a symmetric distribution.

```
#Read in the cars2 dataset and transform weightlbs in different ways
```

```
cars2 <- read.csv(file="c:/cars2.txt", stringsAsFactors=TRUE)
```

```
#Natural Log Transformation
```

```
natlog.weightlbs <- log(cars2$weightlbs)
```

```
natlog.weightlbs
```

```
#Square Root Transformation
```

```
sqrt.weightlbs <- sqrt(cars2$weightlbs)
```

```
sqrt.weightlbs
```

### #Inverse Square Root Transformation

```
invsqrt.weightlbs <- 1/sqrt(cars2$weightlbs)
```

```
invsqrt.weightlbs
```

For you to do – comment on the skewness!!

- Using R calculate the skewness for the car weightlbs attribute
- Using R find the skewness for the natural log transformation of the car weightlbs attribute
- Using R find the skewness for the square root transformation of the car weightlbs attribute
- Using R find the skewness for the inverse square root transformation of the car weightlbs attribute
- Now re-calculate skewness using 3 (normality) transformation methods but with car weightlbs transformed using Min-Max Normalisation
- Interpret your results. What do you concur?

## Binning Data

Let us look at equal frequency , equal- binning and k-means clustering as binning strategies  
Here are some examples

### #Create a vector of values for binning

```
xdata <-c(1,11,2,1,1,2,12,11,44,2,12,1)
```

### #get the sample size of the variable

```
n <- length(xdata)
```

### #Declare number of bins and bin indicator

```
nbins <- 3
```

```
whichbin <- c(rep(0,n))
```

### #Equal Frequency (equal-depth)

```
freq <- n/nbins
```

```
#sort the data
```

```
xsorted <- sort(xdata)
```

```
for(i in 1:nbins) {
```

```
  for(j in 1:n) {
```

```
    if ((i-1)*freq < j && j <= i*freq)
```

```
      whichbin[j] <- i
```

```
  }
```

```
}
```

```
whichbin
```

```
xsorted
```

```

#equal-width
#note bin 2 has 0 values while bin 3 has 1 value
range.xdata <- max(xdata) - min(xdata) +1
binwidth <- range.xdata/nbins
for(i in 1:nbins) {
  for(j in 1:n) {
    if ((i-1)*binwidth < xdata[j] && xdata[j] <= (i)*binwidth)
      whichbin[j] <- i
  }
}
whichbin
xdata

#K-means clustering as a binning strategy where k = 3
kmeansclustering <- kmeans(xdata,centers=nbins)
whichbin <- kmeansclustering$cluster
whichbin

# What is the best binning strategy from above? Explain your reasoning.

```

## How to create a Histogram and a Scatteplot

```

#Read in the cars2 dataset if you have not done so already
cars2 <- read.csv(file="c:/cars2.txt", stringsAsFactors=TRUE)

#Create a histogram
#Set up the plot area
par(mfrow=c(1,1))

#Examine the weightlbs histogram and write down your observations
hist(cars2$weightlbs,
      breaks=30,
      xlim= c(0,5000),
      col="blue",
      border="black",
      ylim=c(0,40),
      xlab="Weight in lbs",
      ylab="Counts",
      main="Histogram of Car Weights")
box(which="plot", lty="solid", col="black")

```

### FOR YOU TO DO!

Create a histogram for the variable MPG

#Use the following command before creating the weight and z score of weight histograms if you want to see them size by side.

```
par(mfrow = c(1,2))
```

#Let us create a ScatterPlot of MPG by Weight. Examine the plot. What do you concur? Are there outliers? Would they be deemed outliers for the individual variables.

```
plot(cars2$weightlbs, cars2$mpg,  
     xlim= c(0,5000),  
     ylim=c(0,600),  
     xlab="Weight",  
     ylab="MPG",  
     main="Scatter Plot of MPG by Weight",  
     type ="p",  
     pch=16,  
     col="blue")  
points(cars2$weightlbs,  
       cars2$mpg,  
       type="p",  
       col="red")
```

### FOR YOU TO DO!

Choose 2 other numeric variables and produce a scatter plot. Can you identify any outliers? Are they correlated?

# Create a Histogram with fitted Normal Distribution and Normal Probability Plot

## # A histogram inverse square root of weight

```
invsqrt.weightlbs <- 1/sqrt(cars2$weightlbs)
invsqrt.weightlbs

x <- rnorm(1000000,
mean=mean(invsqrt.weightlbs),
sd=sd(invsqrt.weightlbs))

par(mfrow=c(1,1))
hist(invsqrt.weightlbs,
breaks=30,
xlim= c(0.0125,0.0275),
col="lightblue",
prob= TRUE,
border="black",
xlab="Inverse Square Root of Weightlbs",
ylab="Counts",
main="Histogram of Inverse Square Root of Car Weightlbs")
lines(density(x),
col="red")

box(which="plot",
lty="solid",
col="black")

lines(density(x),
col="red")
```

## # Normal Probability plot that indicates non-normality

```
par(mfrow= c(1,1))

qqnorm(invsqrt.weightlbs,
datax =TRUE,
col="red",
ylim=c(0.01,0.03),
main ="Normal Q-Q Plot of inverse Square Root of Weightlbs")

qqline(invsqrt.weightlbs,
col="blue",
datax=TRUE)
```

## We can also test for normality is a formal way.

Using the Shapiro-Wilks test, we can examine the p-value. The p-value tells you what the chances are that the sample comes from a normal distribution. The lower this value, the smaller the chance it comes from a normal distribution. Statisticians typically use a value of 0.05 as a cutoff. When the p-value is less than 0.05, one can conclude that the sample deviates from normality.

```
result <- shapiro.test( invsqrt.weightlbs)
```

```
result$p.value
```

## A Note of transforming the data and duplicate records

If you have transformed the data using Z-Score, Min-Max Normalisation, Inverse Square Root Transform etc, you will have to detransform the data at some stage to see the actual values here is an example

### Transforming the data

```
x <- cars2$weightlbs[1]; x
```

```
[1] 4209
```

```
h1 <- head(cars2$weightlbs); h1
```

```
[1] 4209 1925 3449 3761 2051 3900
```

### #Transform x using the inverse sqrt

```
y <- 1/sqrt(x); y
```

```
[1] 0.01541383
```

```
z <- 1/sqrt(h1); z
```

```
[1] 0.01541383 0.02279212 0.01702760 0.01630603 0.02208092
```

```
[6] 0.01601282
```

### # Detransform x using $1/(y)^2$

```
detransx <- 1/y^2; detransx
```

```
[1] 4209
```

```
detransz <- 1/z^2; detransz
```

```
[1] 4209 1925 3449 3761 2051 3900
```

### Finding duplicate Records in a data frame

```
#Finding Duplicate Records in a data frame
```

```
anyDuplicated(cars2)
```

```
[1] 0
```

```
duplicated(cars2)
```

```
# Note: TRUE Indicates a record which is a duplicate of a previous record
```

```
# FALSE indicates a record which is not a duplicate of a previous record
```

```
# Duplicate the first record to make a new dataset
```

```
new.cars2 <- rbind(cars2,cars2[1,])
```

```
anyDuplicated(new.cars2)
```

```
[1] 262
```

```
# The 262nd record is duplicated
```

```
duplicated(new.cars2)
```

```
#TRUE Indicates a record which is a duplicate of a previous record
```

## Hands-On Analysis

Using the churn dataset (churn.txt)

1. Explore whether there are any missing values for any of the variables
2. Use a graph to visually determine whether there are any outliers among the number of calls to customer service
3. Identify the range of customer service calls that should be considered outliers using:
  - i. Z-Score Method
  - ii. IQR
4. Transform the day minutes attribute using Z-Score standardisation
5. Work with skewness as follows.
  - i. Calculate the skewness of day minutes
  - ii. Then calculate the skewness of the Z-score standardised day minutes.  
Comment
  - iii. Based on the skewness value, would you consider day minutes to be skewed or nearly perfectly symmetric