

## R Reference Card

**Help for R studio** Open the cheat sheet for RStudio by selecting the Help menu -> Cheatsheets -> RStudio IDE Cheat Sheet. Note the cheatsheet will usually be downloaded in a web browser as a .pdf.

**Installing Packages** Use the argument `dependencies=TRUE` to load all other packages required by the targeted package.

**Loading Packages** After installing a package, we can load it into the R environment using the `library()` or `require()` functions, which more or less do the same thing.

Use `library(dplyr)` The package `dplyr` contains several easy-to-use data management functions that we will learn to use to manage our data.

**Starting R and RStudio** Launch R Studio, R starts automatically within RStudio

**Quitting R Studio** Use

`q()` Type yes or no to confirm choice

**Loading data set (CSV File)** `my_data <- read.csv("path/to/your/data.csv")` #

Load data from a CSV file **Viewing**

**Data:** After loading your data, you can inspect it using the `head()` function to view the first few rows.

**Building a Linear Regression Model:** You can build a linear regression model using the `lm()` function. For example, let's predict miles per gallon (`mpg`) based on weight (`wt`) and horsepower (`hp`) using the `mtcars` dataset.

```
model <- lm(mpg ~ wt + hp, data = mtcars) # Build a linear model
```

**Model Interpretation** The `summary(model)` function provides key statistics, such as coefficients, p-values, R-squared value, etc., which help assess the model's fit and significance.

**Understanding data**

`str()` Displays the structure of the dataset, including variable types.

`summary()`: Provides summary statistics for each variable (e.g., min, max, mean, median).

`head()`: Displays the first few rows of the data.

**Evaluating a model Residuals**

```
residuals <- resid(model) # Get
```

```
residuals from the model plot(residuals)
# Plot residuals
```

**R-squared:** The **R-squared value** from the model summary indicates how well the independent variables explain the variation in the dependent variable (`mpg`). Higher R-squared values (close to 1) indicate a better fit.

**Interacting with RStudio Console:**

This is where you enter R commands. It executes your commands and displays results. You can type code directly here for quick tests or experimentation.

**Environment/History Pane:** Displays the objects (variables, models, etc.) created in your R session and the history of previously executed commands.

**Plots Pane:** This is where graphical outputs (like plots and charts) are displayed. You can save plots using the "Export" button for further use.

**Packages Pane:** View and manage the libraries and packages installed in R. You can install new packages here and load them for use.

## Interacting with R

### Basic Arithmetic

`3 + 2` # Addition

`10 / 2` # Division

`4 * 5` # Multiplication

`2^3` # Exponentiation

**Creating Variables:** assign values to variables using the `<-`

`x <- 10` # Assign the value 10 to x

`y <- 20` # Assign the value 20 to y

### Frequently used commands

`summary()`: Gives an overview of a dataset or model.

`str()`: Displays the structure of an object, such as a dataset or model.

`mean()` and `sd()`: Used to calculate the mean and standard deviation of a variable.

`lm()`: Fits a linear model to data.

Basic Plotting: Use `plot()`

- `length(collection)` returns the number of elements contained in the variable `collection`
- `c(value1, value2, value3)` creates a vector
- `container[i]` selects the *i*'th element from the variable `container`

<https://swcarpentry.github.io/r-novice-inflammation/reference.html#glossary>

## References

Group, I. S. C. (n.d.-b). Introduction to R.

[https://stats.oarc.ucla.edu/stat/data/intro\\_r/intro\\_r\\_interactive.html#\(15\)](https://stats.oarc.ucla.edu/stat/data/intro_r/intro_r_interactive.html#(15))

FOR THE R CHEAT SHEET

## Basic Operation

- `#` this is a comment in R
- Use `x <- 3` to assign a value, 3, to a variable, `x`
- R counts from 1, unlike many other programming languages (e.g., Python)

### 1.7 The R Zone

#### Getting Started with R

# Comments, indents, and semicolons

```
# Anything prefaced by a pound sign (#) is a comment.
# Comments are not executed by R. Instead, they explain what the code is doing.
# Indented code (that is not a comment) will run in R as if it was on one line
# Code separated by semicolons will run as if the code was on separate lines,
# with the semicolon marking the line break
```

# Open a dataset and display the data

```
# Replace "C:/./" with the exact location of the file you want to open
cars <- read.csv(file = "C:/./cars.txt",
stringsAsFactors = FALSE)
cars # To display the whole dataset, type the dataset name
head(cars) # Display the first few records of a dataset
names(cars) # Display variable names of a data frame, one kind of data in R
cars$weight # Look at only the weight variable within data frame cars
```

# Matrices

```
# Create a matrix with three rows, two columns, and every value equal to 0.0
mat <- matrix(0.0, nrow = 3, ncol = 2); mat
colnames(mat) <- c("Var 1", "Var 2") # Give a matrix variable names
colnames(mat) # Display variable names of a matrix
```

# Subset data and declare new variables

```
cars.rsub <- cars[1:50,] # Subset the data by rows
cars.csub <- cars[,1:3] # Subset by columns
cars.rcsub <- cars[c(1,3,5), c(2,4)] # Subset by specific rows and columns
cars.vsub <- cars[which(cars$mpg > 30),] # Subset by a logical condition
# To declare new variables, type the variable name, a left-arrow, then the value
firstletter <- "a"
weight <- cars$weight
```

# Display more than one figure at a time

```
par(mfrow=c(1,1)) # plots one figure; the default setting
par(mfrow=c(2,3)) # plots six figures: three in the top row, three in the bottom row
# Plots will fill the plot space row by row
```

# Download and install an R Package

```
# Example: ggplot2, from Lesson 3. install.packages("ggplot2")
install.packages("ggplot2")
# Pick any CRAN mirror, as shown
# Open the new package
library(ggplot2)
```



## 2.22 The R Zone

### # Read in the Cars and Cars2 datasets

```
cars <- read.csv("C:/../cars.txt",
stringsAsFactors = FALSE)
cars2 <- read.csv("C:/../cars2.txt",
stringsAsFactors = FALSE)
```

### # Missing data

```
# Look at four variables from cars
cars.4var <- cars[, c(1, 3, 4, 8)]
head(cars.4var)
```

```
head(cars.4var)
mpg cubicinches hp brand
1 14.0 350 165 us
2 11.9 89 71 europe
3 17.0 302 140 us
4 15.0 400 130 us
5 20.1 98 63 us
6 23.0 350 121 us
```

```
# Make certain entries missing
cars.4var[2,2] <- cars.4var[4,4] <- NA
head(cars.4var)
```

```
head(cars.4var)
mpg cubicinches hp brand
1 14.0 350 165 us
2 11.9 NA 71 europe
3 17.0 302 140 us
4 15.0 400 130 <NA>
5 20.1 98 63 us
6 23.0 350 121 us
```

```
# Replace missing values with constants
cars.4var[2,2] <- 0
cars.4var[4,4] <- "Missing"
head(cars.4var)
```

```
head(cars.4var)
mpg cubicinches hp brand
1 14.0 350 165 us
2 11.9 0 71 europe
3 17.0 302 140 us
4 15.0 400 130 Missing
5 20.1 98 63 us
6 23.0 350 121 us
```

```
# Replace values with mean and mode
cars.4var[2,2] <-
mean(na.omit(cars.4var$cubicinches))
our_table <- table(cars.4var$brand)
our_mode <- names(our_table)[our_table ==
max(our_table)]
cars.4var[4,4] <- our_mode
head(cars.4var)
```

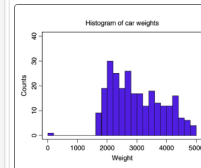
```
head(cars.4var)
mpg cubicinches hp brand
1 14.0 350 165 us
2 11.9 202.1344 71 europe
3 17.0 302 140 us
4 15.0 400 130 us
5 20.1 98 63 us
6 23.0 350 121 us
```

```
# Generate random observations
obs_brand <-
sample(na.omit(cars.4var$brand), 1)
obs_cubicinches <-
sample(na.omit(cars.4var$cubicinches), 1)
cars.4var[2,2] <- obs_cubicinches
cars.4var[4,4] <- obs_brand
head(cars.4var)
```

```
head(cars.4var)
mpg cubicinches hp brand
1 14.0 350 165 us
2 11.9 89 71 europe
3 17.0 302 140 us
4 15.0 400 130 europe
5 20.1 98 63 us
6 23.0 350 121 us
```

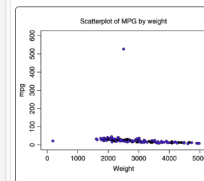
### # Create a Histogram

```
# Set up the plot area
par(mfrow = c(1,1))
# Create the
histogram bars
hist(cars2$weight,
breaks = 30,
xlim = c(0, 5000),
col = "blue",
border = "black",
ylim = c(0, 40),
xlab = "Weight",
ylab = "Counts",
main = "Histogram
of Car Weights")
# Make a box around
# the plot
box(which = "plot",
lty = "solid",
col = "black")
```



### # Create a Scatterplot

```
plot(cars2$weight, cars2$mpg,
xlim = c(0, 5000),
ylim = c(0, 600),
xlab = "Weight",
ylab = "MPG",
main = "Scatterplot
of MPG by
Weight",
type = "p",
pch = 16,
col = "blue")
#Add open black # circles
points(cars2$weight,
cars2$mpg,
type = "p",
col = "black")
```



### # Descriptive Statistics

```
mean(cars$weight) # Mean
median(cars$weight) # Median
length(cars$weight) # Number of observations
sd(cars$weight) # Standard deviation
summary(cars$weight) # Min, Q1, Median, Q3, Max
```

### # Transformations

```
# Min-max normalization
summary(cars$weight)
mi <- min(cars$weight)
ma <- max(cars$weight)
minmax.weight <- (cars$weight - mi)/(ma - mi)
nlimax.weight
```

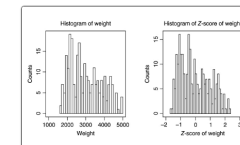
```
# Z-score standardization
s <- sd(cars$weight)
z.weight <- (cars$weight - mi)/s
length(cars$weight)
```

```
# Decimal scaling
d.weight <- cars$weight/(10^4); d.weight
```

### # Side-by-Side Histograms

```
par(mfrow = c(1,2))
# Create two histograms
hist(cars$weight, breaks = 20,
xlim = c(1000, 5000),
main = "Histogram of Weight",
xlab = "Weight",
ylab = "Counts",
box(which = "plot",
lty = "solid",
col = "black"))
```

```
hist(z.weight,
breaks = 20,
xlim = c(-2, 3),
main = "Histogram of Z-
score of Weight",
xlab = "Z-score of Weight",
ylab = "Counts",
box(which = "plot",
lty = "solid",
col = "black"))
```



### # Skewness

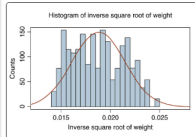
```
(3*(mean(cars$weight) - median(cars$weight))/sd(cars$weight))
(3*(mean(z.weight) - median(z.weight))/sd(z.weight))
```

### # Transformations for Normality

```
sqrt.weight <- sqrt(cars$weight) # Square root
sqrt.weight_skew <- (3*(mean(sqrt.weight) - median(sqrt.weight)) / sd(sqrt.we.
ln.weight <- log(cars$weight) # Natural log
ln.weight_skew <- (3*(mean(ln.weight) - median(ln.weight)) / sd(ln.weight))
invsqrt.weight <- 1 / sqrt(cars$weight) # Inverse square root
invsqrt.weight_skew <- (3*(mean(invsqrt.weight) - median(invsqrt.weight)) / sd
```

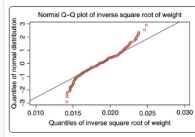
### # Histogram with Normal Distribution Overlay

```
par(mfrow=c(1,1))
x <- rnorm(1000000,
mean = mean(invsqrt.weight),
sd = sd(invsqrt.weight))
hist(invsqrt.weight,
breaks = 30,
xlim = c(0.0125, 0.0275),
col = "lightblue",
prob = TRUE,
border = "black",
xlab = "Inverse Square
Root of Weight",
ylab = "Counts")
main = "Histogram of Inverse Square Root
of Weight")
box(which = "plot",
lty = "solid",
col="black")
# Overlay with Normal density
lines(density(x), col = "red")
```



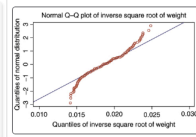
### # Normal Q-Q Plot

```
qqnorm(invsqrt.weight,
datax = TRUE,
col = "red",
ylim = c(0.01, 0.03),
main = "Normal Q-Q Plot of Inverse
Square Root of Weight")
qqline(invsqrt.weight,
col = "blue",
datax = TRUE)
```



### # Normal Q-Q Plot

```
qqnorm(invsqrt.weight,
datax = TRUE,
col = "red",
ylim = c(0.01, 0.03),
main = "Normal Q-Q Plot of Inverse
Square Root of Weight")
qqline(invsqrt.weight,
col = "blue",
datax = TRUE)
```



### # De-transform data

```
# Transform x using y = 1 / sqrt(x)
x <- cars$weight[1]; y <- 1 / sqrt(x)
# Detransform x using x = 1 / (y)^2
detransformedx <- 1 / y^2
x; y; detransformedx
```

```
x: y: detransformedx
[1] 4209
[1] 0.01541363
[1] 4209
```

### # Create indicator variables

```
north_flag <- east_flag <- south_flag <-
c(rep(NA, 10))
region <- c(rep(c("north", "south", "east",
"west"), 2), "north", "south")
# Change the region variable to indicators
for (i in 1:length(region)) {
if(region[i] == "north") north_flag[i] = 1
else north_flag[i] = 0
if(region[i] == "east") east_flag[i] = 1
else east_flag[i] = 0
if(region[i] == "south") south_flag[i] = 1
else south_flag[i] = 0
}
north_flag; east_flag; south_flag
```

```
north_flag east_flag south_flag
[1] 0 0 0 0 0 0 1 0
[2] 0 1 0 0 0 0 0 0
[3] 0 0 0 0 0 0 0 0
```

### # Index fields

```
# Data frames have an index field;
# the left-most column
cars
cars[order(cars$mpg),]
```

```
# For vectors or matrices,
# add a column to act as an index field
x <- c(1,1,3;1,1;4,3); y <- c(9,9;1)
z <- c(2,1;9)
matrix <- t(rbind(x,y,z)); matrix
indexed_m <- cbind(c(1:length(x)), matrix
indexed_m[order(z),]
```

### # Duplicate records

```
er of duplicate records, use anyDuplicated
ted(cars)
ne each record, use Duplicated
(cars)
record is a duplicate,
record is not a duplicate
```

```
# Let's duplicate the first record
new.cars <- rbind(cars, cars[1,])
# Check for duplicates
anyDuplicated(new.cars)
# The 262nd record is a duplicate
duplicated(new.cars)
```

## CHAPTER 3

### 3.13 The R Zone

#### # Read in the Churn data set

```
churn <- read.csv(file =
"C:/../churn.txt",
stringsAsFactors=TRUE)
# Show the first ten records
churn[1:10,]
```

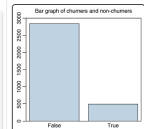
```
churn
  churn  international plan  no international plan
1  No      No               No
2  No      No               No
3  No      No               No
4  No      No               No
5  No      No               No
6  No      No               No
7  No      No               No
8  No      No               No
9  No      No               No
10 No      No               No
```

```
# Summarize the Churn variable
sum.churn <- summary(churn$Churn)
sum.churn
```

```
# Calculate proportion of churners
prop.churn <- sum(churn$Churn ==
"True") / length(churn$Churn)
prop.churn
```

#### # Bar chart of variable Churn

```
barplot(sum.churn,
ylim = c(0, 3000),
main = "Bar Graph of Churners and
Non-Churners",
col = "lightblue")
box(which = "plot",
lty = "solid",
col="black")
```



#### # Create a table of proportions over rows

```
row.margin <- round(prop.table(counts,
margin = 1),
4)*100
row.margin
```

```
row.margin
International plan
Churn:  No  Yes
False  88  57.59
True   71  84.2838
```

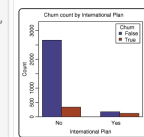
#### # Create a table of proportions over columns

```
col.margin <- round(prop.table(counts,
margin = 2),
4)*100
col.margin
```

```
col.margin
International plan
Churn:  No  Yes
False  88  57.59
True   71  84.2838
```

#### # Clustered Bar Chart, with legend

```
barplot(counts,
col = c("blue", "red"),
ylim = c(0, 3000),
ylab = "Count",
xlab = "International Plan",
main = "Churn Count by
International Plan",
beside = TRUE)
legend("topright",
c(rownames(counts)),
col = c("blue", "red"),
pch = 15,
title = "Churn")
box(which = "plot",
lty = "solid",
col="black")
```



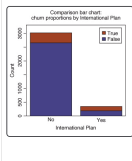
## # Make a table for counts of Churn and International Plan

```
counts <- table(churn$Churn,
  churn$Int.l.Plan,
  dnn=c("Churn", "International Plan"))
counts
```



## # Overlaid bar chart

```
barplot(counts,
  legend = rownames(counts),
  col = c("blue", "red"),
  ylim = c(0, 3300),
  ylab = "Count",
  xlab = "International Plan",
  main = "Comparison Bar Chart:
  Churn Proportions by
  International Plan")
box(which = "plot",
  lty = "solid",
  col="black")
```



## # Create a table with sums for both variables

```
sumentable <- addmargins(counts,
  FUN = sum)
sumentable
```



## # Clustered Bar Chart of Churn and International Plan with legend

```
barplot(t(counts),
  col = c("blue", "green"),
  ylim = c(0, 3300),
  ylab = "Counts",
  xlab = "Churn",
  main = "International Plan Count by
  Churn",
  beside = TRUE)
legend("topright",
  c(rownames(counts)),
  col = c("blue", "green"),
  pch = 15,
  title = "Int'l Plan")
box(which = "plot",
  lty = "solid",
  col="black")
```



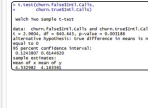
## # Histogram of non-overlaid Customer Service Calls

```
hist(churn$CustServ.Calls,
  xlim = c(0,10),
  col = "lightblue",
  ylab = "Count",
  xlab = "Customer Service Calls",
  main = "Histogram of Customer Service
  Calls")
```



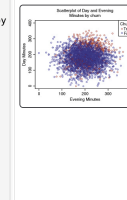
## # Two-sample T-Test for Int'l Calls

```
# Partition data
churn.false <- subset(churn,
  churn$Churn ==
  "False")
churn.true <- subset(churn,
  churn$Churn ==
  "True")
# Run the test
t.test(churn.false$Intl.Calls,
  churn.true$Intl.Calls)
```



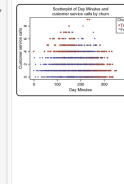
## # Scatterplot of Evening Minutes and Day Minutes, colored by Churn

```
plot(churn$Eve.Mins,
  churn$Day.Mins,
  xlim = c(0, 400),
  ylim = c(0, 400),
  xlab = "Evening Minutes",
  ylab = "Day Minutes",
  main = "Scatterplot of Day
  and Evening Minutes by
  Churn",
  col = ifelse(churn$Churn=="
  "True",
  "red",
  "blue"))
legend("topright",
  c("True",
  "False"),
  col = c("red",
  "blue"),
  pch = 1,
  title = "Churn")
```



## # Scatterplot of Day Minutes and Customer Service Calls, colored by Churn

```
plot(churn$Day.Mins,
  churn$CustServ.Calls,
  xlim = c(0, 400),
  xlab = "Day Minutes",
  ylab = "Customer Service Calls",
  main = "Scatterplot of Day Minutes and
  Customer Service Calls by Churn",
  col = ifelse(churn$Churn=="True",
  "red",
  "blue"),
  pch = ifelse(churn$Churn=="True",
  16, 20))
legend("topright",
  c("True",
  "False"),
  col = c("red",
  "blue"),
  pch = c(16, 20),
  title = "Churn")
```



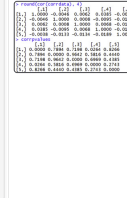
## # Scatterplot matrix

```
pairs(~churn$Day.Mins+
  churn$Day.Calls+
  churn$Day.Charge)
```

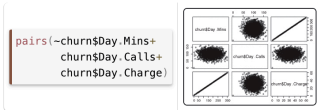


## # Correlation values and p-values in matrix form

```
# Collect variables of interest
corrdata <-
  cbind(churn$Account.Length,
  churn$Mail.Message,
  churn$Day.Mins,
  churn$Day.Calls,
  churn$CustServ.Calls)
# Declare the matrix
corrvalues <- matrix(rep(0, 25),
  ncol = 5)
# Fill the matrix with correlations
for (i in 1:4) {
  for (j in (i+1):5) {
    corrvalues[i,j] <-
      round(cor.test(corrdata[i,i],
        corrdata[j,j])$p.value,
        4)
  }
}
round(cor(corrdata), 4)
corrvalues
```



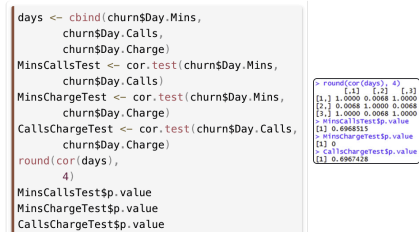
## # Scatterplot matrix



## # Regression of Day Charge vs Day Minutes



## # Correlation values, with p-values



## # Download and install the R Package ggplot2



## # Overlaid bar charts

