

# SHC 798 R-code

Richard Lubega

2025-07-02

## Introduction

### R Markdown Preamble

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

See also the **R Markdown Cookbook**)

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

### Sample Inbuilt Data Set

```
summary(cars)
```

```
##      speed      dist
##  Min.   : 4.0    Min.   : 2.00
## 1st Qu.:12.0    1st Qu.: 26.00
##  Median :15.0    Median : 36.00
##   Mean  :15.4    Mean   : 42.98
## 3rd Qu.:19.0    3rd Qu.: 56.00
##   Max.  :25.0    Max.    :120.00
```

### Including Plots - Example

#### Embed Plot

You can also embed plots, for example:

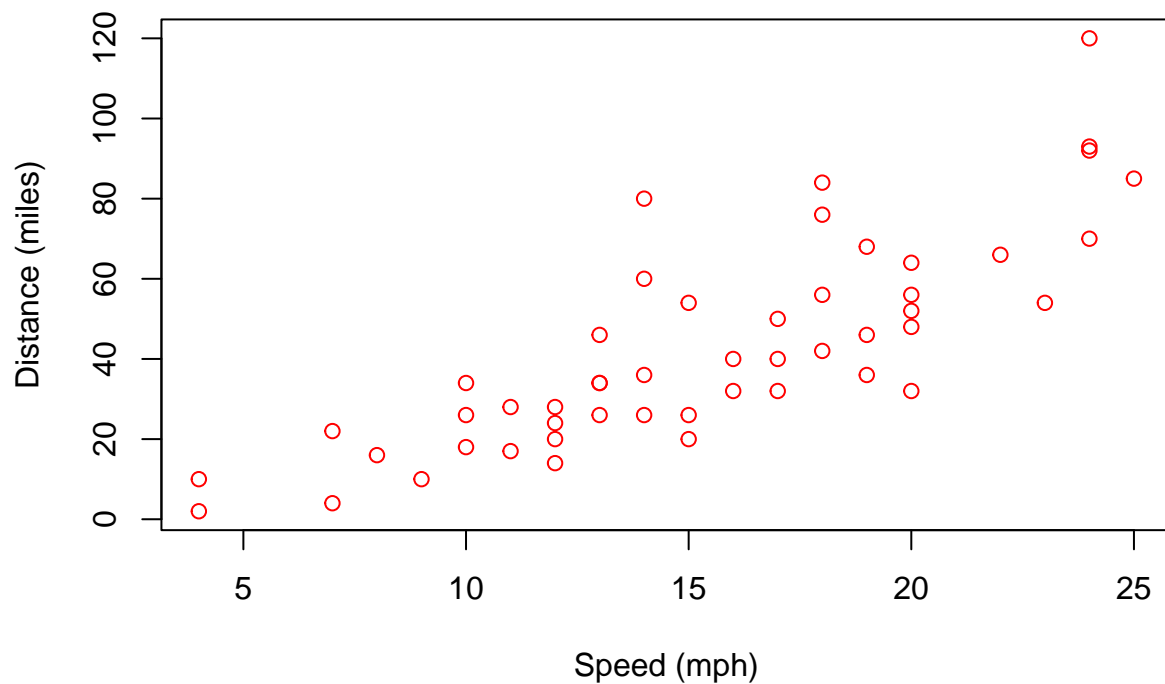


Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.

### Another Plot

Let's create another Test plot

```
plot(cars, col = "red", xlab = "Speed (mph)", ylab = "Distance (miles)" )
```



A summary of the data frame is given below

```
pacman::p_load(knitr)
kable(summary(cars))
```

speed	dist
Min. : 4.0	Min. : 2.00
1st Qu.:12.0	1st Qu.: 26.00
Median :15.0	Median : 36.00
Mean :15.4	Mean : 42.98
3rd Qu.:19.0	3rd Qu.: 56.00
Max. :25.0	Max. :120.00

## Start of My SHC 798 R-code

### R basics

1. Why install **Rtools**?
2. R Shiny **Gallery**: There's Shiny for R and Shiny for Python

# R Operations 1

What can be done?

```
# R as a calculator  
sqrt(12)
```

```
## [1] 3.464102
```

```
x <- 3  
y <- x^2  
x+y
```

```
## [1] 12
```

```
# Creating Vectors  
v1 <- c(1,5, 80)  
v1
```

```
## [1] 1 5 80
```

```
v2 <- 2:11  
v2
```

```
## [1] 2 3 4 5 6 7 8 9 10 11
```

```
a <- c(1,6,10,22,7,13)  
mean(a)
```

```
## [1] 9.833333
```

```
# Incomplete Statement  
3*(4 +  
  +2)
```

```
## [1] 18
```

```
#Assignments and Function Calls  
t.a <- 3*(4+2)  
t.b <- t.a + 2.5  
mn <- mean(c(t.a,t.b))  
mn
```

```
## [1] 19.25
```

## R Help

### Resources

1. Very useful and helpful **R Q&A Website**
2. Paradis 2005 - R for Beginners [Textbook]
3. R Reference Card [in the notes, 6 pages - Material from *R for Beginners*]
4. Base R Cheat Sheet [in the notes, 2 pages]
5. Venables et al 2017 - An introduction to R [Textbook]
6. **R for Data Science, r4ds** or the **r4ds-2e**
7. **W<sup>3</sup>Schools** R Certification

## R Operations 2

### Importing Data sets

```
# Import Data: Website
url <- "https://stat.ethz.ch/Teaching/Datasets/WBL/sport.dat"
d.sport <- read.table(url, header = TRUE)
head(d.sport)
```

```
##           weit kugel hoch  disc stab speer punkte
## OBRIEN      7.57 15.66  207 48.78  500 66.90   8824
## BUSEMANN     8.07 13.60  204 45.04  480 66.86   8706
## DVORAK       7.60 15.82  198 46.28  470 70.16   8664
## FRITZ        7.77 15.31  204 49.84  510 65.70   8644
## HAMALAINEN   7.48 16.32  198 49.62  500 57.66   8613
## NOOL         7.88 14.01  201 42.98  540 65.48   8543
```

```
# Setting the Working Directory
# Use:
getwd() #Prints the current working directory
```

```
## [1] "D:/2025 MEng Transportation/SHC 798 R-Proj"
```

```
# setwd("D:/2025 MEng Transportation/SHC 798 R-Proj") ~ Sets the working directory
# Alternatively, use "Session" → "Set Working Directory" → "Choose Directory..."
# Import Data: Files
# - Different ways depending on the format (csv, txt, xlsx, etc.)
# - Alternative: use the "Import Dataset" tool in RStudio (upper-right panel)
# Save data or write data to a file
# - Text files
# - Excel files: use CSV
```

### R Objects & Indexing

## Statistics

**lm(y ~ x, data=df)**

Linear model.

**glm(y ~ x, data=df)**

Generalised linear model.

**summary**

Get more detailed information out a model.

**t.test(x, y)**

Perform a t-test for difference between means.

**pairwise.t.test**

Perform a t-test for paired data.

**prop.test**

Test for a difference between proportions.

**aov**

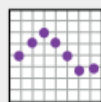
Analysis of variance.

## Distributions

	Random Variates	Density Function	Cumulative Distribution	Quantile
Normal	<b>rnorm</b>	<b>dnorm</b>	<b>pnorm</b>	<b>qnorm</b>
Poisson	<b>rpois</b>	<b>dpois</b>	<b>ppois</b>	<b>qpois</b>
Binomial	<b>rbinom</b>	<b>dbinom</b>	<b>pbinom</b>	<b>qbinom</b>
Uniform	<b>runif</b>	<b>dunif</b>	<b>punif</b>	<b>qunif</b>

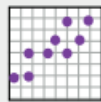
## Plotting

Also see the **ggplot2** package.



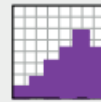
**plot(x)**

Values of x in order.



**plot(x, y)**

Values of x against y.



**hist(x)**

Histogram of x.

## Dates

See the **lubridate** package.

Learn more at [web page](#) or [vignette](#) • package version • Updated: 3/15

Figure 1: *Snippet from Resource 4*

```
# R Objects: Data frames (Most essential)  
str(d.sport)
```

```
## 'data.frame': 15 obs. of 7 variables:  
## $ weit : num 7.57 8.07 7.6 7.77 7.48 7.88 7.64 7.61 7.27 7.49 ...  
## $ kugel : num 15.7 13.6 15.8 15.3 16.3 ...  
## $ hoch : int 207 204 198 204 198 201 195 213 207 204 ...  
## $ disc : num 48.8 45 46.3 49.8 49.6 ...  
## $ stab : int 500 480 470 510 500 540 540 520 470 470 ...  
## $ speer : num 66.9 66.9 70.2 65.7 57.7 ...  
## $ punkte: int 8824 8706 8664 8644 8613 8543 8422 8318 8307 8300 ...
```

```
# R Objects: Vectors  
# (e.g., a column from the data set d.sport)  
kugel <- d.sport$kugel  
str(kugel)
```

```
## num [1:15] 15.7 13.6 15.8 15.3 16.3 ...
```

```
participant <- rownames(d.sport)  
str(participant)
```

```
## chr [1:15] "OBRIEN" "BUSEMANN" "DVORAK" "FRITZ" "HAMALAINEN" "NOOL" ...
```

```
# Select elements  
participant[4]
```

```
## [1] "FRITZ"
```

```
d.sport[c(3,6,4), c(1:3,7)]
```

```
##      weit kugel hoch punkte  
## DVORAK 7.60 15.82 198 8664  
## NOOL 7.88 14.01 201 8543  
## FRITZ 7.77 15.31 204 8644
```

```
d.sport["FRITZ", ]
```

```
##      weit kugel hoch disc stab speer punkte  
## FRITZ 7.77 15.31 204 49.84 510 65.7 8644
```

```
# Accessing Parts of an Object  
# To access only part of an object, use []  
# For vectors: myvector[x]  
# For two-dimensional objects, e.g. data frames or matrices: mydata.frame[x,y]  
d.sport[ , ]
```

```
##      weit kugel hoch  disc stab speer punkte
## OBRIEN      7.57 15.66 207 48.78 500 66.90 8824
## BUSEMANN     8.07 13.60 204 45.04 480 66.86 8706
## DVORAK      7.60 15.82 198 46.28 470 70.16 8664
## FRITZ       7.77 15.31 204 49.84 510 65.70 8644
## HAMALAINEN  7.48 16.32 198 49.62 500 57.66 8613
## NOOL        7.88 14.01 201 42.98 540 65.48 8543
## ZMELIK      7.64 13.53 195 43.44 540 67.20 8422
## GANIYEV     7.61 14.71 213 44.86 520 53.70 8318
## PENALVER    7.27 16.91 207 48.92 470 57.08 8307
## HUFFINS     7.49 15.57 204 48.72 470 60.62 8300
## PLAZIAT     7.82 14.85 204 45.34 490 52.18 8282
## MAGNUSSON   7.28 15.52 195 43.78 480 61.10 8274
## SMITH       7.47 16.97 195 49.54 500 64.34 8271
## MUELLER     7.25 14.69 195 45.90 510 66.10 8253
## CHMARA      7.75 14.51 210 42.60 490 54.84 8249
```

```
c(1,3,7)
```

```
## [1] 1 3 7
```

```
1:10
```

```
## [1] 1 2 3 4 5 6 7 8 9 10
```

```
d.sport[1:10, ]
```

```
##      weit kugel hoch  disc stab speer punkte
## OBRIEN      7.57 15.66 207 48.78 500 66.90 8824
## BUSEMANN     8.07 13.60 204 45.04 480 66.86 8706
## DVORAK      7.60 15.82 198 46.28 470 70.16 8664
## FRITZ       7.77 15.31 204 49.84 510 65.70 8644
## HAMALAINEN  7.48 16.32 198 49.62 500 57.66 8613
## NOOL        7.88 14.01 201 42.98 540 65.48 8543
## ZMELIK      7.64 13.53 195 43.44 540 67.20 8422
## GANIYEV     7.61 14.71 213 44.86 520 53.70 8318
## PENALVER    7.27 16.91 207 48.92 470 57.08 8307
## HUFFINS     7.49 15.57 204 48.72 470 60.62 8300
```

```
d.sport[-c(1, 3,7), ] #negative indices are excluded
```

```
##      weit kugel hoch  disc stab speer punkte
## BUSEMANN     8.07 13.60 204 45.04 480 66.86 8706
## FRITZ       7.77 15.31 204 49.84 510 65.70 8644
## HAMALAINEN  7.48 16.32 198 49.62 500 57.66 8613
## NOOL        7.88 14.01 201 42.98 540 65.48 8543
## GANIYEV     7.61 14.71 213 44.86 520 53.70 8318
## PENALVER    7.27 16.91 207 48.92 470 57.08 8307
## HUFFINS     7.49 15.57 204 48.72 470 60.62 8300
## PLAZIAT     7.82 14.85 204 45.34 490 52.18 8282
## MAGNUSSON   7.28 15.52 195 43.78 480 61.10 8274
```



```
## SMITH      7.47 16.97 195 49.54 500 64.34 8271
## MUELLER    7.25 14.69 195 45.90 510 66.10 8253
## CHMARA     7.75 14.51 210 42.60 490 54.84 8249
```

```
d.sport[, 2:3]
```

```
##           kugel hoch
## O'BRIEN    15.66 207
## BUSEMANN   13.60 204
## DVORAK     15.82 198
## FRITZ      15.31 204
## HAMALAINEN 16.32 198
## NOOL       14.01 201
## ZMELIK     13.53 195
## GANIYEV    14.71 213
## PENALVER   16.91 207
## HUFFINS    15.57 204
## PLAZIAT    14.85 204
## MAGNUSSON  15.52 195
## SMITH      16.97 195
## MUELLER    14.69 195
## CHMARA     14.51 210
```

```
d.sport[c(1,3,6), 2:3]
```

```
##           kugel hoch
## O'BRIEN 15.66 207
## DVORAK 15.82 198
## NOOL   14.01 201
```

```
# Function Calls
```

```
mean(kugel)
```

```
## [1] 15.19867
```

```
quantile(kugel)
```

```
##      0%   25%   50%   75%  100%
## 13.53 14.60 15.31 15.74 16.97
```

```
quantile(kugel,probs = c(.75, 0.9))
```

```
##      75%   90%
## 15.740 16.674
```

```
# Functions consist of mandatory and optional arguments:
```

```
# mean(x, trim = 0, na.rm = FALSE, ...)
```

```
# x: mandatory argument
```

```
# trim: optional argument, default is 0
```

```
# na.rm: optional argument, default is FALSE
```

```
# The arguments of a function have a defined order and each argument has its own unique name
```

```
mean(x = kugel, na.rm = TRUE)
```

```
## [1] 15.19867
```

```
mean(x = kugel, ,TRUE)
```

```
## [1] 15.19867
```

```
# Useful Functions
```

```
nrow(d.sport)
```

```
## [1] 15
```

```
ncol(d.sport)
```

```
## [1] 7
```

```
dim(d.sport)
```

```
## [1] 15 7
```

```
summary(d.sport)
```

```
##      weit      kugel      hoch      disc      stab
##  Min.   :7.250   Min.   :13.53   Min.   :195.0   Min.   :42.60   Min.   :470
## 1st Qu.:7.475   1st Qu.:14.60   1st Qu.:196.5   1st Qu.:44.32   1st Qu.:480
## Median :7.600   Median :15.31   Median :204.0   Median :45.90   Median :500
## Mean   :7.597   Mean   :15.20   Mean   :202.0   Mean   :46.38   Mean   :498
## 3rd Qu.:7.760   3rd Qu.:15.74   3rd Qu.:205.5   3rd Qu.:48.85   3rd Qu.:510
## Max.   :8.070   Max.   :16.97   Max.   :213.0   Max.   :49.84   Max.   :540
##      speer      punkte
##  Min.   :52.18   Min.   :8249
## 1st Qu.:57.37   1st Qu.:8278
## Median :64.34   Median :8318
## Mean   :61.99   Mean   :8445
## 3rd Qu.:66.48   3rd Qu.:8628
## Max.   :70.16   Max.   :8824
```

```
# apply(d.sport)
```

```
head(d.sport)
```

```
##      weit kugel hoch disc stab speer punkte
## OBRIEN  7.57 15.66 207 48.78 500 66.90 8824
## BUSEMANN 8.07 13.60 204 45.04 480 66.86 8706
## DVORAK  7.60 15.82 198 46.28 470 70.16 8664
## FRITZ   7.77 15.31 204 49.84 510 65.70 8644
## HAMALAINEN 7.48 16.32 198 49.62 500 57.66 8613
## NOOL    7.88 14.01 201 42.98 540 65.48 8543
```

```
tail(d.sport)
```

```
##      weit kugel hoch  disc stab speer punkte
## HUFFINS  7.49 15.57  204 48.72  470 60.62  8300
## PLAZIAT  7.82 14.85  204 45.34  490 52.18  8282
## MAGNUSSON 7.28 15.52  195 43.78  480 61.10  8274
## SMITH    7.47 16.97  195 49.54  500 64.34  8271
## MUELLER   7.25 14.69  195 45.90  510 66.10  8253
## CHMARA    7.75 14.51  210 42.60  490 54.84  8249
```

## R Packages

```
#install.packages("MASS")
#require(MASS) # for every R Session
#library(MASS) # or use this
```

Online resources: University of Pretoria | SHC 798 | Introduction to R 36

- o List of all packages: <http://cran.r-project.org/web/packages/>
- o By topic: <http://cran.r-project.org/web/views/>
- o Ask Google / ChatGPT / GroK

## Missing Values

```
d.sport.NA <- d.sport
d.sport.NA[2, "kugel"] <- NA
d.sport.NA[3, "hoch"] <- -999
# missing values are coded as NA (not available) and are treated in a special way, e.g. is.na():
is.na(d.sport.NA) #one logical value per element
```

```
##      weit kugel hoch  disc  stab speer punkte
## OBRIEN  FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## BUSEMANN FALSE  TRUE FALSE FALSE FALSE FALSE FALSE
## DVORAK   FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## FRITZ    FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## HAMALAINEN FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## NOOL     FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## ZMELIK   FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## GANIYEV  FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## PENALVER FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## HUFFINS  FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## PLAZIAT  FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## MAGNUSSON FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## SMITH    FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## MUELLER  FALSE FALSE FALSE FALSE FALSE FALSE FALSE
## CHMARA   FALSE FALSE FALSE FALSE FALSE FALSE FALSE
```

```
sum(is.na(d.sport.NA)) # adds up the TRUE elements
```

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

```
which(is.na(d.sport.NA), arr.ind = TRUE) # where are the NA's
```

```
##           row col  
## BUSEMANN    2    2
```

```
# Specify missing values after reading in the data:
```

```
d.sport.NA[d.sport.NA == -999] <- NA
```

```
# Many functions have an argument to handle missing values, e.g. na.rm, na.omit:
```

```
sum(d.sport.NA$kugel)
```

```
## [1] NA
```

```
sum(d.sport.NA$kugel, na.rm = TRUE)
```

```
## [1] 214.38
```

```
na.omit(d.sport.NA)
```

```
##           weit kugel hoch  disc stab speer punkte  
## OBRIEN      7.57 15.66  207 48.78  500 66.90  8824  
## FRITZ       7.77 15.31  204 49.84  510 65.70  8644  
## HAMALAINEN  7.48 16.32  198 49.62  500 57.66  8613  
## NOOL        7.88 14.01  201 42.98  540 65.48  8543  
## ZMELIK      7.64 13.53  195 43.44  540 67.20  8422  
## GANIYEV     7.61 14.71  213 44.86  520 53.70  8318  
## PENALVER    7.27 16.91  207 48.92  470 57.08  8307  
## HUFFINS     7.49 15.57  204 48.72  470 60.62  8300  
## PLAZIAT     7.82 14.85  204 45.34  490 52.18  8282  
## MAGNUSSON   7.28 15.52  195 43.78  480 61.10  8274  
## SMITH       7.47 16.97  195 49.54  500 64.34  8271  
## MUELLER     7.25 14.69  195 45.90  510 66.10  8253  
## CHMARA      7.75 14.51  210 42.60  490 54.84  8249
```

## Basic Graphics

```
# The Plot Function: only 1 mandatory argument i.e., x. The 2nd most important one is y.
```

```
# many optional arguments [col,pch,main,cex, ...]
```

```
# use function par (?par) to set or query graphical parameters
```

```
data(iris)
```

```
iris
```

##	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
## 1	5.1	3.5	1.4	0.2	setosa
## 2	4.9	3.0	1.4	0.2	setosa
## 3	4.7	3.2	1.3	0.2	setosa
## 4	4.6	3.1	1.5	0.2	setosa
## 5	5.0	3.6	1.4	0.2	setosa
## 6	5.4	3.9	1.7	0.4	setosa
## 7	4.6	3.4	1.4	0.3	setosa
## 8	5.0	3.4	1.5	0.2	setosa
## 9	4.4	2.9	1.4	0.2	setosa
## 10	4.9	3.1	1.5	0.1	setosa
## 11	5.4	3.7	1.5	0.2	setosa
## 12	4.8	3.4	1.6	0.2	setosa
## 13	4.8	3.0	1.4	0.1	setosa
## 14	4.3	3.0	1.1	0.1	setosa
## 15	5.8	4.0	1.2	0.2	setosa
## 16	5.7	4.4	1.5	0.4	setosa
## 17	5.4	3.9	1.3	0.4	setosa
## 18	5.1	3.5	1.4	0.3	setosa
## 19	5.7	3.8	1.7	0.3	setosa
## 20	5.1	3.8	1.5	0.3	setosa
## 21	5.4	3.4	1.7	0.2	setosa
## 22	5.1	3.7	1.5	0.4	setosa
## 23	4.6	3.6	1.0	0.2	setosa
## 24	5.1	3.3	1.7	0.5	setosa
## 25	4.8	3.4	1.9	0.2	setosa
## 26	5.0	3.0	1.6	0.2	setosa
## 27	5.0	3.4	1.6	0.4	setosa
## 28	5.2	3.5	1.5	0.2	setosa
## 29	5.2	3.4	1.4	0.2	setosa
## 30	4.7	3.2	1.6	0.2	setosa
## 31	4.8	3.1	1.6	0.2	setosa
## 32	5.4	3.4	1.5	0.4	setosa
## 33	5.2	4.1	1.5	0.1	setosa
## 34	5.5	4.2	1.4	0.2	setosa
## 35	4.9	3.1	1.5	0.2	setosa
## 36	5.0	3.2	1.2	0.2	setosa
## 37	5.5	3.5	1.3	0.2	setosa
## 38	4.9	3.6	1.4	0.1	setosa
## 39	4.4	3.0	1.3	0.2	setosa
## 40	5.1	3.4	1.5	0.2	setosa
## 41	5.0	3.5	1.3	0.3	setosa
## 42	4.5	2.3	1.3	0.3	setosa
## 43	4.4	3.2	1.3	0.2	setosa
## 44	5.0	3.5	1.6	0.6	setosa
## 45	5.1	3.8	1.9	0.4	setosa
## 46	4.8	3.0	1.4	0.3	setosa
## 47	5.1	3.8	1.6	0.2	setosa
## 48	4.6	3.2	1.4	0.2	setosa
## 49	5.3	3.7	1.5	0.2	setosa
## 50	5.0	3.3	1.4	0.2	setosa
## 51	7.0	3.2	4.7	1.4	versicolor
## 52	6.4	3.2	4.5	1.5	versicolor
## 53	6.9	3.1	4.9	1.5	versicolor

## 54	5.5	2.3	4.0	1.3 versicolor
## 55	6.5	2.8	4.6	1.5 versicolor
## 56	5.7	2.8	4.5	1.3 versicolor
## 57	6.3	3.3	4.7	1.6 versicolor
## 58	4.9	2.4	3.3	1.0 versicolor
## 59	6.6	2.9	4.6	1.3 versicolor
## 60	5.2	2.7	3.9	1.4 versicolor
## 61	5.0	2.0	3.5	1.0 versicolor
## 62	5.9	3.0	4.2	1.5 versicolor
## 63	6.0	2.2	4.0	1.0 versicolor
## 64	6.1	2.9	4.7	1.4 versicolor
## 65	5.6	2.9	3.6	1.3 versicolor
## 66	6.7	3.1	4.4	1.4 versicolor
## 67	5.6	3.0	4.5	1.5 versicolor
## 68	5.8	2.7	4.1	1.0 versicolor
## 69	6.2	2.2	4.5	1.5 versicolor
## 70	5.6	2.5	3.9	1.1 versicolor
## 71	5.9	3.2	4.8	1.8 versicolor
## 72	6.1	2.8	4.0	1.3 versicolor
## 73	6.3	2.5	4.9	1.5 versicolor
## 74	6.1	2.8	4.7	1.2 versicolor
## 75	6.4	2.9	4.3	1.3 versicolor
## 76	6.6	3.0	4.4	1.4 versicolor
## 77	6.8	2.8	4.8	1.4 versicolor
## 78	6.7	3.0	5.0	1.7 versicolor
## 79	6.0	2.9	4.5	1.5 versicolor
## 80	5.7	2.6	3.5	1.0 versicolor
## 81	5.5	2.4	3.8	1.1 versicolor
## 82	5.5	2.4	3.7	1.0 versicolor
## 83	5.8	2.7	3.9	1.2 versicolor
## 84	6.0	2.7	5.1	1.6 versicolor
## 85	5.4	3.0	4.5	1.5 versicolor
## 86	6.0	3.4	4.5	1.6 versicolor
## 87	6.7	3.1	4.7	1.5 versicolor
## 88	6.3	2.3	4.4	1.3 versicolor
## 89	5.6	3.0	4.1	1.3 versicolor
## 90	5.5	2.5	4.0	1.3 versicolor
## 91	5.5	2.6	4.4	1.2 versicolor
## 92	6.1	3.0	4.6	1.4 versicolor
## 93	5.8	2.6	4.0	1.2 versicolor
## 94	5.0	2.3	3.3	1.0 versicolor
## 95	5.6	2.7	4.2	1.3 versicolor
## 96	5.7	3.0	4.2	1.2 versicolor
## 97	5.7	2.9	4.2	1.3 versicolor
## 98	6.2	2.9	4.3	1.3 versicolor
## 99	5.1	2.5	3.0	1.1 versicolor
## 100	5.7	2.8	4.1	1.3 versicolor
## 101	6.3	3.3	6.0	2.5 virginica
## 102	5.8	2.7	5.1	1.9 virginica
## 103	7.1	3.0	5.9	2.1 virginica
## 104	6.3	2.9	5.6	1.8 virginica
## 105	6.5	3.0	5.8	2.2 virginica
## 106	7.6	3.0	6.6	2.1 virginica
## 107	4.9	2.5	4.5	1.7 virginica

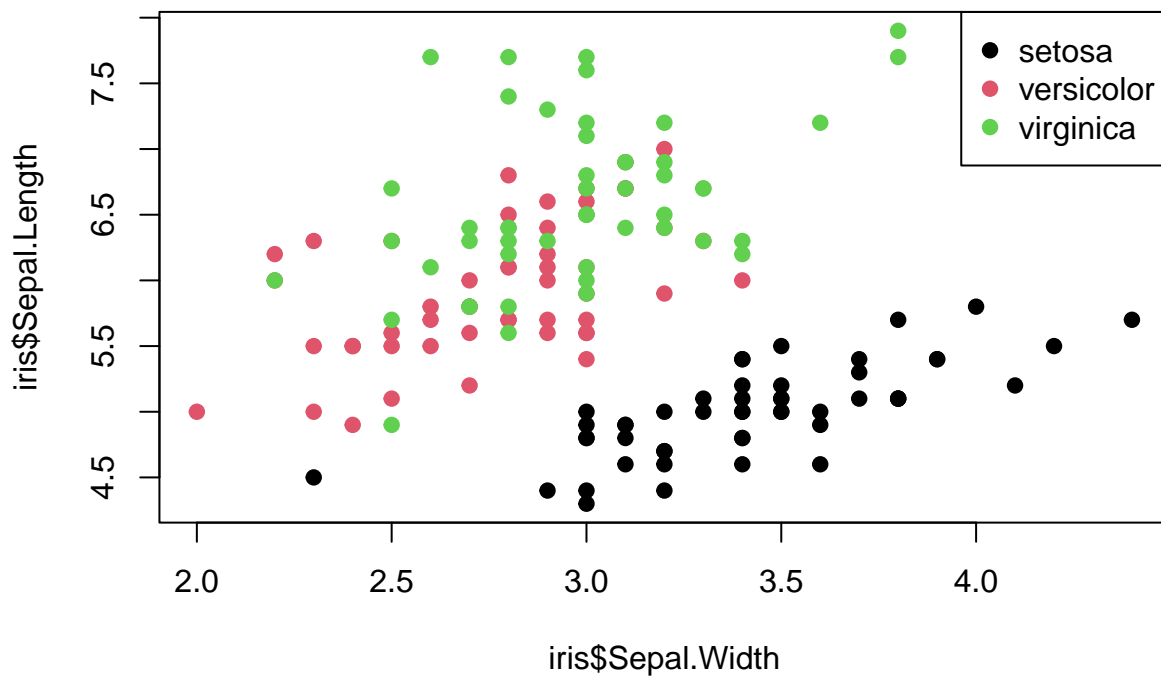
```
## 108      7.3      2.9      6.3      1.8 virginica
## 109      6.7      2.5      5.8      1.8 virginica
## 110      7.2      3.6      6.1      2.5 virginica
## 111      6.5      3.2      5.1      2.0 virginica
## 112      6.4      2.7      5.3      1.9 virginica
## 113      6.8      3.0      5.5      2.1 virginica
## 114      5.7      2.5      5.0      2.0 virginica
## 115      5.8      2.8      5.1      2.4 virginica
## 116      6.4      3.2      5.3      2.3 virginica
## 117      6.5      3.0      5.5      1.8 virginica
## 118      7.7      3.8      6.7      2.2 virginica
## 119      7.7      2.6      6.9      2.3 virginica
## 120      6.0      2.2      5.0      1.5 virginica
## 121      6.9      3.2      5.7      2.3 virginica
## 122      5.6      2.8      4.9      2.0 virginica
## 123      7.7      2.8      6.7      2.0 virginica
## 124      6.3      2.7      4.9      1.8 virginica
## 125      6.7      3.3      5.7      2.1 virginica
## 126      7.2      3.2      6.0      1.8 virginica
## 127      6.2      2.8      4.8      1.8 virginica
## 128      6.1      3.0      4.9      1.8 virginica
## 129      6.4      2.8      5.6      2.1 virginica
## 130      7.2      3.0      5.8      1.6 virginica
## 131      7.4      2.8      6.1      1.9 virginica
## 132      7.9      3.8      6.4      2.0 virginica
## 133      6.4      2.8      5.6      2.2 virginica
## 134      6.3      2.8      5.1      1.5 virginica
## 135      6.1      2.6      5.6      1.4 virginica
## 136      7.7      3.0      6.1      2.3 virginica
## 137      6.3      3.4      5.6      2.4 virginica
## 138      6.4      3.1      5.5      1.8 virginica
## 139      6.0      3.0      4.8      1.8 virginica
## 140      6.9      3.1      5.4      2.1 virginica
## 141      6.7      3.1      5.6      2.4 virginica
## 142      6.9      3.1      5.1      2.3 virginica
## 143      5.8      2.7      5.1      1.9 virginica
## 144      6.8      3.2      5.9      2.3 virginica
## 145      6.7      3.3      5.7      2.5 virginica
## 146      6.7      3.0      5.2      2.3 virginica
## 147      6.3      2.5      5.0      1.9 virginica
## 148      6.5      3.0      5.2      2.0 virginica
## 149      6.2      3.4      5.4      2.3 virginica
## 150      5.9      3.0      5.1      1.8 virginica
```

```
str(iris)
```

```
## 'data.frame':  150 obs. of  5 variables:
## $ Sepal.Length: num  5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
## $ Sepal.Width : num  3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
## $ Petal.Length: num  1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
## $ Petal.Width : num  0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
## $ Species      : Factor w/ 3 levels "setosa","versicolor",...: 1 1 1 1 1 1 1 1 1 1 ...
```

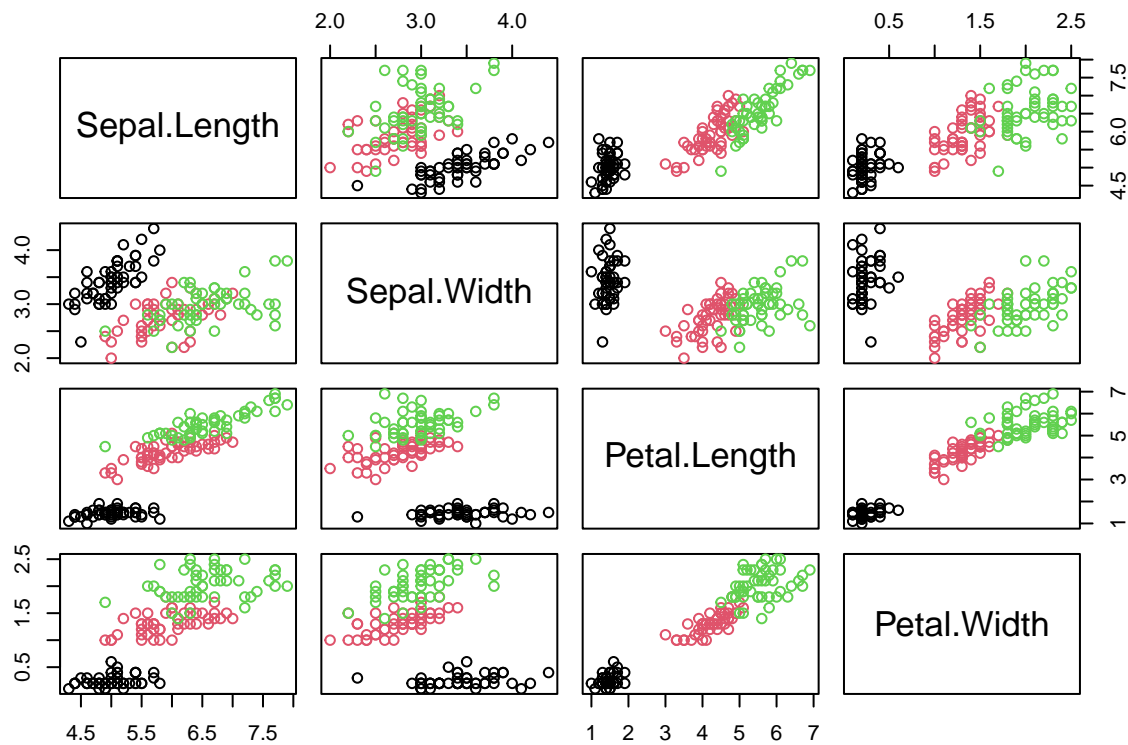
```
# A factor represents categorical values with different "levels"
# High-level plotting function: opens a plot
plot(x = iris$Sepal.Width, y = iris$Sepal.Length, col = iris[, "Species"], pch = 19)

#Low-level plotting functions: add to an existing plot
legend("topright", legend = levels(iris[, "Species"]), pch = 19, col = 1:3) # adds legend
```



```
# In an Rmd file such as this one, you need to call both the plot() & legend() functions at the same time
pairs(iris[, -5], col = iris[, 5])
```





### Arguments of plot

Statement	Meaning
type	Style of drawing (single points, lines etc.)
log	logarithmic scale
xlim	range of x-coordinates
ylim	range of y-coordinates
pch	Plotting character
col	Coloring points
lty	line type
lwd	line width
main	main title (appears above the plot)
xlab	label of x-axis
ylab	label of y-axis

Figure 2: Plot Arguments

### Three categories of R graphics functions:

- High-level plotting functions such as `plot()` to generate a new graphics display.
- Low-level plotting functions such as `legend()` to add further graphical elements to an existing plot.

- Interactive functions such as `identify()` to amend or collect information interactively from a plot.

### Low-level plotting functions

Statement	Meaning
<code>points(x, y, pch = 1)</code>	Draws points pictured as pch.
<code>text(x, y, text)</code>	Writes text at coordinate (x,y).
<code>lines(x, y, lty = 1)</code>	Adds a line to graph.
<code>abline(a, b)</code>	Adds a line with intercept a and slope b.
<code>abline(h = y, v = x)</code>	Horizontal and vertical lines.
<code>legend(x, y, text, lty, pch)</code>	Creates a legend.

Figure 3: Low-level plotting functions

### Useful plot functions

- `plot`, `pairs`, `interaction.plot`
- `boxplot`, `hist`
- `plot3d`

### Graphical Output

```
pdf(file = "iris_plot.pdf") # open the graphics device
plot(Sepal.Length ~ Sepal.Width, data = iris)
# add anything else you want in your plots
dev.off() #close the graphic device
```

```
## pdf
## 2
```

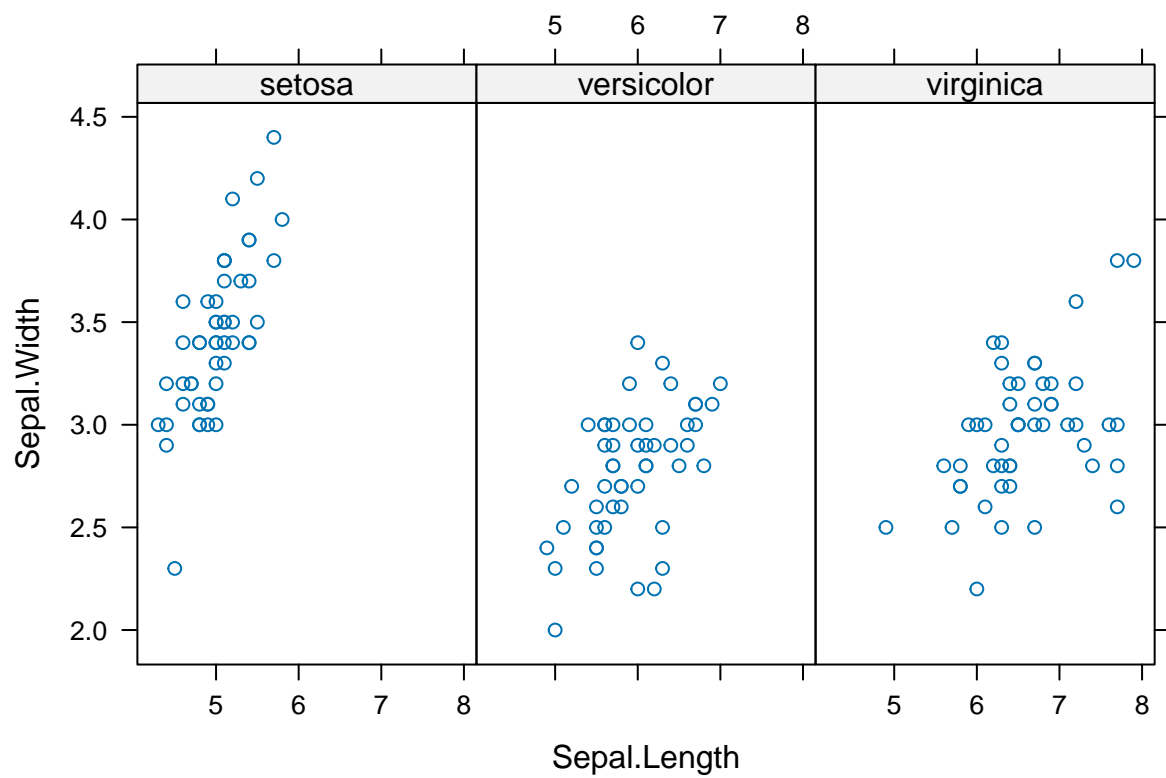
Several plots in one graphical window: splits the graphical window into 3 rows and 2 columns.

```
par(mfrow=c(3,2))
```

### Other Graphics: in lattice

The lattice package functions: good for repeating graphs for various groups. See the **Lattice Graphs in R** (at <http://www.statmethods.net/advgraphs/trellis.html>) for more information.

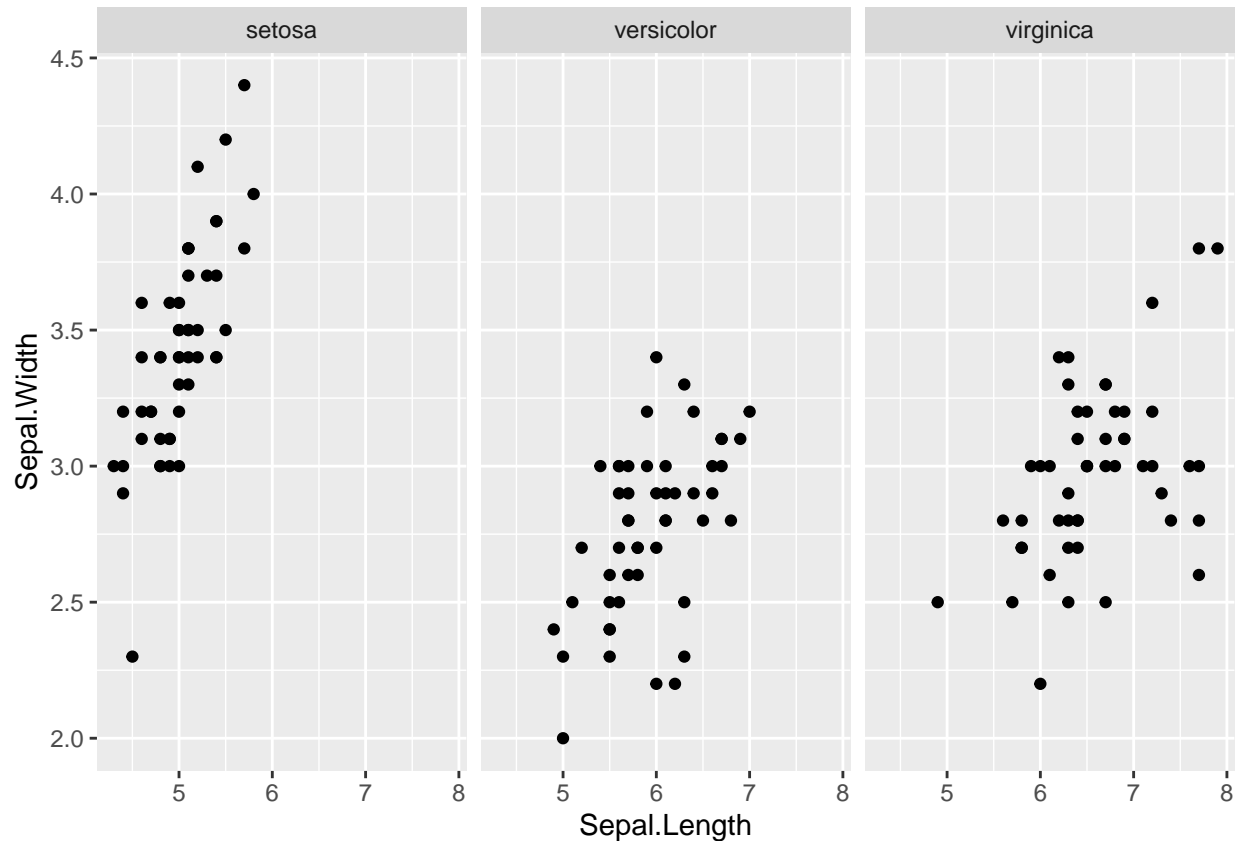
```
pacman::p_load(lattice)
xyplot(Sepal.Width ~ Sepal.Length | Species, data = iris)
```



### Other Graphics: in ggplot2

**ggplot2** package: very flexible, based on grammar of graphics.

```
pacman::p_load(ggplot2)
ggplot(data = iris, aes(x = Sepal.Length, y = Sepal.Width)) + facet_grid(rows = ~ Species) + geom_point
```



## Hypothesis Testing

Approach: Hypothesis testing in 6 steps:

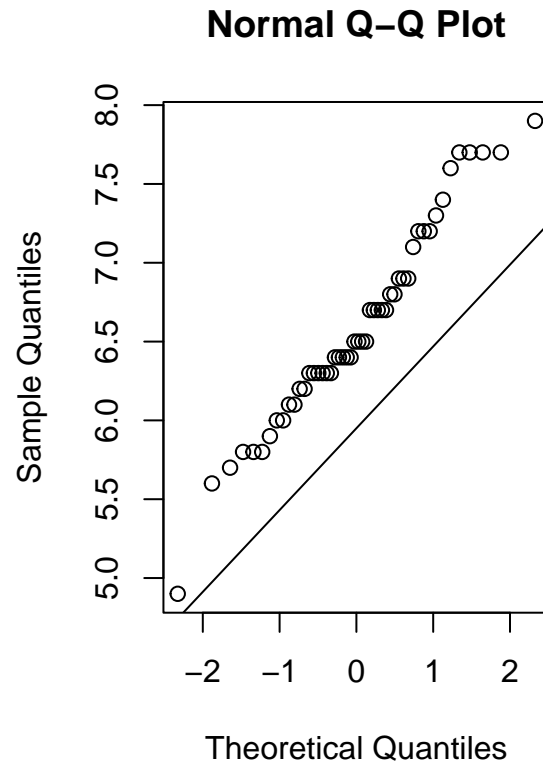
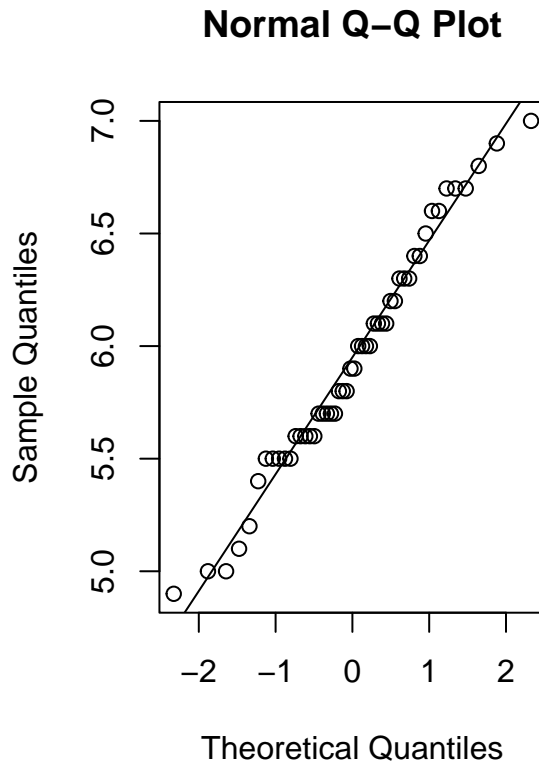
1. Declare *model* by which data were generated (e.g. population is normally distributed, large sample size and  $\sigma$  not known).
2. Define null hypothesis,  $H_0$  and alternative hypothesis,  $H_A$  ; where  $H_0$  is the statement being tested in a test of (statistical) significance and  $H_A$  is the statement that is hoped or expected to be true instead of the null hypothesis
3. Choose the *level of significance*,  $\alpha$ .
4. Determine *critical values* for the  $\alpha$  level of significance and *degrees of freedom*,  $df = (n-1)$
5. Define and calculate **test statistic**, e.g. one-sample test:
6. **Compare** the test statistic to the **critical values** and make **decision** to **reject** or **fail to reject**  $H_0$

### Hypothesis Tests – An Example

```
# Is the sepal length of versicolor different to that of virginica? Let's use a *t-Test and a *Wilcoxon
testdata <- iris[iris$Species != "setosa", c("Sepal.Length", "Species")]
testdata$Species <- droplevels(testdata$Species)
str(testdata) # prepare and check the data

## 'data.frame':   100 obs. of  2 variables:
##  $ Sepal.Length: num  7 6.4 6.9 5.5 6.5 5.7 6.3 4.9 6.6 5.2 ...
##  $ Species      : Factor w/ 2 levels "versicolor","virginica": 1 1 1 1 1 1 1 1 1 1 ...
```

```
# Check normality assumption of t-Test using QQ-Plot:
versi.id <- testdata$Species == "versicolor"
par(mfrow=c(1,2))
qqnorm(testdata$Sepal.Length[versi.id]); qqline(testdata$Sepal.Length[versi.id])
qqnorm(testdata$Sepal.Length[!versi.id]); qqline(testdata$Sepal.Length[!versi.id])
```

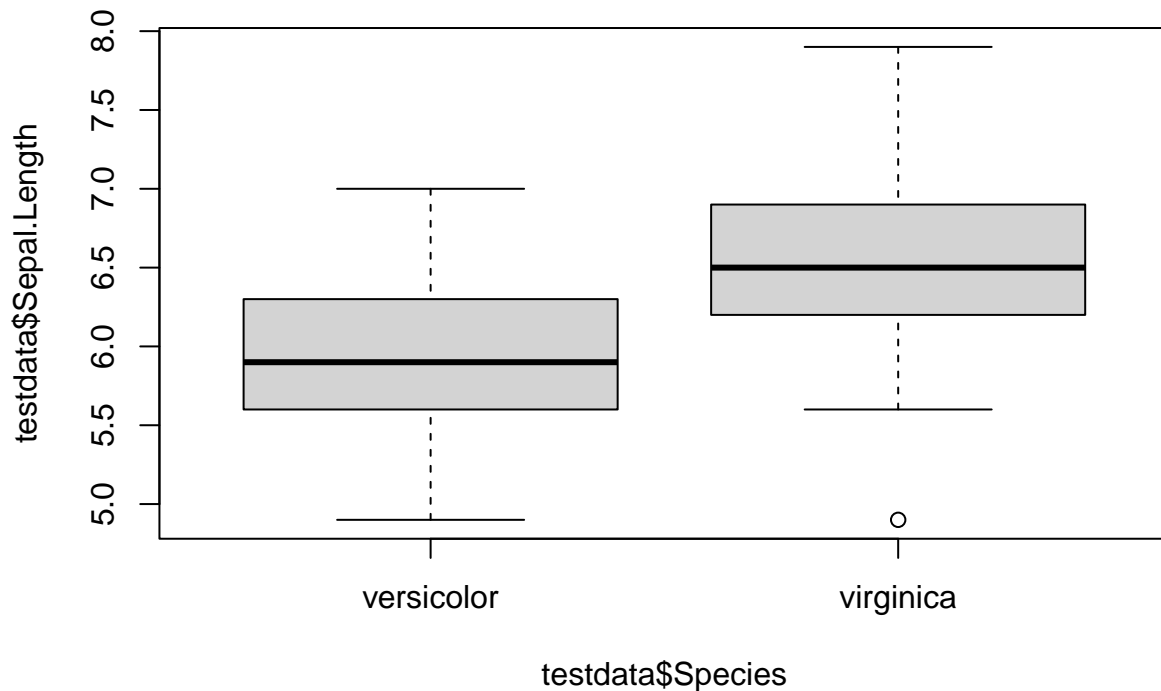


```
# Two-sample t-test:
versi.id <- testdata$Species == "versicolor"
t.test(x =testdata$Sepal.Length[versi.id], y = testdata$Sepal.Length[!versi.id], var.equal = TRUE)
```

```
##
## Two Sample t-test
##
## data: testdata$Sepal.Length[versi.id] and testdata$Sepal.Length[!versi.id]
## t = -5.6292, df = 98, p-value = 1.725e-07
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -0.8818516 -0.4221484
## sample estimates:
## mean of x mean of y
## 5.936 6.588
```

*# T-test rejects the null hypothesis at 5% significance level. Do not forget to visually check the normality assumption.*

```
# Check assumption of \*Wilcoxon Rank Test: same distribution, just a location shift.
boxplot(testdata$Sepal.Length ~ testdata$Species)
```



```
# Now, perform the test:
versi.id <- testdata$Species == "versicolor"
wilcox.test(x =testdata$Sepal.Length[versi.id], y = testdata$Sepal.Length[!versi.id], var.equal = TRUE)
```

```
##
## Wilcoxon rank sum test with continuity correction
##
## data: testdata$Sepal.Length[versi.id] and testdata$Sepal.Length[!versi.id]
## W = 526, p-value = 5.869e-07
## alternative hypothesis: true location shift is not equal to 0
```

```
# Wilcoxon Rank Test also rejects the null hypothesis at 5% significance level.
```

The **Wilcoxon Rank Test** is the *preferred* test for a two-sample statistical test.

## Hypothesis Tests - Summary

How to proceed:

- Formulate the null & alternative hypotheses

- Choose the appropriate test
- Collate data, i.e., do an experiment
- Look at data: plot(), pairs(), hist(), boxplot()
- Validate assumptions of test (e.g., T-test, Wilcoxon test)
- Carry out the test and interpret result

	1 sample / 2 dep. samples	2 indep. samples
parametric	t-Test → normality	t-Test → normality (& equal variance)
non-param.	Wilcoxon Test → symmetric distribution	Wilcoxon Test → location shift

Figure 4: Hypothesis Tests

#### Hypothesis Tests – Chi-squared test of independence

- Hypothesis:  $H_0$ : Independence of education and marriage status
- $H_A$ : Dependence of education and marriage status

```
url <- "https://stat.ethz.ch/Teaching/Datasets/edu.txt"
d.edu <- read.table(url, header = TRUE)

# Cross-tables in R
# Count number of cases with same value:
table(d.edu[, "Married"])
```

```
##
## Married more Married once
##      205      1231
```

```
# Cross-table
table(d.edu[, "Education"], d.edu[, "Married"])
```

```
##
##      Married more Married once
## College      61      550
## No College   144      681
```

```
# Now we perform a Chi-squared test
chisq.test(d.edu[, "Education"], d.edu[, "Married"])
```

```
##
## Pearson's Chi-squared test with Yates' continuity correction
##
## data: d.edu[, "Education"] and d.edu[, "Married"]
## X-squared = 15.405, df = 1, p-value = 8.675e-05
```

```
# Result: Reject  $H_0$ , i.e. education and marriage are dependent.
```

## Correlation

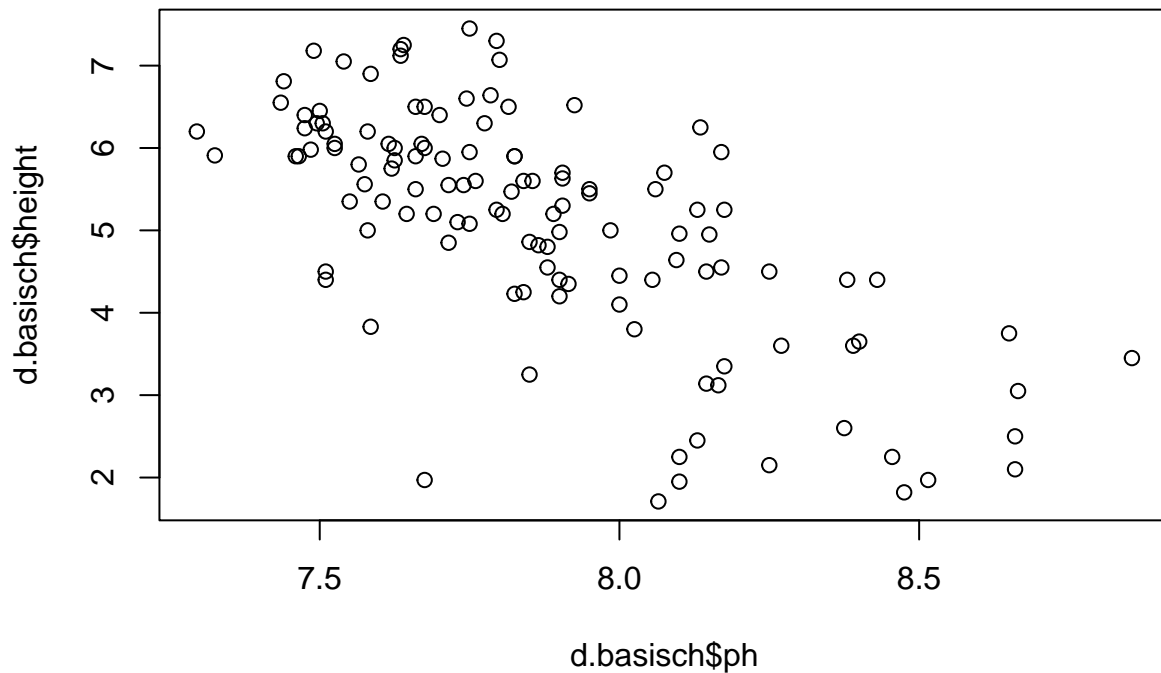
```
# Correlation
url1 <- "https://stat.ethz.ch/Teaching/Datasets/basischOhneNA.dat"
d.basisch <- read.table(url1, header = TRUE)
str(d.basisch)
```

```
## 'data.frame': 123 obs. of 4 variables:
## $ ph : num 7.33 7.69 7.9 8.14 7.62 ...
## $ l.sar : num 0.0969 0.4393 1 1.316 0.0607 ...
## $ height: num 5.91 5.2 4.4 4.5 6.05 6 5.35 5.55 4.95 5.2 ...
## $ h.quad: num 34.9 27 19.4 20.2 36.6 ...
```

```
# Calculate the (Pearson) correlation of ph and height:
cor(d.basisch$ph, d.basisch$height)
```

```
## [1] -0.6925717
```

```
# Corresponding plot:
plot(d.basisch$ph, d.basisch$height)
```



All plots show 2 variables with a correlation of 0.7



- one looks good
- another does not
- outlier(s) influence the result
- ALWAYS FIRST LOOK AT PLOTS

## Regression

### Simple Linear Regression (SLR)

From the `d.basisch()` data,

1. **Response variable:** *height* or *h.quad*: Height of trees or squared height, respectively.
2. **Possible explanatory variables:** *ph*: pH-values of soil and *l.sar*: log(sodium absorption ratio)

The simple linear regression model is:

$$Y_i = \alpha + \beta x_i + E_i \quad \text{with } E_i^{i,i,d} \sim \mathcal{N}(0, \sigma^2)$$

Figure 5: The Simple Linear Regression Model

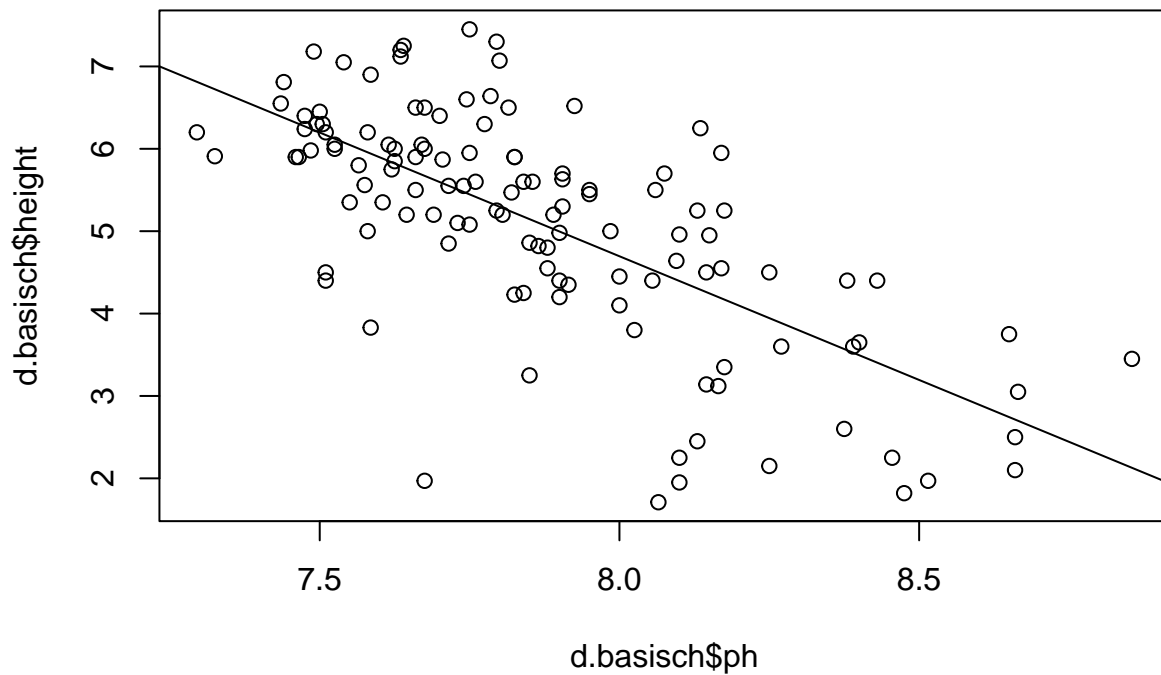
Let us pick the variable *ph* as the explanatory variable.

```
# Fit to data using lm:
fit <- lm(formula = height ~ ph, data = d.basisch)
summary(fit)

##
## Call:
## lm(formula = height ~ ph, data = d.basisch)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.7020 -0.5471  0.0874  0.6663  2.0033
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  28.7227     2.2395   12.82  <2e-16 ***
## ph          -3.0034     0.2844  -10.56  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.008 on 121 degrees of freedom
## Multiple R-squared:  0.4797, Adjusted R-squared:  0.4754
## F-statistic: 111.5 on 1 and 121 DF, p-value: < 2.2e-16

# Estimated equation: height = 28.7 - 3.0pH

# Drawing line into scatterplot:
plot(d.basisch$ph, d.basisch$height) + abline(fit)
```



```
## integer(0)
```

```
# Fit to data using lm:
fit <- lm(formula = height ~ ph, data = d.basisch)
summary(fit)
```

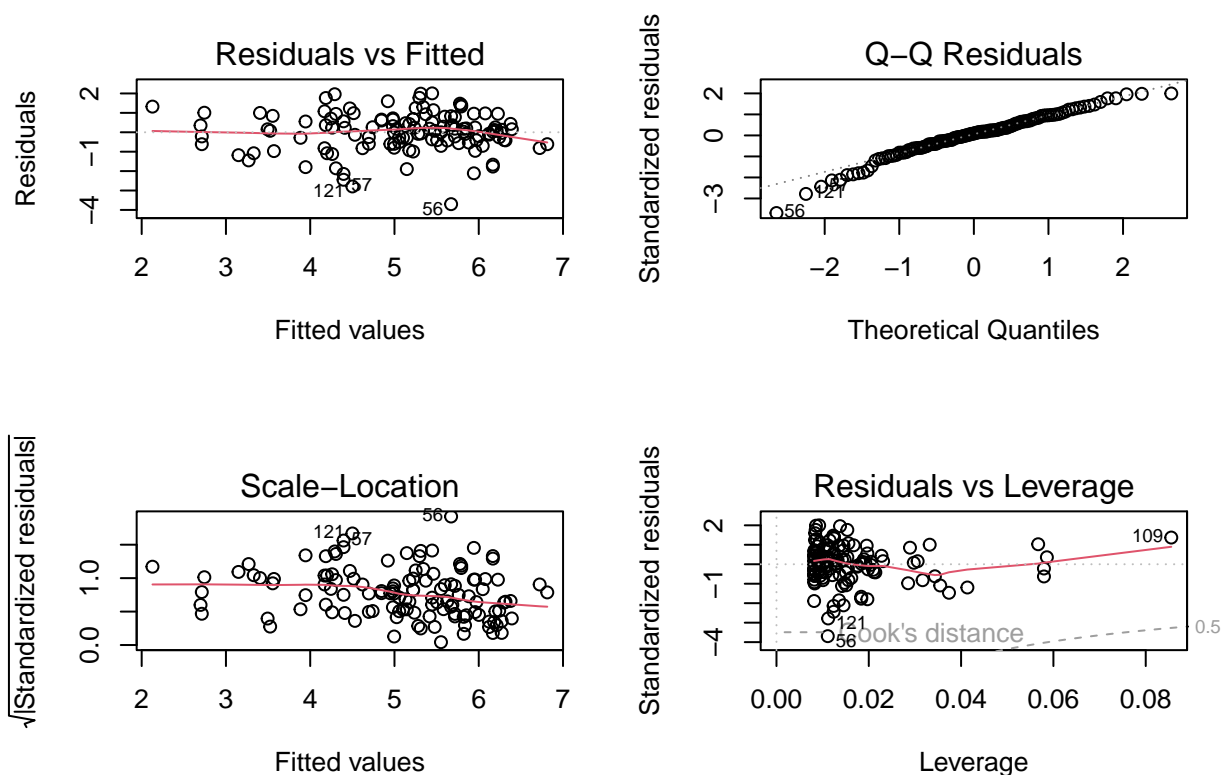
```
##
## Call:
## lm(formula = height ~ ph, data = d.basisch)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.7020 -0.5471  0.0874  0.6663  2.0033
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  28.7227     2.2395   12.82  <2e-16 ***
## ph           -3.0034     0.2844  -10.56  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.008 on 121 degrees of freedom
## Multiple R-squared:  0.4797, Adjusted R-squared:  0.4754
## F-statistic: 111.5 on 1 and 121 DF, p-value: < 2.2e-16
```

```
# Estimated equation: height = 28.7 - 3.0pH
# Drawing line into scatterplot:
plot(d.basisch$ph, d.basisch$height) + abline(fit)
```

```
## integer(0)
```

SLR - Residual Analysis Diagnostics plots are straightforward:

```
par(mfrow = c(2,2))
plot(fit)
```



1. Tukey-Anscombe plot (is the variance of the errors  $E_i$  constant? Is the regression function correct?)
2. Q-Q plot (are the errors  $E_i$  normally distributed?)
3. Scale location plot (similar to Tukey-Anscombe plot)
4. Leverage plot (what points have a strong influence on the fit?)

Residual plots by hand:

```
par(mfrow = c(1,2))
plot(fit$fitted, fit$resid)

#Tukey-Anscombe
```

```
qqnorm(fit$resid) #quantil-Quantil Plot
qqline(fit$resid) # adds the diagonal line
```

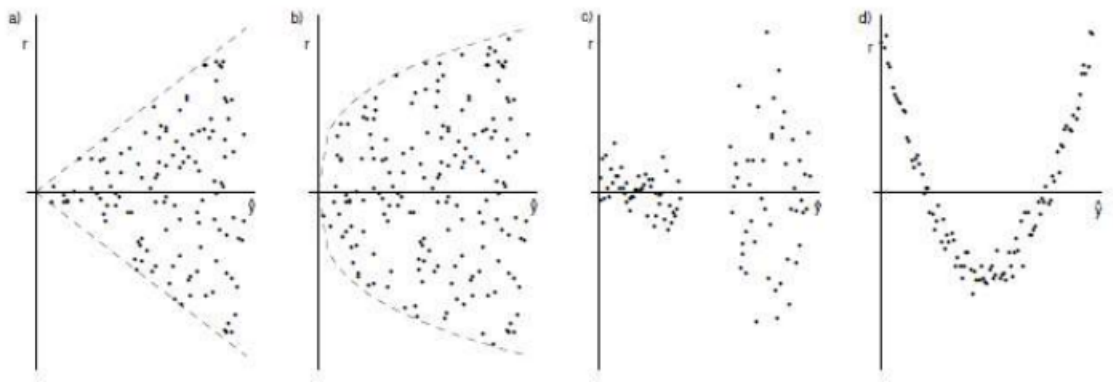
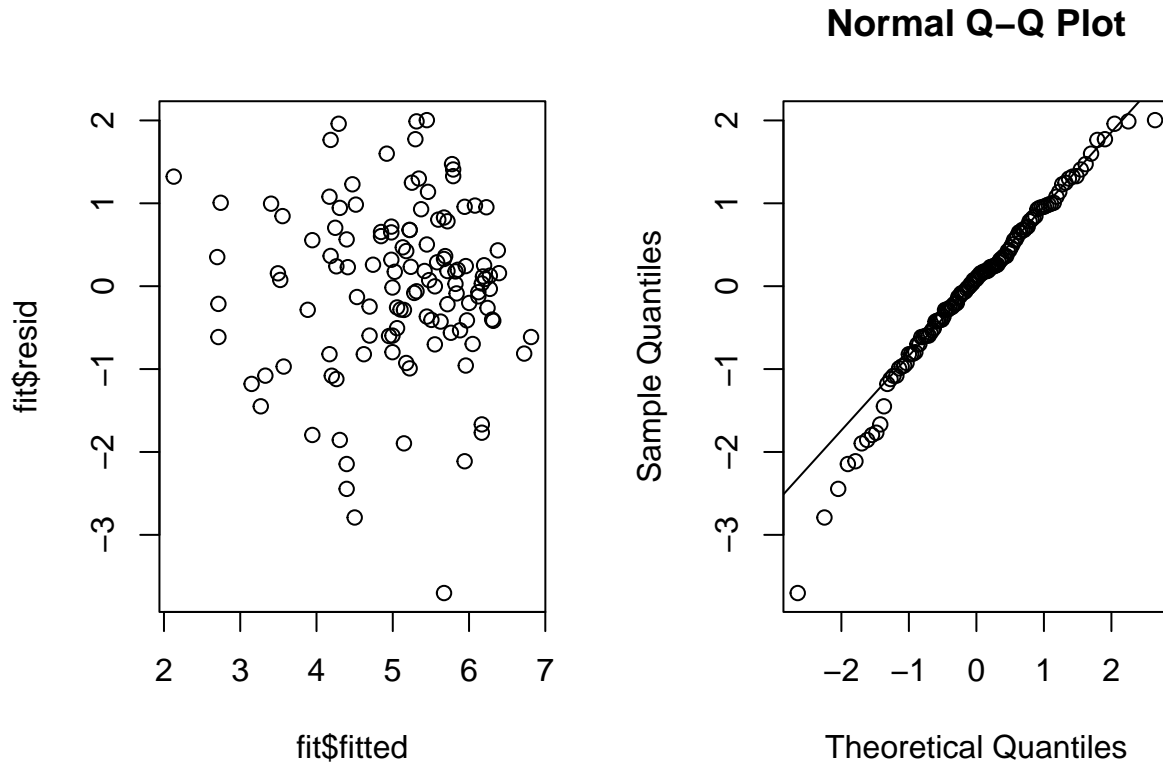


Figure 6: Some bad Tukey-Anscombe plots

## Multiple Linear Regression

Expand the simple linear model to more than one explanatory variable.

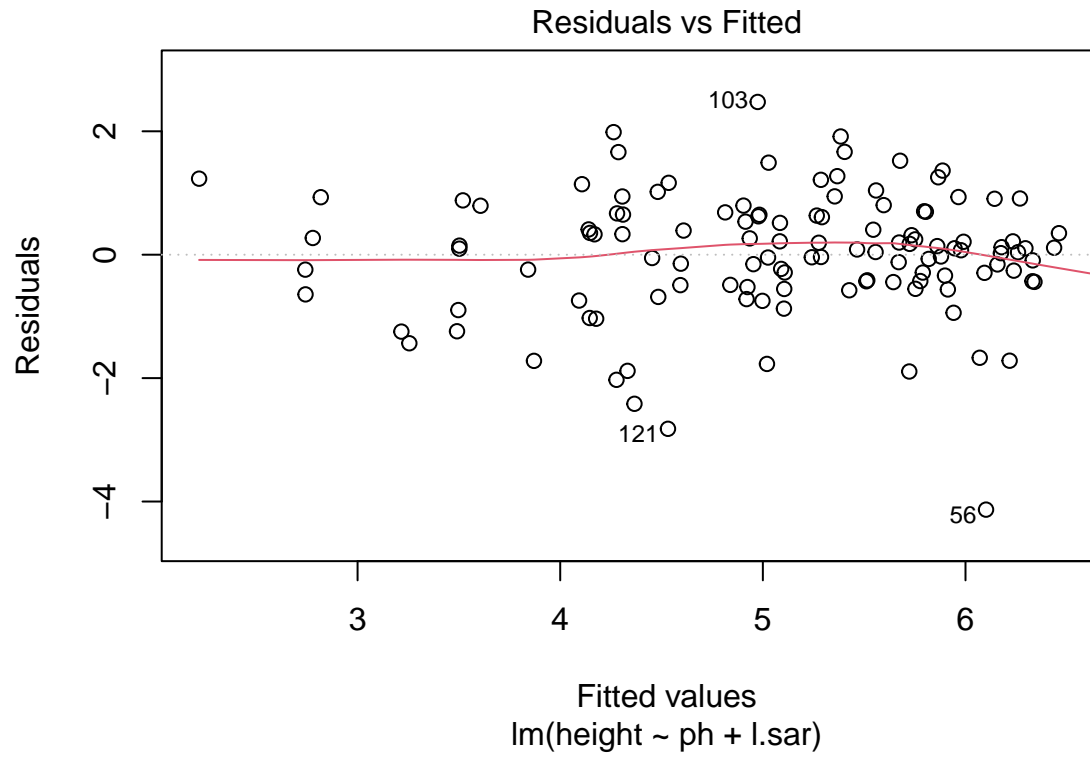
$$Y_i = \alpha + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_n x_n + E_i \quad \text{with } E_i^{i,i,d} \sim \mathcal{N}(0, \sigma^2)$$

Figure 7: The Multiple Linear Regression Model

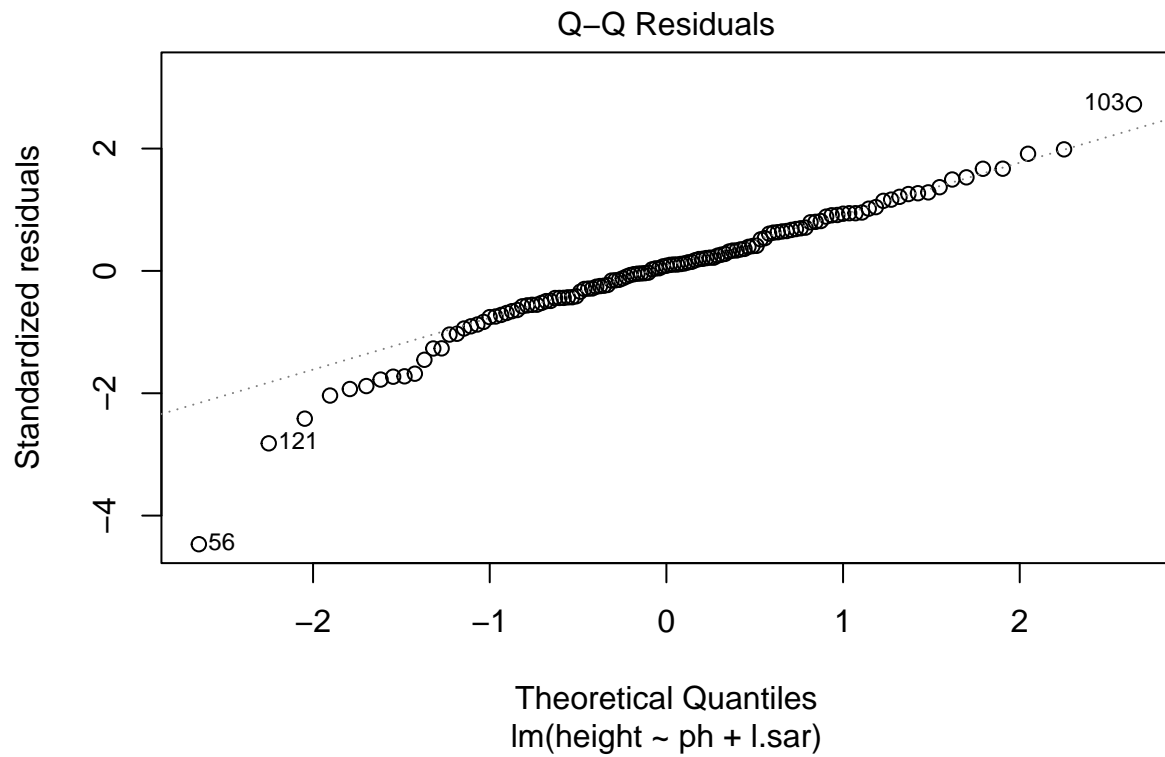
```
# Fit the model with lm
fitm <- lm(height ~ ph + l.sar, data = d.basisch)
summary(fitm)

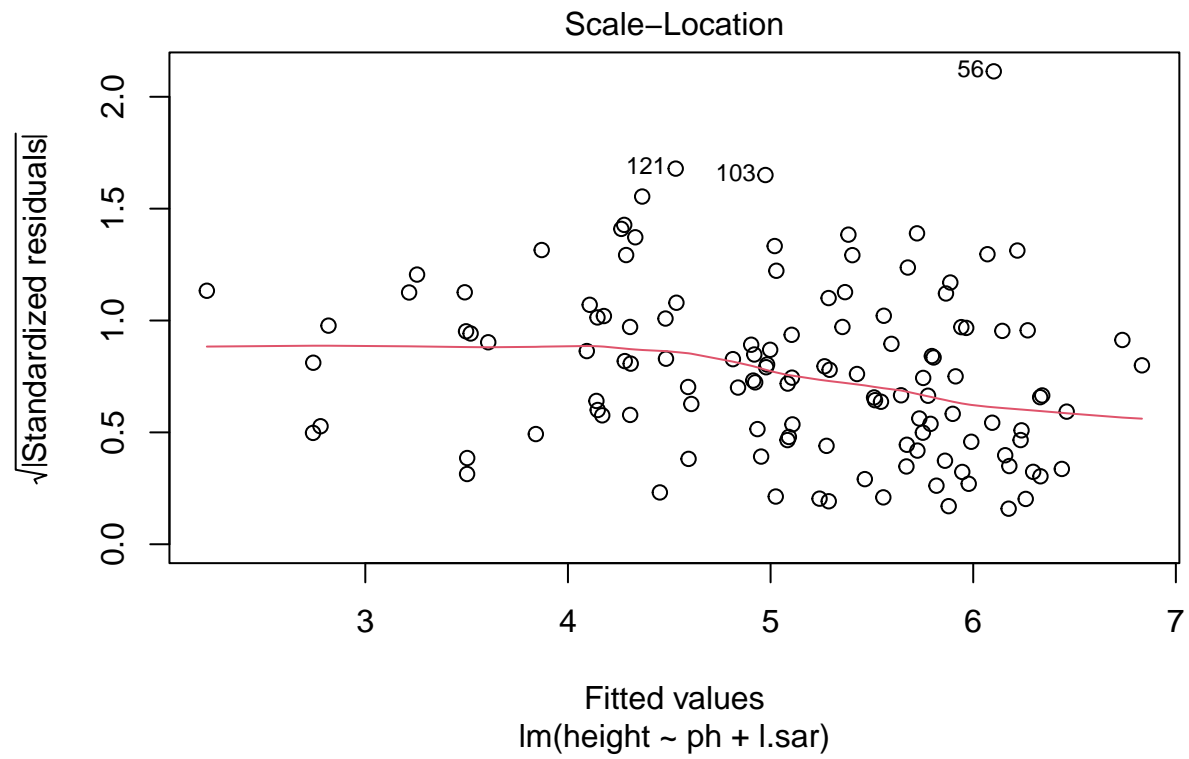
##
## Call:
## lm(formula = height ~ ph + l.sar, data = d.basisch)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.1314 -0.4911  0.0849  0.6488  2.4754
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  26.9466     2.7445   9.818 < 2e-16 ***
## ph          -2.7558     0.3603  -7.649 5.6e-12 ***
## l.sar        -0.2519     0.2255  -1.117  0.266
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.007 on 120 degrees of freedom
## Multiple R-squared:  0.485, Adjusted R-squared:  0.4764
## F-statistic: 56.51 on 2 and 120 DF, p-value: < 2.2e-16

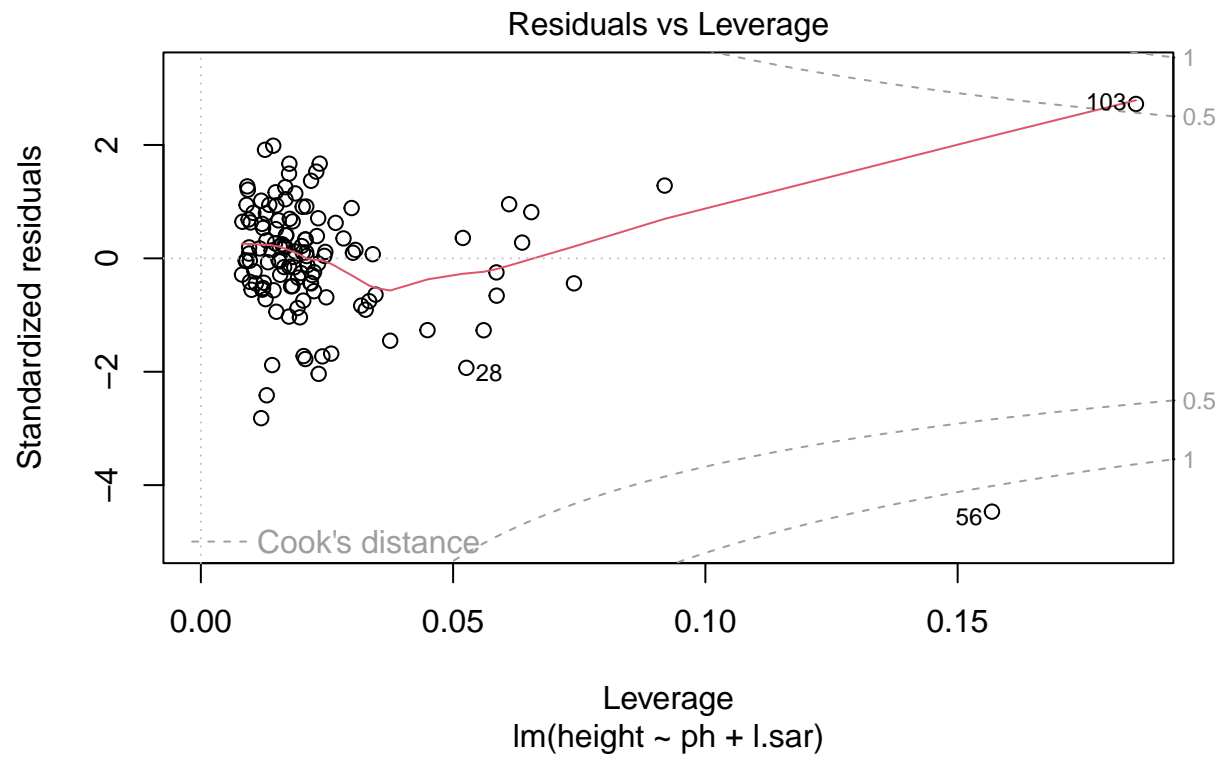
# Look at the same plots as for simple linear regression
plot(fitm)
```



MLR - Residual Analysis

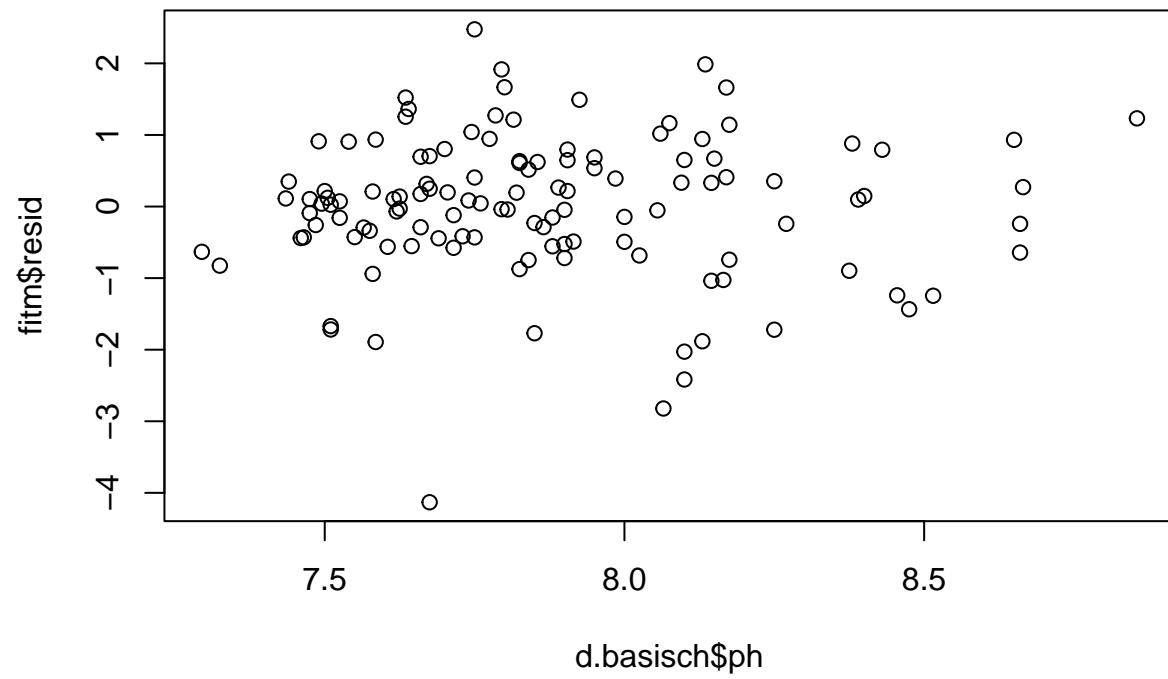






```
# It may help to plot the explanatory variables against the residuals.
plot(d.basisch$ph, fitm$resid)
```





```
plot(d.basisch$l.sar, fitm$resid)
```

