

Week 1 Lab - Introduction to R

Unit Coordinator - Dr Liwan Liyanage

School of Computing, Engineering and Mathematics

What is R?

R is a software environment for statistical computing and graphics. It runs on just about any platform (except iPad!) and is completely free (in the GNU sense).

It is used extensively by academic statisticians for research and teaching and is gaining ground in business.

It has 4634 extension packages available.

Pros

Its free and open source. It has most methods for most things mostly before any other package. It has the best graphics. It extendable.

Cons

It has a steep learning curve. No GUI by default. Poor (but improving) memory management; difficulty with very large data sets.

R Resources

- <http://www.r-project.org> - Main R website.
- CRAN - <http://cran.csiro.au> - Comprehensive R Archive Network - base software and add-on packages.
- RStudio - <http://www.rstudio.com> - is a powerful IDE for R
- R Commander - `install.package(Rcmdr)` - is a partial GUI interface to R - requires TclTk.
- R Graph Gallery - <http://gallery.r-enthusiasts.com/> - loads of pretty pictures.
- <http://cran.csiro.au/doc/contrib/Torfs+ Brauer-Short-R-Intro.pdf> - “A (very) short Introduction to R”
- “Introductory Statistics with R”, Peter Dalgaard, Springer 2008.

Access to R

For students both undergrad and postgrad: When a student enrolls in a subject or course that falls under the old SCEM banner, a CDMS (SCEM) account automatically gets created and an email explaining this is sent to their WSU student email address.

The email has a link to a site generated for each new account that contains their initial password. The web page then explains how they can change this password.

There is an R server available to all CDMS account holders.

Access is via: <https://r.cdms.westernsydney.edu.au>

Installing R

Go to (<https://cran.r-project.org/>)

Download R for your OS and install



CRAN
[Mirrors](#)
[What's new?](#)
[Task Views](#)
[Search](#)

About R
[R Homepage](#)
[The R Journal](#)

Software
[R Sources](#)
[R Binaries](#)
[Packages](#)
[Other](#)

Documentation
[Manuals](#)
[FAQs](#)
[Contributed](#)

The Comprehensive R Archive Network

Download and Install R

Precompiled binary distributions of the base system and contributed packages, **Windows and Mac** users most likely want one of these versions of R:

- [Download R for Linux](#)
- [Download R for Mac/OS X](#)
- [Download R for Windows](#)

R is part of many Linux distributions, you should check with your Linux package management system in addition to the link above.

Source Code for all Platforms

Windows and Mac users most likely want to download the precompiled binaries listed in the upper box, not the source code. The sources have to be compiled before you can use them. If you do not know what this means, you probably do not want to do it!

- The latest release (2020-06-22, Taking Off Again) [R-4.0.2.tar.gz](#), read [what's new](#) in the latest version.
- Sources of [R alpha and beta releases](#) (daily snapshots, created only in time periods before a planned release).
- Daily snapshots of current patched and development versions are [available here](#). Please read about [new features](#) and [bug fixes](#) before filing corresponding feature requests or bug reports.
- Source code of older versions of R is [available here](#).
- Contributed extension [packages](#)

Questions About R

- If you have questions about R like how to download and install the software, or what the license terms are, please read our [answers to frequently asked questions](#) before you send an email.

What are R and CRAN?

R is 'GNU S', a freely available language and environment for statistical computing and graphics which provides a wide variety of statistical and graphical techniques: linear and nonlinear modelling, statistical tests, time series analysis, classification, clustering, etc. Please consult the [R project homepage](#) for further information.



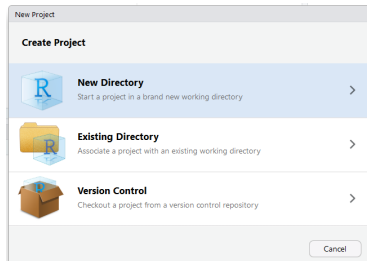
##Installing RStudio

R projects

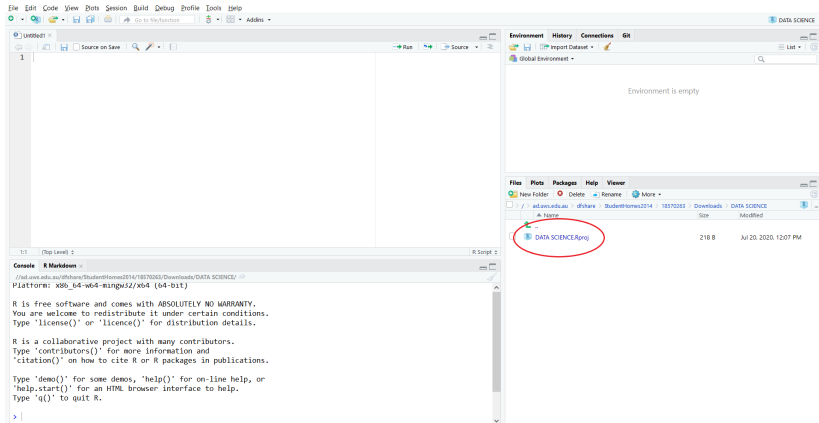
Projects are very helpful in organizing your codes. They keep all your related files together. When you create a new project, it creates a folder where all files will be kept and assign that folder as the working directory.

To create a project, go to

File->New Project->New Directory->New Project



R projects



R Commands

R can be used as a basic calculator.

```
1+1
```

```
## [1] 2
```

```
sqrt(3)
```

```
## [1] 1.732051
```

```
2^3
```

```
## [1] 8
```

R Commands ctd. . .

It can store and print variables.

```
x=10  
print(x)
```

```
## [1] 10
```

R Commands ctd...

It understands vectors and matrices.

```
x <- c(1,2)
print(x)
```

```
## [1] 1 2
```

```
m <- matrix(c(1,2,3,4), ncol=2, byrow=TRUE)
print(m)
```

```
##      [,1] [,2]
## [1,]    1    2
## [2,]    3    4
```

R Commands ctd...

It has functions, and you can write them.

```
x <- sqrt(2)
sqr <- function(x) x^2
sqr(x)
```

```
## [1] 2
```

Uploading Data Into R

When data set iris.csv is stored in a folder with path
C:/LIWAN/R/2016/Intro to Data Science, use

```
iris<- read.csv("C:/LIWAN/R/2016/Intro to Data Science/iris.csv")  
attach(iris)
```

When the data set is uploaded to the same folder where R project is
saved, use

```
iris<- read.csv("iris.csv")  
attach(iris)
```

Alternatively use "Import Dataset" function in the Environment

Data in R

Tables are stored in data.frames

```
head(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

```
sapply(iris,class)
```

##	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
##	"numeric"	"numeric"	"numeric"	"numeric"	"factor"

Summary

```
names(iris)
```

```
## [1] "Sepal.Length" "Sepal.Width" "Petal.Length" "Petal.Width"  
## [5] "Species"
```

```
dim(iris)
```

```
## [1] 150 5
```

```
summary(iris)
```

```
##      Sepal.Length      Sepal.Width      Petal.Length      Petal.Width  
## Min.      :4.300    Min.      :2.000    Min.      :1.000    Min.      :0.100  
## 1st Qu.:5.100    1st Qu.:2.800    1st Qu.:1.600    1st Qu.:0.100  
## Median :5.800    Median :3.000    Median :4.350    Median :1.300  
## Mean   :5.843    Mean   :3.057    Mean   :3.758    Mean   :1.326  
## 3rd Qu.:6.400    3rd Qu.:3.300    3rd Qu.:5.100    3rd Qu.:1.800
```

Basic Statistics

```
x<-rnorm(100)
mean(x)
```

```
## [1] -0.001880122
```

```
var(x)
```

```
## [1] 0.8723308
```

```
sd(x)
```

```
## [1] 0.9339865
```

```
fivenum(x)
```

```
## [1] -1.91539426 -0.73196551 0.03462727 0.69070222 1.8332
```

```
minimum lower quartile median upper quartile maximum
```


Basic Statistics

```
t.test(x)
```

```
##  
## One Sample t-test  
##  
## data: x  
## t = -0.02013, df = 99, p-value = 0.984  
## alternative hypothesis: true mean is not equal to 0  
## 95 percent confidence interval:  
## -0.1872033 0.1834431  
## sample estimates:  
## mean of x  
## -0.001880122
```

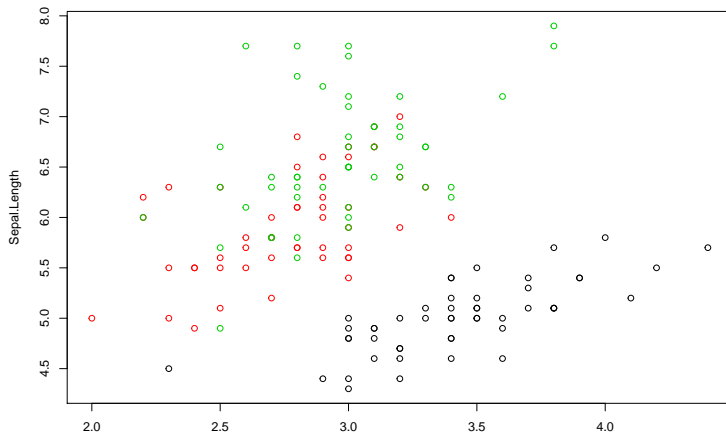
Basic Statistics

```
t.test(iris$Sepal.Length)
```

```
##  
## One Sample t-test  
##  
## data: iris$Sepal.Length  
## t = 86.425, df = 149, p-value < 2.2e-16  
## alternative hypothesis: true mean is not equal to 0  
## 95 percent confidence interval:  
## 5.709732 5.976934  
## sample estimates:  
## mean of x  
## 5.843333
```

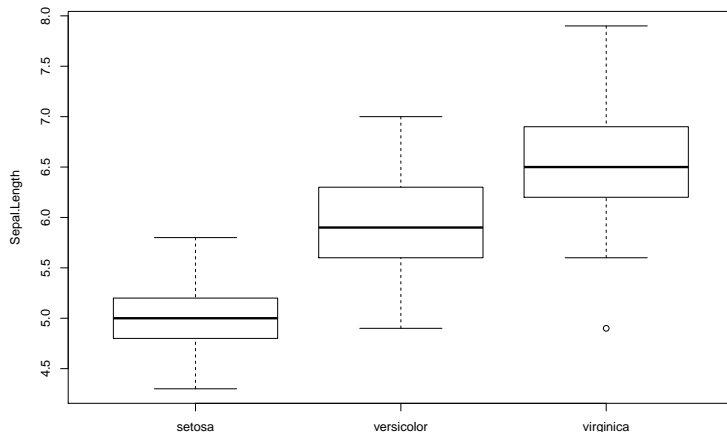
R has extensive plotting

```
plot(Sepal.Length~Sepal.Width, col=Species, data=iris)
```



R has extensive plotting

```
boxplot(Sepal.Length~Species, data=iris)
```



Help in R

Everything in R has a help file.

```
help(t.test)
```

Or see the help pane in RStudio

Will illustrate R further within Regression Analysis

Getting Ready for the Data Analysis covered in Lectures

Data Import

```
install.packages("ISLR")
```

```
install.packages("MASS")
```

```
library(ISLR)
```

```
library(MASS)
```

```
library(class)
```

```
library(DMwR)
```

```
attach(Smarket)
```

```
attach(Boston)
```

```
attach(Carseats)
```

```
attach(iris)
```

Data Sets

Supervised Learning:

- Advertising
- Income
- Heart
- Smarket
- Caravan (Insurance Data)

Unsupervised Learning:

- USAarrests
- groceries

View Advertising Data and Discuss How to Initiate Knowledge Discovery

Exercise: List Possible Research Questions?

```
Advertising<-read.csv("Advertising.csv")
attach(Advertising)
head(Advertising)
```

##		TV	Radio	Newspaper	Sales
## 1		230.1	37.8	69.2	22.1
## 2		44.5	39.3	45.1	10.4
## 3		17.2	45.9	69.3	9.3
## 4		151.5	41.3	58.5	18.5
## 5		180.8	10.8	58.4	12.9
## 6		8.7	48.9	75.0	7.2

View heart Data and Discuss How to Initiate Knowledge Discovery

Exercise: List Possible Research Questions?

```
Heart<-read.csv("heart.csv")
attach(Heart)
head(Heart)
```

##	X	Age	Sex	ChestPain	RestBP	Chol	Fbs	RestECG	MaxHR	ExAr
## 1	1	63	1	typical	145	233	1	2	150	
## 2	2	67	1	asymptomatic	160	286	0	2	108	
## 3	3	67	1	asymptomatic	120	229	0	2	129	
## 4	4	37	1	nonanginal	130	250	0	0	187	
## 5	5	41	0	nontypical	130	204	0	2	172	
## 6	6	56	1	nontypical	120	236	0	0	178	
##	Ca			Thal	AHD					
## 1	0			fixed	0					
## 2	0			normal	1					

View groceries Data and Discuss How to Initiate Knowledge Discovery

Exercise: List Possible Research Questions?

```
Groceries<-read.csv("groceries.csv")  
attach(Groceries)  
head(Groceries)
```

```
## frankfurter sausage liver.loaf ham meat finished.products  
## 1 0 0 0 0 0  
## 2 0 0 0 0 0  
## 3 0 0 0 0 0  
## 4 0 0 0 0 0  
## 5 0 0 0 0 0  
## 6 0 0 0 0 0  
## organic.sausage chicken turkey pork beef hamburger meat  
## 1 0 0 0 0 0  
## 2 0 0 0 0 0
```

Explore Default Data set from the ISLR Library

```
#install.packages("ISLR")
```

```
library(ISLR)
```

```
attach(Default)
```

```
View(Default)
```

```
dim(Default)
```

```
## [1] 10000      4
```

```
head(Default)
```

```
##      default student   balance   income
## 1         No      No  729.5265 44361.625
## 2         No     Yes  817.1804 12106.135
## 3         No      No 1073.5492 31767.139
## 4         No      No  529.2506 35704.494
## 5         No      No  785.6559 38463.496
```

Save and read datasets within R

Save a data set downloaded from a library within R as a csv file

```
write.csv(Default,file="Default.csv")
```

To read csv files

```
read.csv("Default.csv",header=TRUE)
```

To read any other files read.table("file",header=False)

How to change a factor variable to a numeric variable

Add another variable named Defcode to table Default check the levels of the new variable (It will be same class as the original variable)

```
Defcode = Default$default  
levels(Defcode)
```

```
## [1] "No"  "Yes"
```

Change the levels as 1 for Yes and 0 for No

```
levels(Defcode) [levels(Defcode)=="No"]=0  
levels(Defcode) [levels(Defcode)=="Yes"]=1  
levels(Defcode)
```

```
## [1] "0" "1"
```

Still Defcode variable is a factor variable and cannot use as a numeric variable in regression setting.

To summarise a factor variable

```
Defcode = as.character(Default$default)
table(Defcode)
```

```
## Defcode
##      No   Yes
## 9667  333
```

```
Defcode[Defcode=="No"]=0
Defcode[Defcode=="Yes"]=1
table(Defcode)
```

```
## Defcode
##      0    1
## 9667  333
```

Change a factor variable to a numeric variable

```
Defcode = as.numeric(Defcode)  
class(Defcode)
```

```
## [1] "numeric"
```


Exercises

For each of the three data sets, iris, heart and groceries

- Explore the variables
- List the quantitative variables and qualitative variables
- State a Research question and identify the target variable if applicable
- Comment if they are supervised learning or unsupervised learning.