

# **Data Sciences Course**

## **(Day X - Session X)**

### **R Simplified**

**Covers**

**R basics and startup for data summarization**

## ✓ Getting Data into R:

- Several ways – lets have overview of some in day-to-day functions

### ✓ Entering data

- “c” function – combine or concatenates terms
- Ex: XYZ <- c(10, 15, 12, 18)
- Ex: matrix(c(10, 15, 12, 18), nrow=2, ncol=2)
- Ex: matrix(c(10, 15, 12, 18), nrow=2, ncol=2, byrow = TRUE)
- Ex: array(c(10, 15, 12, 18), dim = c(2, 2))
- Ex: data.entry()

### ✓ Reading external data

- icities <- read.table("India\_Cities.csv", header = T, sep = ",", dec = ".", nrows = -1, comment.char = "", strip.white = T, stringsAsFactors = F)
- will read the external file as a data.frame

### ✓ Reading/Loading data

- load("/dir/icities.RData")
- library(xx); data(xx);
- recollect command;
- try(data(package = "ggplot2"))
- data(cancer, package = "survival")

### ✓ Generating data

- Ex: a1 <- rnorm(100)
- Ex: a2 <- seq(1, 10, by = 0.01)
- Ex: a3 <- rep(1, 100)
- Ex: a4 <- runif(100)
- Ex: cbind(a1, a2, a3, a4)
- Ex: rbind(a1, a2, a3, a4)

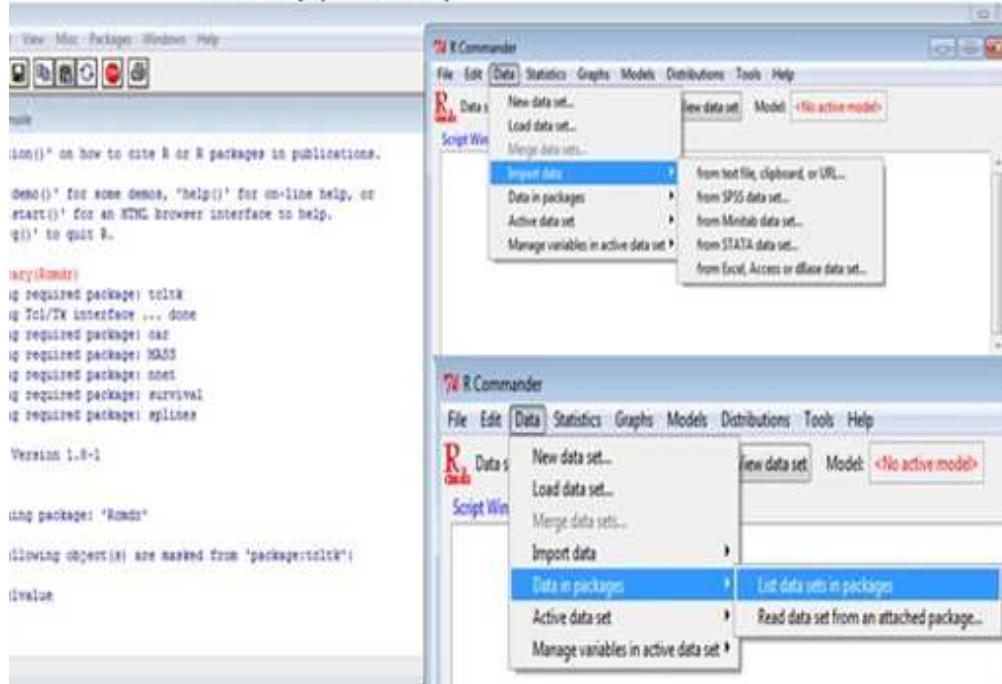
### ✓ Reading External data

- read.table.url(url) ;
- download.file(url); url.show(url)
- read.csv(f, header=TRUE)
- read.csv2(f, header=TRUE, sep="; quote="\\"", dec=",")
- read.delim(f, header=TRUE, sep="\t", quote="\\"", dec=".")
- read.delim2(f, header=TRUE, sep="\t", quote="\\"", dec=",")
- read.fwf(file, widths=c(3,5,3), header=FALSE, sep="", as.is=FALSE)
- as.is=TRUE; # not to be converted into a factor
- na.strings=c(".", "NA", "", "#") # characters for missing
- save(icities, file = "icities.Rdata")

# Understanding R Easily – Using GUI's

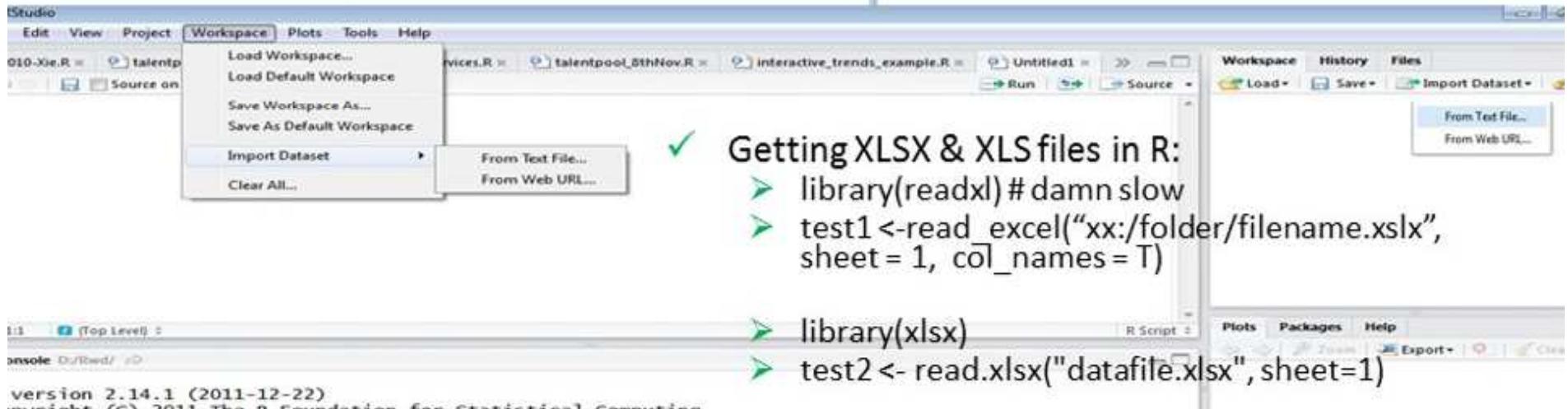
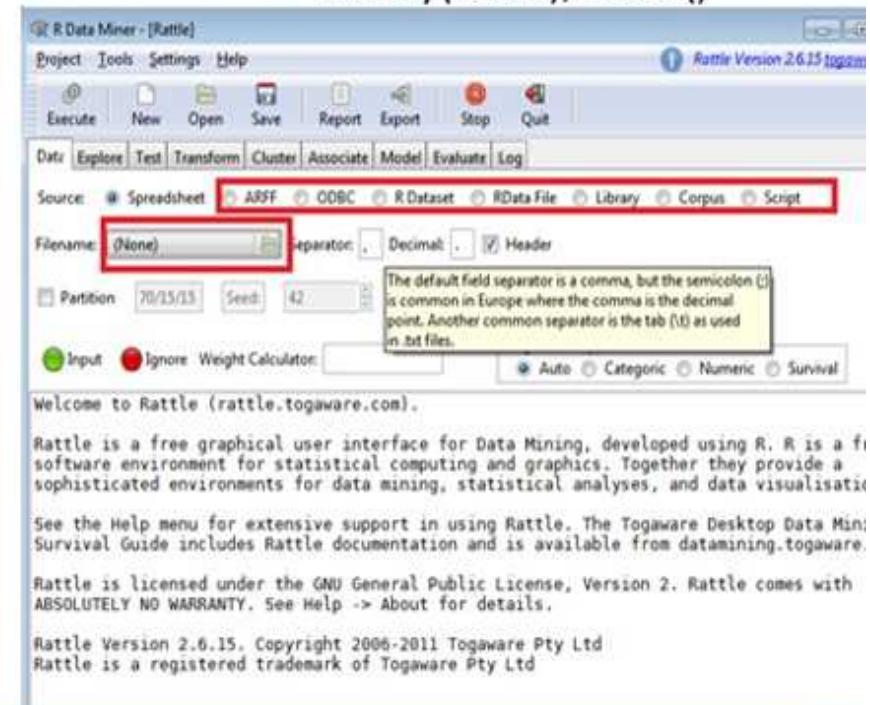
## ✓ Getting Data in R:

- library(Rcmdr)



## ✓ Getting Data in R:

- library(rattle); rattle()



## ✓ Getting XLSX & XLS files in R:

- library(readxl) # damn slow
- test1 <- read\_excel("xx:/folder/filename.xlsx", sheet = 1, col\_names = T)
- library(xlsx)
- test2 <- read.xlsx("datafile.xlsx", sheet=1)

## Understanding R Easily – Now lets us know what is in read data

- ✓ `str()` # compactly displays the structure of an arbitrary R object
- ✓ **variable info:**
  - `names(df)`
  - `length(df)` # number of elements
  - `attributes(df)`; # more details than "str()"
  - `nrow(df); ncol(df)`
  - `is.data.frame(df)`
  - `is.character(df$x1)`
  - `is.numeric(df$x1)`
  - `is.matrix(df)`
  - `dim(df)`
  - `dimnames(df)`
- ✓ Best practice is to have a data summary sheet.



### Useful functions and operators for understanding and analyzing the data

Mathematical & Statistical		Arithmetic		Logical	
<code>log</code>	Natural Logarithm	<code>+</code>	Addition	<code>!X</code>	Logical negation (NOT)
<code>exp</code>	Exponential	<code>-</code>	Subtraction	<code>X &amp; Y</code>	Logical And
<code>sin</code>	SINE	<code>*</code>	Multiplication	<code>X   Y</code>	Logical OR
<code>cos</code>	COSINE	<code>/</code>	Division	<code>xor</code>	indicates elementwise exclusive OR
<code>tan</code>	Tangent	<code>^</code>	Power	Matrices	
<code>sqrt</code>	Square Root	Comparison		<code>t(x)</code>	Transpose of x
<code>min</code>	Returns minimum among numbers	<code>&lt;</code>	Lesser than	<code>diag(x)</code>	Diagonal of x
<code>max</code>	Returns maximum among numbers	<code>&gt;</code>	Greater than	<code>%*%</code>	Matrix Multiplication
<code>sum</code>	Sum of numbers	<code>&lt;=</code>	Lesser than or equal to	<code>solve(x)</code>	Matrix inverse of x
<code>mean</code>	Averages of numbers	<code>&gt;=</code>	Greater than or equal to	<code>rowSums(x)</code>	Sum of rows for a matrix-like object
<code>var</code>	Variance of numbers	<code>==</code>	Equal	<code>colSums(x)</code>	Sum of columns for a matrix-like object
<code>sd</code>	Standard Deviation of numbers	<code>!=</code>	Not equal to	<code>rowMeans(x)</code>	Average of rows for a matrix-like object

# Understanding R Easily – Data Summarization

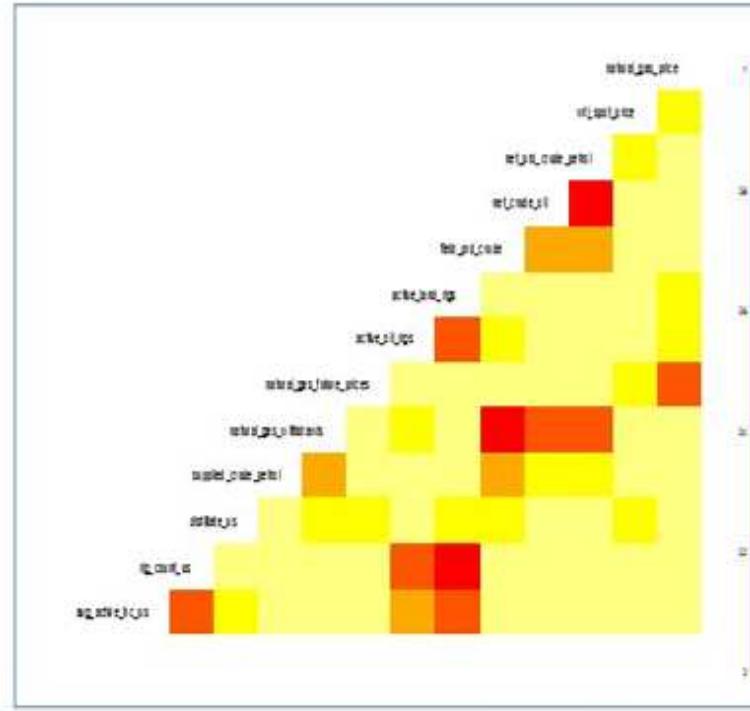
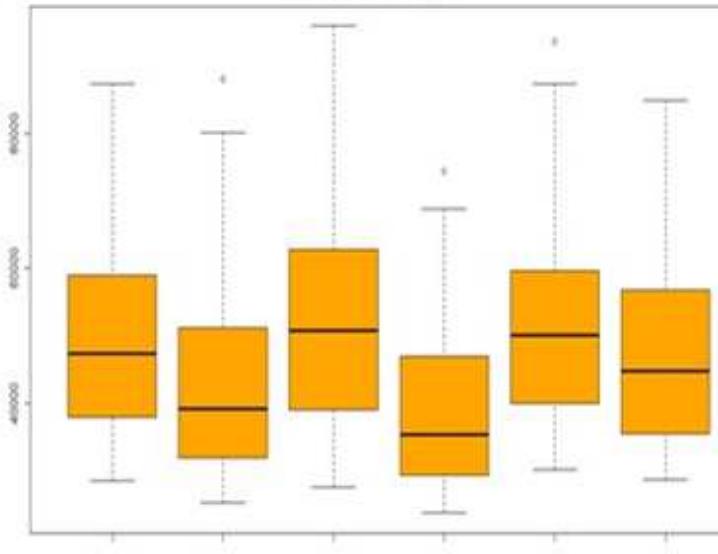
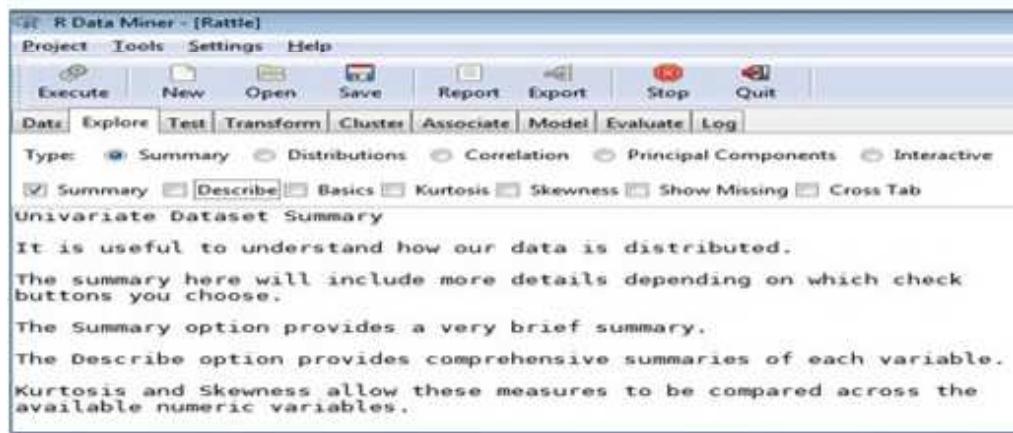
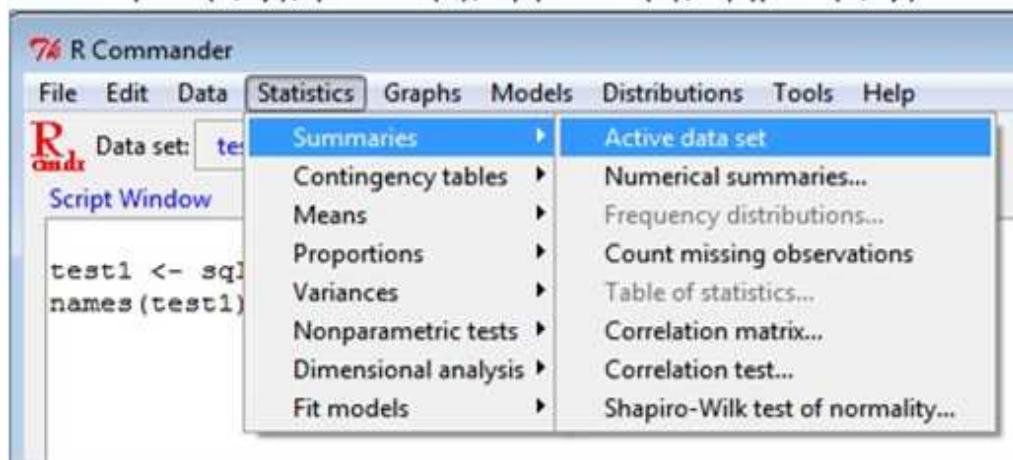
## Box Plot

### ✓ Data summarization/descriptive functions

- `summary(df)`; `fivenum(numeric)`; `cor(df)`;
- `library(Hmisc)`; `describe(df)`;
- `mean()`; `median()`; `max()`; `min()`; `std()`; `var()`;
- `library(arm)`; `corrplot(df, col = T)`;

### ✓ Descriptive plots

- `boxplot(df)`; `dotchart(df)`; `paris(df)`
- `hist(df$x1)`; `barplot(df$x1)`
- `plot(x, y)`; `plot.ts(x)`; `qqnorm(x)`; `qqplot(x, y)`



## ✓ Has built-in help pages?

- `help(t.test)` or `?t.test`
- `help.start()` # starts help page (both in R GUI & Rstudio)
- `search()` # prints the list of default loaded/attached packages
- `help.search("linear models")` # searches the help systems and provides all matching string results

R R Information - Help for "t.test"

t.test package:ctest R Documentation

**Student's t-Test**

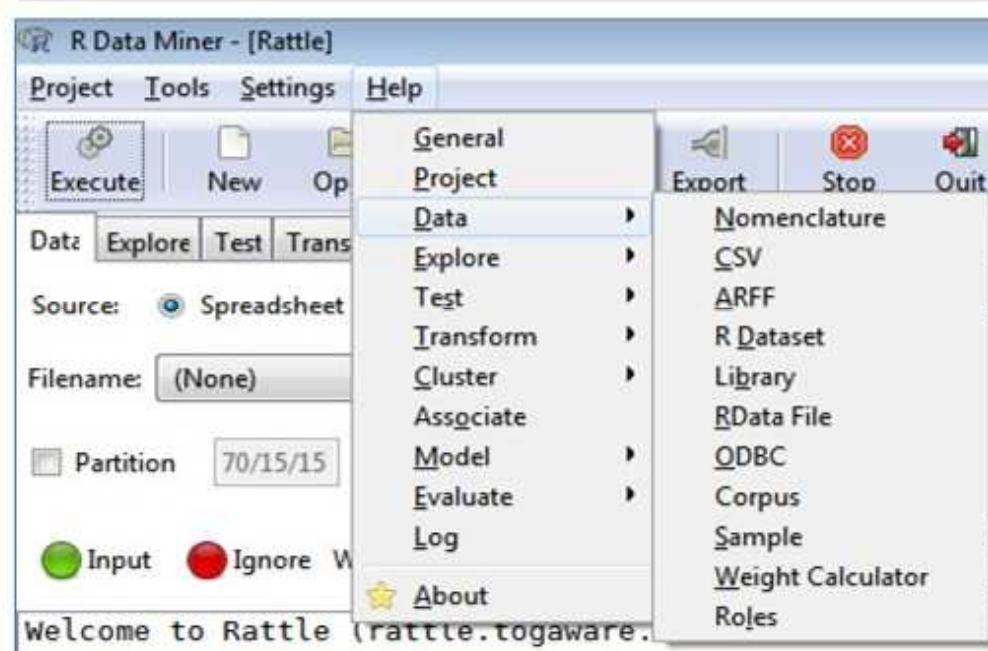
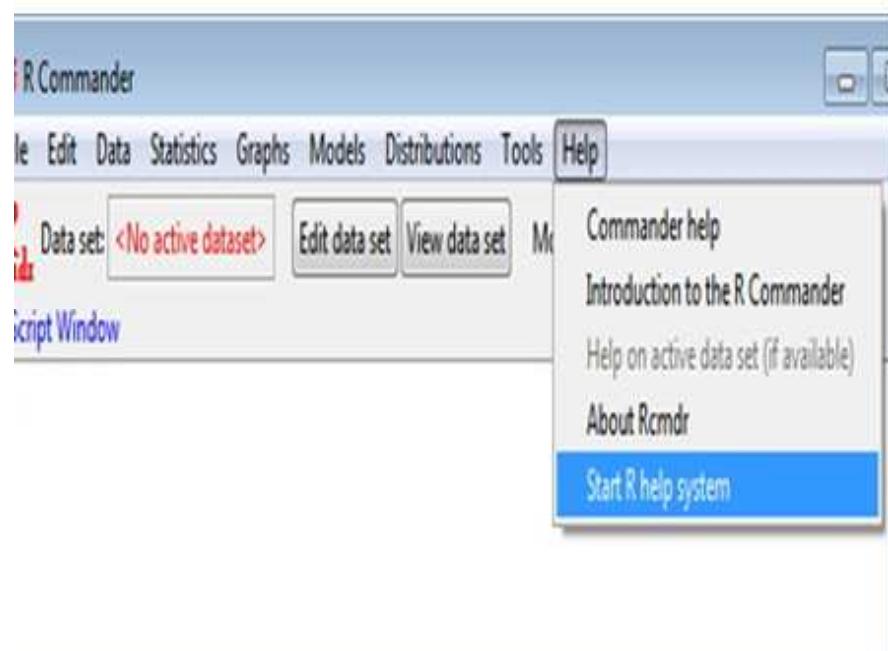
**Description:**  
Performs one and two sample t-tests on vectors of data.

**Usage:**

```
t.test(x, y = NULL, alternative = c("two.sided", "less", "greater"),
       mu = 0, paired = FALSE, var.equal = FALSE,
       conf.level = 0.95, ...)
t.test(formula, data, subset, na.action, ...)
```

**Arguments:**

- `x`: a numeric vector of data values.
- `y`: an optional numeric vector data values.
- `alternative`: a character string specifying the alternative hypothesis, must be one of "two.sided" (default), "greater" or "less". You can specify just the initial letter.
- `mu`: a number indicating the true value of the mean (or difference in means if you are performing a two sample test).
- `paired`: a logical indicating whether you want a paired t-test.
- `var.equal`: a logical variable indicating whether to treat the two



## ✓ Dates & Times

- 'as.Date()' – built-in function – to handle dates.
- 'chron' package has few better options to handle dates, but cannot handle time zones.
- 'POSIXct' and 'POSIXlt' – built-in functions – for handle/control time zones.
- 'Sys.Date()' prints today's date and 'Sys.time()' prints both today's date and current time.
- R also follows – default date of origin – 1970-01-01 - programming languages - it is day zero.
- Except for - POSIXlt class - all dates are stored internally - as the number of days or seconds from reference date.
- Thus, dates will have - a numeric mode - however - class function – explains - format.
- Default format/type of dates – "YYYY-MM-DD" which in R takes the code value as "%Y-%m-%d".
- Coming to POSIXlt class - it stores - date/time values - as a list of components (hour, min, sec, mo etc.) - making them easy for extract time and date components.
- e.g. as.Date('1995-6-16') prints "1995-06-16"
- e.g. as.Date('2/14/2002',format='%m/%d/%Y') prints "2002-02-14"
- e.g. as.Date('April 27, 2009',format='%B %d, %Y') prints "2009-04-27"
- e.g. as.Date('21JUL11',format='%d%b%y') prints "2011-07-21" # better not to use 2-digit year

Time Unit	R-symbol & meaning	Illustration
Day	(%d) Day as a number	01-31
Days of the year	(%j) Days of year as decimal number	001-366
Hours (24)	(%H) Hours as decimal number	00-23
Hours (12)	(%I) Hours as decimal number	01-12
Abbreviated Weekday	(%a) Abbreviated Weekday	Mon
Unabbreviated Weekday	(%A) Unabbreviated Weekday	Monday
Weekday	(%w) Weekday as decimal number	0-6 (0=Sunday)
Week of the year	(%W) Week of the year as decimal	First Monday as day 1 of week 1
Week of the year	(%U) Week of the year as decimal	First Sunday as day 1 of week 1
Month	(%m) Month as a number	01-12
Abbreviated Month	(%b) Abbreviated Month	JUL

## Understanding R Easily – Date & Time Formats (contd.)

Time Unit	R-symbol & meaning	Illustration
Unabbreviated Month	(%B) Unabbreviated Month	JULY
Year	(%Y) Four-Digit Year	2011
Year	(%y) Two-Digit Year	11
Minute	(%M) Minute as decimal number	00-59
AM/PM	(%p) AM, PM indicator in locale (used with Hours)	AM or PM
Seconds	(%S) Second as decimal number (allows for up-to 2 leap seconds)	00-61
Time Zone	(%z) Abbreviated Time Zone	IST
Locale Options	(%x) Date, Locale Specific Option	Locale Specific
Locale Options	(%X) Time, Locale Specific Option	Locale Specific
Locale Options	(%c) Date and Time, Locale Specific Option	Locale Specific

✓ Don't get confused with date & time classes with data classes and also try 'lubridate' package

Date & Time Class	Data Class
Date or as.Date	ts
POSIXct or as.POSIXct	data.frame
POSIXct or as.POSIXct	vector
numeric	matrix
character	fts
chron	timeSeries
yearmon	mts
yearqtr	its
timeDate	irts
force_tz	zoo



Session open for questions

Thank You