

ECON 172 - Section 1 : Introduction to R

Gwyneth Miner and Jose Azcunaga

Spring 2022

1 Welcome to Econ 172 Section

Hi! We are your GSIs for Econ 172 this semester, Gwyneth and Jose. Our goal is to make these sections as interactive as possible. The more you and your classmates participate, the more everyone will get out of section. We will mainly focus on reviewing the econometric methods discussed in lecture and applying them using R.

Gwyneth gwyneth_miner@berkeley.edu

- Tuesday 10-11 AM (102)
- Thursday 9-10 AM (104)
- Friday 9-10 AM (106)

Jose jedazcunaga@berkeley.edu

- Monday 2-3 PM (101)
- Wednesday 3-4 PM (103)
- Thursday 5-6 PM (105)

All of the Zoom links for section are on bCourses. We will also post all materials on bCourses. As in lecture, we'd appreciate if you keep your camera on for section, but understand that this may not always be possible for all students. We understand that things are still difficult this semester and we will try to be as accommodating as possible, so please communicate any issues to us and we will try our best to help.

Feel free to email with any questions you may have that can be answered over email, for example, administrative questions. Please include [ECON 172] in the subject line of any email communications. We will do our best to respond to all emails within 3 business days. For in-depth questions that will require a longer response, stick around after section or email us asking for some time to meet outside of class. Often an in-person discussion can help clear things up more effectively than communication over email.

2 Welcome to R

R is a programming language and environment for statistical computing and graphics which is free to use and runs on Windows, Unix and MacOS. To do analysis in R, one can either write one's own commands or take advantage of the huge number of pre-existing packages written by R users and freely available. We are not assuming any prior knowledge of R. If you are not familiar with programming, it can be a bit tricky at first but we are here to make sure you get the hang of it. R is used throughout the social sciences and in industry, so while this course will provide a basic overview, you may find it worthwhile to dive a bit deeper and really master it.

2.1 Getting Started

2.1.1 Using R on your local machine (optional)

R can be downloaded from <https://cran.rstudio.com/>. We also recommend downloading RStudio, a user interface which makes R easier to use. It can be downloaded from <https://www.rstudio.com/products/rstudio/download/>.

If you run R on your local computer, you will need to install any packages you want to use (just one time), then load them into R's library for the current session. After that, it's easy to access any of the commands in the package (which are also generally extensively documented in a sort of user guide for the package).

2.1.2 Using R on your UC Berkeley Datahub (recommended)

RStudio also runs on the UC Berkeley Datahub, so you don't need to install it locally on your computer. To start RStudio on Datahub, simply go to: <https://r.datahub.berkeley.edu/>

In our sections (and for the problem sets), we are going to show you this way of using RStudio. However, anything that we show you should also work on your local machine.

2.2 R Code

When we use R, we want to write and save our commands in a script. This will let us reproduce everything we do without having to retype our commands every time we work on a project. To do this, you can write and save a file as a .R file, which typically contains only code and code comments. You can find a separate .R file containing the code used to conduct the analysis below in the folder with the material for this section.

2.3 RMarkdown

What's even nicer about R is a file format called RMarkdown, such as the one you are currently working in. This file type allows you to combine text, mathematical equations (in a format called LaTeX), code as well as dynamic output from this code, and automatically compile / output everything into a pdf document. (In RMarkdown language, this process is called "knitting" the document.)

2.4 Setting up the Environment

First, we need to set up the R environment, including installing all the required packages, loading all the required packages. This is done using R-code. The way to let RMarkdown know that the following output is a 'code chunk' is by enclosing the code chunk in triple quotations, as below. The code chunk indicator is followed by curly brackets, and r (indicating this code chunk is in the R language), a name for the code chunk (here "setup"), and some additional parameters. The parameter "results='hide'" and "message=FALSE" lets RMarkdown know that we will not want to include the output of or any messages generated by this code chunk in the pdf document (but the code will still run).

In your code, you can add comments, using # (beginning of the line) and ## (end of the line). These comments are not run as code, but will help your later self (and others reading your code) to understand what you were doing.

```
knitr::opts_chunk$set(echo = TRUE) ## sets global option to include code chunks in-line
#install.packages("haven") ## only run once to install
library(haven) #this library allows you to load datasets in Stata format

#install.packages("tidyverse") ## only run once to install
library(tidyverse) #provides _many_ functions useful for data analysis

#install.packages("summarytools") ## only run once to install
library(summarytools) #provides tools to summarize data
```

```
## Warning in fun(libname, pkgname): couldn't connect to display ":0"
#install.packages("stargazer") ## only run once to install
library(stargazer) ##This package is great for making tables in .html, .tex and many other formats.
```

The code chunk specifies that, by default, code chunks are part of the pdf output. Next, it installs a few packages. You only need to run this installation once. Afterwards, the package will remain installed on your datahub. (Remove the # before the install.packages commands to install each package for the first time, then “comment out” these lines using #, as above.) Next, the “library” command lets R know that you will want to use these packages in your code.

2.5 Loading Data

The first step in using R for statistical analysis is to store data in what R calls a “data frame,” a matrix whose columns have different modes (like numeric, words, etc.) and each row is a unit of observation (i.e: a person or a country, etc.)

R can load data from many types of files, including .csv and .txt. R has its own type of data file, with a .rds or .Rdata extension, and can also load data from Stata, SPSS, and SAS (other statistics computing programs) with the appropriate commands. The .rds format allows for a single object to be saved at a time, like a single dataset. The .RData format allows for multiple objects to be saved at once. Either one is a great way to save your data in R format.

For this section, we will be using data from Robinson, Acemoglu and Johnson (2001). For your convenience, the following link imports this data directly into your Datahub repository: <https://r.datahub.berkeley.edu/hub/user-redirect/git-pull?repo=https%3A%2F%2Fgithub.com%2Fwynethminer%2FECON172-SP22&urlpath=rstudio%2F&branch=main> The folder should now appear as “ECON172-SP22” on your workspace (bottom right). Click on the folder, then “Section_1” and then click on the .Rmd file to open the file that produced the PDF you’re reading.

When we load a dataframe in R, we give it a name for reference, such as `mydata` or something more specific, like `prop_rights`. We also need to tell R where to find the data. We can do that by specifying a full file path to its exact location of the data. For the purposes of this section, we’ll use the dataset `prop_right` which is an example dataset that contains data on a property rights index (protection against expropriation risk), GDP per capita and other variables for countries that were European colonies in the past.

```
## Load in data in csv form, from the folder just loaded into your directory.
## Important: Paths are always relative to the RMarkdown file location
##           Alternatively, specify your working directory using
##           knitr::opts_chunk$set(root.dir = "xxx")
prop_rights <- read.csv("prop_rights.csv")
```

Next we can visualize the database we just imported, and take a look at what variables are included. To do so, it’s easiest to click directly on the dataset in your Environment pane on the top right in RStudio.

```
## Let's take a look on the dataset.
prop_rights ## print the dataframe to screen
head(prop_rights) ## or print first few observations

## Let's take a look the variables in this dataset.
names(prop_rights) ## just print
prop_rights_variables <- names(prop_rights) ## store list, assigned "prop_rights_variables"
prop_rights_variables ## take a look at this list
```

2.6 Manipulating Data Frames

Often, you might want to augment your data frame, perhaps by keeping a subset of observations or variables, creating new variables, or cleaning data. The following are some basic commands that you can use to play

with your data.

```
#create a new df of only Africa countries
africa <- filter(prop_rights, africa==1)

#Create a new variable equal to GDP (rather than log)
africa$levelGDP <- exp(africa$logGDP)

#Keep only the variables "Country, Protection, Level GDP, and Log Settler Mortality
africa <- select(africa, country, protection, levelGDP, logsettlermortality )
```

Let's break this down a bit. The first command creates a new dataframe called `africa` that keeps only the observations of `prop_rights` where `africa` equals 1. Then, we create a new variable containing the actual GDP (recall `exp` is the inverse of `log`). The `$` instructs R to refer to the specified column (e.g. `logGDP`) of the dataframe (e.g. `africa`). Since there previously was no variable `levelGDP`, this gets added to the dataframe as a new column. Finally, we use the `select` function to keep only the subset of variables we specify. Notice that since we call this object `africa`, it overwrites our existing data frame. If we had wanted to keep the dataframe with the old variables, we would have to call the new dataframe something else.

3 Analysis in R

3.1 Summary Statistics

Once a dataframe is loaded in R, it's very easy to run basic summary statistics. Unlike many other statistics programs, R provides fairly minimal output, and it's possible to store the results of a command without ever displaying them. For example, we may be interested in summary statistics for the GDP per capita of countries in the sample, which we can get by typing `summary(prop_rights$gdppc)`. Note, the `results = 'markup'` bit tells RMarkdown to display the output of the code exactly as displayed in the R console.

```
## Let's look at some summary statistics.
## Here is the min, 25%-ile, median, mean, 75%-ile, and max for the variable gdppc
summary(prop_rights$gdppc) ## just print
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      450   1480   2835   5445   6968   27330
```

```
prop_rights_sumstats <- summary(prop_rights$gdppc) ## store summary, "prop_rights_variables"
prop_rights_sumstats ## take a look at this summary
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      450   1480   2835   5445   6968   27330
```

Or, we may be interested specifically in the **mean** value of GDP per capita for countries in the sample specifically, which we can get by typing `mean(prop_rights$gdppc)`. As before, this command instructs R to look in the `prop_rights` dataframe and calculate the mean for the variable `gdppc`. We can store value resulting from this calculation, by using the `<-` symbol to store the results in a new object called `meangdp`, which can be called by typing `meangdp` into the console, or used in another command or context later on.

```
## Let's look specifically at mean GDP
mean(prop_rights$gdppc) ## just print
```

```
## [1] 5445.458
```

```
meangdp <- mean(prop_rights$gdppc) ## save as a variable, called "meangdp"
meangdp ## take a look at the value stored in of "meangdp"
```

```
## [1] 5445.458
```

Similar commands exist for standard deviation (**sd**), variance (**var**), minimum (**min**), maximum (**max**), median (**med**), range (**range**) and quantile (**quantile**).

3.2 Regression

Linear regression is also very straightforward to apply in R. The `lm` command, which comes preloaded, can be used for univariate or multivariate regression:

$$Y_i = \alpha + \beta X_{1,i} + \gamma X_{2,i} + \dots + e_i$$

(Note: The above mathematical expression was written in a typesetting environment called LaTeX, and can easily be included in RMarkdown. For examples, see the Script.) Examples of how to run a regression of log GDP per capita on protection against the property rights index (protection) and latitude are below:

```
## Univariate regression of log gdp per capita on property rights index.
lm(logGDP ~ protection, data=prop_rights) # display regression results

##
## Call:
## lm(formula = logGDP ~ protection, data = prop_rights)
##
## Coefficients:
## (Intercept)    protection
##      4.6604         0.5221

## Multivariate, adding absolute latitude.
lm(logGDP ~ protection + lat_abst, data=prop_rights) # display regression results

##
## Call:
## lm(formula = logGDP ~ protection + lat_abst, data = prop_rights)
##
## Coefficients:
## (Intercept)    protection    lat_abst
##      4.7281         0.4679         1.5769

## Storing results
reg1 <- lm(logGDP ~ protection, data=prop_rights) # save results, assign to "reg1"
reg2 <- lm(logGDP ~ protection + lat_abst, data=prop_rights) # save results, assign to "reg2"
```

In the last example, we store the univariate regression results in the object `reg1` and the multivariate regression results in the object `reg2`. If you try this in R, you'll notice the results don't display on the screen when you store the results, but you can access them by naming the object.

```
## Now access regression results for the multivariate regression
summary(reg2) # access complete regression results
```

```
##
## Call:
## lm(formula = logGDP ~ protection + lat_abst, data = prop_rights)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.6845 -0.4233  0.1408  0.4584  1.1858
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  4.72808    0.39732  11.900  < 2e-16 ***
```

```
## protection    0.46789    0.06416    7.292 7.29e-10 ***
## lat_abst     1.57688    0.71031    2.220 0.0301 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.6917 on 61 degrees of freedom
## Multiple R-squared:  0.5745, Adjusted R-squared:  0.5605
## F-statistic: 41.18 on 2 and 61 DF,  p-value: 4.805e-12
```

3.3 Tables

There are many packages in R that let one create very nice tables. We recommend **stargazer**. Take the two regressions we ran, stored as **reg1** and **reg2**, and say we want to export both to a table to be used in word or LaTeX.

To include the final table in the knitted pdf created by RMarkdown, you can use the following template. (Note, we want to use **results='asis'** and **header=FALSE**, in order for RMarkdown to compile the LaTeX table into a pretty format when knitting the pdf):

```
stargazer(reg1, reg2,
  out="Table 1",type="latex",header=FALSE, table.placement = "!h",
  title="Property Rights and Development",align=TRUE,
  report = "vc*st", omit.stat=c("LL","ser","f","rsq","adj.rsq"),no.space=TRUE)
```

Table 1: Property Rights and Development

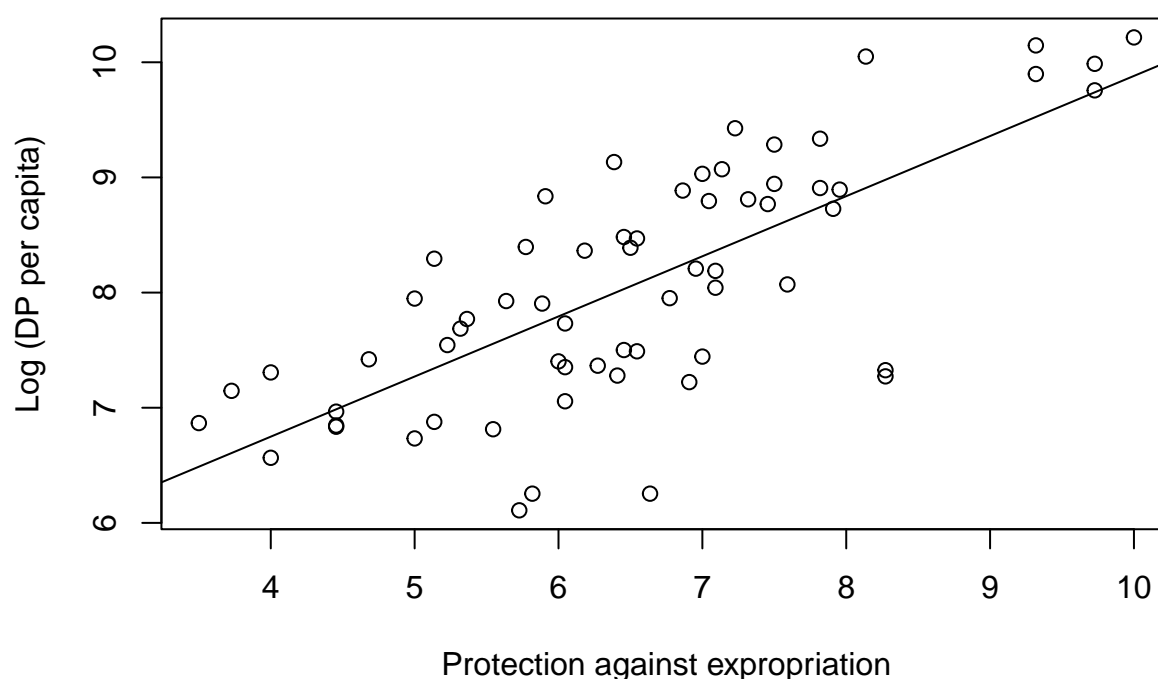
	<i>Dependent variable:</i>	
	logGDP	
	(1)	(2)
protection	0.522*** (0.061) $t = 8.533$	0.468*** (0.064) $t = 7.292$
lat_abst		1.577** (0.710) $t = 2.220$
Constant	4.660*** (0.409) $t = 11.408$	4.728*** (0.397) $t = 11.900$
Observations	64	64
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01	

The above command combines regression output saved in **reg1** and **reg2** into one table. Stargazer lets you specify title, alignment of coefficients inside columns and rows, which statistics to report (here, the variable names **v**, coefficients **c**, significance stars *****, standard errors **s** and t-statistics **t**). Moreover, you can tell it to omit statistics such as the log-likelihood **LL**, the R^2 **rsq**, etc. You can also set the table placement in the final pdf (here, **table.placement="H"** means that the table will be placed directly after the code chunk). Stargazer has extensive documentation (try **??stargazer** in R) and can also be used to make great summary statistics tables - try it out! Note that **type="latex"** produces table code that can be interpreted by LaTeX to produce nicely-formatted tables when you knit your RMarkdown file to PDF. Another alternative is to use the **type="text"** option to have a nicely formatted version of the table print to screen (for example, to later copy/paste into a word document).

3.4 Plotting

R also makes it easy to do basic plots. For example, one might want to create a scatterplot from two variables, and graph the regression line - or line of best fit - on the plot. This is easy to do with the `plot` command, which as a default plots the first variable listed on the horizontal axis, and the second on the vertical axis. Below we show a plot of the property rights index (protection against expropriation risk) and log of GDP per capita with titles and labels. We also use the `abline` command to add a regression line to the graph.

```
## A scatterplot with titles and line of best fit from univariate regression
plot(prop_rights$protection, prop_rights$logGDP,
     xlab="Protection against expropriation", ylab="Log (DP per capita)")
abline(reg1) # line of best fit using regression results as per "reg1"
```



4 Resources for further study

- UCLA has a number of excellent resources to help you learn more about how to use R, available at: <http://www.ats.ucla.edu/stat/r/>.
- You can also find interactive lessons at Try R <http://tryr.codeschool.com/> for practice writing code.
- We also recommend the text on the syllabus, Hanck, Christoph, Martin Arnold, Alexander Gerber and Martin Schmelzer. (2018). Introduction to Econometrics with R, <https://www.econometrics-with-r.org/>.
- You may also find <https://www.r-bloggers.com/> to be a useful resource for specific questions.
- Finally, we include two useful resources that summarize useful R commands and common RMarkdown features.

Base R

Cheat Sheet

Getting Help

Accessing the help files

?mean

Get help of a particular function.

help.search('weighted mean')

Search the help files for a word or phrase.

help(package = 'dplyr')

Find help for a package.

More about an object

str(iris)

Get a summary of an object's structure.

class(iris)

Find the class an object belongs to.

Using Libraries

install.packages('dplyr')

Download and install a package from CRAN.

library(dplyr)

Load the package into the session, making all its functions available to use.

dplyr::select

Use a particular function from a package.

data(iris)

Load a built-in dataset into the environment.

Working Directory

getwd()

Find the current working directory (where inputs are found and outputs are sent).

setwd('C://file/path')

Change the current working directory.

Use projects in RStudio to set the working directory to the folder you are working in.

Vectors

Creating Vectors

c(2, 4, 6)	2 4 6	Join elements into a vector
2:6	2 3 4 5 6	An integer sequence
seq(2, 3, by=0.5)	2.0 2.5 3.0	A complex sequence
rep(1:2, times=3)	1 2 1 2 1 2	Repeat a vector
rep(1:2, each=3)	1 1 1 2 2 2	Repeat elements of a vector

Vector Functions

sort(x)

Return x sorted.

rev(x)

Return x reversed.

unique(x)

See counts of values.

See unique values.

Selecting Vector Elements

By Position

x[4]

The fourth element.

x[-4]

All but the fourth.

x[2:4]

Elements two to four.

x[-(2:4)]

All elements except two to four.

x[c(1, 5)]

Elements one and five.

By Value

x[x == 10]

Elements which are equal to 10.

x[x < 0]

All elements less than zero.

x[x %in% c(1, 2, 5)]

Elements in the set 1, 2, 5.

Named Vectors

x['apple']

Element with name 'apple'.

Programming

For Loop

```
for (variable in sequence){  
  Do something  
}
```

Example

```
for (i in 1:4){  
  j <- i + 10  
  print(j)  
}
```

While Loop

```
while (condition){  
  Do something  
}
```

Example

```
while (i < 5){  
  print(i)  
  i <- i + 1  
}
```

If Statements

```
if (condition){  
  Do something  
} else {  
  Do something different  
}
```

Example

```
if (i > 3){  
  print('Yes')  
} else {  
  print('No')  
}
```

Functions

```
function_name <- function(var){  
  Do something  
  return(new_variable)  
}
```

Example

```
square <- function(x){  
  squared <- x*x  
  return(squared)  
}
```

Reading and Writing Data

Input	Output	Description
df <- read.table('file.txt')	write.table(df, 'file.txt')	Read and write a delimited text file.
df <- read.csv('file.csv')	write.csv(df, 'file.csv')	Read and write a comma separated value file. This is a special case of read.table/write.table.
load('file.Rdata')	save(df, file = 'file.Rdata')	Read and write an R data file, a file type special for R.

Conditions

a == b	Are equal	a > b	Greater than	a >= b	Greater than or equal to	is.na(a)	Is missing
a != b	Not equal	a < b	Less than	a <= b	Less than or equal to	is.null(a)	Is null

Types

Converting between common data types in R. Can always go from a higher value in the table to a lower value.

as.logical	TRUE, FALSE, TRUE	Boolean values (TRUE or FALSE).
as.numeric	1, 0, 1	Integers or floating point numbers.
as.character	'1', '0', '1'	Character strings. Generally preferred to factors.
as.factor	'1', '0', '1', levels: '1', '0'	Character strings with preset levels. Needed for some statistical models.

Maths Functions

log(x)	Natural log.	sum(x)	Sum.
exp(x)	Exponential.	mean(x)	Mean.
max(x)	Largest element.	median(x)	Median.
min(x)	Smallest element.	quantile(x)	Percentage quantiles.
round(x, n)	Round to n decimal places.	rank(x)	Rank of elements.
signif(x, n)	Round to n significant figures.	var(x)	The variance.
cor(x, y)	Correlation.	sd(x)	The standard deviation.

Variable Assignment

```
> a <- 'apple'
> a
[1] 'apple'
```

The Environment

ls()	List all variables in the environment.
rm(x)	Remove x from the environment.
rm(list = ls())	Remove all variables from the environment.

You can use the **environment** panel in **RStudio** to **browse variables in your environment**.

RStudio® is a trademark of RStudio, Inc. • [CC BY](#) Mhairi McNeill • mhairi.mcnell@gmail.com • 844-448-1212 • [rstudio.com](#)

Matrixes

```
m <- matrix(x, nrow = 3, ncol = 3)
# Create a matrix from x.
```

	m[2,] - Select a row	t(m) Transpose 
	m[, 1] - Select a column	m %*% n Matrix Multiplication 
	m[2, 3] - Select an element	solve(m, n) Find x in: m * x = n

Lists

```
l <- list(x = 1:5, y = c('a', 'b'))
# A list is collection of elements which can be of different types.
```

l[[2]]	l[[1]]	l\$x	l['y']
Second element of l.	New list with only the first element.	Element named x.	New list with only element named y.

Data Frames

Also see the **dplyr** library.

```
df <- data.frame(x = 1:3, y = c('a', 'b', 'c'))
# A special case of a list where all elements are the same length.
```

x	y
1	a
2	b
3	c

df\$x	df[[2]]
1 2 3	a b c

List subsetting

View(df)	head(df)
See the full data frame.	See the first 6 rows.

Matrix subsetting

df[, 2]	df[2,]	df[2, 2]
Number of rows.	Number of columns.	Number of columns and rows.

cbind	rbind
- Bind columns.	- Bind rows.

Strings

Also see the **stringr** library.

```
paste(x, y, sep = ' ')
# Join multiple vectors together.

paste(x, collapse = ' ')
# Join elements of a vector together.

grep(pattern, x)
# Find regular expression matches in x.

gsub(pattern, replace, x)
# Replace matches in x with a string.

toupper(x)
# Convert to uppercase.

tolower(x)
# Convert to lowercase.

nchar(x)
# Number of characters in a string.
```

Factors

```
factor(x)
# Turn a vector into a factor. Can set the levels of the factor and the order.

cut(x, breaks = 4)
# Turn a numeric vector into a factor but 'cutting' into sections.
```

Statistics

```
lm(x ~ y, data=df)
# Linear model.

glm(x ~ y, data=df)
# Generalised linear model.

summary
# Get more detailed information out a model.

prop.test
# Test for a difference between proportions.

pairwise.t.test
# Perform a t-test for paired data.

aov
# Analysis of variance.
```

Distributions

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

Plotting

Also see the **ggplot2** library.

	plot(x) Values of x in order.		plot(x, y) Values of x against y.		hist(x) Histogram of x.
---	---	---	---	---	-----------------------------------

Dates

See the **lubridate** library.

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

learn more at rmarkdown.rstudio.com



Rmd files

Reproducible Research

A diagram illustrating the output options for an R Markdown (Rmd) file. At the bottom is a document icon labeled 'Rmd'. Four arrows point upwards from this icon to four different output formats: a PDF icon, an HTML icon, a Word document icon, and a LaTeX icon.

Dynamic Documents

Interactive Documents

- Turn your report into an interactive Shiny document in 4 steps

 - 1 Add runtime: shiny to the YAML header.
 - 2 Call shiny input functions to embed input objects.
 - 3 Call shiny render functions to embed reactive output.
 - 4 Render with rmarkdown::run or click Run Document in RStudio IDE

```
## ---- output: html_document  
runtime: shiny  
  
... {r, echo = FALSE}  
numericInput("n",  
  "How many cars?", 5)  
  
renderTable({  
  head(cars, input$n)  
})  
...}
```

	speed	dist
1	4.00	2.00
2	4.00	10.00
3	7.00	4.00
4	7.00	22.00
5	8.00	16.00

How many cars?
5

Enbed a complete app into your document with
shiny::shinyAppDir()

* Your report will be rendered as a Shiny app, which means you must choose an html output format, like **HTML Document**, and save it with an active R session.

Parameters

1

Parameterize your documents to reuse with different inputs (e.g., data sets, values, etc.)

Add parameters

Create and set parameters in the header as sub-values of `<params>`

2

Call parameters

Call parameter values in code as `<params>name</params>`

3

Set parameters

Set values with `Knit` or with parameters or the `params` argument of `rrender()`:

Learn more at [markdown.rstudio.com](https://www.markdown.rstudio.com) • RStudio IDE 0.99.879 • Updated: 02/16

[illegible]