

Introduction to



Master course
October 5th, 2016

Sarah Bonnin
CRG Bioinformatics Core Facility

Today - outline

- R Basics
- R Studio
- Input and Output
- Library / packages
- Graphing in R: ggplot2
- R Markdown
- Writing functions in R

R Basics

What is R?

- Used for **data manipulation, calculation and graphical display.**

- **Open source !**

<https://www.r-project.org/>

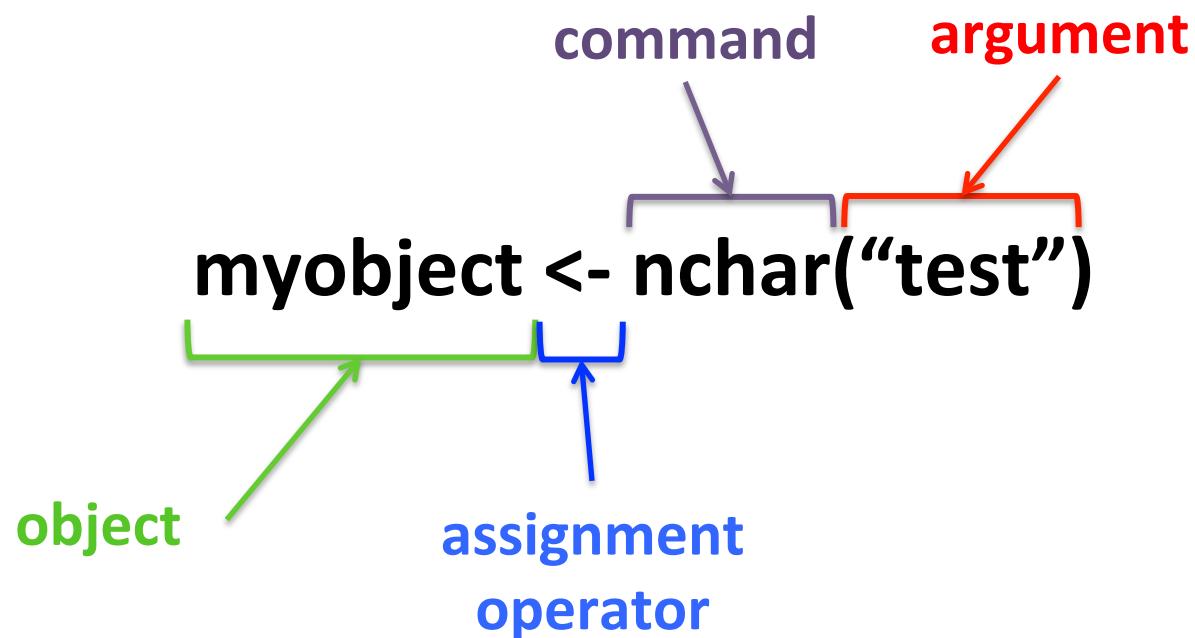
- **Interactive, flexible**

- **Very active community of developers and users!**

R base syntax

a <- function(arguments)

example:



Trying a few commands in the terminal

```
getwd()
```

```
[1] "/Users/sbonnin"
```

```
setwd("Desktop/")
```

```
my_directory <- "R_introduction"
```

```
dir.create(my_directory)
```

```
setwd("R_introduction/")
```

```
getwd()
```

```
[1] "/Users/sbonnin/Desktop/R_introduction"
```

R syntax

- Case sensitive
- Comment lines start with #
- Commands separated by ; or a new line

```
> my_directory  
[1] "R_introduction"  
>  
>  
> my_Directory  
Error: object 'my_Directory' not found  
>  
>  
> # my_Directory  
>  
>  
> getwd(); 2+4  
[1] "/Users/sbonnin/R_introduction"  
[1] 6
```

Getting help on a function

- `help(getwd)`
- `?getwd`

R Studio

R studio?

- **Free and open source** Integrated Development Environment for R
- Available for Windows, Mac OS and Linux



Screen: 4 windows

The screenshot shows the R Studio interface with four open windows:

- 1. Console**: Located at the bottom left, it displays R command-line interactions. It shows commands like `getwd()`, `library("ggplot2")`, and `list.files()` along with their results.
- 2. R script**: Located in the top-left quadrant, it shows an R script named "Untitled1.R" with code for setting the working directory, loading packages, creating a matrix, and creating a list.
- 3. Environment and history**: Located in the top-right quadrant, it shows the Global Environment pane with objects "a" and "b" defined, and the History pane below it.
- 4. Files, plots, packages, help**: Located in the bottom-right quadrant, it shows the Files pane displaying the contents of the current directory, including subfolders like Applications, Desktop, Documents, Library, Movies, Music, Pictures, Projects, and Public.

Data Structures

Vectors

Matrices

Data frames

List

Data Types

- Every object has a **data type** that tells what sort of value it is:
 - Numeric (Numbers)
 - Character (Text)
 - Logical (True / False)

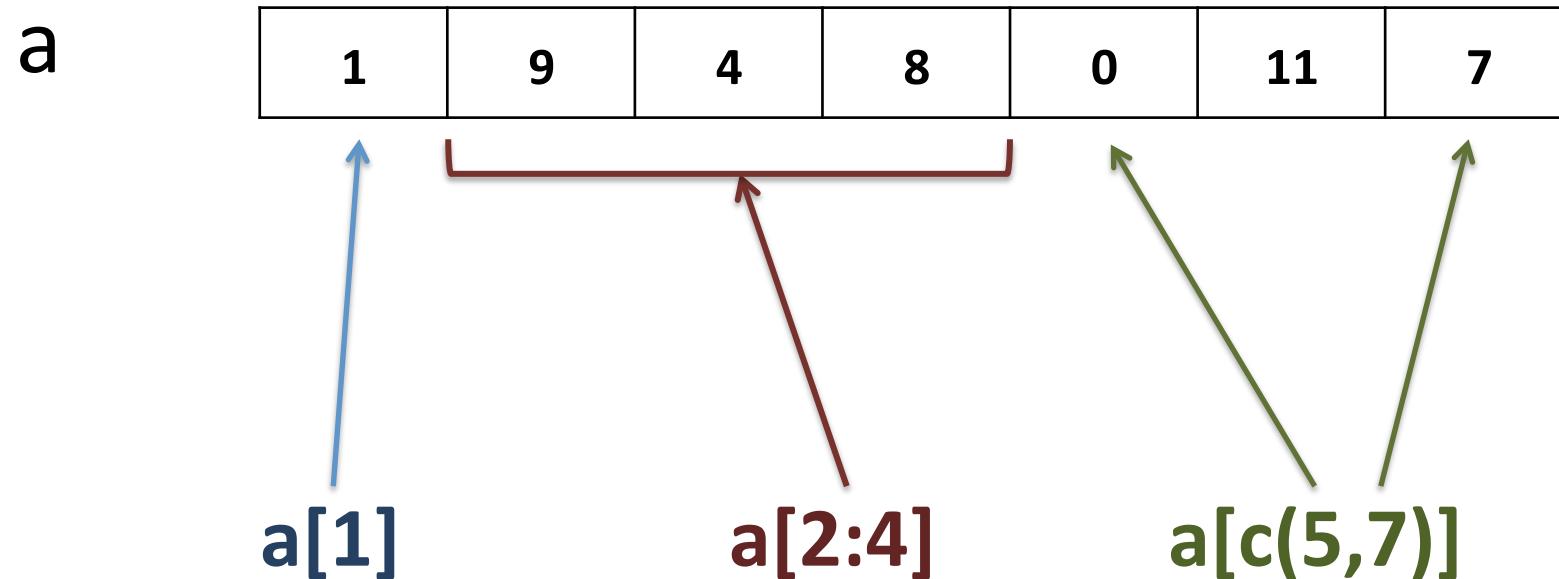
Vectors

- Sequence of data elements of the same type
- Assignment of values to vector using the **c** command (combining elements)

```
a <- c(1, 9, 4, 8, 0, 11, 7)
```

Vectors

- Fetch elements of a vector:



Exercise 1:

Vector manipulation

Exercise 1

Vector manipulation

- Create vector **x** of 1000 random numbers.
 - Use **rnorm** function (random generation of numbers from a normal distribution)
- Compute the mean, median, minimum and maximum of this vector.
 - Functions: **mean**, **median**, **min**, **max**.
- Create **y**, vector containing the first 100 elements of **x**.
 - **x[?]**
- Compute summary statistics of **y** using the **summary** function.

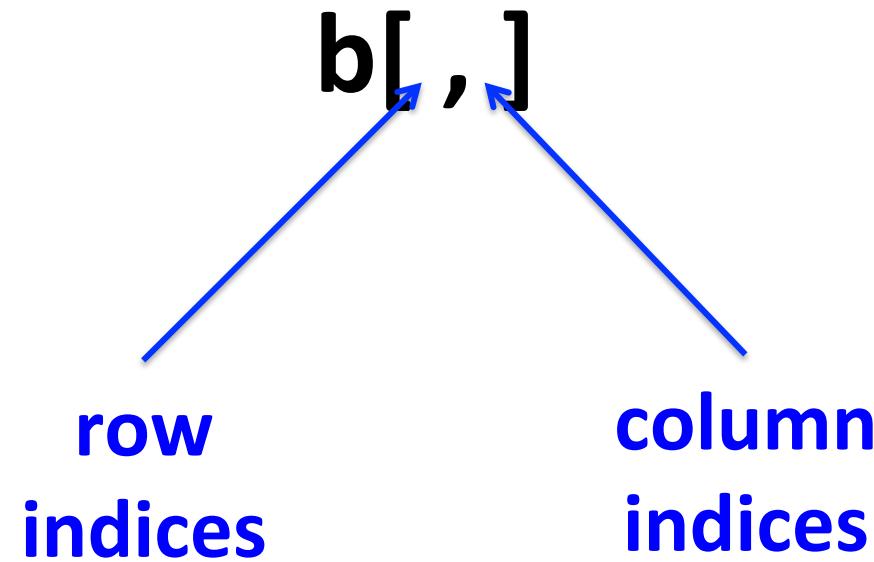
Matrices

- A matrix is a vector of **2 dimensions**
- All columns in a matrix must have :
 - the same **type** (numeric, character, logical)
 - the same **length**

```
b <- matrix(c(1, 0, 34, 44, 12, 4),  
            nrow=3,  
            ncol=2)
```

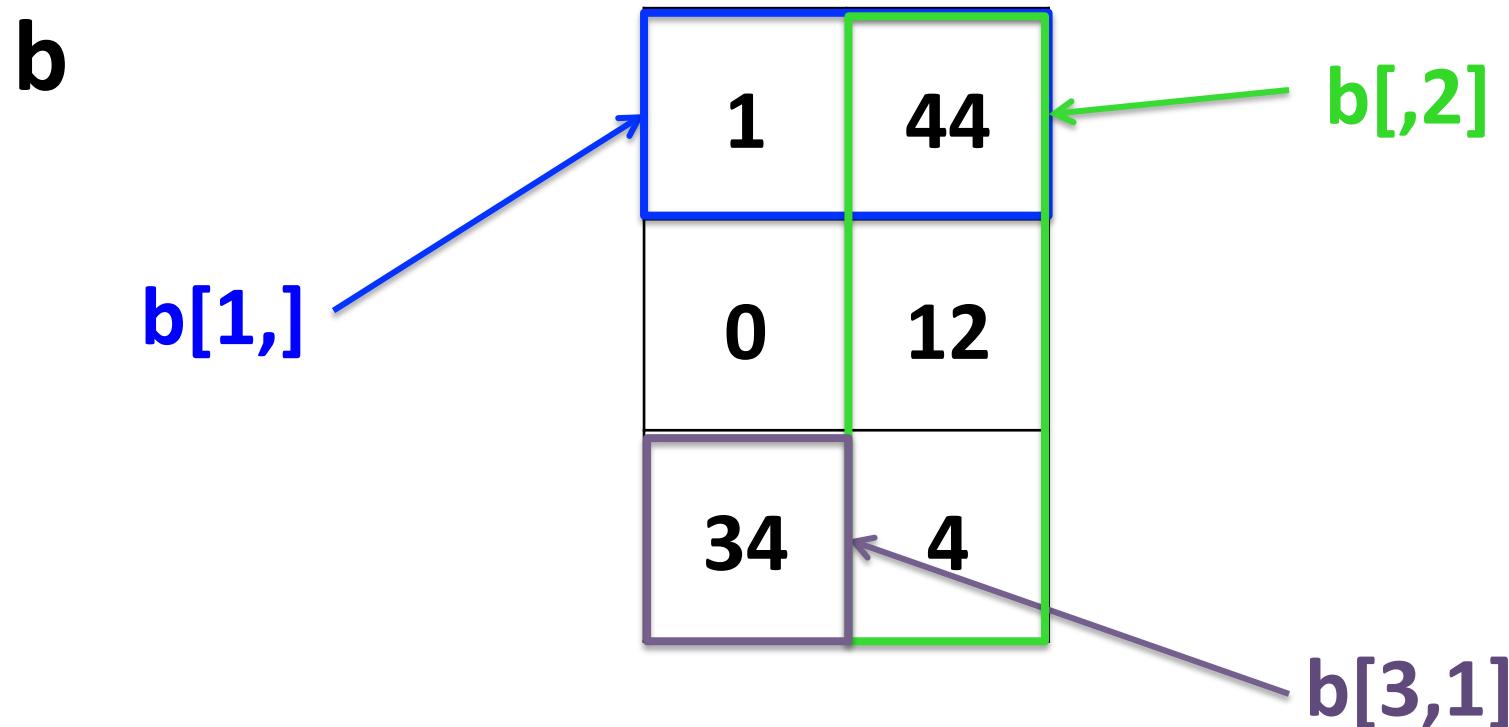
Matrices

- Fetch rows, columns or single elements of a matrix:



Matrices

- Fetch rows, columns or single elements of a matrix:



Data frames

- Structures of 2 dimensions.
- More general than matrices.
- Different columns can have **different types** but must have the **same length**.

```
d <- data.frame(c(1, 3, 34),  
                  c("no", "ok", "yes"),  
                  c(TRUE, TRUE, FALSE))
```

Data frames

- Fetch rows, columns or single elements of a data frame:

d

d[1,]	→	1	no	TRUE
		3	ok	TRUE
d[5,1]	→	34	yes	FALSE

Dimension names

- Column and / or row names can be added to matrices and data frames:

```
colnames(d) <- c("numb", "charac", "logic")
```

```
rownames(d) <- c("row1", "row2", "row3")
```

Lists

- Linear structures.
- A component of a list can be **any data structure**:
→ matrix, vector, data frame, list...

Lists

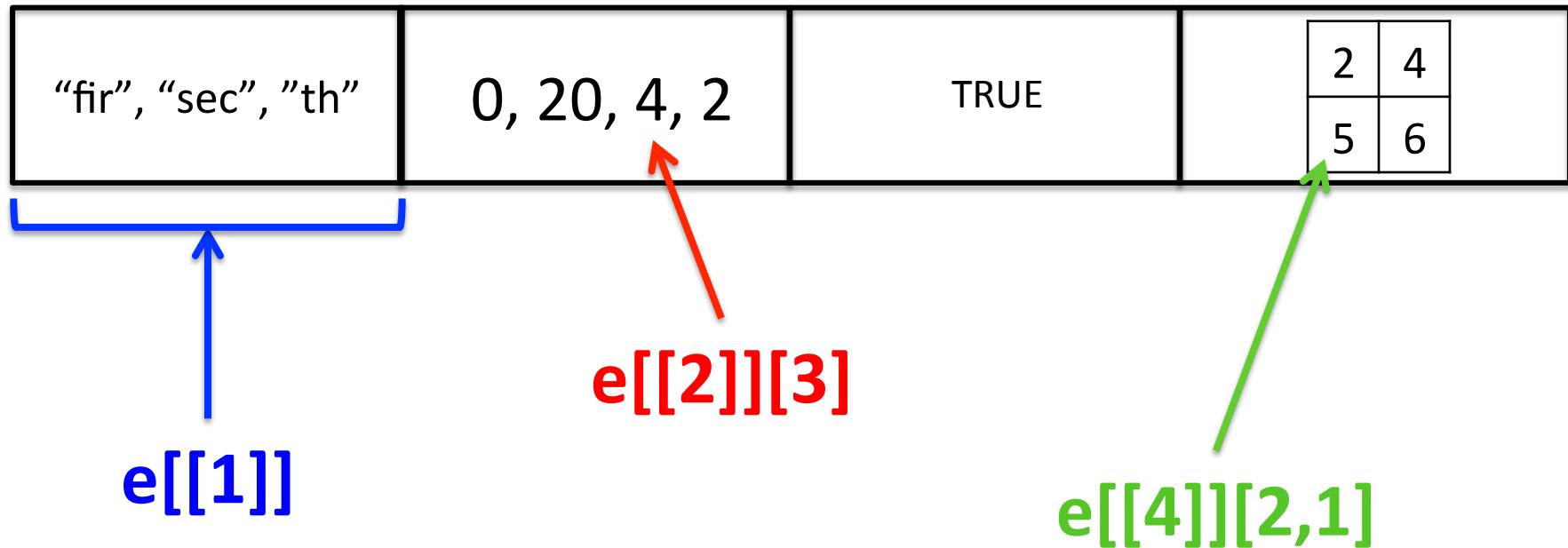
- Create a list:

```
e <- list() # Empty list that can be filled later on
```

```
e <- list(c("fir", "sec", "th"),      # 1st element of the list  
          c(0, 20, 4, 2),      # 2d element of the list  
          TRUE,      # 3d element of the list  
          matrix(c(2, 5, 4, 6),      # 4th element of the list  
                  nrow=2, ncol=2))
```

Lists

- Accessing elements of a list **e**:



Input and Output

File manipulation

Writing and reading a text file

- Let's start from a simple matrix

```
b <- matrix(c(1:6), nrow=3)
```

- Write this matrix to a file

```
write.table(b, "matrixb.txt")
```

... and read it back!

```
b <- read.table("matrixb.txt")
```

Producing graphs

- Graphs can be saved into many formats, including: **pdf**, **jpeg**, **bmp**, **tiff**.
- Open the file:
pdf("my_graph.pdf")
- Produce a graph:
plot(1:10)
- And close it:
dev.off()

Exercise 2:

data structures + I/O

Exercise 2

Data structures + I/O

- Create a data frame from this shopping list  (save it in object **shop**):

Yogurt	6
Sugar	1
Apple	4

- Add column names to the data frame: respectively **item** and **quantity**.
- Write this data frame to file **shopping_list.txt**:
 - without quotes
 - without row names
 - with column names
 - fields separated with tab

Libraries / packages

Library/packages

- **Packages** are collections of R functions, data, and compiled code in a well-defined format.
- The directory where packages are stored is called the **library**.

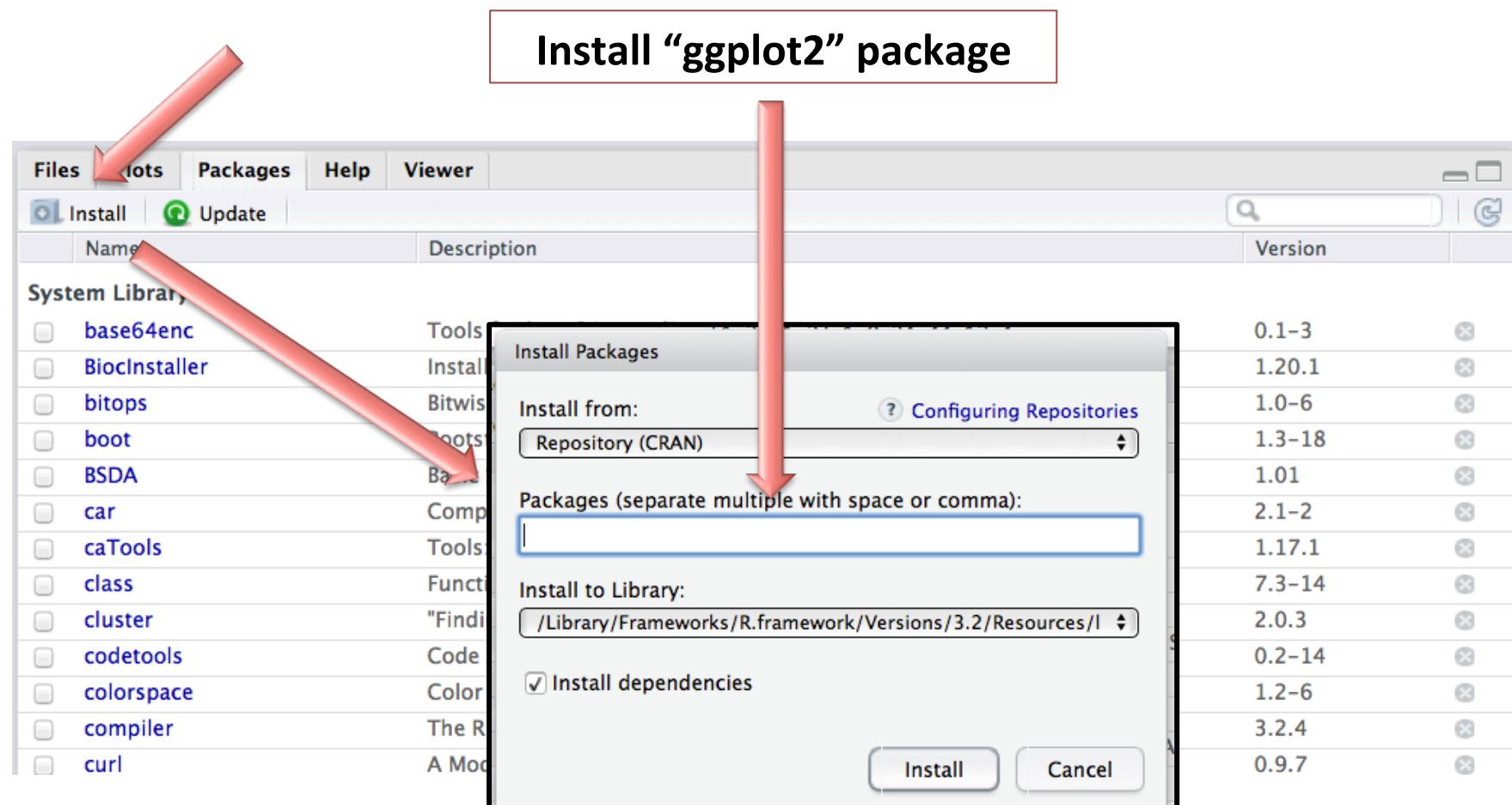
Definitions from <http://www.statmethods.net/interface/packages.html>

Standard packages

- About 25 standard packages are supplied with R by default (example: base, stats, graphics).
- On October 3d 2016, **9278 packages** were available!
- Anyone can contribute to the R package repository

Library / packages

Install a package with R studio



Listing functions from packages

- **ls("package:ggplot2")**

```
[1] "%+%"           "%+replace%"  
[3] "Coord"         "CoordCartesian"  
[5] "CoordFixed"    "CoordFlip"  
[7] "CoordMap"      "CoordPolar"  
[9] "CoordQuickmap" "CoordTrans"  
[11] "Geom"          "GeomAbline"  
[13] "GeomAnnotationMap" "GeomArea"  
[15] "GeomBar"       "GeomBlank"
```

...

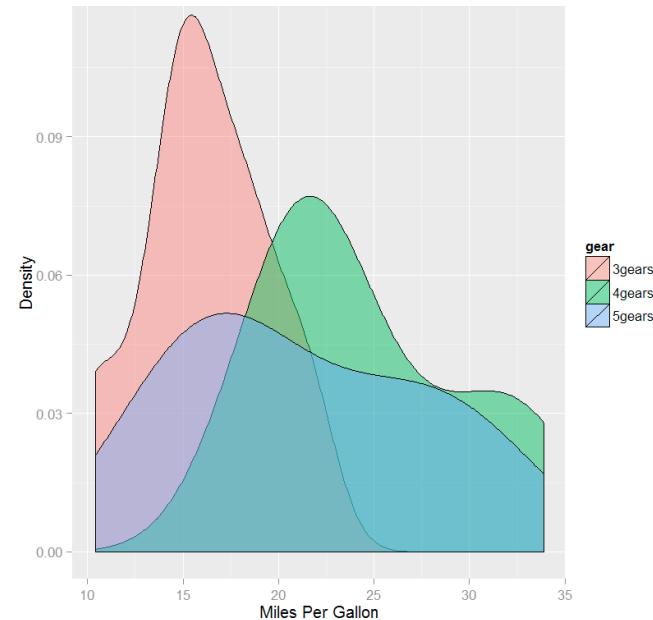
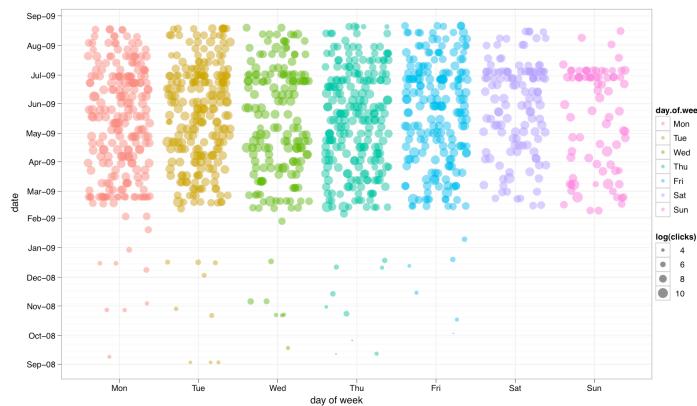
Graphing in R: short introduction to ggplot2

ggplot2

- Graphing package inspired by the Grammar of Graphics seminal work of Leland Wilkinson.
- A tool that enables to concisely describe the components of a graphic.

ggplot2

- Flexible
- Customizable
- Pretty
- Well documented



ggplot2

- Base layer to plot two variables

ggplot(dataframe, aes(x, y))

- Add layer specifying what kind of plot you want.

Example of a scatter plot:

ggplot(dataframe, aes(x, y)) + geom_point()

Exercise 3:

ggplot2

Exercise 3

ggplot2

- Create a data frame containing columns **x** and **y**, each a random generation of 100 normal distribution values.
 - Remember **rnorm** function.
- Produce a scatter plot using variable **x** and **y** from the data frame.
- Save plot in a **jpeg** file.

Dataset

- **diamonds** data set in ggplot2 package
- **dim(diamonds); colnames(diamonds); head(diamonds)**
- subset diamonds:
`diams2 <- diamonds[sample(rownames(diamonds), 1000),]`

Scatter plots

A scatter plot is a type of plot or mathematical diagram using Cartesian coordinates to display values for typically two variables for a set of data.

[Wikipedia]

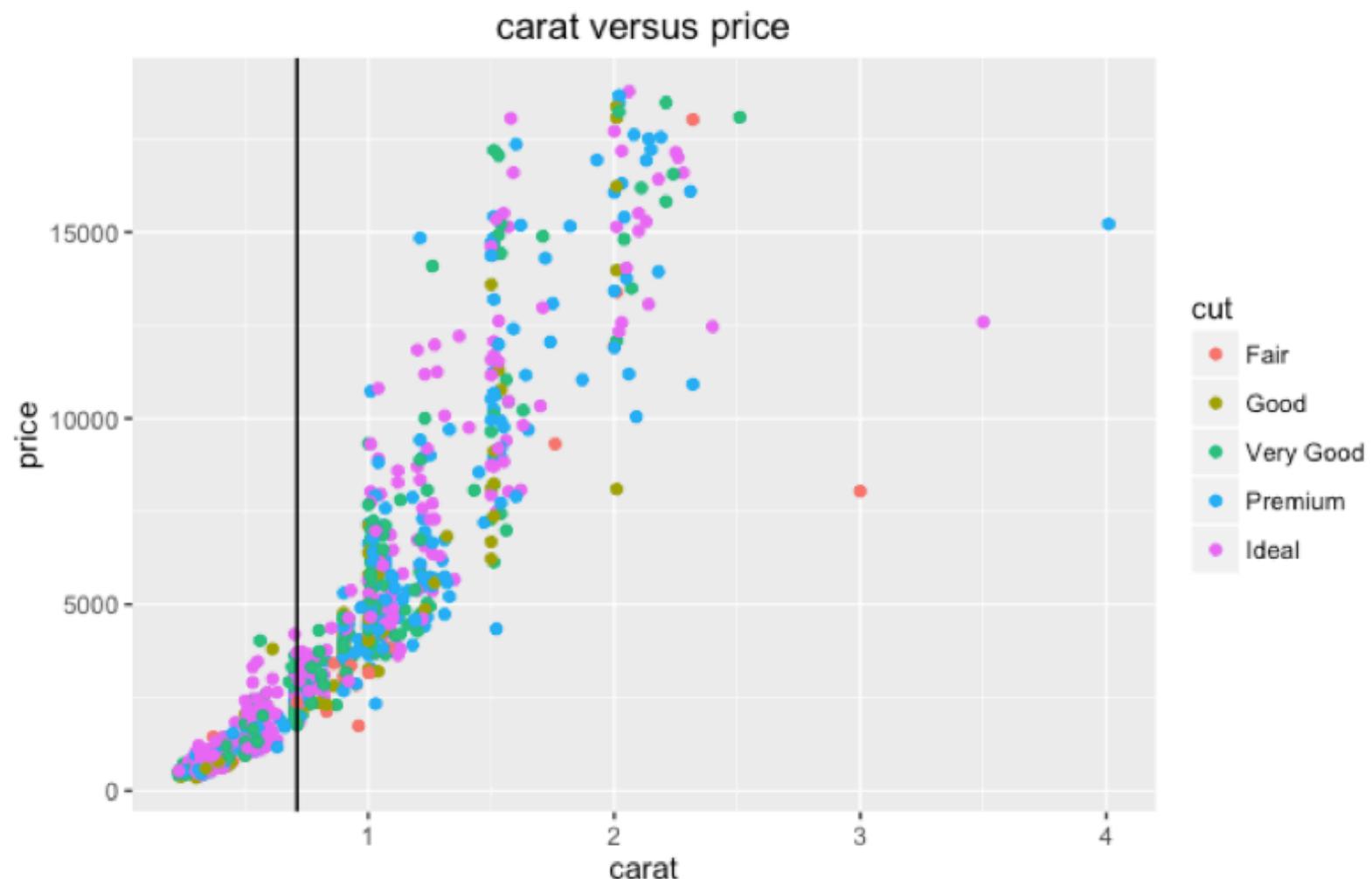
Scatter plots

- ```
p <- ggplot(diams2, aes(carat, price, color=cut)) +
 geom_point()
```
- add title with **ggtitle**  

```
p2 <- p + ggtitle("carat versus price")
```
- draw a vertical line at the median carat value with  
**geom\_vline**  

```
p3 <- p2 + geom_vline(xintercept=median(diams2$carat))
```

# Scatter plots



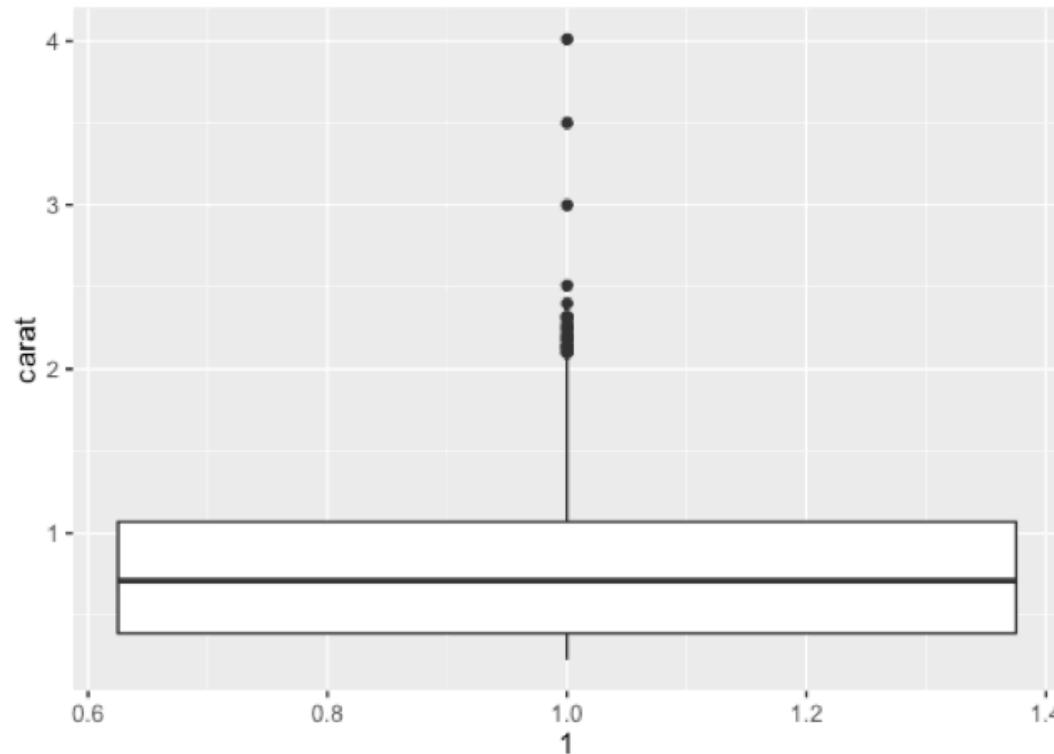
# Boxplots

*In descriptive statistics, a boxplot is a convenient way of graphically depicting groups of numerical data through their quartiles.*

[Wikipedia]

# Boxplots

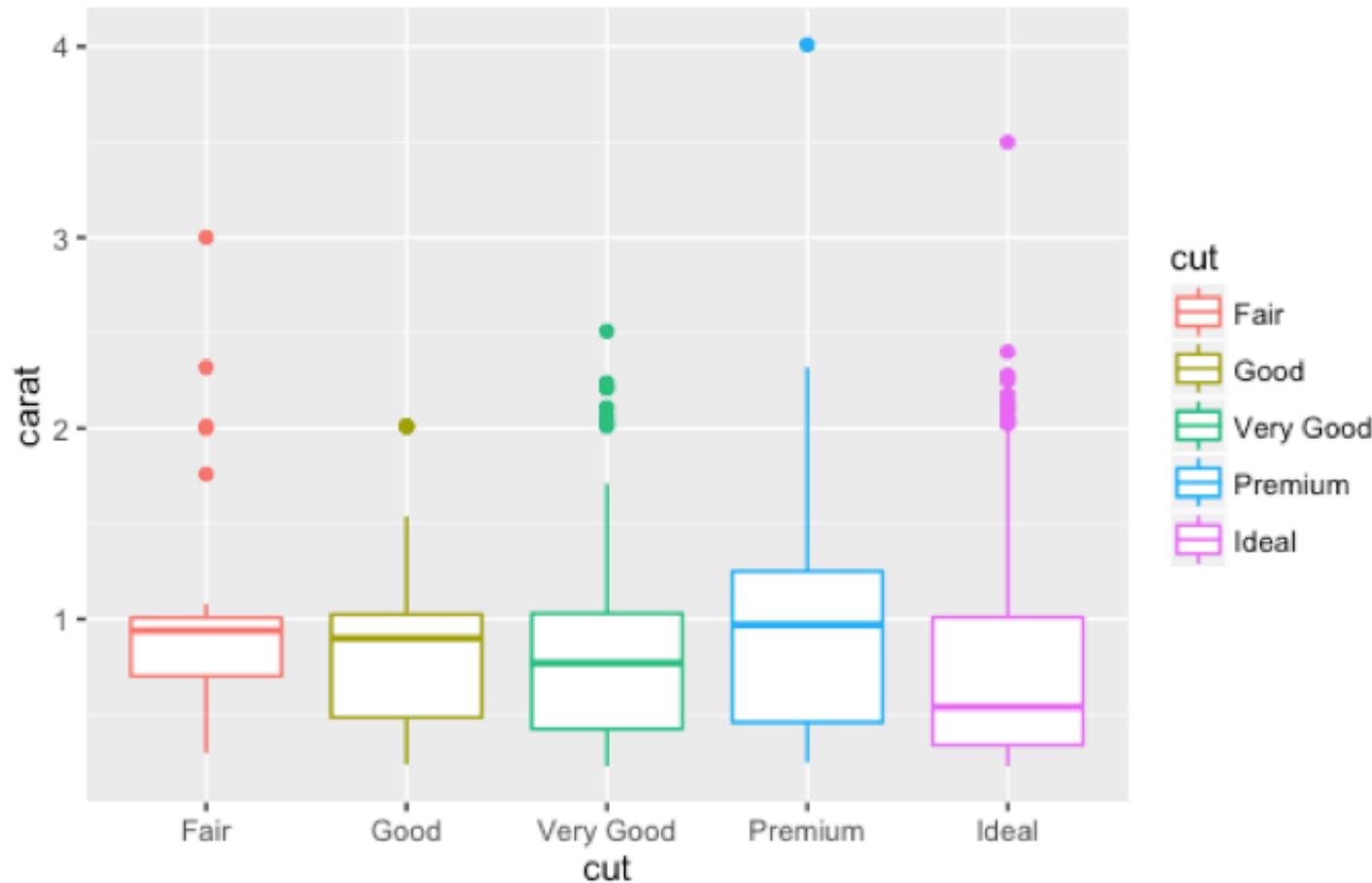
- **ggplot(diams2, aes(x=1, y=carat)) + geom\_boxplot()**



# Boxplots

- ```
p <- ggplot(diams2, aes(x=cut, y=carat)) +  
  geom_boxplot()
```
- add colors to the boxes lines with **color** in aes
 - ```
p <- ggplot(diams2, aes(x=cut, y=carat, color=cut)) +
 geom_boxplot()
```

# Boxplots



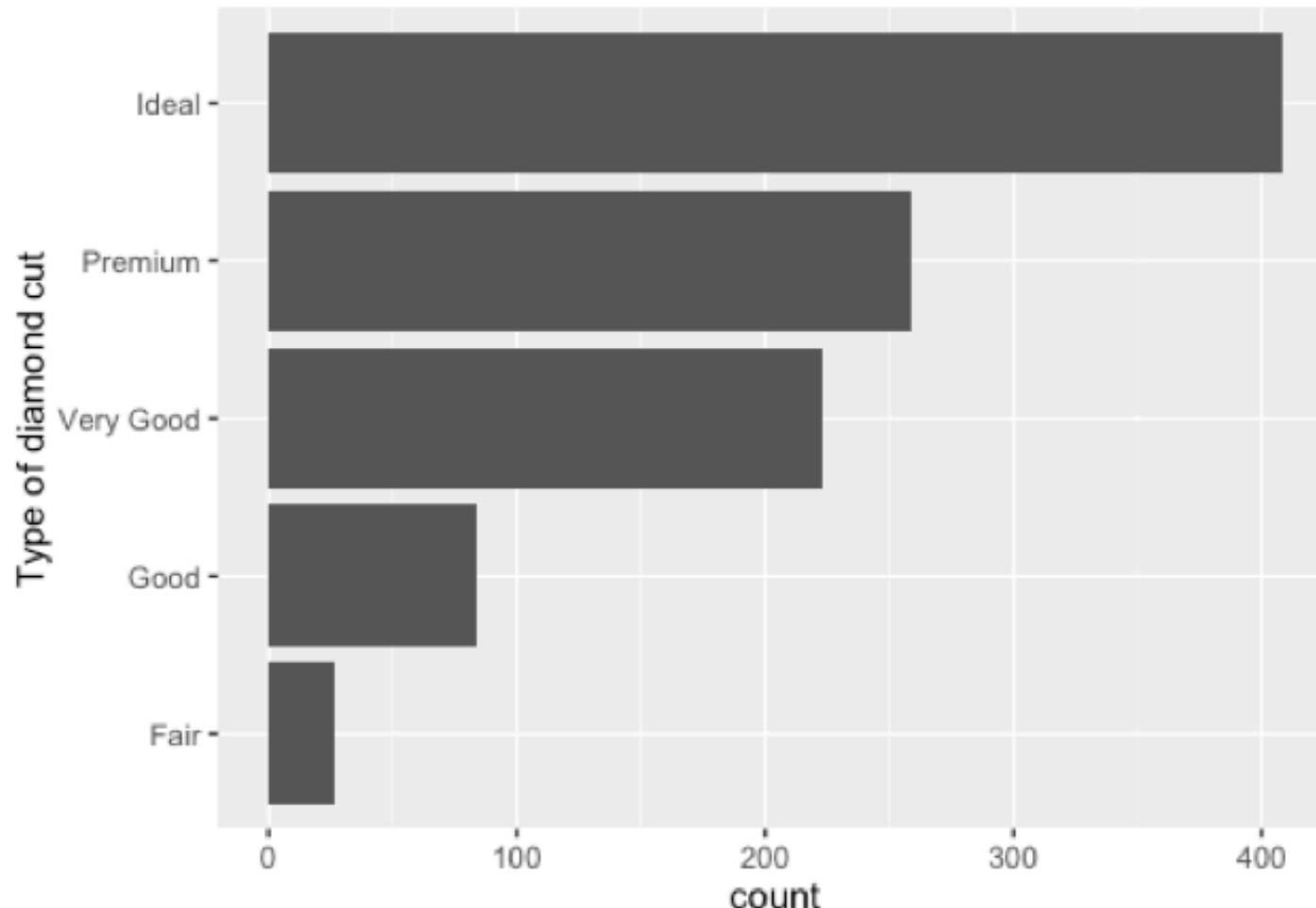
# Barplots

*A bar chart or bar graph is a chart that presents grouped data with rectangular bars with lengths proportional to the values that they represent.*  
[Wikipedia]

# Barplots

- `p <- ggplot(diams2, aes(x=cut)) + geom_bar()`
- Change x axis label with `scale_x_discrete`:  
`p2 <- p + scale_x_discrete(name="Type of diamond cut")`
- Swapping x and y axis with `coord_flip()`:  
`p3 <- p2 + coord_flip()`

# Barplots



# Histograms

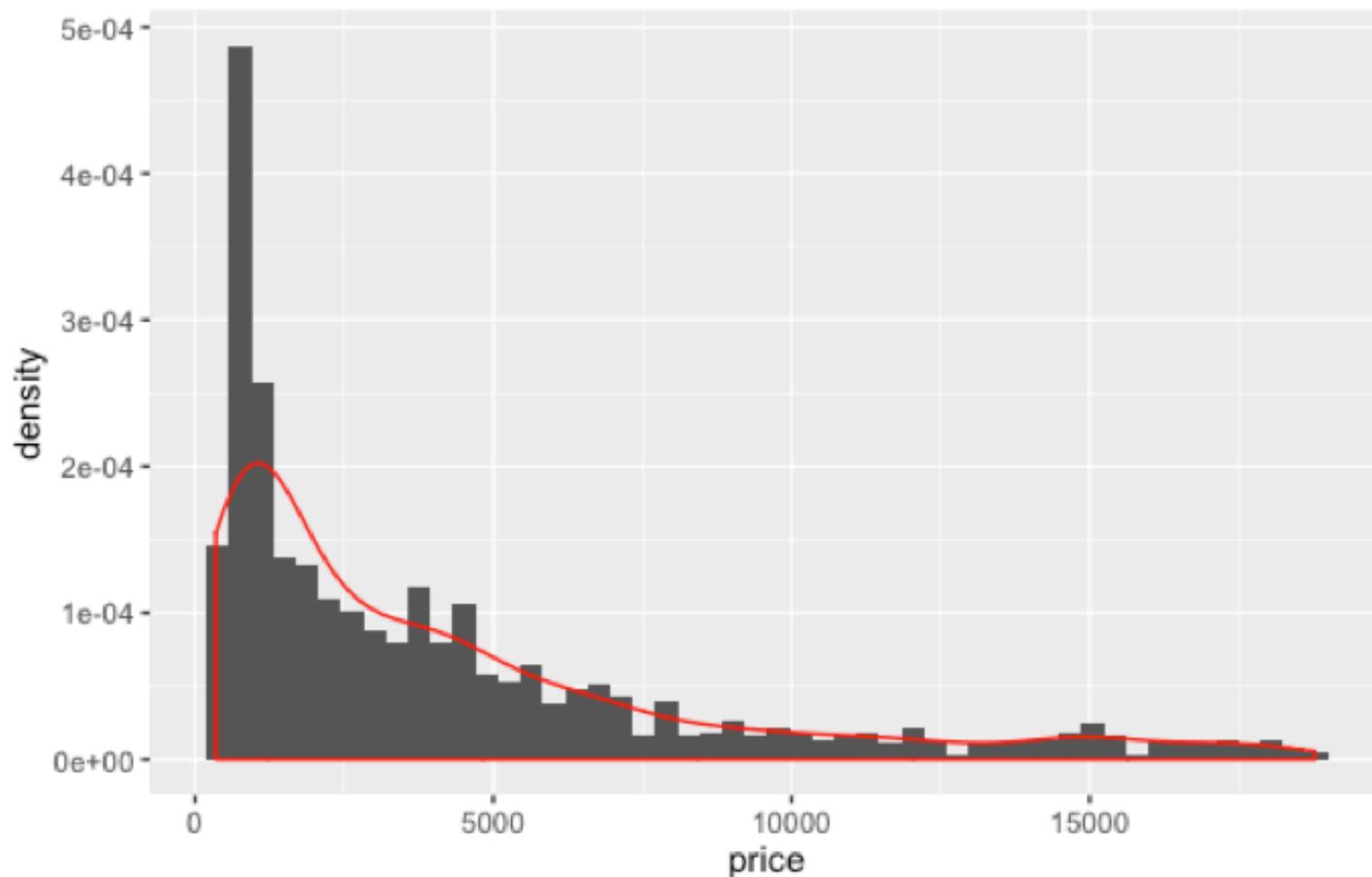
*A histogram is a plot that lets you discover, and show, the underlying frequency distribution of a set of continuous data.*

[Wikipedia]

# Histograms

- `ggplot(diams2, aes(x=price)) + geom_histogram()`
- Increase resolution with **bins** option:  
`ggplot(diams2, aes(x=price)) + geom_histogram(bins=50)`
- Add a density curve:  
`ggplot(diams2, aes(x=price))`
  - + `geom_histogram(aes(y=..density..), bins=50)`
  - + `geom_density(color="red")`

# Histograms



# RMD Markdown

# RMD: R MarkDown

- Format for writing **reproducible, dynamic reports** with R.
- Work directly inserted into formatted documents (HTML, PDF and Word)
- Easy to use!

# RMD: R MarkDown

The screenshot shows the RStudio interface with an R Markdown file open. The file contains the following content:

```
1 ---
2 title: "Test_rmd"
3 author: "Sarah Bonnin"
4 date: "14 Apr 2016"
5 output: word_document
6 ---
7
8 ```{r setup, include=FALSE}
9 knitr::opts_chunk$set(echo = TRUE)
10 ```
11
12 ## R Markdown
13
14 This is an R Markdown document. Markdown is a simple way to write
details on using R Markdown see <http://rmarkdown.rstudio.com>
15
16 When you click the **Knit** button a document will be generated that includes
code chunks within the document. You can embed an R chunk like
17
18 ```{r, echo=TRUE}
19 a <- c(1:10, nrow=2, ncol=5)
20 a
21 ```
22
23 ## Including Plots
24
25 You can also embed plots, for example:
26
27 ```{r, echo=FALSE}
28 plot(pressure)
29 plot(1:10, col=1:10, pch=1:10)
30 ```
31
32 Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.
```

A callout box highlights the code chunk starting at line 8, specifically the line `knitr::opts_chunk$set(echo = TRUE)`. The RStudio interface includes tabs for cmd.R\*, R packages available, and Test\_rmd.Rmd, along with standard toolbar icons.

# RMD: R MarkDown

The screenshot shows the RStudio interface with an R Markdown file named "Test\_rmd.Rmd". The code editor on the left contains the following R Markdown code:

```
1 ---
2 title: "Test_rmd"
3 author: "Sarah Bonnin"
4 date: "14 Apr 2016"
5 output: word_document
6 ---
7
8 ```{r setup, include=FALSE}
9 knitr::opts_chunk$set(echo = TRUE)
10 ```
11 ## R Markdown
12
13 This is an R Markdown document. Markdown is a simple way to write natural looking documentation.
14 details on using R Markdown see <http://rmarkdown.rstudio.com>
15
16 When you click the **Knit** button a document will be generated that includes both the
17 code chunks within the document. You can embed an R chunk as follows:
18 ```{r, echo=TRUE}
19 a <- c(1:10, nrow=2, ncol=5)
20 a
21```
22
23 ## Including Plots
24
25 You can also embed plots, for example:
26
27 ```{r, echo=FALSE}
28 plot(pressure)
29 plot(1:10, col=1:10, pch=1:10)
30```
31
32 Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.
```

A callout box highlights the R code chunk starting at line 8. The box contains the following text:

## R Markdown  
This is an R Markdown document. Markdown is a simple way to write natural looking documentation.  
details on using R Markdown see <<http://rmarkdown.rstudio.com>>

When you click the \*\*Knit\*\* button a document will be generated that includes both the code chunks within the document. You can embed an R chunk as follows:

```
```{r, echo=TRUE}  
a <- c(1:10, nrow=2, ncol=5)  
a  
```
```

# RMD: R MarkDown

The screenshot shows an RStudio interface with the following code in the 'Test\_rmd.Rmd' file:

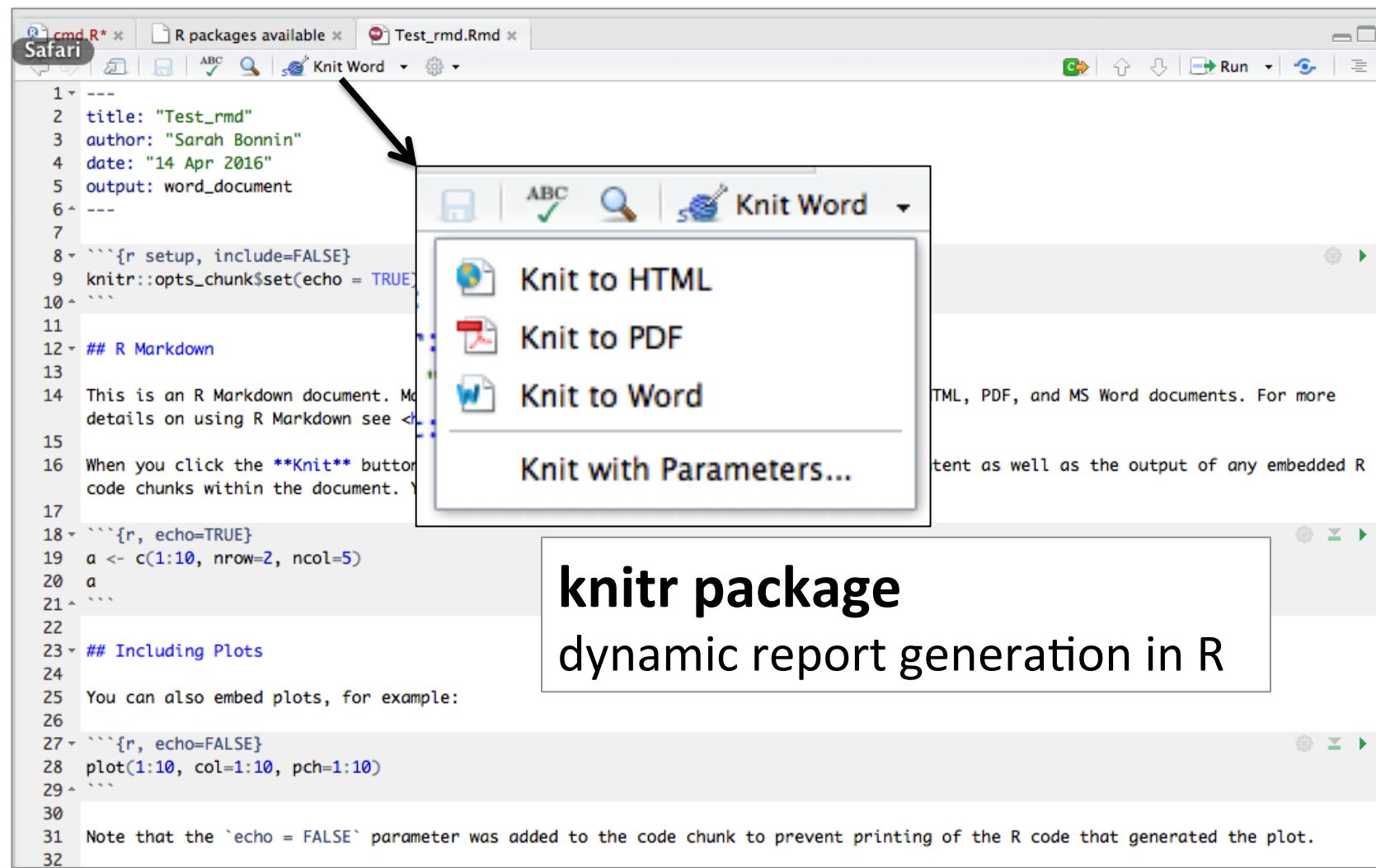
```
1 ---
2 title: "Test_rmd"
3 author: "Sarah Bonnin"
4 date: "14 Apr 2016"
5 output: word_document
6 ---
7
8 ```{r setup, include=FALSE}
9 knitr::opts_chunk$set(echo = TRUE)
10 ```
11
12 ## R Markdown
13
14 This is an R Markdown document. Markdown is a simple formalism for representing structured documents using plain text. For details on using R Markdown see <http://rmarkdown.rstudio.com>.
15
16 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:
17
18 ```{r, echo=TRUE}
19 a <- c(1:10, nrow=2, ncol=5)
20 a
21```
22
23 ## Including Plots
24
25 You can also embed plots, for example:
26
27 ```{r, echo=FALSE}
28 plot(pressure)
29 plot(1:10, col=1:10, pch=1:10)
30```
31
32 Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.
```

A callout box highlights two code chunks:

- echo=TRUE eval=TRUE**: code + output in report
- echo=FALSE eval=TRUE**: only output in report

Arrows point from the text descriptions to the corresponding code blocks in the RMD file.

# RMD: R MarkDown



The screenshot shows the RStudio interface with a code editor window containing an R Markdown document named 'Test\_rmd.Rmd'. The code includes sections for metadata, R code chunks, and comments about the knitr package. A black arrow points from the 'Knit Word' button in the toolbar to a dropdown menu titled 'Knit Word' which lists options: 'Knit to HTML', 'Knit to PDF', 'Knit to Word', and 'Knit with Parameters...'. Below the code editor, a large text box displays the title 'knitr package' and subtitle 'dynamic report generation in R'.

```
1 ---
2 title: "Test_rmd"
3 author: "Sarah Bonnin"
4 date: "14 Apr 2016"
5 output: word_document
6 ---
7
8 ```{r setup, include=FALSE}
9 knitr::opts_chunk$set(echo = TRUE)
10```
11
12## R Markdown
13
14 This is an R Markdown document. More details on using R Markdown see <http://rmarkdown.rstudio.com>.
15
16 When you click the **Knit** button, R will render all code chunks as
17 content as well as the output of any embedded R
18```{r, echo=TRUE}
19 a <- c(1:10, nrow=2, ncol=5)
20 a
21```
22
23## Including Plots
24
25 You can also embed plots, for example:
26
27```{r, echo=FALSE}
28 plot(1:10, col=1:10, pch=1:10)
29```
30
31 Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.
32
```

# RMD: R MarkDown

exercise3.html | [Open in Browser](#) | [Find](#)

CRG Bioinformatics Unit, [sarah.bonnin@crg.eu](mailto:sarah.bonnin@crg.eu)

October 5th, 2016

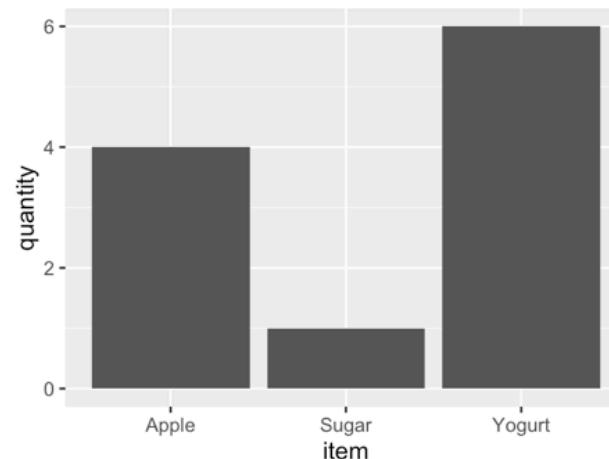
**Let's see how things works**

Read in the shopping list from exercise 2:

Load ggplot2

Produce a simple barplot

```
ggplot(a, aes(x=item, y=quantity)) + geom_bar(stat="identity")
```



**Knit to HTML**  
→ Creates a .html file  
(knitr package)

# **Exercise 4:**

# **first RMD script**

# RMD file

- Open file [Exercise4.Rmd](#) in R studio: it will serve as a template.  
→ File in  
[https://biocore.crg.eu/wiki/CRG\\_PhD\\_%26\\_Masters\\_Course\\_2016\\_Introduction\\_to\\_Statistics\\_in\\_R#MODUL\\_E\\_I.\\_Introduction\\_to\\_R\\_.Basic\\_functions\\_and\\_plotting.\\_Oct\\_5.2C\\_2016](https://biocore.crg.eu/wiki/CRG_PhD_%26_Masters_Course_2016_Introduction_to_Statistics_in_R#MODUL_E_I._Introduction_to_R_.Basic_functions_and_plotting._Oct_5.2C_2016)
- In the RMD file:
  - Try to implement some of the plots produced before.
  - Play with parameters (echo, eval, fig.height, fig.width)
  - Produce a pdf file.

# Writing functions in R

# User-written functions

- Functions are pieces of code written **to carry out specified tasks** and be able to repeat them easily.
- R allows you to create your own functions.

# Functions' structure

```
myfunction <- function(arg1, arg2, ...){
 commands } what myfunction does
 return(something)
}
```

Function's name

Creating the function

myfunction's list of arguments

Value / object myfunction returns

# Objects in functions

```
myfunction <- function(arg1){
 a <- arg1
 return(a+1)
}
```

```
> myfunction(10)
[1] 11
> a
Error: object 'a' not found
```

```
> a <- 12
> myfunction(10)
[1] 11
> a
[1] 12
```

# **Exercise 5:**

# **Functions**

## *Exercise 5*

# Functions

- Create a new file: `myfunctions.R`
- Write a function that outputs the **minimum** and a **maximum** of a numerical vector.
- Save `myfunctions.R`
- Source `myfunctions.R` with the **source** function.
- Test your function on a simple vector, for example:  
`ve <- c(0, 3, -2, 1.5, 8, 4.3, -2.1)`

# **More useful basic commands**

# Size/dimensions of objects

number of elements  
in vector, factor or list

**length(x)**

---

dimensions  
of matrix or data frame

**dim(x)**

---

number of rows  
of matrix or data frame

**nrow(x)**

---

number of columns  
of matrix or data frame

**ncol(x)**

# Elementary arithmetic operators

addition                    +

---

subtraction                -

---

division                    /

---

multiplication            \*

---

exponentiation            ^ or \*\*

# Logical operators

inferior <

inferior or equal <=

superior >

superior or equal >=

equality ==

inequality !=

intersection (“and”) &

union (“or”) |

# Obtaining summary statistics

|              |                  |
|--------------|------------------|
| average/mean | <b>mean(x)</b>   |
| median       | <b>median(x)</b> |
| minimum      | <b>min(x)</b>    |
| maximum      | <b>max(x)</b>    |
| variance     | <b>var(x)</b>    |
| correlation  | <b>cor(x)</b>    |

# Some common arithmetic functions

|                               |                |
|-------------------------------|----------------|
| natural logarithm             | <b>log(x)</b>  |
| exponential function<br>$e^x$ | <b>exp(x)</b>  |
| sine                          | <b>sin(x)</b>  |
| cosine                        | <b>cos(x)</b>  |
| tangent                       | <b>tan(x)</b>  |
| absolute value                | <b>abs(x)</b>  |
| square root                   | <b>sqrt(x)</b> |

# Objects stored in the global environment

- Listing:  
**ls()** or **objects()**
- Removing one object from environment:  
**rm(x)**
- Removing several objects:  
**rm(x, y)**
- Removing all object from environment:  
**rm(list=ls())**

# Saving objects or session

- Save objects x and y into “myobjects.RData” file
  - **save(x, y, file=“myobjects.RData”)**
- Load objects x and y into current directory:
  - **load(“myobjects.RData”)**
- Save the current workspace (all objects):
  - **save.image(file=“.RData”)**

# Command history

- Last 25 commands:
  - **history()**
- All previous commands
  - **history(max.show=Inf)**
- Save command history:
  - **savehistory()**
- Load command history:
  - **loadhistory()**

# Information about the current session

**sessionInfo()**

R version 3.2.4 (2016-03-10)

Platform: x86\_64-apple-darwin13.4.0 (64-bit)

Running under: OS X 10.9.5 (Mavericks)

**R version**

locale:

[1] C/UTF-8/C/C/C/C

attached base packages:

[1] stats graphics grDevices utils

datasets methods base

other attached packages:

[1] ggplot2\_2.1.0

**Packages attached**

loaded via a namespace (and not attached):

[1] colorspace\_1.2-6 scales\_0.4.0 plyr\_1.8.3 gtable\_0.2.0  
[5] Rcpp\_0.12.4 grid\_3.2.4 munsell\_0.4.3