

Basics of R

Workshop: Analysis of Longitudinal Data

12th Nov 2024

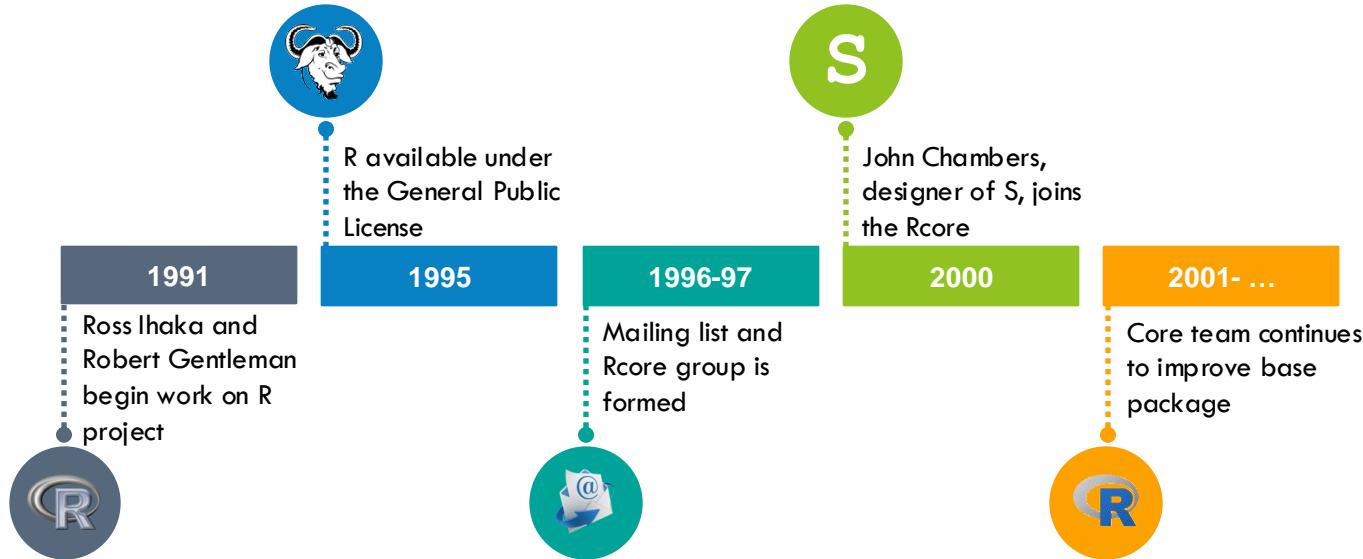
Jaroslaw Harezlak
Armando Teixeira-Pinto



Outline

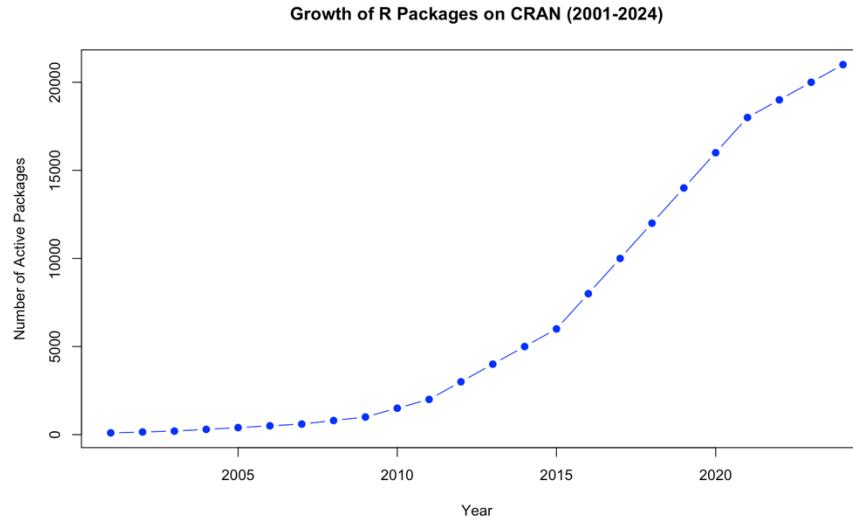
- R and RStudio
- Objects and functions
- Managing data
- Packages/Libraries

Some history



Current state

- Version 4.4.2
- Almost 22,000 packages available in CRAN



Statistical software



STATA

SPSS® SAS

Open source	Expensive	Availability	Expensive	Free Uni edition
Programming	Point & click menus	User-friendly	Point & click menus	Programming
Quickly added	Only in new versions	Updates & New methods	Only in new versions	Only in new versions
Advanced	Limited	Graphical capabilities	Limited	Limited
Many examples and large community	Technical support	Help	Technical support	Technical support

Should you learn R?

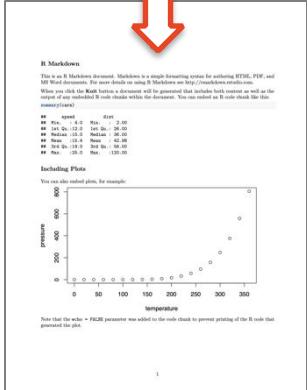
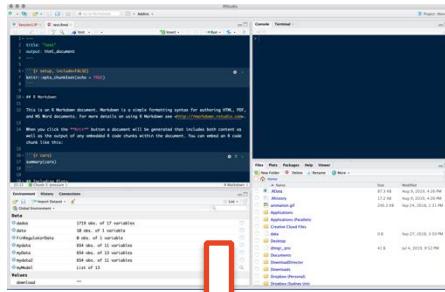
- Free
- Nice graphics
- Use advance methods
- Programming (e.g., implement new methods or run simulations)
- Reproducible research
- Look smart!

- But a steeper learning curve!

R is a programming language for statistical computing and graphics

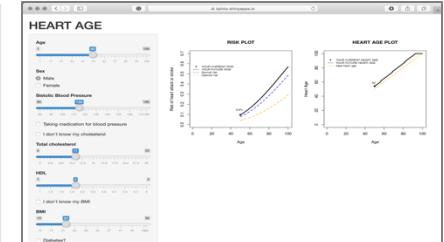
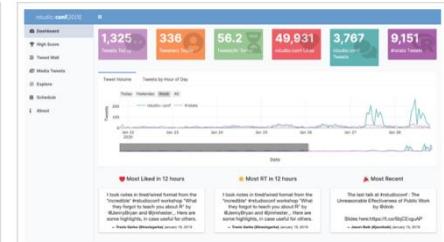
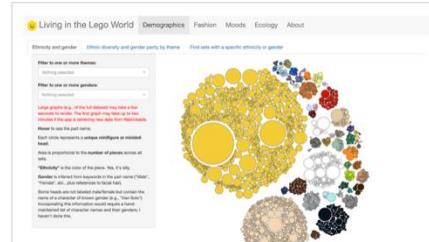
R “extensions”

– Rmarkdown / quarto



– R shiny

(<https://tpinto.shinyapps.io/shinyapp/>)

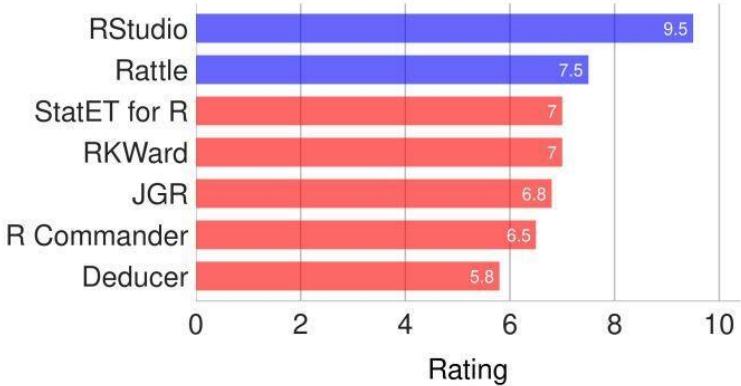


RStudio

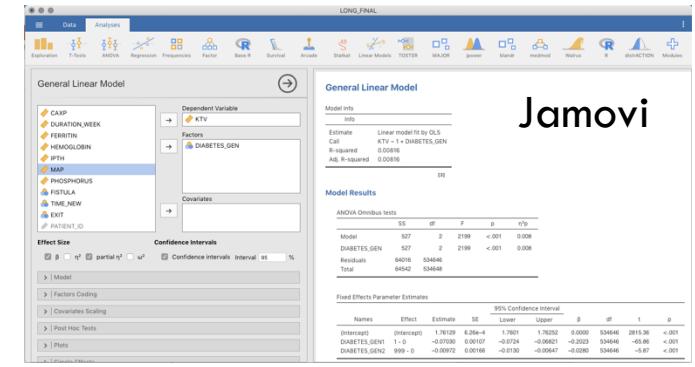


Best Free GUIs for R

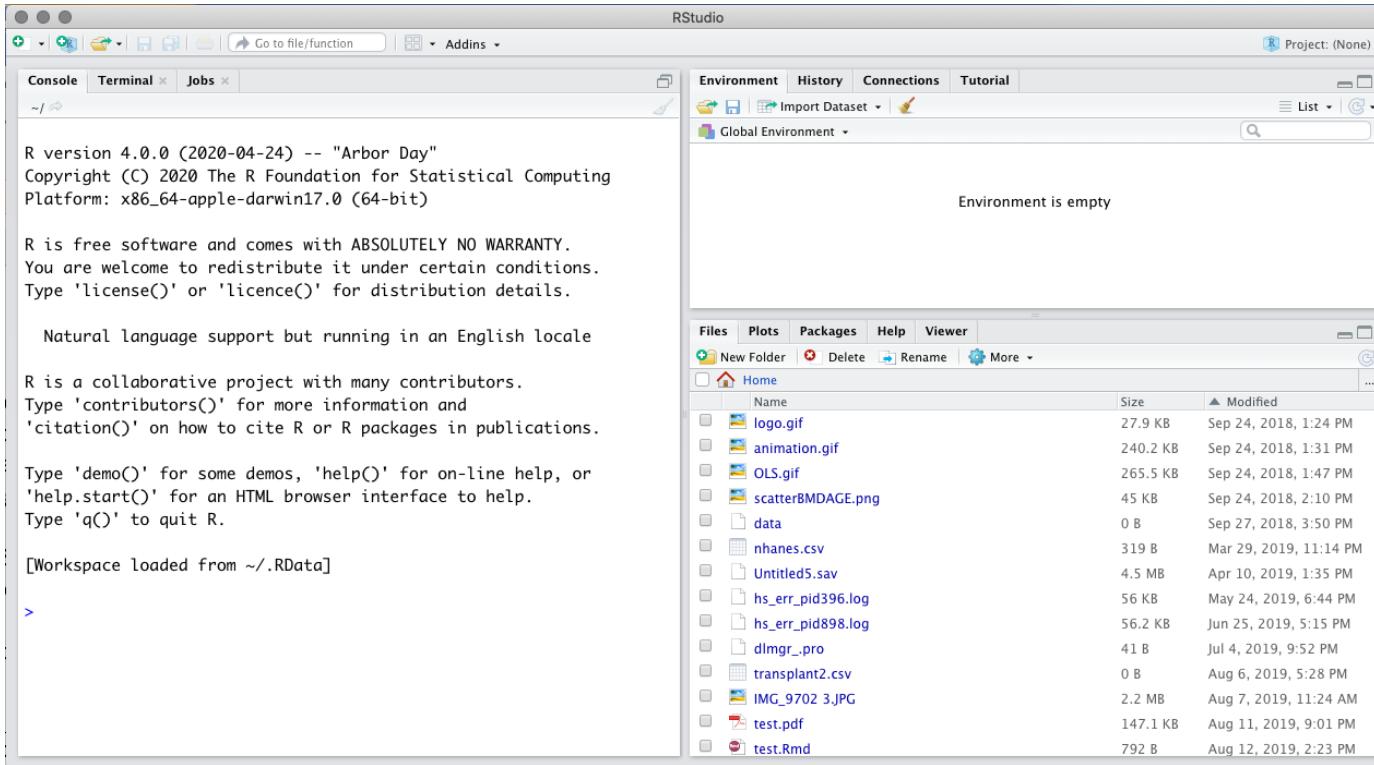
■ Recommended ■ Good



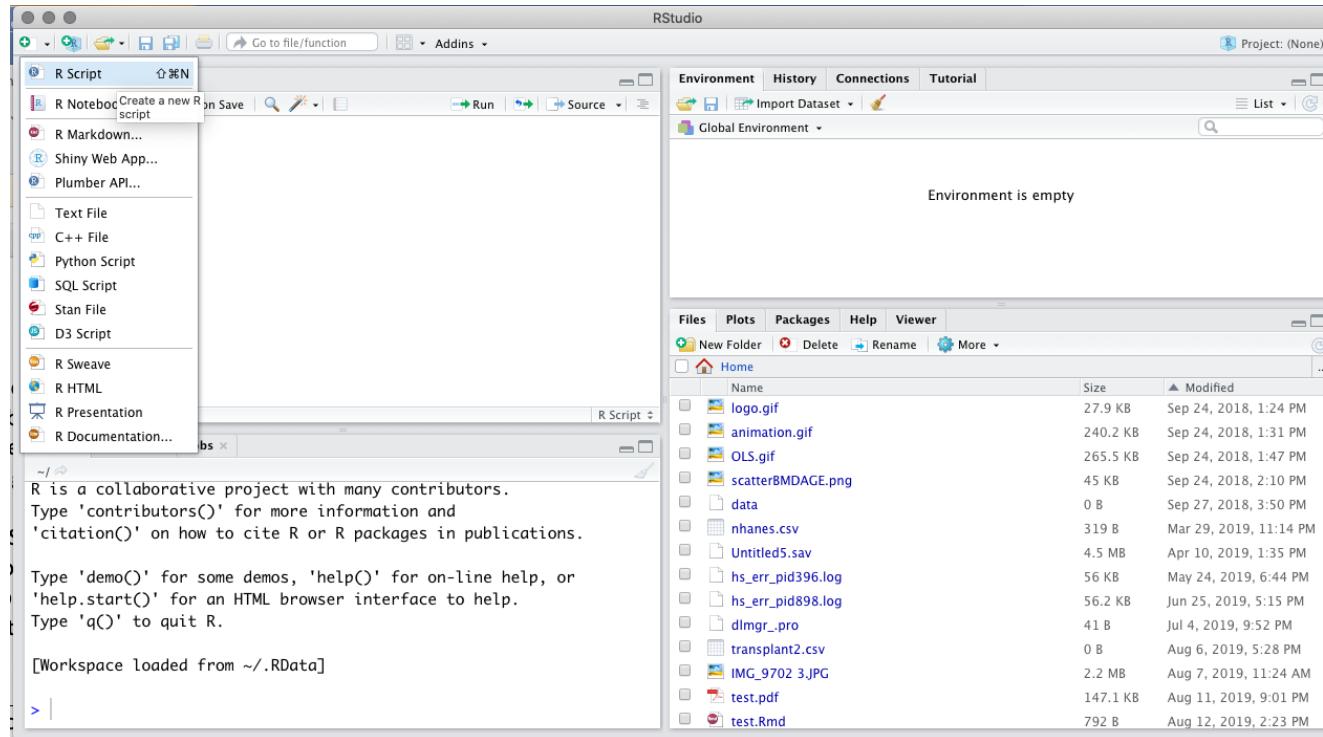
LINUX
LINKS



RStudio



RStudio



Objects

- R is an **object-oriented** programming language
- Examples of “objects”:

- **number, string, logic value**



} Data
values

- **vector (variable), matrix, data frame, list**

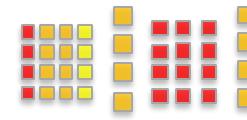
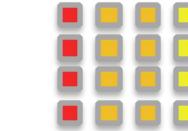
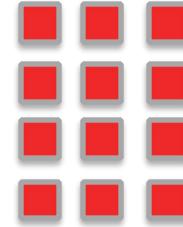


} Data
structures

- **results of a model, plot, table, ...**

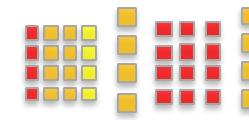
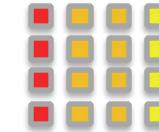
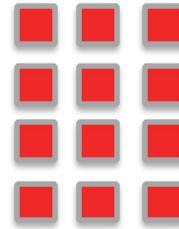
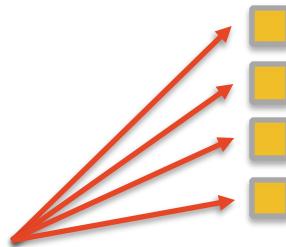
Objects

Vector	Matrix	Data frame	list
--------	--------	------------	------



Objects

Vector	Matrix	Data frame	list
--------	--------	------------	------



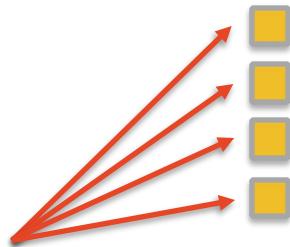
Same type of
data values

Age

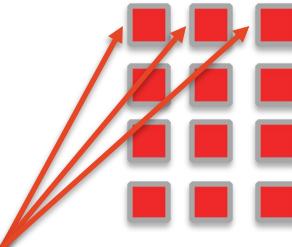
32
54
87
54
37
...

Objects

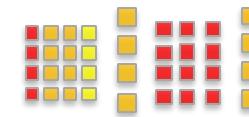
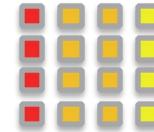
Vector	Matrix	Data frame	list
--------	--------	------------	------



Same type of
data values



Same type
of vectors



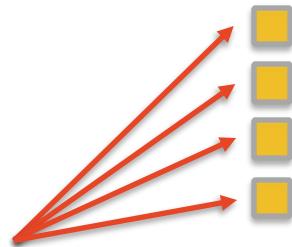
Age Correlation matrix

32
54
87
54
37
...

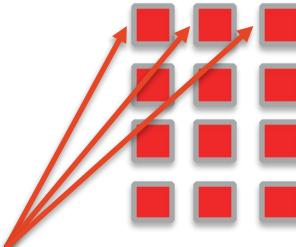
$$\begin{pmatrix} 1 & 0.32 & 0.43 \\ 0.32 & 1 & 0.27 \\ 0.43 & 0.27 & 1 \end{pmatrix}$$

Objects

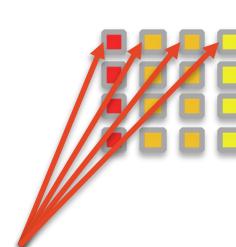
Vector	Matrix	Data frame	list
--------	--------	------------	------



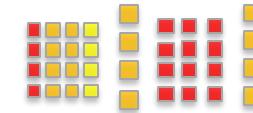
Same type of
data values



Same type
of vectors



Possibly different
type of vectors



Age Correlation matrix Dataset

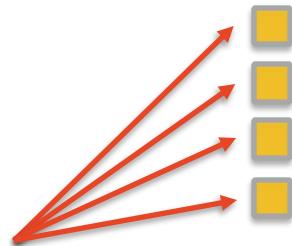
32
54
87
54
37
...

$$\begin{pmatrix} 1 & 0.32 & 0.43 \\ 0.32 & 1 & 0.27 \\ 0.43 & 0.27 & 1 \end{pmatrix}$$

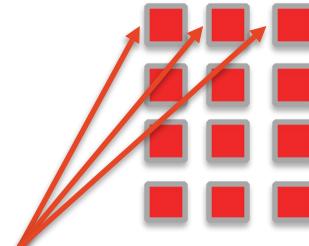
id	sex	scim	egfr_lm	cig	lung	coronary	...
1	M	134.74573	52.120399	N	N	N	
2	F	98.27387	49.141222	N	N	N	
3	F	85.16724	68.361577	N	Y	Y	
4	M	NA	NA	N	N	Y	
5	M	476.46724	9.899353	F	N	Y	
6	F	NA	NA	F	N	N	
7	F	162.64212	31.572882	N	N	N	
8	F	377.51957	13.208584	F	N	N	
9	M	100.00430	80.077657	C	N	N	

Objects

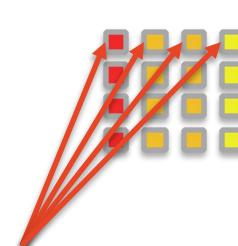
Vector	Matrix	Data frame	list
--------	--------	------------	------



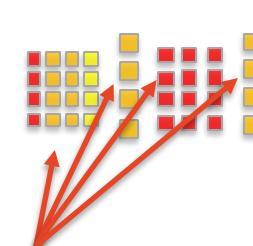
Same type of
data values



Same type
of vectors



Possibly different
type of vectors



Possibly different
type of objects

Age

Correlation matrix

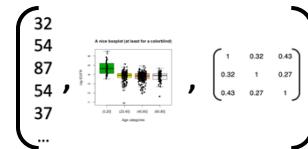
32
54
87
54
37
...

$$\begin{pmatrix} 1 & 0.32 & 0.43 \\ 0.32 & 1 & 0.27 \\ 0.43 & 0.27 & 1 \end{pmatrix}$$

Dataset

id	sex	scim	egfr_lm	cig	lung	coronary	...
1	M	134.74573	52.120399	N	N	N	
2	F	98.27387	49.141222	N	N	N	
3	F	85.16724	68.361577	N	Y	Y	
4	M	NA	NA	N	N	Y	
5	M	476.46724	9.899353	F	N	Y	
6	F	NA	NA	F	N	N	
7	F	162.64212	31.572802	N	N	N	
8	F	377.51957	13.208504	F	N	N	
9	M	100.00430	80.077657	C	N	N	

List



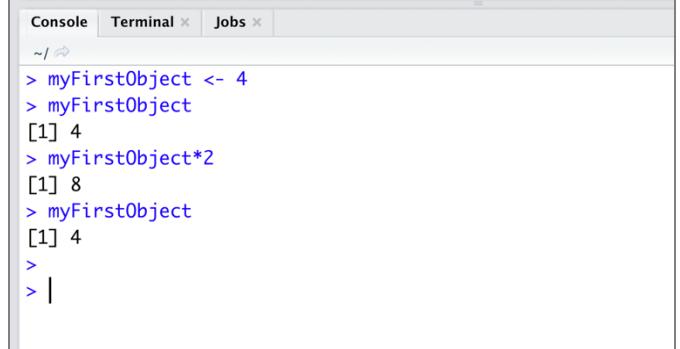
Objects

- You can give a name to the “objects”
- For example,

```
myFirstObject <- 4
```

- This creates an object (numeric) with value 4 and name `myFirstObject` (the “`<-`” (“less than”, “minus”) makes suggests an arrow)

```
# Try in R:  
  
myFirstObject <- 4  
myFirstObject  
myFirstObject*2  
myFirstObject
```



The screenshot shows an R console interface with three tabs: "Console", "Terminal", and "Jobs". The "Console" tab is active. The history pane at the top shows the command `> myFirstObject <- 4`. The main pane displays the results of the assignment and subsequent operations:

```
> myFirstObject <- 4
> myFirstObject
[1] 4
> myFirstObject*2
[1] 8
> myFirstObject
[1] 4
>
> |
```

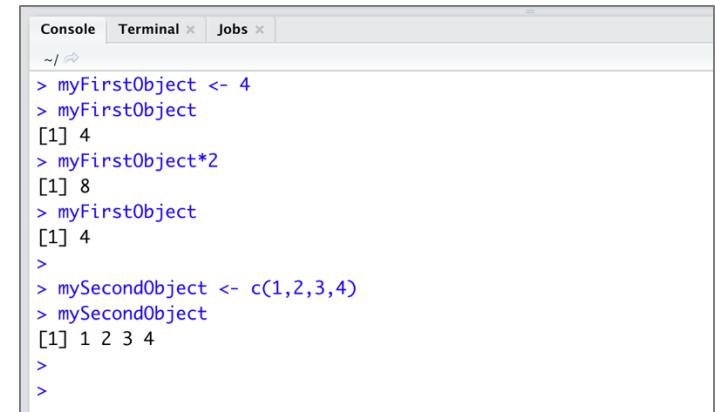
Functions

- Now, a vector

```
mySecondObject <- c(1,2,3,4)
```

```
# Try in R:
```

```
mySecondObject <- c(1,2,3,4)  
mySecondObject
```



The screenshot shows an R console interface with three tabs: 'Console' (selected), 'Terminal', and 'Jobs'. The console area displays the following R session:

```
Console Terminal Jobs  
~/  
> myFirstObject <- 4  
> myFirstObject  
[1] 4  
> myFirstObject*2  
[1] 8  
> myFirstObject  
[1] 4  
>  
> mySecondObject <- c(1,2,3,4)  
> mySecondObject  
[1] 1 2 3 4  
>  
>
```

Functions

- Now, a vector

```
mySecondObject <- c(1,2,3,4)
```

- And we can now apply functions to the objects

function(object)

- For example, we can compute the mean

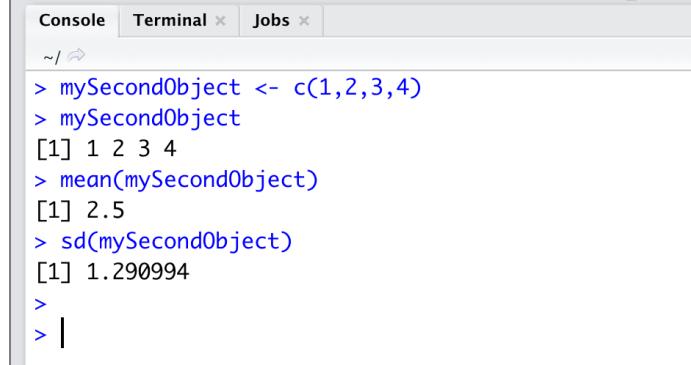
mean(
 ^{function}mySecondObject
 ^{object})

- And the standard deviation

sd(mySecondObject)

```
#Try in R:
```

```
mySecondObject <- c(1,2,3,4)
mySecondObject
mean(mySecondObject)
sd(mySecondObject)
```



The screenshot shows an R console window with three tabs: 'Console' (selected), 'Terminal', and 'Jobs'. The console output is as follows:

```
> mySecondObject <- c(1,2,3,4)
> mySecondObject
[1] 1 2 3 4
> mean(mySecondObject)
[1] 2.5
> sd(mySecondObject)
[1] 1.290994
>
> |
```

Data management

- Let's now load a dataset `transplant.csv`

id	sex	sc1m	egfr_1m	cig	lung	coronary	cvd	age	ftime	died	ID	M
2	1 M	134.74573	52.120399	N	N	N	N	67.740600	2	0	1	
3	2 F	98.7387	49.141221	N	N	N	N	59.660398	2	0	2	
4	3 F	85.167236	68.361577	N	Y	N	N	33.152289	2	1	3	
5	4 M	NA	NA	N	N	Y	Y	49.114292	2	0	4	
6	5 M	476.467232	9.899353	F	N	N	Y	25.778478	2	1	5	
7	6 F	162.642124	31.572801	N	N	N	N	53.835584	2	0	6	
8	7 F	162.642124	31.5728018	N	N	N	N	51.800447	2	0	11	
9	8 F	377.519574	13.208509	F	N	N	N	53.8545713	3	1	14	
10	9 M	100.00429	80.077656	C	N	N	N	43.5955863	3	0	15	
11	10 M	94.73901	51.708822	N	N	N	N	60.8375893	3	0	16	
12	11 M	67.112341	44.229434	N	N	N	N	58.835586	3	0	17	
13	12 M	82.481882	123.66017	N	N	N	N	15.3639781	4	0	18	
14	13 F	71.8358818	94.2294346	N	Y	N	N	35.4501034	4	0	19	
15	14 F	102.58364	58.9585379	N	S	N	N	46.726535	4	0	20	
16	15 M	124.179883	53.535461	N	N	N	N	67.2874862	4	0	21	
17	16 M	69.907095	53.848312	N	N	N	N	60.837589	4	0	22	
18	17 M	NA	NA	C	S	S	Y	52.0488498	4	0	23	
19	18 M	203.144742	36.2273846	N	S	N	N	49.0262819	4	1	26	
20	19 M	142.545616	46.5371895	N	N	N	N	70.6254839	5	0	27	
21	20 M	262.960022	24.2487807	C	N	N	N	46.6359351	5	1	31	
22	21 M	148.387411	50.0391627	C	N	N	N	28.2398711	5	1	32	

#Try in R:

```
#Usually, the file would be stored locally  
#myData <- read.csv("c:\\yourpath\\transplant.csv")  
  
myData <-  
read.csv("https://www.dropbox.com/scl/fi/9qidag3d5odyvlqcrixln/transplant.cs  
v?rlkey=livvep6peuz2xpvmr2qcqt84b&st=wjoqpefr&dl=1")  
myData
```

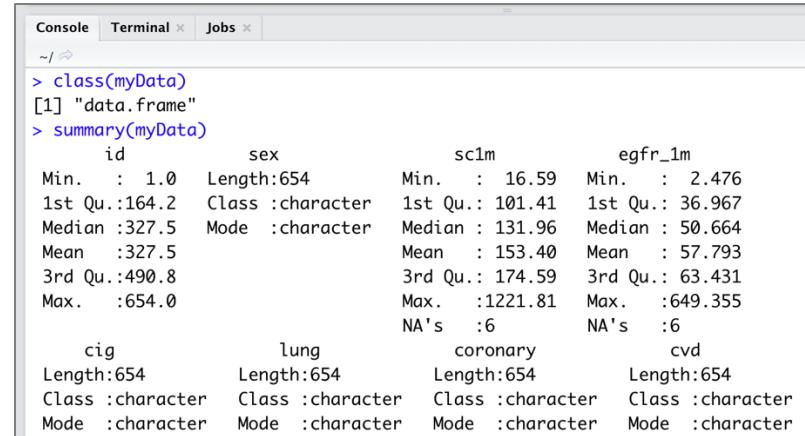
	id	sex	sc1m	egfr_1m	cig	lung	coronary	cvd	age	ftime	died	
1	1	M	134.74573	52.120399	N	N	N	N	67.740600	2	0	
2	2	F	98.27387	49.141222	N	N	N	N	59.660398	2	0	
3	3	F	85.16724	68.361577	N	Y	Y	Y	33.152289	2	1	
4	4	M	NA	NA	N	N	Y	Y	49.114293	2	0	
5	5	M	476.46724	9.899353	F	N	N	Y	25.778478	2	1	
6	6	F	162.64212	31.572802	N	N	N	N	53.835584	2	0	
7	7	F	162.642124	31.572802	N	N	N	N	51.800447	2	0	
8	8	F	377.51957	13.208504	F	N	N	N	53.854571	3	1	
9	9	M	100.00430	80.077657	C	N	N	N	43.595586	3	0	
10	10	F	94.73901	51.708823	N	N	N	N	60.837589	3	0	
11	11	M	67.1151512	121.944604	F	N	N	N	28.885556	3	0	
12	12	M	82.48188	122.606717	N	N	N	N	15.363978	4	0	

Data management

- We can also apply *functions* to this new “object” myData
- For example,

summary(myData)

```
#Try in R:  
#let's first check that myData is a data frame  
  
class(myData)  
summary(myData)
```



The screenshot shows the RStudio interface with the 'Console' tab selected. The command > summary(myData) is entered, and the resulting output is displayed. The output shows the summary statistics for various variables in the myData data frame. The variables include id, sex, sc1m, egfr_1m, cig, lung, coronary, and cvd. For each variable, the output provides the minimum value, the first quartile, the median, the mean, the third quartile, and the maximum value. It also indicates the data type (character or mode) and the number of NA values.

Variable	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	Class / Mode	NA's
id	1.0	164.2	327.5	490.8	654.0	1221.81	character	6
sex	Length:654	Class :character	Mode :character	Length:654	Class :character	Mode :character	character	6
sc1m	16.59	101.41	131.96	153.40	174.59	221.81	character	6
egfr_1m	2.476	36.967	50.664	57.793	63.431	649.355	character	6
cig	Length:654	Class :character	Mode :character	Length:654	Class :character	Mode :character	character	6
lung	Length:654	Class :character	Mode :character	Length:654	Class :character	Mode :character	character	6
coronary	Length:654	Class :character	Mode :character	Length:654	Class :character	Mode :character	character	6
cvd	Length:654	Class :character	Mode :character	Length:654	Class :character	Mode :character	character	6

Data management

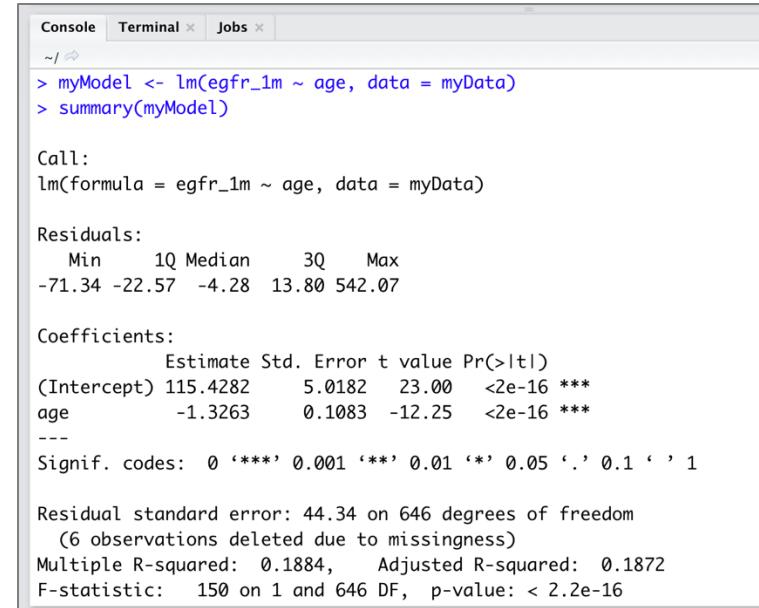
- We can also apply *functions* to this new “object” `myData`
- For example,

`summary(myData)`

Notice that if we apply the function `summary` to a different type of object, e.g. a model, it will produce a different output

R (usually) understands the context of the object and function

```
#Try in R:  
#let's first check that myData is a data frame  
  
class(myData)  
summary(myData)  
  
# "summary" for a model will produce a different output  
  
myModel <- lm(egfr_1m ~ age, data = myData)  
summary(myModel)
```



The screenshot shows the RStudio interface with the 'Console' tab selected. The console window displays the following R session:

```
Console Terminal × Jobs ×  
~/  
> myModel <- lm(egfr_1m ~ age, data = myData)  
> summary(myModel)
```

Call:

```
lm(formula = egfr_1m ~ age, data = myData)
```

Residuals:

Min	1Q	Median	3Q	Max
-71.34	-22.57	-4.28	13.80	542.07

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	115.4282	5.0182	23.00	<2e-16 ***
age	-1.3263	0.1083	-12.25	<2e-16 ***

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 44.34 on 646 degrees of freedom
(6 observations deleted due to missingness)

Multiple R-squared: 0.1884, Adjusted R-squared: 0.1872
F-statistic: 150 on 1 and 646 DF, p-value: < 2.2e-16

Data management

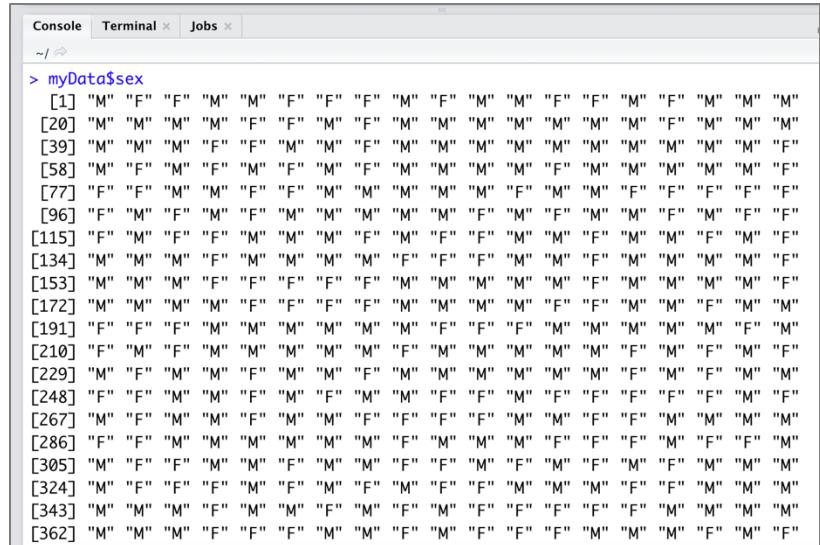
- Each variable in `myData` is also an object

id, sex, sc1m, egfr_1m,...

- The variable `sex`, e.g., is identified as

`myData$sex`

```
#Try in R:  
  
myData$sex
```



A screenshot of the RStudio interface showing the console tab. The command `> myData$sex` is entered, and the output is a large vector of character strings representing gender values ('M' or 'F') for 362 observations.

```
Console Terminal × Jobs ×  
~ /  
> myData$sex  
[1] "M" "F" "M" "M" "F" "F" "M" "F" "M" "M" "F" "F" "M" "M" "M" "M" "M"  
[20] "M" "M" "M" "M" "F" "F" "M" "F" "M" "M" "M" "M" "M" "M" "M" "M" "M"  
[39] "M" "M" "M" "F" "F" "M" "M" "F" "M"  
[58] "M" "F" "M" "F" "M" "F" "M" "F" "M" "M" "M" "M" "F" "M" "M" "M" "M" "M"  
[77] "F" "F" "M" "F" "M" "M" "M" "M" "M" "M" "M" "F" "M" "M" "F" "F" "F" "F"  
[96] "F" "M" "F" "F" "M" "M" "M" "M" "F" "M" "F" "M" "M" "F" "M" "M" "F" "F"  
[115] "F" "M" "F" "M" "M" "M" "F" "M" "F" "F" "M" "M" "F" "M" "F" "M" "F" "F"  
[134] "M" "M" "F" "M" "M" "M" "M" "F" "F" "F" "M" "M" "F" "M" "M" "M" "M" "F"  
[153] "M" "M" "F" "F" "F" "F" "M" "M" "M" "M" "M" "F" "M" "M" "M" "M" "M" "F"  
[172] "M" "M" "M" "F" "F" "F" "M" "M" "M" "M" "M" "F" "F" "M" "M" "F" "M" "M"  
[191] "F" "F" "M" "M" "M" "M" "M" "F" "F" "F" "F" "M" "M" "M" "M" "M" "F" "M"  
[210] "F" "M" "F" "M" "M" "M" "M" "F" "M" "M" "M" "M" "M" "F" "M" "F" "M" "F"  
[229] "M" "F" "M" "F" "M" "M" "F" "M" "M" "M" "M" "M" "F" "M" "F" "M" "F" "M"  
[248] "F" "F" "M" "F" "M" "F" "M" "F" "F" "M" "F" "F" "F" "F" "F" "M" "F" "F"  
[267] "M" "F" "M" "M" "F" "M" "M" "F" "F" "F" "M" "M" "F" "F" "M" "M" "M" "M"  
[286] "F" "F" "M" "M" "M" "M" "M" "F" "M" "M" "M" "F" "F" "M" "F" "M" "F" "F"  
[305] "M" "F" "F" "M" "M" "F" "M" "F" "F" "M" "F" "M" "F" "M" "F" "M" "F" "M"  
[324] "M" "F" "F" "M" "M" "F" "M" "F" "F" "M" "M" "F" "M" "F" "M" "F" "M" "M"  
[343] "M" "M" "M" "F" "F" "M" "M" "F" "M" "F" "M" "F" "F" "M" "M" "F" "M" "M"  
[362] "M" "M" "M" "F" "F" "M" "M" "F" "M" "F" "M" "F" "F" "M" "M" "F" "M" "F"
```

Data management

- **Exercise:** Identify the variable age and get the mean age
- Let's use the function `table` to get frequencies of each category of sex
`table(myData$sex)`
`function(object)`
- **Exercise:** How many patients died (variable `died`)?

```
#Try in R:  
mean(myData$age)  
  
myData$sex  
table(myData$sex)
```



The screenshot shows the RStudio interface with the 'Console' tab selected. The command `> table(myData$sex)` is entered, and the resulting frequency table is displayed:

	F	M
255	399	

Data management

- Specific observation (values) in the variable can be identified using []
- If you want the first four observations of the variable `sex`

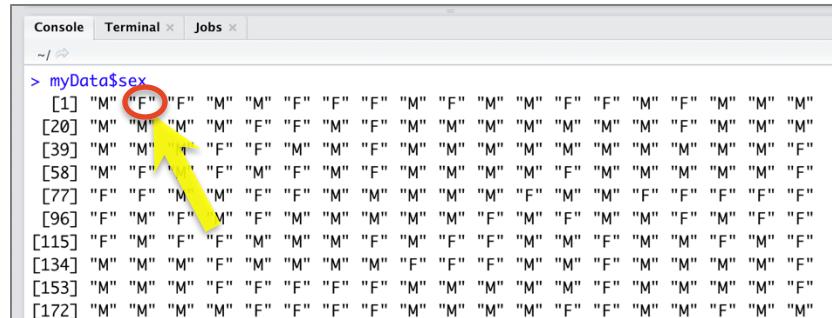
`myData$sex[2]`

- If you want the first four observations of the variable `sex`

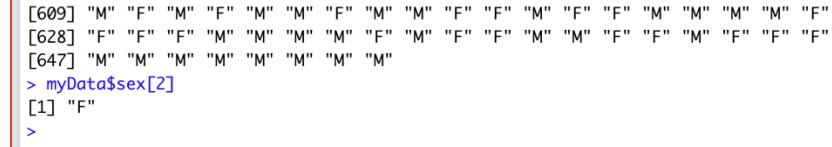
`myData$sex[c(1, 2, 3, 4)]`

#Try in R:

```
myData$sex[2]  
myData$sex[c(1, 2, 3, 4)]
```



```
Console Terminal × Jobs ×  
~/  
> myData$sex  
[1] "M" "F" "F" "M" "M" "F" "F" "M" "F" "M" "M" "F" "F" "M" "F" "M" "M" "M"  
[20] "M" "M" "M" "M" "F" "F" "M" "F" "M" "M" "M" "M" "M" "M" "M" "F" "M" "M" "M"  
[39] "M" "M" "F" "F" "M" "M" "F" "M" "F"  
[58] "M" "F" "F" "M" "E" "M" "F" "M" "M" "M" "F" "M" "M" "M" "M" "M" "M" "M" "F"  
[77] "F" "F" "M" "M" "F" "F" "M" "M" "M" "M" "F" "F" "M" "M" "F" "F" "F" "F" "F"  
[96] "F" "M" "F" "M" "F" "M" "M" "M" "F" "M" "F" "M" "F" "M" "F" "M" "F" "M" "F"  
[115] "F" "M" "F" "M" "M" "M" "F" "M" "F" "F" "M" "F" "M" "F" "M" "F" "M" "F" "F"  
[134] "M" "M" "F" "M" "M" "M" "F" "F" "F" "M" "M" "F" "M" "M" "F" "M" "M" "M" "F"  
[153] "M" "M" "M" "F" "F" "F" "F" "M" "M" "M" "M" "M" "F" "M" "M" "M" "M" "M" "F"  
[172] "M" "M" "M" "M" "F" "F" "F" "F" "M" "M" "M" "M" "F" "F" "M" "M" "F" "M" "M"
```

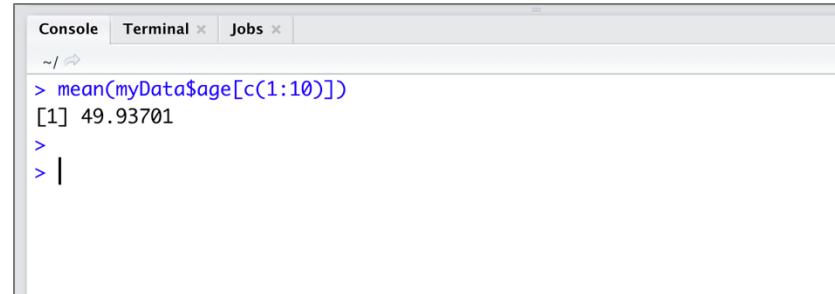


```
[609] "M" "F" "M" "F" "M" "M" "F" "M" "M" "F" "F" "M" "F" "F" "M" "M" "M" "M" "F"  
[628] "F" "F" "F" "M" "M" "M" "F" "M" "F" "F" "M" "M" "F" "F" "M" "F" "F" "F"  
[647] "M" "F" "F"  
> myData$sex[2]  
[1] "F"  
>
```

Data management

- **Exercise:** Get the mean age of the 10 first individuals

```
#Try in R:  
  
#c(1:10) will give all 1,2, ..., 10  
mean(myData$age[c(1:10)])
```

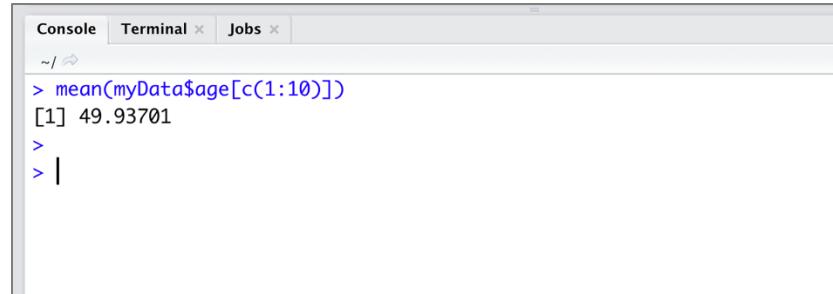


The screenshot shows an R console interface with three tabs: 'Console' (selected), 'Terminal', and 'Jobs'. The current working directory is indicated as '~/'. In the console area, the command `> mean(myData$age[c(1:10)])` is entered, followed by its output: `[1] 49.93701`. Below the command line, there are two empty lines starting with a greater than sign (>).

Data management

- **Exercise:** Get the mean age of the 10 first individuals
- Now, let's try something useful! Let's compute the mean age of men
- First, let's see how do we get the ages of every men

```
#Try in R:  
  
#c(1:10) will give all 1,2, ..., 10  
mean(myData$age[c(1:10)])
```



The screenshot shows an R console interface with three tabs: 'Console' (selected), 'Terminal', and 'Jobs'. The current working directory is indicated as '~/'. In the console area, the user has typed the command `> mean(myData$age[c(1:10)])`. The console then displays the output: `[1] 49.93701`. Below the output, there are two additional prompt lines: `>` and `> |`.

Data management

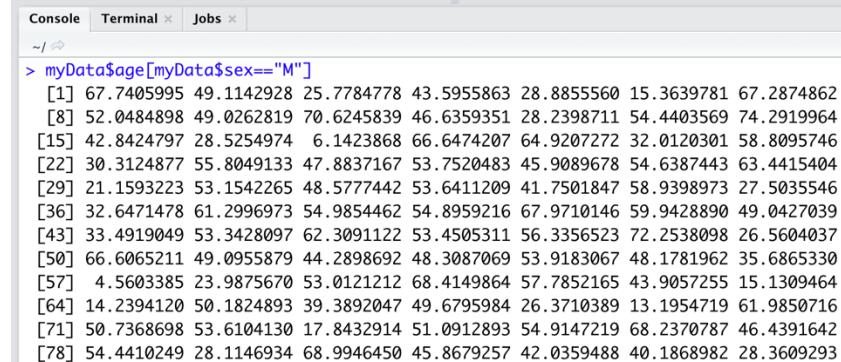
- **Exercise:** Get the mean age of the 10 first individuals
- Now, let's try something useful! Let's compute the mean age of men
- First, let's see how do we get the ages of every men

myData\$age [myData\$sex == "M"]

Values of age in the positions where the variable sex is M

equal
↓

```
#Try in R:  
  
#c(1:10) will give all 1,2, ..., 10  
mean(myData$age[c(1:10)])  
  
myData$age[myData$sex=="M"]
```



The screenshot shows the RStudio interface with the 'Console' tab selected. The command `> myData$age[myData$sex=="M"]` is entered, and its output is displayed below. The output consists of 78 numerical values representing the ages of men from the dataset.

```
[1] 67.7405995 49.1142928 25.7784778 43.5955863 28.8855560 15.3639781 67.2874862  
[8] 52.0484898 49.0262819 70.6245839 46.6359351 28.2398711 54.4403569 74.2919964  
[15] 42.8424797 28.5254974 6.1423868 66.6474207 64.9207272 32.0120301 58.8095746  
[22] 30.3124877 55.8049133 47.8837167 53.7520483 45.9089678 54.6387443 63.4415404  
[29] 21.1593223 53.1542265 48.5777442 53.6411209 41.7501847 58.9398973 27.5035546  
[36] 32.6471478 61.2996973 54.9854462 54.8959216 67.9710146 59.9428890 49.0427039  
[43] 33.4919049 53.3428097 62.3091122 53.4505311 56.3356523 72.2538098 26.5604037  
[50] 66.6065211 49.0955879 44.2898692 48.3087069 53.9183067 48.1781962 35.6865330  
[57] 4.5603385 23.9875670 53.0121212 68.4149864 57.7852165 43.9057255 15.1309464  
[64] 14.2394120 50.1824893 39.3892047 49.6795984 26.3710389 13.1954719 61.9850716  
[71] 50.7368698 53.6104130 17.8432914 51.0912893 54.9147219 68.2370787 46.4391642  
[78] 54.4410249 28.1146934 68.9946450 45.8679257 42.0359488 40.1868982 28.3609293
```

Data management

- **Exercise:** Get the mean age of the 10 first individuals
- Now, let's try something useful! Let's compute the mean age of men
- First, let's see how do we get the ages of every men

myData\$age [myData\$sex == "M"]

Values of age in the positions where the variable sex is M

- So, the only thing we need is to apply the mean function to the above

mean (myData\$age [myData\$sex=="M"])

```
#Try in R:  
  
#c(1:10) will give all 1,2, ..., 10  
mean(myData$age[c(1:10)])  
  
myData$age[myData$sex=="M"]  
mean (myData$age [myData$sex=="M" ] )
```



The screenshot shows an R console interface with three tabs: 'Console', 'Terminal', and 'Jobs'. The 'Console' tab is active, displaying the following session:

```
Console Terminal Jobs  
~/  
> mean(myData$age[myData$sex=="M"])  
[1] 44.10213  
>  
> |
```

Data management

- **Exercise:** get the number of men and women, older than 50 years old

Data management

- **Exercise:** get the number of men and women, older than 50 years old

```
table(myData$sex[myData$age>50])
```

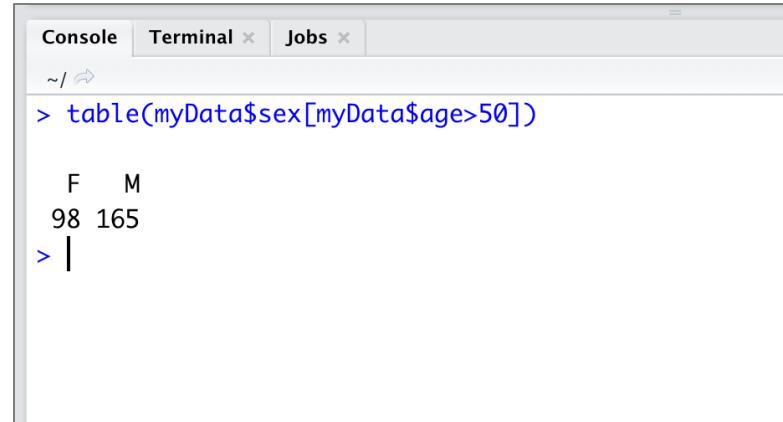
The function with () simplifies the writing of the expression above:

```
with(myData, table(sex[age>50]))
```

For data management there is another way of writing commands (different grammar) that is very popular. If you want to know more about this, <https://dplyr.tidyverse.org>

```
#Try in R:
```

```
table(myData$sex[myData$age>50])
```



```
Console Terminal × Jobs ×
~/
> table(myData$sex[myData$age>50])
   F   M
 98 165
> |
```

	F	M
98	165	

Data management

- **Exercise:** get the number of men and women, older than 50 years old

```
table(myData$sex[myData$age>50])
```

The function with () simplifies the writing of the expression above:

```
with(myData, table(sex[age>50]))
```

For data management
writing commands
very popular
this, <https://>

```
#Try in R:
```

```
table(myData$sex[myData$age>50])
```

```
Console Terminal × Jobs ×
~/
> table(myData$sex[myData$age>50])
   F   M
 98 165
> |
```

Data management

- Data frames use the same indexing idea but now with two dimensions, row number and column number

myData[1, 3]

Row 1 Column 3
(3rd variable)

Instead of the column's number, you can use the name of the variable instead

myData[1, "sc1m"]

```
#Try in R:
```

```
myData  
myData[1, 3]
```

Console Terminal × Jobs ×

~/

```
> myData
```

	id	sex	sc1m	egfr_1m	cig	lung	coronary	cvd	age	ftime	died
1	1	M	134.74573	52.120399	N	N	N	N	67.740600	2	0
2	2	F	98.27387	49.141222	N	N	N	N	59.660398	2	0
3	3	F	85.16724	68.361577	N	Y	Y	N	33.152289	2	1
4	4	M	NA	NA	N	N	Y	Y	49.114293	2	0
5	5	M	476.46724	9.899353	F	N	Y	N	25.778478	2	1
6	6	F	NA	NA	F	N	N	N	53.835844	2	0
7	7	F	162.64212	31.572802	N	N	N	N	51.800447	2	0
8	8	F	377.51957	13.208504	F	N	N	N	53.854571	3	1

```
89 89 M 122.33592 51.994959 F N N N 57.785216 21 0  
90 90 M 151.80705 45.462498 N N N N 43.905726 21 0  
[ reached 'max' / getOption("max.print") -- omitted 564 rows ]  
> myData[1,3]  
[1] 134.7457  
>  
> |
```

Data management

- If we want to subset the data, for example, only women, we can use

myData [myData\$sex=="F",]

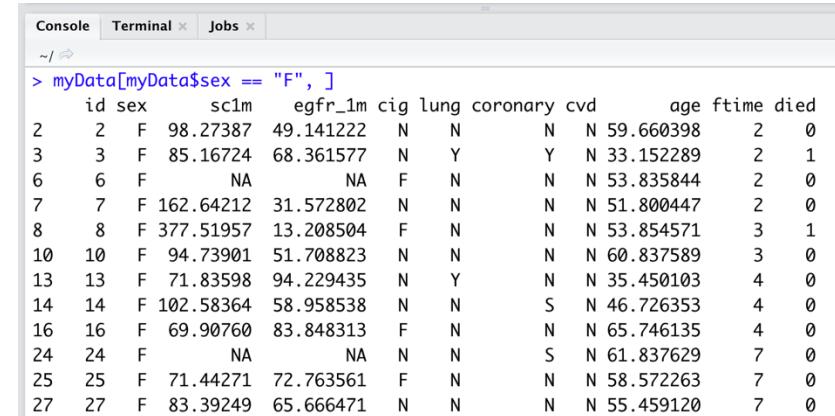
All the rows
where sex is F



No number
indicates all
the variables

#Try in R:

```
myData[myData$sex == "F", ]  
myData.Females <- myData[myData$sex == "F", ]
```



The screenshot shows the RStudio interface with the 'Console' tab selected. The command `> myData[myData$sex == "F",]` is entered, and the resulting data frame is displayed below. The data frame contains 27 rows and 12 columns, with the first few rows shown:

	id	sex	sc1m	egfr_1m	cig	lung	coronary	cvd	age	ftime	died
2	2	F	98.27387	49.141222	N	N	N	N	59.660398	2	0
3	3	F	85.16724	68.361577	N	Y	Y	N	33.152289	2	1
6	6	F	NA	NA	F	N	N	N	53.835844	2	0
7	7	F	162.64212	31.572802	N	N	N	N	51.800447	2	0
8	8	F	377.51957	13.208504	F	N	N	N	53.854571	3	1
10	10	F	94.73901	51.708823	N	N	N	N	60.837589	3	0
13	13	F	71.83598	94.229435	N	Y	N	N	35.450103	4	0
14	14	F	102.58364	58.958538	N	N	S	N	46.726353	4	0
16	16	F	69.90760	83.848313	F	N	N	N	65.746135	4	0
24	24	F	NA	NA	N	N	S	N	61.837629	7	0
25	25	F	71.44271	72.763561	F	N	N	N	58.572263	7	0
27	27	F	83.39249	65.666471	N	N	N	N	55.459120	7	0

Data management

- If we want to subset the data, for example, only women, we can use

myData [myData\$sex=="F",]

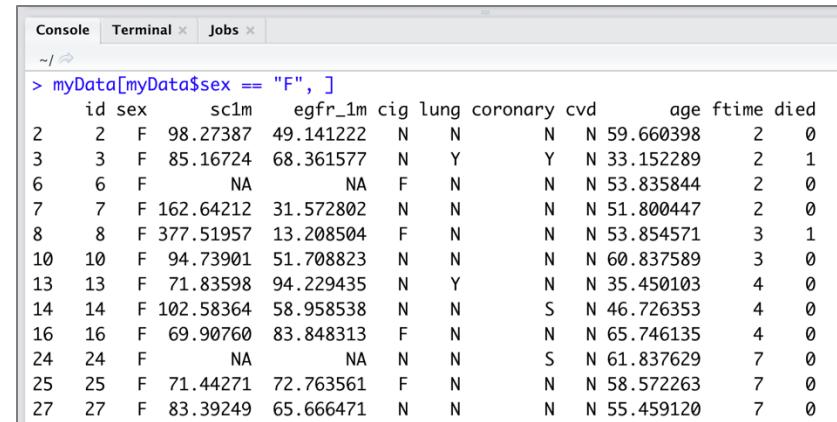
All the rows
where sex is F



No number
indicates all
the variables

#Try in R:

```
myData[myData$sex == "F", ]  
myData.Females <- myData[myData$sex == "F", ]
```



		id	sex	sc1m	egfr_1m	cig	lung	coronary	cvd	age	ftime	died
2	2	2	F	98.27387	49.141222	N	N	N	N	59.660398	2	0
3	3	3	F	85.16724	68.361577	N	Y	Y	N	33.152289	2	1
6	6	6	F	NA	NA	F	N	N	N	53.835844	2	0
7	7	7	F	162.64212	31.572802	N	N	N	N	51.800447	2	0
8	8	8	F	377.51957	13.208504	F	N	N	N	53.854571	3	1
10	10	10	F	94.73901	51.708823	N	N	N	N	60.837589	3	0
13	13	13	F	71.83598	94.229435	N	Y	N	N	35.450103	4	0
14	14	14	F	102.58364	58.958538	N	N	S	N	46.726353	4	0
16	16	16	F	69.90760	83.848313	F	N	N	N	65.746135	4	0
24	24	24	F	NA	NA	N	N	S	N	61.837629	7	0
25	25	25	F	71.44271	72.763561	F	N	N	N	58.572263	7	0
27	27	27	F	83.39249	65.666471	N	N	N	N	55.459120	7	0

Data management

- If we want to subset the data, for example, only women, we can use

```
myData [myData$sex=="F", ]
```

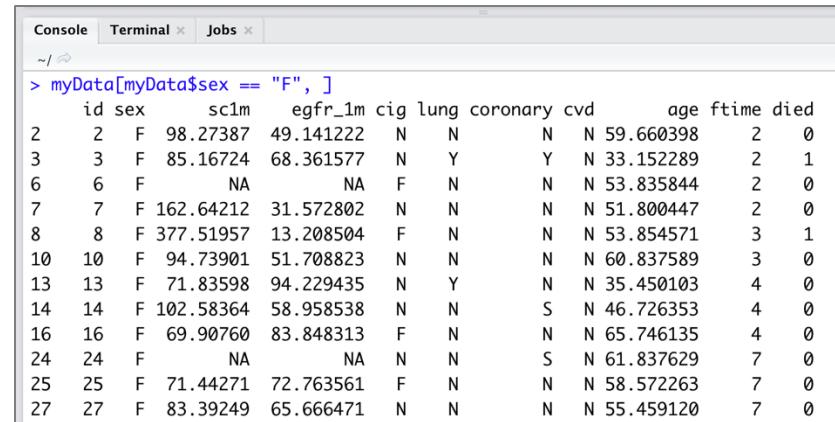
All the rows
where sex is F

No number
indicates all
the variables

- **Exercise:** Store the dataset with only females, in a new object

#Try in R:

```
myData[myData$sex == "F", ]  
myData.Females <- myData[myData$sex == "F", ]
```



The screenshot shows the RStudio interface with the 'Console' tab selected. The command `> myData[myData$sex == "F",]` is entered, and its output is displayed below. The output is a data frame with 27 rows and 12 columns. The columns are labeled: id, sex, sc1m, egfr_1m, cig, lung, coronary, cvd, age, ftime, and died. The 'sex' column shows values 'F' for all rows. The 'id' column ranges from 2 to 27. Other columns contain various numerical and categorical values.

		id	sex	sc1m	egfr_1m	cig	lung	coronary	cvd	age	ftime	died
2	2	2	F	98.27387	49.141222	N	N	N	N	59.660398	2	0
3	3	3	F	85.16724	68.361577	N	Y	Y	N	33.152289	2	1
6	6	6	F	NA	NA	F	N	N	N	53.835844	2	0
7	7	7	F	162.64212	31.572802	N	N	N	N	51.800447	2	0
8	8	8	F	377.51957	13.208504	F	N	N	N	53.854571	3	1
10	10	10	F	94.73901	51.708823	N	N	N	N	60.837589	3	0
13	13	13	F	71.83598	94.229435	N	Y	N	N	35.450103	4	0
14	14	14	F	102.58364	58.958538	N	N	S	N	46.726353	4	0
16	16	16	F	69.90760	83.848313	F	N	N	N	65.746135	4	0
24	24	24	F	NA	NA	N	N	S	N	61.837629	7	0
25	25	25	F	71.44271	72.763561	F	N	N	N	58.572263	7	0
27	27	27	F	83.39249	65.666471	N	N	N	N	55.459120	7	0

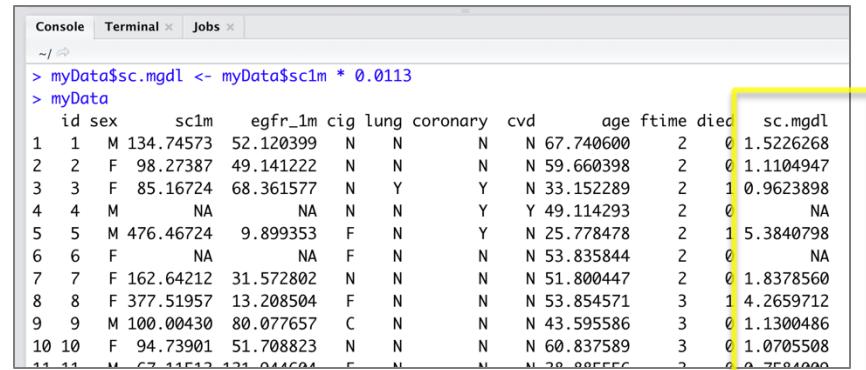
Data management

- How to compute new variables from existing ones?
- Serum creatinine seems to be measure in $\mu\text{mol/l}$.
- To convert it to mg/dl we need to multiply the original values by 0.0113

myData\$sc.mgdl <- myData\$sc1m * 0.0113

#Try in R:

```
myData$sc.mgdl <- myData$sc1m * 0.0113  
myData
```



The screenshot shows an R console window with three tabs: 'Console', 'Terminal', and 'Jobs'. The 'Console' tab is active and displays the following R code and output:

```
> myData$sc.mgdl <- myData$sc1m * 0.0113
> myData
```

	id	sex	sc1m	egfr_1m	cig	lung	coronary	cvd	age	ftime	died	sc.mgdl
1	1	M	134.74573	52.120399	N	N	N	N	67.740600	2	0	1.5226268
2	2	F	98.27387	49.141222	N	N	N	N	59.660398	2	0	1.1104947
3	3	F	85.16724	68.361577	N	Y	Y	N	33.152289	2	1	0.9623898
4	4	M	NA	NA	N	N	Y	Y	49.114293	2	0	NA
5	5	M	476.46724	9.899353	F	N	Y	N	25.778478	2	1	5.3840798
6	6	F	NA	NA	F	N	N	N	53.835844	2	0	NA
7	7	F	162.64212	31.572802	N	N	N	N	51.800447	2	0	1.8378560
8	8	F	377.51957	13.208504	F	N	N	N	53.854571	3	1	4.2659712
9	9	M	100.00430	80.077657	C	N	N	N	43.595586	3	0	1.1300486
10	10	F	94.73901	51.708823	N	N	N	N	60.837589	3	0	1.0705508
11	11	N	67.11517	171.044004	F	N	N	N	20.000000	2	0	0.7504000

Data management

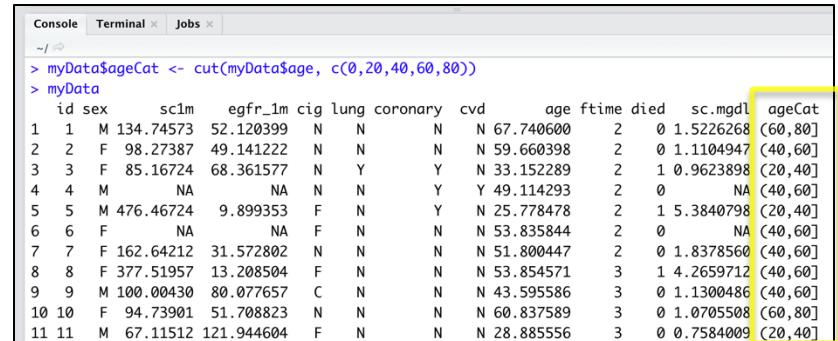
- Another common operation is to categorise continuous variables
- Let's create the following categories for age: <=20, (20-40], (40-60], >60
- We will use the function `cut()`

```
myData$ageCat <- cut(myData$age, c(0, 20, 40, 60, 80))
```

- **Exercise:** Get the frequencies for the age categories

```
#Try in R:
```

```
myData$ageCat <- cut(myData$age, c(0, 20, 40, 60, 80))  
myData
```



	id	sex	sc1m	egfr_1m	cig	lung	coronary	cvd	age	ftime	died	sc.mgdl	ageCat
1	1	M	134.74573	52.120399	N	N	N	N	67.740600	2	0	1.5226268	(60,80]
2	2	F	98.27387	49.141222	N	N	N	N	59.660398	2	0	1.1104947	(40,60]
3	3	F	85.16724	68.361577	N	Y	Y	N	33.152289	2	1	0.9623898	(20,40]
4	4	M	NA	NA	N	N	Y	Y	49.114293	2	0	NA	(40,60]
5	5	M	476.46724	9.899353	F	N	Y	N	25.778478	2	1	5.3840798	(20,40]
6	6	F	NA	NA	F	N	N	N	53.835844	2	0	NA	(40,60]
7	7	F	162.64212	31.572802	N	N	N	N	51.800447	2	0	1.8378560	(40,60]
8	8	F	377.51957	13.208504	F	N	N	N	53.854571	3	1	4.2659712	(40,60]
9	9	M	100.00430	80.077657	C	N	N	N	43.595586	3	0	1.1300486	(40,60]
10	10	F	94.73901	51.708823	N	N	N	N	60.837589	3	0	1.0705508	(60,80]
11	11	M	67.11512	121.944604	F	N	N	N	28.885556	3	0	0.7584009	(20,40]

Data management

- **Exercise:** Get the frequencies for the age categories

```
table(myData$ageCat)
```

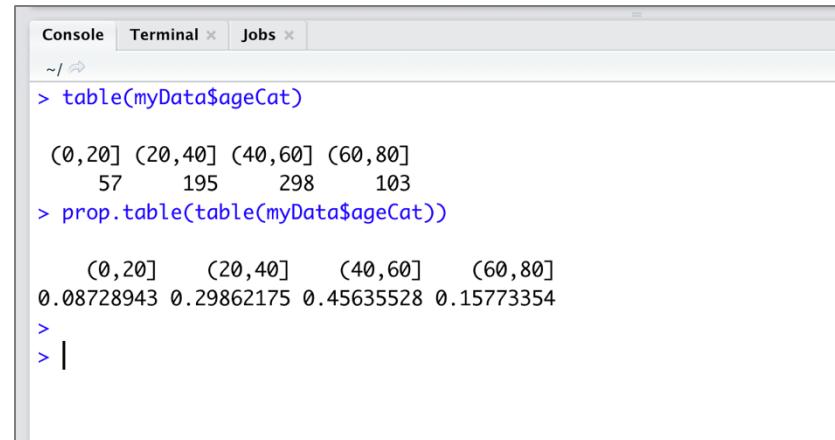
- There is no direct function for the percentages of each category. We need to do:

```
prop.table(table(myData$ageCat))
```

- A bit clunky...

```
#Try in R:
```

```
table(myData$ageCat)  
prop.table(table(myData$ageCat))
```



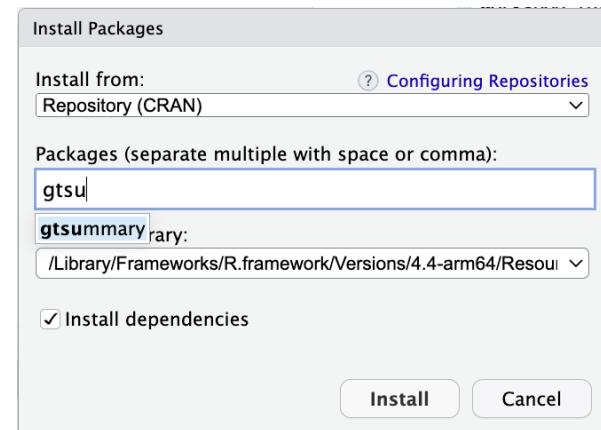
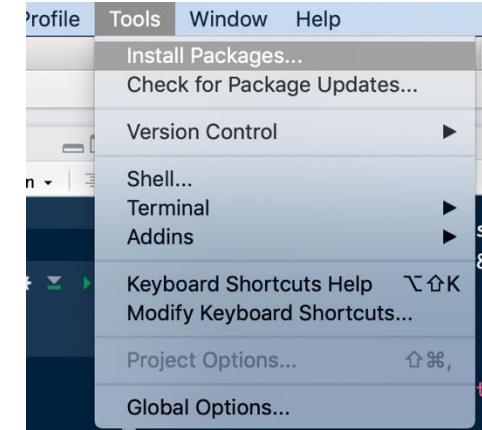
The screenshot shows an R console window with three tabs: 'Console' (selected), 'Terminal', and 'Jobs'. The console area contains the following R session:

```
Console Terminal × Jobs ×  
~ /  
> table(myData$ageCat)  
(0,20] (20,40] (40,60] (60,80]  
57      195     298     103  
> prop.table(table(myData$ageCat))  
(0,20] (20,40] (40,60] (60,80]  
0.08728943 0.29862175 0.45635528 0.15773354  
>  
> |
```

Libraries

- We are not restricted to the functions in the base version of R
- Having all community contributing with new functions is one of the main advantages of R
- Let's use the library `gtsummary`
- First we need to download it
(can also use `install.packages("gtsummary")`)
- And then **load** it in R:

`library(gtsummary)`
- You'll need to load it every time you open R



Libraries

- There are thousands of libraries available
- Each library usually have multiple functions
- The documentation is (most of the times) helpful
- So, if you type `?gtsummary`, you will have a brief description of the library and the main functions
- You can also use `?tbl_summary` for functions, e.g.,
`?tbl_summary`

```
#Try in R:  
#After downloading "gtsummary"  
  
library(gtsummary)  
?gtsummary  
  
# tbl_summary is one of the functions in the package  
?tbl_summary
```

gtsummary-package (gtsummary) R Documentation

gtsummary: Presentation-Ready Data Summary and Analytic Result Tables

Description

Creates presentation-ready tables summarizing data sets, regression models, and more. The code to create the tables is concise and highly customizable. Data frames can be summarized with any function, e.g. `mean()`, `median()`, even user-written functions. Regression models are summarized and include the reference rows for categorical variables. Common regression models, such as logistic regression and Cox proportional hazards regression, are automatically identified and the tables are pre-filled with appropriate column headers.

Author(s)

Maintainer: Daniel D. Sjoberg daniel.sjoberg@gmail.com ([ORCID](#))

Authors:

- Joseph Larmanange ([ORCID](#))
- Michael Curry ([ORCID](#))
- Jessica Lavery ([ORCID](#))
- Karissa Whiting ([ORCID](#))
- Emily C. Zabor ([ORCID](#))

Other contributors:



Libraries

- Let's used this function with the entire dataset

```
tbl_summary(data=myData)
```



When the function accepts many arguments (options), it is better to tell R what is the argument that you are passing to the function.

```
tbl_summary( data, by = NULL, label = NULL, statistic =  
list(all_continuous() ~ "{median} ({p25}, {p75})", all_categorical() ~  
"({n} ({p}%)", digits = NULL, type = NULL, value = NULL, missing =  
c("ifany", "no", "always"), missing_text = "Unknown", missing_stat =  
"{N_miss}", sort = all_categorical(FALSE) ~ "alphanumeric", percent =  
c("column", "row", "cell"), include = everything() )
```

```
#Try in R:
```

```
tbl_summary(data=myData)
```

Characteristic	N = 654 ¹
id	328 (164, 491)
sex	
F	255 (39%)
M	399 (61%)
sc1m	132 (101, 175)
Unknown	6
egfr_1m	51 (37, 63)
Unknown	6
cig	
C	68 (10%)
F	190 (29%)

Libraries

– And by survival status

```
tbl_summary(data=myData, by='died')
```

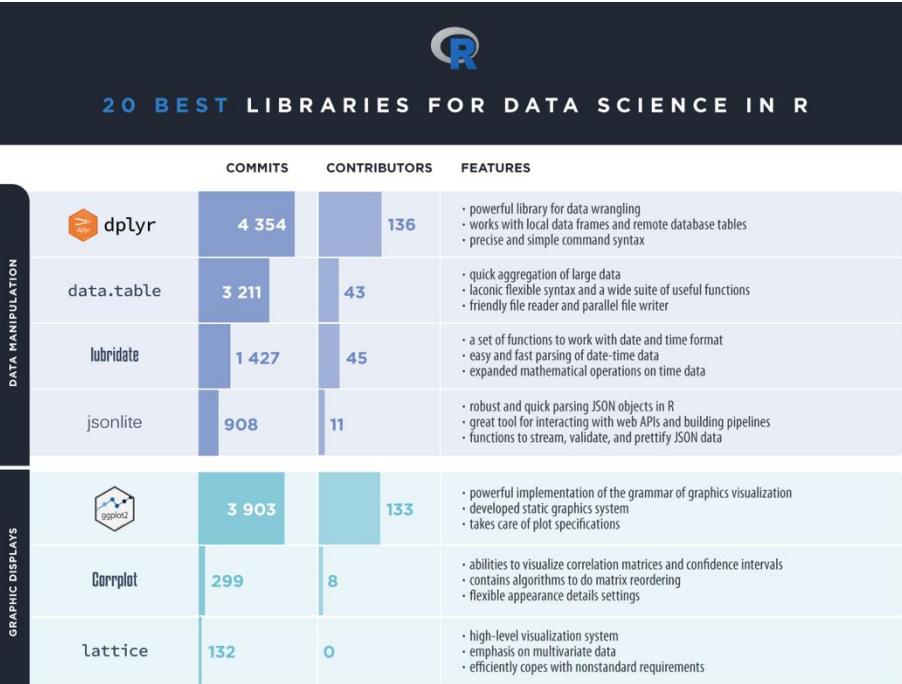
#Try in R:

```
tbl_summary(data=myData)
```

```
#now, stratified by survival status  
tbl_summary(data=myData, by='died')
```

Characteristic	0 N = 516 ¹	1 N = 138 ¹
id	329 (173, 507)	323 (147, 461)
sex		
F	193 (37%)	62 (45%)
M	323 (63%)	76 (55%)
sc1m	129 (101, 167)	149 (103, 205)
Unknown	6	0
egfr_1m	52 (39, 64)	44 (29, 60)
Unknown	6	0
cig		
C	45 (8.7%)	23 (17%)
F	153 (30%)	37 (27%)
N	318 (62%)	78 (57%)

Libraries



HTML WIDGETS	plotly	2 989	26	<ul style="list-style-type: none"> rich features and plenty of available charts web-based toolbox for building visualizations abilities to make ggplot2 graphics interactive
	ggvis	2 159	21	<ul style="list-style-type: none"> implementation of an interactive grammar of graphic incorporates shiny reactive programming model and dplyr grammar of data transformation
	DT DataTables	1 919	21	<ul style="list-style-type: none"> displays R matrices and data frames as interactive HTML tables creates sortable tables with a minimum of code many useful features and styling options for tables
	rCharts	638	11	<ul style="list-style-type: none"> interactive JS charts from R tools for creation, customization, and sharing
REPRODUCIBLE RESEARCH	kable	5 467	96	<ul style="list-style-type: none"> transparent tool for easy dynamic report generation in R enables integration of R code into LaTeX, LyX, HTML, Markdown, AsciiDoc, and reStructuredText documents
	markdown	2 297	56	<ul style="list-style-type: none"> next generation implementation of R Markdown based on pandoc many static and dynamic output formats abilities to define new formats for custom publishing requirements
	slidify	302	7	<ul style="list-style-type: none"> generates reproducible HTML5 slides from r markdown allows embedded code chunks and mathematical formulas rich sharing and customizing opportunities
	mlr	3 915	55	<ul style="list-style-type: none"> extensible framework for classification, regression, survival analysis, and clustering easy extension mechanism through S3 inheritance
MACHINE LEARNING	XGBoost	3 188	259	<ul style="list-style-type: none"> implementation of the Gradient Boosted Decision Trees algorithm reach tools for regression, classification, and ranking problems high speed and performance
	caret	1 659	59	<ul style="list-style-type: none"> many models for classification and regression powerful tools and algorithms for creating predictive models
	gbm	731	26	<ul style="list-style-type: none"> represents Generalized Boosted Regression Models includes plenty of regression methods tools variable selection and final stage precision modeling
	Prophet	190	20	<ul style="list-style-type: none"> high-quality forecasts for time series data manages data that has multiple seasonality with linear or non-linear growth robust to missing data, shifts in the trend, and large outliers
randomForest	randomForest	56	0	<ul style="list-style-type: none"> implements Breiman's random forest algorithm for classification and regression builds multiple decision trees and gives back the mean prediction of the individual trees

Updated: December 2017

Created by ActiveWizards

Libraries

The screenshot shows a blog post titled "Quick list of useful R packages" by Garrett Grolemund. The post includes code snippets for installing and loading packages, and a sidebar with related articles.

Quick list of useful R packages

Garrett Grolemund
Today at 03:54

Recommended Packages

Many useful R functions come in packages, free libraries of code written by R's active user community. To install an R package, open an R session and type at the command line

```
install.packages("<the package's name>")
```

R will download the package from CRAN, so you'll need to be connected to the internet. Once you have a package installed, you can make its contents available to use in your current R session by running

```
library("<the package's name>")
```

Premium email support Sign in

Related articles

- Importing Data with RStudio
- Working Directories and Workspaces
- Getting Started with R
- How to run R scripts from the command line
- Installing older versions of packages

<https://support.rstudio.com/hc/en-us/articles/201057987-Quick-list-of-useful-R-packages>

Example of a statistical analysis

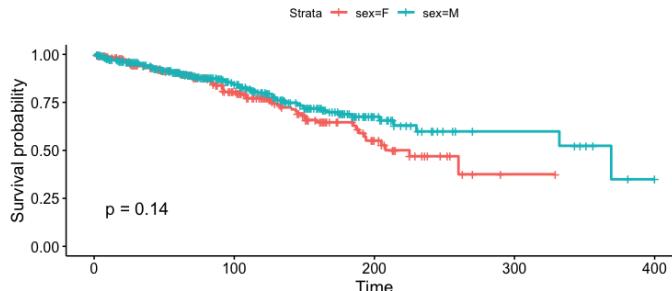
- Let's look at differences on survival time after transplant between sexes
- We could start by looking at the Kaplan-Meier curve
- Perform a log-rank test
- And then move to a multivariable Cox regression

```
#####Survival analysis#####
library(survival)
library(survminer)

#Kaplan meier and log-rank test comparing sexes
KMsex <- survfit(Surv(ftime,died) ~ sex,
                   data = myData)
ggsurvplot(KMsex, pval=TRUE)

#Cox model
myCoxModel <- coxph(Surv(ftime,died) ~ sex + age + egfr_lm,
                      data = myData)
summary(myCoxModel)
ggforest(myCoxModel, data = myData)

#checking the PH assumption
test.ph <- cox.zph(myCoxModel)
test.ph
plot(test.ph)
```



Links for learning (more) about R

Edx.org (online courses)

- Statistics and R (Harvard School of Public Health) (<https://www.edx.org/course/statistics-and-r-3>)

datacamp.com (courses and tutorials for data science)

- Survival Analysis in R For Beginners (<https://www.datacamp.com/community/tutorials/survival-analysis-R>)
- Logistic Regression in R Tutorial (<https://www.datacamp.com/community/tutorials/logistic-regression-R>)

Books

- An Introduction to Biostatistics Using R - Waveland Press (<http://waveland.com/Glover-Mitchell/r-guide.pdf>)

News

- <https://www.r-bloggers.com>
- <https://awesome-r.com>
- <https://www.statmethods.net> (Quick-R)