

Introduction to R

Department of Epidemiology and Data Science

8 - 12 June 2026

Outline

Importing, saving and managing data

Data manipulation

Programming structures in R

Graphics

Extra Documentation and help

Workspace management

- R uses **working directories**
- Working directory is a file path setting the location of any files you will use
- You can only have one active working directory at the time
- `getwd()` shows what your working directory is
- `setwd()` allows you to define your working directory
 - > `setwd("C:\\Users\\P073787\\Documents\\")`
 - > `getwd()`

Exporting Data

- Export to text file: `write.table()`
- Export to SPSS, Stata and SAS: package **haven**
(`write_sav()`, `write_dta()`, `write_sas()`)
- Export to xlsx: package **writexl** (`write_xlsx()`)

EXERCISES

Let's do Questions 7, 8, 9

Outline

Importing, saving and managing data

Data manipulation

Programming structures in R

Graphics

Extra Documentation and help

Sorting data (I)

Three functions to help you sort data:

- `sort()`: sorts a vector in an ascending order (default) or descending order (set `decreasing = TRUE`)
- `rank()`: gives the respective rank of the numbers present in the vector, the smallest number receiving the rank 1.
- `order()`: returns the indices of the vector in a sorted order (ascending or descending)

```
> x <- c(3, 1, 2, 5, 4)
> sort(x)
```

```
[1] 1 2 3 4 5
```

```
> rank(x)
```

```
[1] 3 1 2 5 4
```

```
> order(x)
```


Sorting data (II)

To sort the rows of a data frame according to the values of a specific column, use `order()`

```
> titanic_sorted <- titanic4[order(titanic4$age), ]
> head(titanic_sorted[, c(3, 5, 9)])
```

| | | name | age | fare |
|------|---------------------------------|------|--------|---------|
| 764 | Dean, Miss. Elizabeth Gladys | \"M | 0.1667 | 20.5750 |
| 748 | Danbom, Master. Gilbert Sigvard | | 0.3333 | 14.4000 |
| 1241 | Thomas, Master. Assad Alexander | | 0.4167 | 8.5167 |
| 428 | Hamalainen, Master. Viljo | | 0.6667 | 14.5000 |
| 658 | Baclini, Miss. Eugenie | | 0.7500 | 19.2583 |
| 659 | Baclini, Miss. Helene Barbara | | 0.7500 | 19.2583 |

Summarizing data

- Data summary: `summary()`, `rowMeans()`, `colMeans()`
- Contingency tables: `table()`, `xtabs()`, `CrossTable` from **descr** package
- Summary by subgroups:
 - `aggregate()` and `tapply()`, `sapply()`, or `lapply()`
 - **doBy**, **Hmisc**, **compareGroups**, **dplyr**
- Graphical summary of data frames: `dfSummary()` in **summarytools**

EXERCISES

Let's do questions 10-16

Outline

Importing, saving and managing data

Data manipulation

Programming structures in R

Graphics

Extra Documentation and help

If - else (I)

- R has a **conditional construct**: `if(){} else(){}`
- Depending on the outcome of a test, execute one or another statement

```
if(logical statement){  
  do this  
} else {  
  do that  
}
```

- Think carefully when making this construct. Using too many will make your code slow!
- Multiple if - else constructs can be nested within each other

For - loop (III)

```
> results <- rep(0, 2)
> results

[1] 0 0

> for (i in 1:length(results)) {
+ results[i] <- mean(A[i, ])
+ }
> results

[1] 5 6
```

Apply (I)

Functions from the apply family are convenient shorthands for for loops

```
apply(X, MARGIN, FUN, ...)
```

Arguments

X an array, including a matrix

MARGIN for a matrix 1 indicates rows, 2 indicates columns

FUN the function to be applied

Apply (II)

```
> apply(A, 1, mean)
```

```
gene1 gene2
     5     6
```

Apply (III)

- `lapply()`: apply a function over a list or vector
- `sapply()`: similar to `apply` but more user-friendly if output can be coerced into a vector
- `tapply()`: can be used to split a vector in subgroups and apply a function for each subgroup
- `replicate()`: simpler version of `sapply` for the repeated evaluation of an expression
- `aggregate()`: extension of `tapply` for data frames

EXERCISES

Let's do Questions 17, 18, and 19

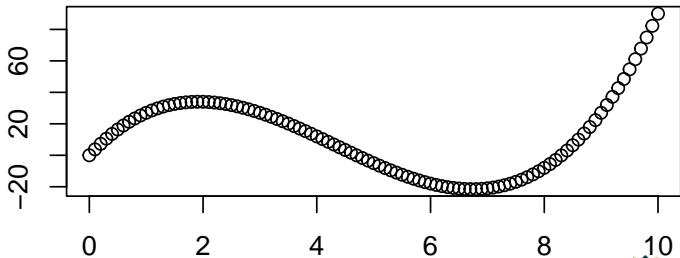
Basic plot (I)

It is straightforward to make a simple plot using functions from the **graphics** package (loaded by default):

```
> x <- (0:100) / 10
```

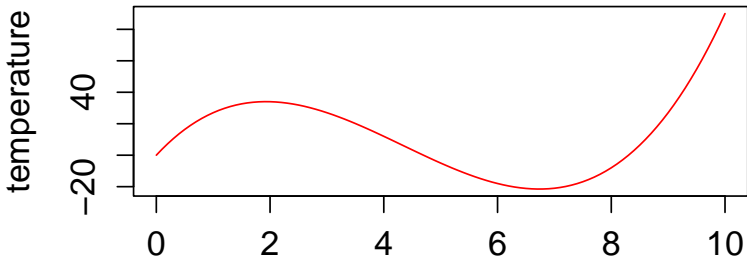
```
> plot(x, x^3 - 13 * x^2 + 39 * x)
```

$x^3 - 13 * x^2 + 39 * x$



Basic plot (IV)

Enhanced plot

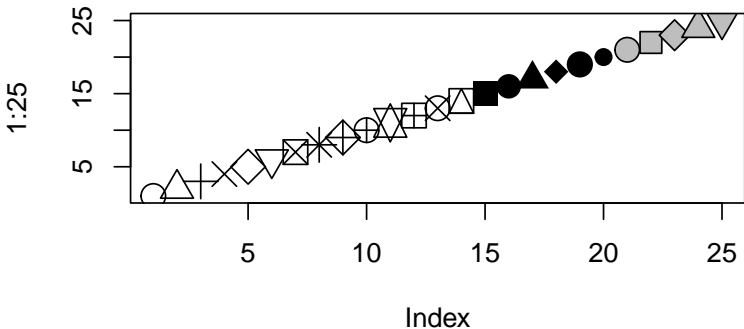


time (hours)
Find the differences

Basic plot (V)

There are 25 different plot symbols, see ?points

```
> plot(1:25, pch = 1:25, cex = 2, bg = "grey")
```



Basic plot (VI)

The most commonly used options of `par`:

- `lwd` sets the line width
- `mfrow` and `mfcol` enable multiple plots in one figure
- `las` rotates axis symbols
- `mar` changes the margins of the figure
- `bg` changes the background colour

For more options see `?par`

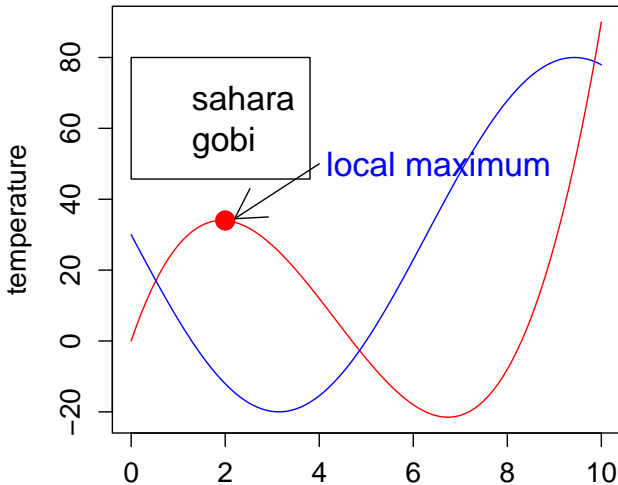
Enhancing a plot (I)

With these functions, you can add extra elements to your plots:

- `points()`: add points
- `lines()`: add line
- `abline()`: add horizontal or vertical lines
- `arrows()`: add arrows
- `curve()`: add a curved line
- `rect()`: add a rectangle
- `text()`: add text
- `legend()`: add a legend
- `axis()`: add an axis

Enhancing a plot (II)

Enhanced plot



Histogram (I)

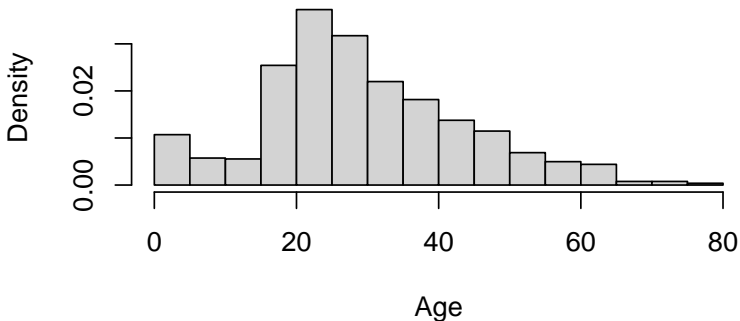
Use `hist()` for plotting histograms.

As always, see `?hist` for the many arguments of this function

```
> hist(titanic3$age,  
+ breaks = 15, freq = FALSE,  
+ main = "Histogram",  
+ xlab = "Age"  
+ )
```

Histogram (II)

Histogram



Boxplot (I)

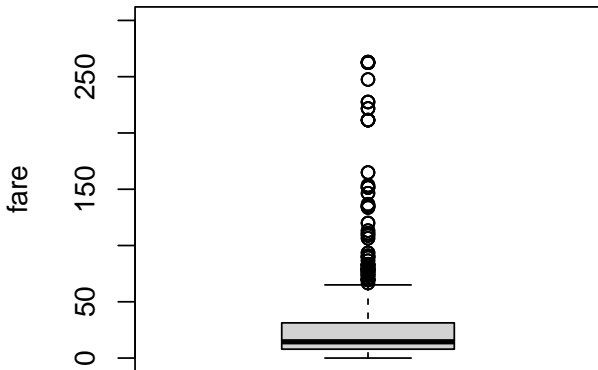
- Use `boxplot()` for plotting boxplots.
- For one group (a single vector):

```
> boxplot(titanic3$fare,  
+ ylim = c(0, 300), ylab = "fare"  
+ )
```

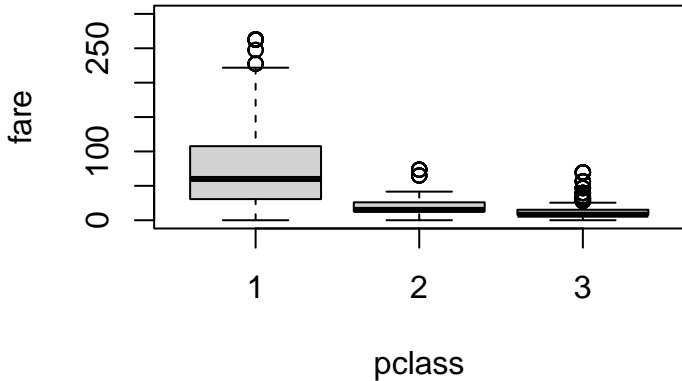
- For multiple groups:

```
> boxplot(fare ~ pclass,  
+ data = titanic3, ylim = c(0, 300), ylab = "fare"  
+ )
```

Boxplot (II)



Boxplot (III)



Barplot (I)

- Use `barplot()` for plotting barplots
- See `?barplot` for the many arguments of this function

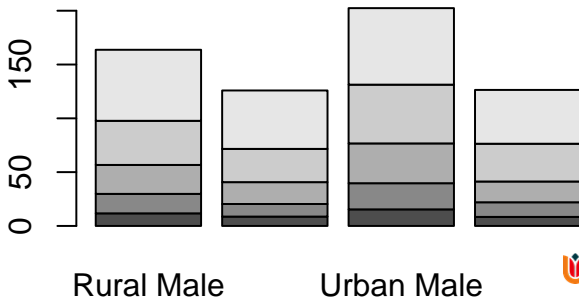
```
> head(VADeaths)
```

| | Rural Male | Rural Female | Urban Male | Urban Female |
|-------|------------|--------------|------------|--------------|
| 50-54 | 11.7 | 8.7 | 15.4 | 8.4 |
| 55-59 | 18.1 | 11.7 | 24.3 | 13.6 |
| 60-64 | 26.9 | 20.3 | 37.0 | 19.3 |
| 65-69 | 41.0 | 30.9 | 54.6 | 35.1 |
| 70-74 | 66.0 | 54.3 | 71.1 | 50.0 |

Barplot (II)

- Use `barplot()` for plotting barplots
- See `?barplot` for the many arguments of this function

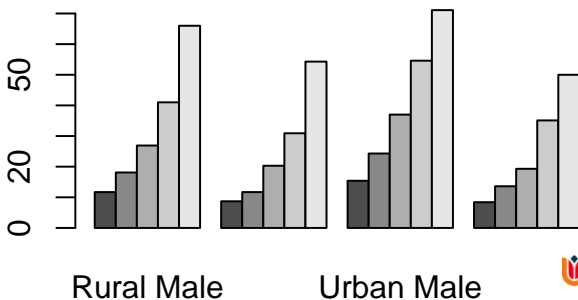
```
> barplot(VADeaths)
```



Barplot (III)

- Use `barplot()` for plotting barplots
- See `?barplot` for the many arguments of this function

```
> barplot(VADeaths, beside = TRUE)
```



Advanced R graphics

- Chapter 12 of “An Introduction to R” gives an introduction to **base** graphics
- **lattice**: very powerful for multipanel conditioning needs to be loaded first; `xypLOT()` is the main function
- **ggplot2**: based on “the grammar of graphics” (see <https://rstudio.github.io/cheatsheets/html/data-visualization.html>)
- **ggvis**, **plotly**, **rCharts**, **Shiny**: interactive visualizations
- and in many more packages (**gplots**, **plotrix**, ...)

Exporting figures

Two types of figure formats:

- Vector format (*pdf, eps, wmf, emf*)
 - digital image consisting of independent geometric objects (segments, polygons, curves, etc.)
 - can be enlarged without losing resolution
- Raster (*png, jpeg, tiff*)
 - rectangular grid of pixels, possibly with color
 - resolution impaired if image is enlarged

Graphics can be saved via the menu in the graphics/plots window, or a specific graphics file type can be created directly (`pdf()`, `win.metafile()`, `png()`) followed by `dev.off()`

Outline

Importing, saving and managing data

Data manipulation

Programming structures in R

Graphics

Extra Documentation and help

Finding information

- `help()`
- `help.search()`: search the help system for documentation matching your search term
- `RSiteSearch()`: search CRAN (online) for documentation matching your search term
- package **sos**
- packages have sometimes *vignettes*, which introduce the package and their functions
- Google / Claude / ...
- [https://bioinformaticslaboratory.eu/gs-computing-in-r:](https://bioinformaticslaboratory.eu/gs-computing-in-r) Section Information on R

