



Modelización Estadística: ¿Qué hacemos con una base de datos?

Introducción a R Estadística Descriptiva

Mathias Bourel
mathias.bourel@fceia.edu.uy

Diciembre 2025



FACULTAD DE
CIENCIAS ECONÓMICAS
Y DE ADMINISTRACIÓN

IESTA INSTITUTO
DE ESTADÍSTICA



UNIVERSIDAD
DE LA REPÚBLICA
URUGUAY

IRL-2030 IFUMI



Plan

- 1 ¿Qué es R?
- 2 Distribuciones
- 3 Variable cuantitativas
- 4 Variables cualitativas

¿Qué es R?

¿Qué es R?



- R es uno de los software estadístico más potente y más ampliamente utilizado.

¿Qué es R?



- R es uno de los software estadístico más potente y más ampliamente utilizado.
- Es un lenguaje y un entorno para el cálculo estadístico y la visualización gráfica.



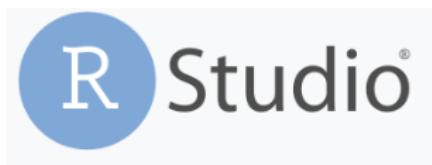
- R es uno de los software estadístico más potente y más ampliamente utilizado.
- Es un lenguaje y un entorno para el cálculo estadístico y la visualización gráfica.
- Es un conjunto integrado de herramientas para la manipulación de datos, el cálculo y la representación gráfica.



- R es uno de los software estadístico más potente y más ampliamente utilizado.
- Es un lenguaje y un entorno para el cálculo estadístico y la visualización gráfica.
- Es un conjunto integrado de herramientas para la manipulación de datos, el cálculo y la representación gráfica.
- Proporciona una amplia variedad de técnicas estadísticas clásicas y modernas.
- Es un proyecto colaborativo, donde todos los usuarios pueden subir sus propios códigos y paquetes.



- R es uno de los software estadístico más potente y más ampliamente utilizado.
- Es un lenguaje y un entorno para el cálculo estadístico y la visualización gráfica.
- Es un conjunto integrado de herramientas para la manipulación de datos, el cálculo y la representación gráfica.
- Proporciona una amplia variedad de técnicas estadísticas clásicas y modernas.
- Es un proyecto colaborativo, donde todos los usuarios pueden subir sus propios códigos y paquetes.
- Es **gratuito y libre**.



- R es uno de los software estadístico más potente y más ampliamente utilizado.
- Es un lenguaje y un entorno para el cálculo estadístico y la visualización gráfica.
- Es un conjunto integrado de herramientas para la manipulación de datos, el cálculo y la representación gráfica.
- Proporciona una amplia variedad de técnicas estadísticas clásicas y modernas.
- Es un proyecto colaborativo, donde todos los usuarios pueden subir sus propios códigos y paquetes.
- Es **gratuito y libre**.
- Funciona en Windows, Mac y Linux.



- R es uno de los software estadístico más potente y más ampliamente utilizado.
- Es un lenguaje y un entorno para el cálculo estadístico y la visualización gráfica.
- Es un conjunto integrado de herramientas para la manipulación de datos, el cálculo y la representación gráfica.
- Proporciona una amplia variedad de técnicas estadísticas clásicas y modernas.
- Es un proyecto colaborativo, donde todos los usuarios pueden subir sus propios códigos y paquetes.
- Es **gratuito y libre**.
- Funciona en Windows, Mac y Linux.
- Dispone de numerosas herramientas para la visualización de datos.



- R es uno de los software estadístico más potente y más ampliamente utilizado.
- Es un lenguaje y un entorno para el cálculo estadístico y la visualización gráfica.
- Es un conjunto integrado de herramientas para la manipulación de datos, el cálculo y la representación gráfica.
- Proporciona una amplia variedad de técnicas estadísticas clásicas y modernas.
- Es un proyecto colaborativo, donde todos los usuarios pueden subir sus propios códigos y paquetes.
- Es **gratuito y libre**.
- Funciona en Windows, Mac y Linux.
- Dispone de numerosas herramientas para la visualización de datos.

Se puede descargar RStudio que tiene una interfase mucho más amigable en
<https://posit.co/downloads/> o usar en la nube <https://posit.cloud>



- R es uno de los software estadístico más potente y más ampliamente utilizado.
- Es un lenguaje y un entorno para el cálculo estadístico y la visualización gráfica.
- Es un conjunto integrado de herramientas para la manipulación de datos, el cálculo y la representación gráfica.
- Proporciona una amplia variedad de técnicas estadísticas clásicas y modernas.
- Es un proyecto colaborativo, donde todos los usuarios pueden subir sus propios códigos y paquetes.
- Es **gratuito y libre**.
- Funciona en Windows, Mac y Linux.
- Dispone de numerosas herramientas para la visualización de datos.

Se puede descargar RStudio que tiene una interfase mucho más amigable en
<https://posit.co/downloads/> o usar en la nube <https://posit.cloud>

Existen muchos tutoriales disponibles: Introduction to R, R RefCard, etc.



Paquetes y librerías

R se basa en paquetes y librerías donde hay funciones, conjuntos de datos, documentaciones,...

- La lista de las librerías disponibles

```
library()
```

Paquetes y librerías

R se basa en paquetes y librerías donde hay funciones, conjuntos de datos, documentaciones,...

- La lista de las librerías disponibles

```
library()
```

- Descargar e instalar un paquete desde el CRAN (Comprehensive R Archive Network)

```
install.packages (...)
```

Paquetes y librerías

R se basa en paquetes y librerías donde hay funciones, conjuntos de datos, documentaciones,...

- La lista de las librerías disponibles

```
library()
```

- Descargar e instalar un paquete desde el CRAN (Comprehensive R Archive Network)

```
install.packages (...)
```

- Cargar paquetes

```
library(MASS) # funciones y bases del libro
```

"Modern Applied Statistics with S" para hacer análisis discriminante, regresiones, multidimensional

```
library(rpart) #para hacer arboles de clasificación y regresión
```

```
library(e1071) #para Support Vector Machines
```

```
library(ks) #estimación de densidad por núcleo
```

Paquetes y librerías

R se basa en paquetes y librerías donde hay funciones, conjuntos de datos, documentaciones,...

- La lista de las librerías disponibles

```
library()
```

- Descargar e instalar un paquete desde el CRAN (Comprehensive R Archive Network)

```
install.packages (...)
```

- Cargar paquetes

```
library(MASS) # funciones y bases del libro  
"Modern Applied Statistics with S" para hacer análisis discriminante,  
regresiones, multidimensional  
library(rpart) #para hacer arboles de clasificación y regresión  
library(e1071) #para Support Vector Machines  
library(ks) #estimación de densidad por núcleo
```

- Documentación sobre R

```
help.start()
```

Paquetes y librerías

R se basa en paquetes y librerías donde hay funciones, conjuntos de datos, documentaciones,...

- La lista de las librerías disponibles

```
library()
```

- Descargar e instalar un paquete desde el CRAN (Comprehensive R Archive Network)

```
install.packages (...)
```

- Cargar paquetes

```
library(MASS) # funciones y bases del libro
```

"Modern Applied Statistics with S" para hacer análisis discriminante, regresiones, multidimensional

```
library(rpart) #para hacer arboles de clasificación y regresión
```

```
library(e1071) #para Support Vector Machines
```

```
library(ks) #estimación de densidad por núcleo
```

- Documentación sobre R

```
help.start()
```

- Documentación sobre librería

```
library(help = rpart)
```

Importación y exportación de bases de datos

- Ayuda sobre funciones:

```
help(mvrnorm) #ayuda para la función que genera datos normales multivariados  
?lm #ayuda sobre la función lm para modelos lineales
```

Importación y exportación de bases de datos

- Ayuda sobre funciones:

```
help(mvrnorm) #ayuda para la función que genera datos normales multivariados  
?lm #ayuda sobre la función lm para modelos lineales
```

- Cargar bases de datos:

- Algunas están en R en las paquetes bases, por ejemplo iris, airquality, cars
- Otras hay que cargarlas con por ejemplo estas líneas de código:

```
datos1=read.table("miejemplo.txt", sep= ",", header=T)  
datos2=read.csv("miejemplo2.csv", sep= ",", header=T)  
datos3=read.csv("/home/mathiasbourel/Escritorio/Investigacion/DATA/  
Historica2009_2019.csv",header=T, sep=";")  
base=read.csv("Desgargas/spam.csv",sep=",",header=T)
```

Importación y exportación de bases de datos

- Ayuda sobre funciones:

```
help(mvrnorm) #ayuda para la función que genera datos normales multivariados  
?lm #ayuda sobre la función lm para modelos lineales
```

- Cargar bases de datos:

- Algunas están en R en las paquetes bases, por ejemplo iris, airquality, cars
- Otras hay que cargarlas con por ejemplo estas líneas de código:

```
datos1=read.table("miejemplo.txt", sep= ",", header=T)  
datos2=read.csv("miejemplo2.csv", sep= ",", header=T)  
datos3=read.csv("/home/mathiasbourel/Escritorio/Investigacion/DATA/  
Historica2009_2019.csv",header=T, sep=";")  
base=read.csv("Desgargas/spam.csv",sep=",",header=T)
```

- Escribir tablas y bases de datos

```
write.table()      write.csv()
```

Importación y exportación de bases de datos

- Ayuda sobre funciones:

```
help(mvrnorm) #ayuda para la función que genera datos normales multivariados  
?lm #ayuda sobre la función lm para modelos lineales
```

- Cargar bases de datos:

- Algunas están en R en las paquetes bases, por ejemplo iris, airquality, cars
- Otras hay que cargarlas con por ejemplo estas líneas de código:

```
datos1=read.table("miejemplo.txt", sep= ",", header=T)  
datos2=read.csv("miejemplo2.csv", sep= ",", header=T)  
datos3=read.csv("/home/mathiasbourel/Escritorio/Investigacion/DATA/  
Historica2009_2019.csv",header=T, sep=";")  
base=read.csv("Desgargas/spam.csv",sep=",",header=T)
```

- Escribir tablas y bases de datos

```
write.table()      write.csv()
```

- Ejemplo completo

```
miejemplo <- data.frame(cbind(matrix(round(runif(24)*100,0),6,4)),  
CAT=rep(c('A','B'),3))  
attr(miejemplo,'names') <- paste('Var',1:5,sep='')  
row.names(miejemplo) <- c('I1','I2','I3','I4','I5','I6')  
miejemplo <- cbind(miejemplo,suma13=colSums(miejemplo[,1:3]),  
logica2=miejemplo[,2]>50)  
ej2 <- miejemplo[-1,-ncol(miejemplo)]; summary(ej2)  
write.table(miejemplo,file='miejemplo.txt',sep='\t')
```

Las Posit Cheatsheets son guías rápidas (hojas resumen) creadas por Posit (antes RStudio) que sirven como material de referencia práctico para trabajar con R, RStudio y el ecosistema tidyverse, entre otras herramientas.

Las Posit Cheatsheets son guías rápidas (hojas resumen) creadas por Posit (antes RStudio) que sirven como material de referencia práctico para trabajar con R, RStudio y el ecosistema tidyverse, entre otras herramientas. ¿Para qué sirven?

- Recordar funciones clave y su sintaxis sin tener que ir a la documentación completa.
- Ver de un vistazo los flujos de trabajo más comunes.
- Aprender buenas prácticas de uso.
- Acompañar clases, talleres o trabajo cotidiano con R.

Las Posit Cheatsheets son guías rápidas (hojas resumen) creadas por Posit (antes RStudio) que sirven como material de referencia práctico para trabajar con R, RStudio y el ecosistema tidyverse, entre otras herramientas. ¿Para qué sirven?

- Recordar funciones clave y su sintaxis sin tener que ir a la documentación completa.
- Ver de un vistazo los flujos de trabajo más comunes.
- Aprender buenas prácticas de uso.
- Acompañar clases, talleres o trabajo cotidiano con R.

Cada cheatsheet suele incluir:

- Funciones más usadas (con ejemplos breves).
- Esquemas visuales (pipes, gramática de gráficos, manipulación de datos).
- Tablas resumen de argumentos importantes.
- Recordatorios de atajos y conceptos clave.

Las Posit Cheatsheets son guías rápidas (hojas resumen) creadas por Posit (antes RStudio) que sirven como material de referencia práctico para trabajar con R, RStudio y el ecosistema tidyverse, entre otras herramientas. ¿Para qué sirven?

- Recordar funciones clave y su sintaxis sin tener que ir a la documentación completa.
- Ver de un vistazo los flujos de trabajo más comunes.
- Aprender buenas prácticas de uso.
- Acompañar clases, talleres o trabajo cotidiano con R.

Cada cheatsheet suele incluir:

- Funciones más usadas (con ejemplos breves).
- Esquemas visuales (pipes, gramática de gráficos, manipulación de datos).
- Tablas resumen de argumentos importantes.
- Recordatorios de atajos y conceptos clave.

Se consiguen en

<https://rstudio.github.io/cheatsheets/contributed-cheatsheets.html> y van algunas a continuación:

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 Packages

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			Programming			While Loop		
Creating Vectors			For Loop			Example		
c(2, 4, 6)	2 4 6	Join elements into a vector	for (variable in sequence){ Do something }			while (condition){ Do something }		
2:6	2 3 4 5 6	An integer sequence	for (i in 1:4){ j <- i + 10 print(j) }			while (i < 5){ print(i) i <- i + 1 }		
seq(2, 3, by=0.5)	2.0 2.5 3.0	A complex sequence	Vector Functions			Example		
rep(1:2, times=3)	1 2 1 2 1 2	Repeat a vector	sort(x)			If Statements		
rep(1:2, each=3)	1 1 1 2 2 2	Repeat elements of a vector	rev(x)			if (condition){ Do something } else { Do something different }		
			unique(x)			Example		
			See counts of values.			functions		
			Selecting Vector Elements			function_name <- function(var){ Do something return(new_variable) }		
			By Position			Example		
			x[4]			Example		
			The fourth element.			square <- function(x){ squared <- x*x return(squared) }		
			x[-4]			Reading and Writing Data		
			All but the fourth.			Also see the <code>readr</code> package.		
			x[2:4]			Input		
			Elements two to four.			Output		
			x[!(2:4)]			Description		
			All elements except two to four.			df <- read.table('file.txt')		
			x[c(1, 5)]			Read and write a delimited text file.		
			Elements one and five.			df <- read.csv('file.csv')		
			By Value			Read and write a comma separated value file. This is a special case of read.table/write.table.		
			x[x == 10]			load('file.RData')		
			Elements which are equal to 10.			save(df, file = 'file.Rdata')		
			x[x < 0]			Conditions		
			All elements less than zero.			a == b		
			x[x %%in% c(1, 2, 5)]			Are equal		
			Elements in the set 1, 2, 5.			a > b		
			Named Vectors			Greater than		
			x['apple']			a >= b		
			Element with name 'apple'.			a < b		
						Less than		
						a <= b		
						Greater than or equal to		
						a != b		
						Less than or equal to		
						is.na(a)		
						is.null(a)		
						is.null(b)		

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.

Matrices

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

<code>m[2,]</code>	Select a row
<code>m[, 1]</code>	Select a column
<code>m[2, 3]</code>	Select an element

`t(m)`

Transpose

`m %*% n`

Matrix Multiplication

`solve(m, n)`

Find x in: $m^{-1} \cdot x = n$

Lists

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

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

Also see the [dplyr package](#).

Data Frames

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

Understanding a data frame

`View(df)`

See the full data frame.

`head(df)`

See the first 6 rows.

Matrix subsetting

<code>df[, 2]</code>	<code>nrow(df)</code>	<code>cbind</code> - Bind columns.
<code>df[2,]</code>	<code>ncol(df)</code>	<code>rbind</code> - Bind rows.
<code>df[2, 2]</code>	<code>dim(df)</code>	<code>bind_rows</code> - Bind rows.

Strings

Also see the [string package](#).

<code>paste(x, y, sep = ' ')</code>	Join multiple vectors together.
<code>paste(x, collapse = ' ')</code>	Join elements of a vector together.
<code>grep(pattern, x)</code>	Find regular expression matches in x.
<code>gsub(pattern, replace, x)</code>	Replace matches in x with a string.
<code>toupper(x)</code>	Convert to uppercase.
<code>tolower(x)</code>	Convert to lowercase.
<code>nchar(x)</code>	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 by 'cutting' into sections.

Statistics

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

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

`summary`
Get more detailed information out a model.

`t.test(x, y)`
Perform a t-test for difference between means.

`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	<code>rnorm</code>	<code>dnorm</code>	<code>pnorm</code>	<code>qnorm</code>
Poisson	<code>rpois</code>	<code>dpois</code>	<code>ppois</code>	<code>qpois</code>
Binomial	<code>rbinom</code>	<code>dbinom</code>	<code>pbinom</code>	<code>qbinom</code>
Uniform	<code>runif</code>	<code>unif</code>	<code>runif</code>	<code>qunif</code>

Plotting

Also see the [ggplot2 package](#).

<code>plot(x)</code>	Values of x in order.
<code>plot(x, y)</code>	Values of x against y.
<code>hist(x)</code>	Histogram of x.

Dates

See the [lubridate package](#).

Using Git and GitHub with RStudio: : CHEATSHEET



Version control control, also known as **source control**, is the practice of tracking and managing changes to software code.

Version control systems are software tools that help software teams manage changes to source code over time.

Git is an **open-source** software for version control, originally developed in 2005 by Linus Torvalds, the creator of the Linux operating system kernel.

Git is a version control tool to track the changes in the source code of a project.

Github is the most popular hosting service for collaborating on code using Git.

Requirements

1. R and RStudio installed
2. Git installed
3. Register a free GitHub account



Check that Git is installed

In the Terminal of RStudio, enter `which git` to request the path to your Git executable:

```
which git  
## /usr/bin/git
```

and `git --version` to see its version:

```
git --version  
## git version 2.34.1
```

Introduce yourself to Git

Open a shell from RStudio Tools > Shell and type each line separately by substituting your name and the email associated with your GitHub account:

```
git config --global user.name 'Jane Doe'  
git config --global user.email  
'jane@example.com'
```

Github Glossary

This [glossary](#) introduces common Git and GitHub terminology.

Basics

<code>git init <directory></code>	Create empty Git repository in specified directory.
<code>git clone <repository></code>	Clone a repository located at <code><repository></code> on your local machine.
<code>git config user.name <username></code>	Define author name to be used for all commits in current repository.
<code>git add <directory></code>	Stage all changes in <code><directory></code> for the next commit.
<code>git commit -m <"message"></code>	Commit the staged snapshot, but instead of launching a text editor, use <code><"message"></code> as the commit message.
<code>git status</code>	List which files are staged, unstaged, and untracked.
<code>git log</code>	Display the entire commit history using the default format.
<code>git diff</code>	Show unstaged changes between your index and working directory.

Remote Repositories

<code>git remote add <name> <url></code>	Create a new connection to a remote repository. After adding a remote, you can use <code><name></code> as a shortcut for <code><url></code> in other commands.
<code>git fetch <remote> <branch></code>	Fetches a specific <code><branch></code> , from the repository. Leave off <code><branch></code> to fetch all remote refs.
<code>git pull <remote></code>	Fetch the specified remote's copy of current branch and immediately merge it into the local copy.
<code>git push <remote> <branch></code>	Push the branch to <code><remote></code> , along with necessary commits and objects. Creates named branch in the remote repository if it doesn't exist.

Undoing Changes

<code>git revert <commit></code>	Create new commit that undoes all of the changes made in <code><commit></code> , then apply it to the current branch.
<code>git reset <file></code>	Remove <code><file></code> from the staging area but leave the working directory unchanged. This unstages a file without overwriting any changes.
<code>git clean -n</code>	Shows which files would be removed from working directory. Use the <code>-f</code> flag in place of the <code>-n</code> flag to execute the clean.

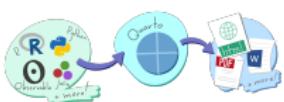
Rewriting Git History

<code>git commit --amend</code>	Replace the last commit with the staged changes and last commit combined. Use with nothing staged to edit the last commit's message.
<code>git rebase <base></code>	Rebase the current branch onto <code><base></code> . <code><base></code> can be a commit ID, branch name, a tag, or a relative reference to HEAD.
<code>git reflog</code>	Show a log of changes to the local repository's HEAD. Add <code>--relative-date</code> flag to show date info or <code>--all</code> to show all refs.

Git Branches

<code>git branch</code>	List all of the branches in your repo. Add a <code><branch></code> argument to create a new branch with the name <code><branch></code> .
<code>git checkout -b <branch></code>	Create and check out a new named <code><branch></code> . Drop the <code>-b</code> flag to checkout an existing branch.
<code>git merge <branch></code>	Merge <code><branch></code> into the current branch.

Publish and Share with Quarto :: CHEATSHEET



Author

WRITE AND CODE IN PLAIN TEXT
Author documents as .qmd files or Jupyter notebooks. Write in a rich Markdown syntax.

Author

SOURCE FILE: hello.qmd

```
---
```

```
title: "Hello, Penguins"
format: html
execute:
  echo: false
```

Meet the penguins

The 'penguins' data contains from three islands in the R

The three species of penguin distributions of physical dimensions (#fig-penguins).

```
##(r)
#| label: fig-penguins
#| fig-cap: "Dimensions of penguins"
#| warning: false
library(tidyverse, quietly = TRUE)
library(palmerpenguins)
```

```
penguins |>
  ggplot(aes(x = flipper_length_mm, y = bill_length_mm)) +
  geom_point(aes(color = species)) +
  scale_color_manual(
    values = c("darkorange", "purple", "cyan4")) +
```

Set format(s) and options Use YAML Syntax

Write with **Markdown**
RStudio Help & Markdown Quick Reference

R P X Use Visual Editor

Include code R, Python, Julia, Observable, or any language with a Jupyter kernel

USE A TOOL WITH A RICH EDITING EXPERIENCE

R RStudio P Posit X VS Code + Quarto extension

Run code cells as you write

Render with a button or keyboard shortcut

Edit Quarto documents with a Visual Editor

OR ANY TEXT EDITOR

Quarto documents (.qmd) can be edited in any tool that edits text.

Apply formatting in Visual Editor. Saved as Markdown in source.
Insert elements like code cells, cross references, and more.

Render → Publish

GENERATE DOCUMENTS, PRESENTATIONS AND MORE

Produce HTML, PDF, MS Word reveal.js, MS Powerpoint, Beamer Websites, blogs, books...

SHARE YOUR WORK WITH THE WORLD

Quickly deploy to GitHub Pages, Netlify, Quarto Pub, Post Cloud, or Posit Connect

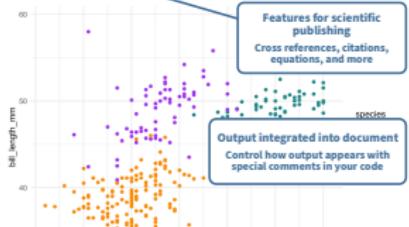
Render

RENDERED OUTPUT: hello.html

Hello, Penguins

Meet the penguins

The three species of penguins have quite distinct distributions of physical dimensions (Figure 1).



Save, then Render to preview the document output.

Terminal
quarto preview hello.qmd

R Use Render button P X Use Review button

The resulting HTML/PDF/MS Word/etc. document will be created and saved in the same directory as the source .qmd file.

BEHIND THE SCENES

When you render a document, Quarto:

1. Runs the code and embeds results and any output associated with: Knitr, if any [r] cells or Jupyter, if any other cells.
2. Converts the .md file into the output format with Pandoc.

GET QUARTO

<https://quarto.org/docs/download/>
Or use version bundled with Positron or RStudio
GET STARTED
<https://quarto.org/docs/get-started/>

Publish

Terminal
quarto publish {venue} hello.qmd
{venue}: quarto-pub, connect, gh-pages, netlify, conference

Quarto Pub Free publishing service for Quarto content.

Posit Connect Org-hosted, control access, schedule updates.

R Use Publish button P X Use Posit Publisher extension

Quarto Projects

CREATE WEBSITES, BOOKS, AND MORE

A directory of Quarto documents + a configuration file (quarto.yml)
See examples at <https://quarto.org/docs/gallery/>
Get started from the command line

Terminal
quarto create project {type}

{type}: website, blog, book, confluence, manuscript

R Use File > New Project

P X Use Quarto: Create Project command

Artwork from "Hello, Quarto" keynote by Julia Lemos and Mire Cetinkaya-Rundel, presented at RStudio Conference 2022. Illustrated by Allison Horst.



Include Code

CODE CELLS

Code cells start with `{{language}}` and end with `{{language}}`.



Use [Insert Code Chunk/Cell](#)

```
```(r)
#| label: chunk-id
library(tidyverse)
...```
```(python)
#| label: chunk-id
import pandas as pd
...```

```

Other languages: {julia}, {ojs}

Add code cell options with #| comments.

Cell options control **execution**, **figures**, **tables**, **layout** and more. See them all at: <https://quarto.org/docs/reference/cells>

EXECUTION OPTIONS

OPTION DEFAULT EFFECTS

echo	true	false: hide code fenced: include code cell syntax
-------------	------	--

eval	true	false: don't run code
-------------	------	-----------------------

include	true	false: don't include code or results
----------------	------	--------------------------------------

output	true	false: don't include results asis: treat results as raw markdown
---------------	------	---

warning	true	false: don't include warnings in output
----------------	------	---

error	false	true: include error in output and continue with render
--------------	-------	--

Set execution options at the **cell level**:

```
```(r)
#| echo: false
...```
```(python)
#| echo: false
...```

```

Or, **globally** in the YAML header with the `execute` option:

```
---
```

Set options in code cells with #| comments and YAML syntax:
key: value

```
---
```

INLINE CODE

Use computed values directly in text sections.
Code is evaluated at render and results appear as text.

KNITR

JUPYTER

OUTPUT

Value is `r^2 + 2` . Value is `{{python}} 2 + 2` . Value is 4.

Set Format and Options

SET FORMAT OPTIONS

```
---
```

title: "My Document"
format:
 html:
 code-fold: true
 toc: true

Indent options 4 spaces
Indent format 2 spaces

MULTIPLE FORMATS

```
---
```

title: "My Document"
toc: true
format:
 html:
 code-fold: true
 pdf: default

Top-level options apply to all formats

OPTION

	html/revealjs pdf/beamer docx/pptx	DESCRIPTION
toc	X X X	Add a table of contents (true or false)
toc-depth	X X	Lowest level of headings to add to table of contents (e.g. 2, 3)
anchor-sections	X	Show section anchors on mouse hover (true or false)
highlight-style	X X X	Syntax highlighting theme (e.g. arrow, pygments, kate, zenburn)
mainfont, monofont	X X	Font name. HTML sets CSS font family; LaTeX via fontspec package
theme	X	Bootstrap theme name (e.g. cosmo, darkly, solar etc.)
css	X	CSS or SCSS file to use to style the document (e.g. "style.css")
reference-doc	X	docx/pptx file containing template styles (e.g. file.docx, file.pptx)
include-in-header	X X	Files of content to include in header of output document, also <code>include-before-body</code> , <code>include-after-body</code>
keep-md	X X X	Keep intermediate markdown (true or false), also <code>keep-ipynb</code> , <code>keep-te</code>
documentclass	X	LaTeX document class, set document class options with <code>classoption</code>
pdf-engine	X	LaTeX engine to produce PDF output (xelatex, pdflatex, lualatex)
cite-method	X	Method used to format citations (citetproc, natbib, biblatex)

LaTeX

code-fold	X	Let readers toggle the display of R code (false, true, or show)
code-tools	X	Add menu for hiding, showing, and downloading code (true or false)
code-overflow	X	Display of wide code (scroll, or wrap)
fig-align	X X X	Alignment of figures (default, left, right, or center)
fig-width, fig-height	X X X	Default width and height for figures in inches
fig-format	X X X	Format for Matplotlib or R figures (retina, png, jpeg, svg, or pdf)

Visit <https://quarto.org/docs/reference/> to see all options by format

Add Content

FIGURES

MARKDOWN

```
![[CAP]](image.png){#fig-LABEL fig-alt="ALT"}
```

COMPUTATION

```
```(python)
#| label: fig-LABEL
#| fig-cap: CAP
#| fig-alt: ALT
{ plot code here }```
Or(r)
```

## CROSS REFERENCES

1. Add a bibliography file to the YAML header:

```

```

bibliography: references.bib

2. Add citations: `[@citation]`, `or @citation`

You can see in `#fig-scatterplot`,  
that...

## TABLES

### MARKDOWN

User @knitr::kable() to produce

Markdown:

```
|object | radius|
|-----:|-----:|
|Sun | 696000|
|Earth | 6371|
```

: CAPTION (#tbl-LABEL)



Insert Table in Visual Editor

Also see the R packages: gt, flexible, kableExtra.

## COMPUTATION

### KNITR

User @knitr::kable() to produce

Markdown:

```
... (python)
#| label: tbl-LABEL
#|tbl-cap: CAPTION
import pandas as pd, tabulate
from IPython.display import Markdown
df = pd.DataFrame([{"A": [1, 2],
"B": [1, 2]})
```

Markdown(df.to\_markdown(index=False))

... (Julia)

## JUPYTER

Add Markown() to Markdown output:

```
... (python)
#| label:tbl-LABEL
#|tbl-cap: CAPTION
import pandas as pd, tabulate
from IPython.display import Markdown
df = pd.DataFrame([{"A": [1, 2],
"B": [1, 2]})
```

Markdown(df.to\_markdown(index=False))

## CITATIONS

1. Add a bibliography file to the YAML header:

```

```

bibliography: references.bib

2. Add citations: `[@citation]`, `or @citation`



Use [Insert Citations](#) dialog in the Visual Editor

Build your bibliography file from your Zotero library,

DOI, Crossref, DataCite, or PubMed

## CALLOUTS

### callout-tip

Instead of `tip` use one of `note`, `warning`, or `important`.

### note

⚠ warning

### Text

⚠ caution

⚠ important

⚠ note

⚠ warning

⚠ important

## SHORTCODES

```
{< include_file.qmd >}
```

```
{< embed_file.ipynb#id >}
```

```
{< video video.mp4 >}
```



# Data visualization with ggplot2 :: CHEATSHEET



## Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same components: **a data layer**, **a coordinate system**, and **geoms**—visual marks that represent data points.



To display values, map variables in the data to visual properties of the geom (**aesthetics**) like **size**, **color**, and **x** and **y** locations.



Complete the template below to build a graph.

```
ggplot(data = <DATA>) +
 <GEOM_FUNCTION>(mapping = aes <MAPPINGS>,
 stat = <STAT>, position = <POSITION>) +
 <COORDINATE_FUNCTION>+
 <FACET_FUNCTION>+
 <SCALE_FUNCTION>+
 <THEME_FUNCTION>
```

required  
Not required, sensible defaults supplied

ggplot(data = mpg, aes(x = cyl, y = hwy)) Begins a plot that you finish by adding layers to. Add one geom function per layer.

last\_plot() Returns the last plot.

ggsave("plot.png", width = 5, height = 5) Saves last plot as 5" x 5" file named "plot.png" in working directory. Matches file type to file extension.

## Aes Common aesthetic values.

color and fill - string ("red", "#RRGGBB")

linetype - integer or string [0 = "blank", 1 = "solid", 2 = "dashed", 3 = "dotted", 4 = "dotdash", 5 = "longdash", 6 = "twodash"]

size - integer (in mm for size of points and text)

linewidth - integer (in mm for widths of lines)

shape - integer/shape name or a single character ("a")  
a circle, square, triangle, diamond, etc.



## Geoms

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

### GRAPHICAL PRIMITIVES

```
a <- ggplot(economics, aes(date, unemploy))
b <- ggplot(seals, aes(x = long, y = lat))
```

**a + geom\_blank()** and **a + expand\_limits()**  
Ensure limits include values across all plots.

**b + geom\_curve**(aes(yend = lat + 1,  
xend = long + 1), curvature = 1) - x, y, end, yend,  
alpha, angle, color, curvature, linetype, size

**a + geom\_path**(linend = "butt",  
linejoin = "round", linemtpe = 1)  
x, y, alpha, color, group, linetype, size

**a + geom\_rect**(aes(alpha = .5)) - x, y, alpha,  
color, fill, group, subgroup, linetype, size

**b + geom\_ribbons**(aes(ymin = long, ymax = lat,  
xmax = long + 1, ymax = lat + 1)) - xmax, xmin,  
ymin, ymax, alpha, color, fill, linetype, size

**a + geom\_rug**(aes(ymin = unemploy - 900,  
xymax = unemploy + 900)) - x, ymax, ymin,  
alpha, color, fill, group, linetype, size

### LINE SEGMENTS

common aesthetics: x, y, alpha, color, linetype, size

**b + geom\_abline**(aes(intercept = 0, slope = 1))

**b + geom\_hline**(aes(intercept = lat))

**b + geom\_vline**(aes(intercept = long))

**b + geom\_segment**(aes(yend = lat + 1, xend = long + 1))

**b + geom\_spoke**(aes(angle = 1:1155, radius = 1))

### ONE VARIABLE continuous

```
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)
```

**c + geom\_area**(stat = "bin")  
x, y, alpha, color, fill, linetype, size

**c + geom\_density**(kernel = "gaussian")  
x, y, alpha, color, fill, group, linetype, size, weight

**c + geom\_dotplot**  
x, y, alpha, color, fill

**c + geom\_freqpoly**  
x, y, alpha, color, group, linetype, size

**c + geom\_histogram**(binwidth = 5)  
x, y, alpha, color, fill, linetype, size, weight

**c2 + geom\_qq**(aes(sample = hwy))  
x, y, alpha, color, fill, linetype, size, weight

### discrete

**c <- ggplot(mpg, aes(ffill))**

**d + geom\_bar**  
x, alpha, color, fill, linetype, size, weight

### continuous bivariate distribution

**h <- ggplot(diamonds, aes(carat, price))**

**h + geom\_bin2d**(binwidth = c(0.25, 500))  
x, y, alpha, color, fill, linetype, size, weight

**h + geom\_hex2d**(x, y, alpha, color, group, linetype, size)

**h + geom\_hex**(x, y, alpha, color, fill, size)

### continuous function

**i <- ggplot(economics, aes(date, unemploy))**

**i + geom\_area**  
x, y, alpha, color, fill, linetype, size

**i + geom\_line**  
x, y, alpha, color, group, linetype, size

**i + geom\_step**(direction = "hv")  
x, y, alpha, color, group, linetype, size

### visualizing error

**df <- data.frame(grp = c("A", "B"))**, fit = 4.5, se = 1.2)

**j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))**

**j + geom\_crossbar**(atten = 2) - x, y, ymax,  
x, y, alpha, color, fill, group, linetype, size

**j + geom\_errorbar**(x, y, ymax, ymin,  
alpha, color, group, linetype, size, width

Also **geom\_errorbarh**.

**j + geom\_linerange**(x, y, ymax, alpha, color, group, linetype, size)

**j + geom\_pointrange**(x, y, ymin, ymax,  
alpha, color, fill, group, linetype, shape, size)

### both discrete

**g <- ggplot(diamonds, aes(cut, color))**

**g + geom\_count**  
x, y, alpha, color, fill, shape, stroke

**e + geom\_jitter**(height = 2, width = 2)  
x, y, alpha, color, fill, shape, size

Draw the appropriate geometric object depending on the simple features passed in the `...args` arguments:  
map\_id, alpha, color, fill, linetype, linewidth,

nc <- st\_sf(st\_read(system.file("shape/nc.shp", package = "sf")))

**ggplot(nc) +**

**geom\_sf**(aes(fill = AREA))

### THREE VARIABLES

**sealsSz <- with(seals, sqrt(delta\_long^2 + delta\_lat^2))**

**i <- ggplot(seals, aes(long, lat))**

**i + geom\_contour**(aes(z = z))  
x, y, alpha, color, group, linetype, size, weight

**i + geom\_contour\_filled**(aes(fill = z))  
x, y, alpha, color, fill, group, linetype, size, subgroup

**i + geom\_raster**(aes(fill = z), hijst = 0.5,  
yjst = 0.5, interpolate = FALSE)  
x, y, alpha, color, fill

**i + geom\_tile**(aes(fill = z))  
x, y, alpha, color, fill, linetype, size, width



## Stats

An alternative way to build a layer.

A stat builds new variables to plot (e.g., count, prop).



Visualizes a stat by changing the default stat of a geom function, `geom_bar(stat = "count")` or by using a stat function, `stat_count(geom = "bar")`, which calls a default geom to make a layer (equivalent to a geom function). Use `stat[stat$name]` syntax to map the stat variable `name` to an aesthetic.

`geom to use`   `stat function`   `geom mappings`  
`i + stat_density_2d(aes(..., after_stat(level)), geom = "polygon")`   `variable created by stat`

`c + stat_bin(binwidth = 1, boundary = 10)`

`x, y, fill | count, density`

`c + stat_count(width = 1) x, y | count, prop`

`x, y, count, density`

`c + stat_density(level = 1, kernel = "gaussian")`

`x, y, density, level`

`e + stat_bin(bins = 30, drop = T)`

`x, y, fill | count, density`

`e + stat_bin_hex(bins = 30) x, y, fill | count, density`

`e + stat_density_2d(contour = TRUE, n = 100)`

`x, y, color, size | level`

`e + stat_ellipse(level = 0.95, segments = 51, type = "t")`

`I + stat_contour(aes(z = z)) x, y, order | level`

`I + stat_summary_hex(aes(z = z), bins = 30, fun = max)`

`x, y, z, fill | value`

`I + stat_summary_2d(aes(z = z), bins = 30, fun = mean)`

`x, y, z, fill | value`

`f + stat_boxplot(coef = 1.5)`

`x, y | lower, middle, upper, width, ymin, ymax`

`f + stat_ydensity(kernel = "gaussian", scale = "area") x, y |`

`density, scaled, count, n, violinwidth, width`

`e + stat_ecdf(n = 40) x, y | x`

`e + stat_quantile(quartiles = c(0.1, 0.9), formula = y ~ log(x), method = "rq") x, y | quantile`

`e + stat_smooth(method = "lm", formula = y ~ x, se = T, level = 0.95) x, y | x, y, ymin, ymax`

`ggplot() + xlim(-5, 5) + stat_function(fun = dnorm, n = 20, geom = "point") x, y`

`ggplot() + stat_qq(aes(sample = rnorm(100)))`

`x, y, sample | sample, theoretical`

`e + stat_sum(x, y, size | n, prop`

`e + stat_summary(fun.data = "mean_cl_boot")`

`h + stat_summary_bin(fun = "mean", geom = "bar")`

`e + stat_identity()`

`e + stat_unique()`

## Scales

Override defaults with `scales` package.

**Scales** map data values to the visual values of an aesthetic. To change a mapping, add a new scale.

```
n <- d + geom_bar()
r + coord_cartesian(slim = c(0, 5)) -> xlim, ylim
The default cartesian coordinate system.

n + scale_geom(aes(fill = fl))
scale_# aesthetic to adjust
scale_fill_manual(values = c("skyblue", "royalblue", "blue", "navy", "darkblue", "black", "grey", "brown", "red", "purple", "pink", "yellow", "orange", "lightgreen", "lightblue", "lightgrey", "white"))
scale_*_manual(name = "text", labels = c("D", "E", "P", "R"))
scale_*_identity()
scale_*_log10()
scale_*_log10(breaks = c(10^0, 10^1, 10^2, 10^3, 10^4))
scale_*_log10(trans = "sqrt") -> x, y, xlim, ylim
Transformed cartesian coordinates. Set xtrans and ytrans to the name of a window function.
scale_*_date()
scale_*_date(date_labels = "%m/%d/%y", date_breaks = "2 weeks") -> Treat data values as dates.
scale_*_date_time()
Treat data values as date times. Same as scale_*_date(). See ?strptime for label formats.
```

### GENERAL PURPOSE SCALES

Use with most aesthetics

- `scale_*_continuous()` - Map cont. values to visual ones.
- `scale_*_discrete()` - Map discrete values to visual ones.
- `scale_*_manual()` - Map discrete values to discrete bins.
- `scale_*_identity()` - Use data values as visual ones.
- `scale_*_manual(values = c(...))` - Map discrete values to manually chosen visual ones.
- `scale_*_date(date_labels = "%m/%d/%y", date_breaks = "2 weeks")` - Treat data values as dates.
- `scale_*_date_time()` - Treat data values as date times. Same as scale\_\*\_date(). See ?strptime for label formats.

### X & Y LOCATION SCALES

Use with x or y aesthetics (x shown here)

```
scale_x_log10() -> Plot x on log10 scale.
scale_x_reverse() -> Reverse the direction of the x axis.
scale_x_sqrt() -> Plot x on square root scale.
```

### COLOR AND FILL SCALES (DISCRETE)

```
n + scale_fill_brewer(palette = "Blues")
For palette choices:
RColorBrewer::display.brewer.all()
n + scale_fill_grey(start = 0.2, end = 0.8, na.value = "red")
```

### COLOR AND FILL SCALES (CONTINUOUS)

`o + c + geom_dotplot(aes(fill = x))`

`o + scale_fill_distiller(palette = "Blues")`

`o + scale_fill_gradient(low = "red", high = "yellow")`

`o + scale_fill_gradient2(low = "red", high = "blue", mid = "white", midpoint = 25)`

`o + scale_fill_gradientn(colors = topo.colors(6))  
# Also: rainbow(), heat.colors(), terrain.colors(), cm.colors(), RColorBrewer::brewer.pal()`

### SHAPE AND SIZE SCALES

`p <- e + geom_point(aes(shape = fl, size = cyl))`

`p + scale_shape() + scale_size()`

`p + scale_shape_manual(values = c(3:7))`

`p + scale_radius(range = c(1, 6))`

`p + scale_size_area(max.size = 6)`

## Coordinate Systems

`r <- d + geom_bar()`

`r + coord_cartesian(slim = c(0, 5)) -> xlim, ylim`  
The default cartesian coordinate system.

`r + coord_fixed(ratio = 1/2)`  
ratio, xlim, ylim - Cartesian coordinates with fixed aspect ratio between x and y units.

`r + coord_flip()`  
Flip cartesian coordinates by switching x and y aesthetic mappings.

`r + coord_polar(theta = "x", direction = 1)`  
theta, start, direction - Polar coordinates.

`r + coord_trans(y = "sqrt") -> x, y, xlim, ylim`  
Transformed cartesian coordinates. Set xtrans and ytrans to the name of a window function.

`r + coord_sf(..., crs = "llcsp")`  
xlim, ylim, crs. Ensures all layers use a common Coordinate Reference System.

## Position Adjustments

Position adjustments determine how to arrange geoms that would otherwise occupy the same space.

`s <- ggplot(mpg, aes(fl, fill = drv))`

`s + geom_bar(position = "dodge")`  
Arrange elements side by side.

`s + geom_bar(position = "fill")`  
Stack elements on top of one another, normalize height.

`s + geom_point(position = "jitter")`  
Add random noise to X and Y position of each element to avoid overlapping.

`s + geom_label(position = "nudge")`  
Nudge labels away from points.

`s + geom_bar(position = "stack")`  
Stack elements on top of one another.

Each position adjustment can be recast as a function with manual width and height arguments:

`s + geom_bar(position = position_stack(width = 1))`

## Themes

`r + theme_bw()`  
White background with grid lines.

`r + theme_gray()`  
Grey background (default theme).

`r + theme_minimal()`  
Minimal theme.

`r + theme_dark()`  
Dark for contrast.

`r + theme_void()`  
Empty theme.

`r + theme()` Customize aspects of the theme such as axis, legend, panel, and facet properties.

`r + labs(title = "Title")` and `theme(plot.title.position = "plot")`

`r + theme(panel.background = element_rect(fill = "blue"))`

## Faceting

Facets divide a plot into subplots based on the values of one or more discrete variables.

`t <- ggplot(mpg, aes(cty, hwy)) + geom_point()`

`t + facet_grid(~ - fl)`  
Facet into columns based on fl.

`t + facet_grid(~year - )`  
Facet into rows based on year.

`t + facet_grid(~year - fl)`  
Facet into both rows and columns.

`t + facet_wrap(~ fl)`  
Wrap facets into a rectangular layout.

Set scales to let axis limits across vary across facets.

`t + facet_grid(~fl, scales = "free")`  
x and y axis limits adjust to individual facets:  
"free\_x" - x axis limits adjust  
"free\_y" - y axis limits adjust

Set labeller to adjust facet label:

`t + facet_grid(~ fl, labeller = label_both)`

`# c - fl; d - fl; e - fl; f - fl; g - fl;`

`t + facet_grid(~ fl, labeller = label_bquote(alpha ^ .(fl)))`

`# a^f # a^d # a^e # a^g # a^c`

Use `label()` to label the elements of your plot.

`t + labels(label = "New")` Add a label below the plot.

`t + subtitle(subtitle = "New")` Add a title above the plot.

`t + caption(caption = "New")` Add a caption below plot.

`t + alt_text(alt = "New")` Add alt text to the plot.

`(new) = New` `(alt) = alt` `(caption) = New` `(label) = New` `(subtitle) = New`

`t + annotate(gem = "text", x = 5, y = 5, label = "A")` Places a geom with manually selected aesthetics.

`p + guides(x = guide_axis_nudge(d = 2))` Avoid crowded or overlapping labels with labels with `guide_axis_nudge(d, angle)`.

`p + guides(fill = "none")` Set legend type for each aesthetic: color/bar, legend, or none (no legend).

`t + theme(legend.position = "bottom")` Place legend at bottom, "top", "left", or "right".

`n + scale_fill_discrete(name = "Title")`

Set legend title and labels with a scale function.

## Zooming

Without clipping (preferred):

`t + coord_cartesian(xlim = c(0, 100), ylim = c(10, 20))`

With clipping (removes unseen data points):

`t + xlim(0, 100) + ylim(10, 20)`

`t + scale_x_continuous(limits = c(0, 100)) + scale_y_continuous(limits = c(0, 100))`

# Machine Learning Modelling in R :: CHEAT SHEET

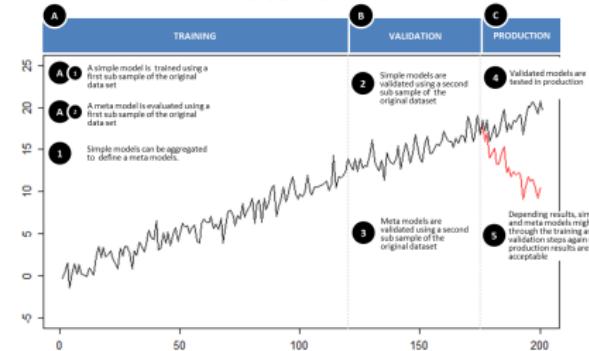
## Supervised & Unsupervised Learning

ALGORITHM	DESCRIPTION	R PACKAGE::FUNCTION	SAMPLE CODE		ALGORITHM	DESCRIPTION	R PACKAGE::FUNCTION	SAMPLE CODE
Naïve Bayes classifier	A classification technique based on Bayes' Theorem with an assumption of independence between predictors. For a given test, a Naïve Bayes classifier assumes that the presence or absence of a particular feature is independent of the presence or absence of any other feature.	e1071::naiveBayes	naiveBayes(class ~ ... , data = x)		Regulation	Regulation edits a generally overfitted model by reducing the number of nodes in the tree.	tree::regTree	regTree(tree, x, y, ntree=100, cp=0.01)
k-Nearest Neighbours	A non-parametric method used for classification and regression. In both cases, the input consists of the k closest training samples in the feature space. The output is a prediction on whether k-NN is used for classification or regression.	class::knn	knn(braai, test, k, l, k=1, l=0, prob = FALSE, use.all = TRUE)		Boosting	A process of iteratively refitting, e.g., by reweighting, of estimated regression and classification functions (though it can also refer to a process of iterative improvement applied to a tree), in order to improve predictive ability.	gbm::gbm	gbm(x ~ ., data = cars, n.trees=100)
Linear Regression	Model the linear relationship between a scalar dependent variable Y and one or more explanatory variables (or independent variables) denoted X.	stats::lm	lm(dish ~ speed, data=cars)		Bagging	Bagging is a way to increase the power of a predictive model by averaging the predictions of multiple separate models (bagging). It involves randomly selecting a subset of the training data set, and using each of them to construct a separate model and separate predictions.	rpart::rpart	rpart(y ~ family = binomial)(link = "logit"), data = X)
Logistic Regression	Used to predict a binary outcome (1 / 0, Yes / No, True / False) given a set of independent variables.	stats::glm	glm(Y ~ ..., family = binomial(link = "logit"), data = X)		Pruning	Pruning is a technique that reduces the size of decision trees by removing sections of the tree that provide little information to classify the data. This improves the accuracy of the final classifier and hence improves predictive accuracy by reducing overfitting.	rpart::prune	prune(rpart, cp = 0.1)
Tree-Based Models	The idea is to consecutively divide (branch) the training dataset based on the input features until an assignment function (leaf node) is reached for the target variable (the "data bucket"). (left) recursive partitioning (right) classification trees.	rpart::rpart	rpart(y ~ hypothesis   Age + Number + Start, data = hypothesis)		Random Forest	Random sampling of observations for training and testing a model can be an issue when faced with a time dimension. This is where random forests come in, as they use the correlation properties in the data which we would like to exploit.	randomForest::randomForest	randomForest(X ~ ., data = mySub)
Artificial Neural Network	Neural networks are built from units called perceptrons, which operate on weighted inputs, an activation function and an output. An ANN model is built up by combining perceptrons in structured layers.	neuralnet::neuralnet	neuralnet(data ~ trite, hidden = 2, stepwise = T)					
Support Vector Machine	A data classification method that separates data using hyperplanes.	e1071::svm	svm(formula, data = NULL, ..., subset, na.action = na.omit, scale = TRUE)					
PCA	A procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components.	stats::princomp	stats::princomp(princomp, data = NULL, subset, na.action = ..., ...)					
Principal Component Analysis		factoextra::PCA	factoextra::PCA(data, quantiles = 11, ...)					
K-Means Clustering	Aims at partitioning a observations into k clusters in which each observation belongs to the cluster with the nearest mean	stats::kmeans	kmeans(cents, centers, iter.max = 10, nstart = 1, algorithm = c("Hartigan-Wong", "K-means++", "Forgy", "MacQueen"), trace = FALSE)					
Hierarchical Clustering	An approach which builds a hierarchy from the bottom-up, and doesn't require the number of clusters to be specified beforehand.	stats::hclust	hclust(d, method = "complete", members = NLL)					

## Standard Modelling Workflow



## Time Series View



CC BY SA Arnaud Amsellem • [thertrader@gmail.com](mailto:thertrader@gmail.com) • [www.thertrader.com](http://www.thertrader.com) • Updated: 2018-03

# Distribuciones

Si  $(\Omega, \mathcal{A}, \mathbb{P})$  es un espacio de probabilidad,  $X$  una variable aleatoria entonces se puede calcular y simular numeros aleatorios que tienen la misma distribución que  $X$ .

- $P(X \in A) = \begin{cases} \sum_{x \in R_X} \mathbb{P}(X = x) & \text{si } X \text{ es discreta} \\ \int_A f_X(x) dx & \text{si } X \text{ es continua} \end{cases}$
- la función de distribución de  $X$  es  $F_X : \mathbb{R} \rightarrow \mathbb{R}$  definida por

$$F_X(x) = \mathbb{P}(X \leq x) \quad \forall x \in \mathbb{R}$$

- la densidad  $\begin{cases} \mathbb{P}(X = x) & \text{si } X \text{ es discreta} \\ f_X(x) & \text{si } X \text{ es continua} \end{cases}$
- los cuantiles  $F^-(p) = \inf\{x : F(x) \geq p\}$ <sup>1</sup>

---

<sup>1</sup>Si  $F$  es biyectiva,  $F^- = F^{-1}$

## ① Distribuciones discretas

- binomial ( $n, p$ )
- hypergeométrica ( $N, n, k$ )
- poisson  $\lambda$
- geométrica  $p$
- ley con soporte finito  $\{(a_i, p_i), i = 1 \dots m\}$

binom  
hyper  
pois  
geom  
sample

## ② Distribuciones continuas

- normal ( $\mu, \sigma^2$ )
- uniforme en  $[a, b]$
- t-student con  $p$  grados de libertad
- chi cuadrado con  $p$  grados de libertad

norm  
unif  
t  
chisq

Si zzz representa el nombre de alguna distribución entonces:

- dzzz( $x, \dots$ ) calcula la densidad de la distribución zzz en  $x$
- pzzz( $x, \dots$ ) calcula la función de distribución en  $x$
- qzzz( $p, \dots$ ) calcula el cuántil de orden  $p$
- rzzz( $x, \dots$ ) da una muestra de tamaño  $n$

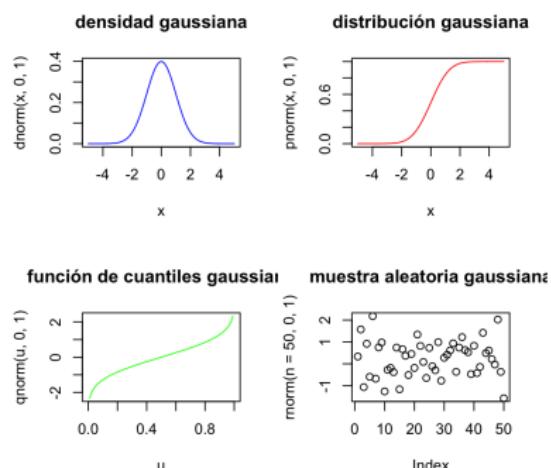
donde .... representan los parámetros específicos de cada distribución

# con la gaussiana $\mathcal{N}(0, 1)$

- $dnorm(x, 0, 1)$  densidad en el punto  $x$
- $pnorm(x, 0, 1)$  función de distribución en el punto  $x$
- $qnorm(p, 0, 1)$  cuantil  $p+$
- $rnorm(n, 0, 1)$  muestra de tamaño  $n$

```
x=seq(-5,5,.01)
u=seq(0,1,.01)
```

```
par(mfrow=c(2,2))
plot(x,dnorm(x,0,1),type='l', col='blue',main='densidad gaussiana')
plot(x,pnorm(x,0,1),type='l', col='red',main='distribución gaussiana')
plot(u,qnorm(u,0,1),type='l', col='green',main='función de cuantiles gaussiana')
plot(rnorm(n=50,0,1), main='muestra aleatoria gaussiana')
```



# Simulación de una distribución discreta con soporte finito

Sea  $X$  una variable aleatoria con soporte en  $\{a_1, \dots, a_L\}$  tal que  $\mathbb{P}(X = a_\ell) = p_k$  para todo  $\ell = 1, \dots, L$

- La función `sample` es un generador de numeros aleatorios
- `sample(x, n, replace = T, prob)` es el comando

```
muestra=sample(x=c(1:6),size=50,
replace=T)
```

```
muestra
```

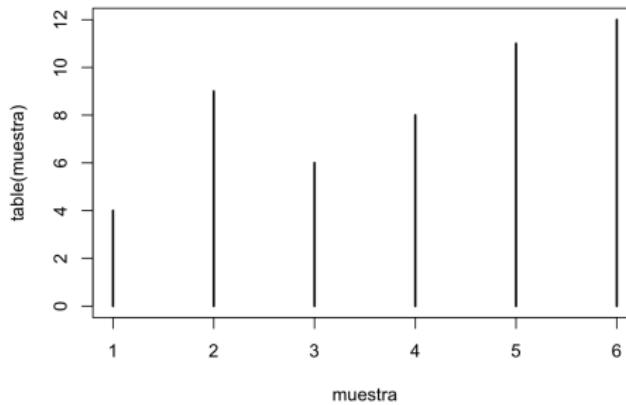
```
[1] 5 1 5 6 5 2 3 5 6 4 3 5 6
[14] 4 2 3 6 5 2 4 3 2 5 3 6 5
[27] 6 1 2 6 2 6 6 2 6 4 4 5 4
[40] 4 5 5 1 4 3 2 1 6 6 2
```

```
table(muestra)
```

```
muestra
```

```
1 2 3 4 5 6
4 9 6 8 11 12
```

```
plot(table(muestra))
```



# Escribir una función

R tiene muchas funciones, pero podemos y deberemos implementar otras.

```
cuentas = function(x, y){
 a = x+2*y
 b=x-y
 c=x/y
 final=c(mean(a),sd(a),mean(b),sd(b),
 mean(c),sd(c))
 salida=list(x=x,y=y,resultados=result)
 return(salida)}
x <- runif(50,0,1)
y <- rnorm(50,0,3)
cuentas(x,y)
```

```
> cuentas(x,y)
$x
[1] 0.02489827 0.00484174 0.13137825 0.02481710 0.09212945 0.47555822 0.19665676
[8] 0.81467395 0.24405923 0.48880764 0.64888248 0.15871647 0.83697791 0.43659653
[15] 0.23349553 0.29885599 0.15942134 0.58556710 0.14877798 0.17905956 0.33922150
[22] 0.18390194 0.40554152 0.61355962 0.68567133 0.82352560 0.37398211 0.95284100
[29] 0.01883216 0.61237497 0.46357376 0.35472962 0.88242303 0.07580286 0.88329119
[36] 0.18822184 0.21506675 0.86953229 0.06083806 0.23489960 0.99301551 0.67819773
[43] 0.43548680 0.37730994 0.46118402 0.80185187 0.66296698 0.46600366 0.05325944
[50] 0.23719792

$y
[1] -2.184657334 -4.621327216 -2.079283842 0.356548299 -4.094128374 1.769948038
[7] 0.868832086 -2.712645077 0.678974827 2.244243487 3.183285760 -0.638544836
[13] -0.280910382 -0.260142405 4.324385268 3.375215677 2.503204703 -0.862022400
[19] 1.119724302 1.209878994 -3.125019883 -5.184913546 1.925490083 -4.587931594
[25] 0.005051065 0.750743464 1.691602171 0.568278713 -2.198561419 2.959097580
[31] 5.215901301 2.643536427 -5.830952704 4.198728556 -0.168167839 1.574742838
[37] 1.866099708 -0.290058220 -0.225789593 3.057471208 2.134885767 2.970786738
[43] 7.148780086 1.993247591 0.622143472 -6.631899332 8.075142010 -1.448030465
[49] 7.124204145 1.123930704

$resultados
[1] 0.52406972 5.94214664 0.45731884 2.94956123 -0.08869252 1.75717340
```

# **Variable cuantitativas**

## Serie estadística continua

Si  $X$  es el puntaje de 120 estudiantes en una prueba de probabilidad:

97	63	81	9	27	65	2	59	25	94	21	85	38	32	36
90	79	57	15	76	44	53	34	51	74	31	88	29	28	27
97	81	3	58	66	4	84	59	50	77	50	14	83	72	11
65	39	84	41	1	40	81	70	40	25	41	26	29	23	55
73	36	87	94	73	66	91	20	21	86	59	95	91	88	64
64	85	80	82	7	25	71	58	32	90	92	47	1	64	5
85	26	92	60	5	74	61	52	64	33	77	67	69	89	26
24	52	51	92	87	70	6	61	43	40	36	59	62	92	61

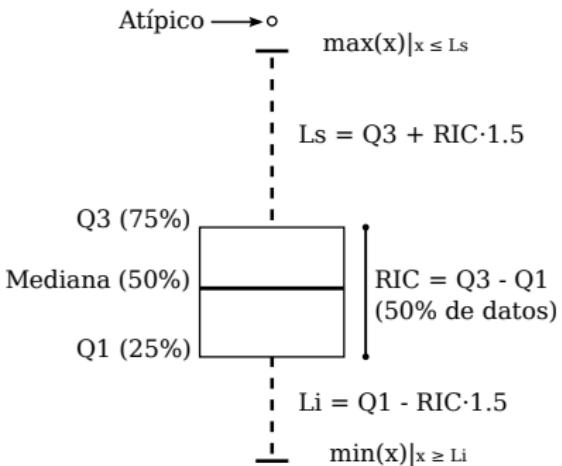
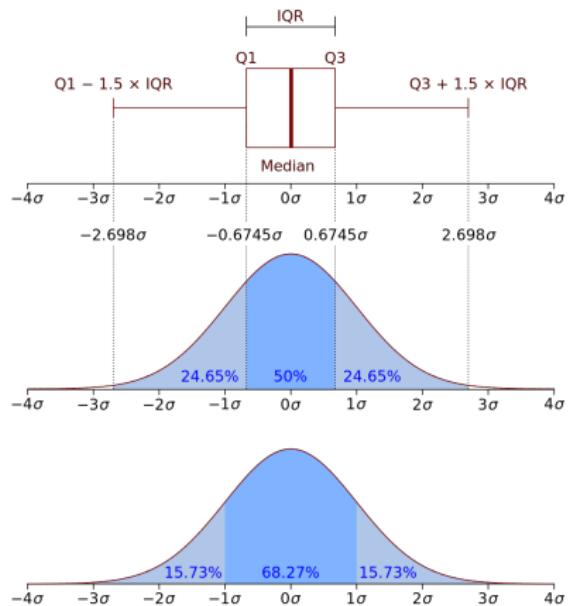
Se pueden identificar algunas cantidades, en particular los 5 números importantes que nos servirán para hacer el diagrama de caja: el mínimo  $m$ , el primer cuartil  $Q_1$ , la mediana  $m$ , el tercer cuartil  $Q_3$ , el máximo  $M$  (se incluye además el promedio pr):

$x_m$	$Q_1$	$m$	$\bar{x}$	$Q_3$	$x_M$
1.00	32	59.00	54.52	79	97.00

Si  $\hat{F}(x)$  es la proporción de individuos en la muestra con un valor menor o igual a  $x$  entonces:

- La mediana es el valor  $m$  de la muestra tal que  $\hat{F}(m) = 0.5$
- Cuartiles:  $\hat{F}(Q_1) = 0.25$ ,  $\hat{F}(Q_2) = 0.5$ ,  $\hat{F}(Q_3) = 0.75$
- Cuantil  $p$ :  $\hat{F}(x_p) = p$

# Boxplot

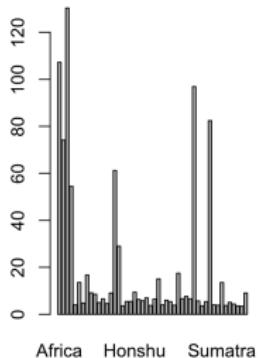


- El valor indicado por la raya en la caja es la mediana
- Las extremidades de la caja corresponden a los cuartiles  $Q_1$  y  $Q_3$
- El rango intercuartílico es  $R = Q_3 - Q_1$
- Las extremidades de los bigotes corresponden al valor más grande menor a  $Q_3 + 1.5R$  y al valor más chico mayor a  $Q_1 - 1.5R$
- Los puntos alineados con los bigotes pero fuera de ellos corresponden a los puntos atípicos o outliers

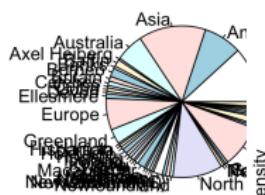
# Estimación densidad

El set de datos Island contiene las áreas, en miles de millas cuadradas, de las masas terrestres que superan las 10 000 millas cuadradas. Como las unidades son grandes tomamos la raíz cuadrada.

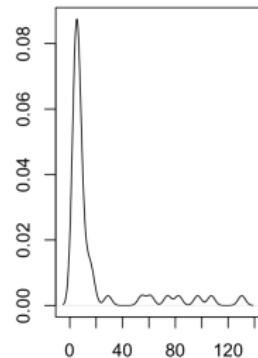
Diagrama en barra



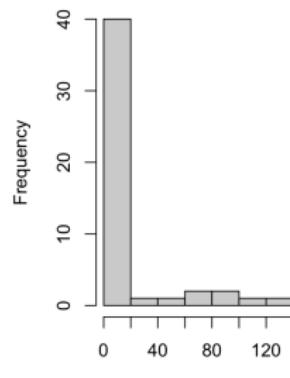
Reparticion



density( $x = x$ )



Histogram of  $x$



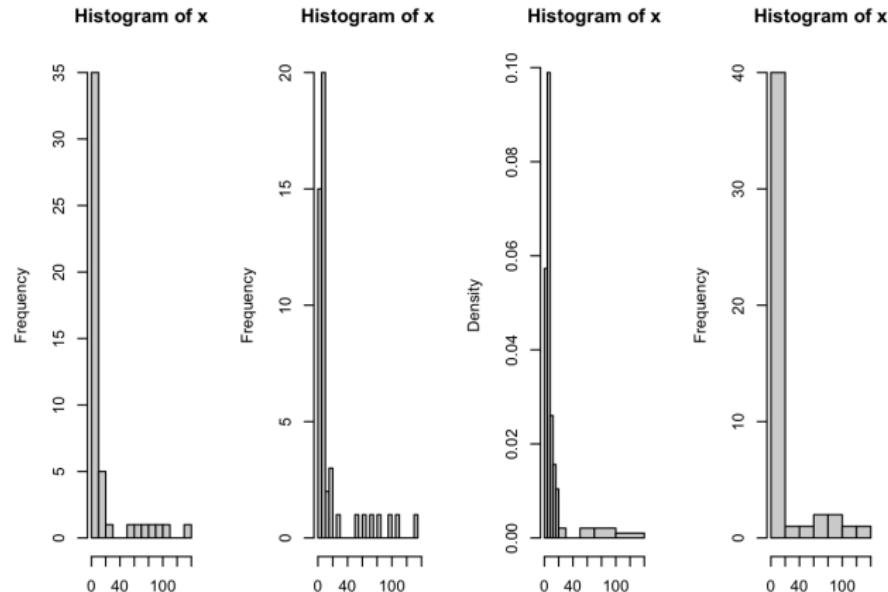
```
x=sqrt(islands)
par(mfrow=c(1,2))
barplot(x,main='Diagrama en barra')
pie(x,main='Reparticion')
plot(density(x))
hist(x)
```

N = 48 Bandwidth = 2.79

# Histogramas

Se puede elegir la cantidad de celdas, o dar un vector, equiespaciado o no, u optimizar según una regla.

```
hist(x, breaks=10)
hist(x, breaks=20)
hist(x, breaks=c(4 *0:5,30, 45, 50 ,70 , 100,140))
hist(x, breaks='Scott')
```

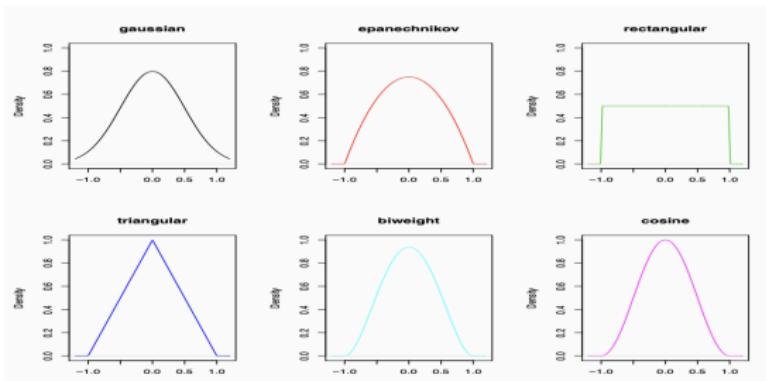


# El estimador por núcleo

$X_1, \dots, X_n$  es una muestra aleatoria con densidad  $f$ . La función density calcula un estimador de la densidad  $f$  de la forma (Kernel Density Estimator):

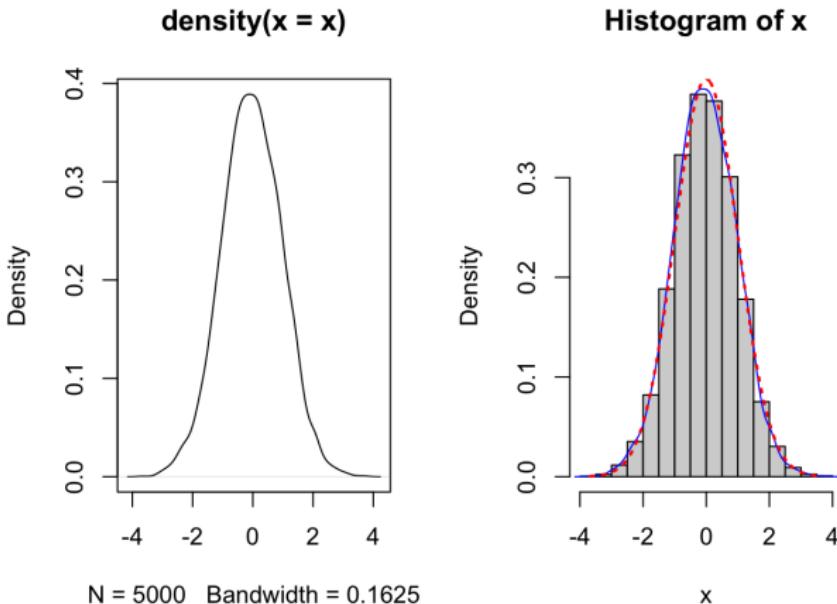
$$f_{KDE}(x) = \frac{1}{nb_w} \sum_{i=1}^n K\left(\frac{x - X_i}{b_w}\right)$$

- La elección del núcleo  $K$  ( $\int K(u) du = 1$ ,  $K(u) = K(-u) \forall u$ )



- La elección de la ventana (bandwidth  $b_w$ ). Este parámetro se puede optimizar.

# Estimador KDE



```
par(mfrow=c(1,2))
x=rnorm (5000)
y=density(x)
plot(y)
hist(x,proba=T)
lines(y,col='blue')
z=seq(min(x),max(x),0.01)
lines(z, dnorm(z, 0, 1), lwd=2, lty =3, col='red')
```

## Estimador KDE

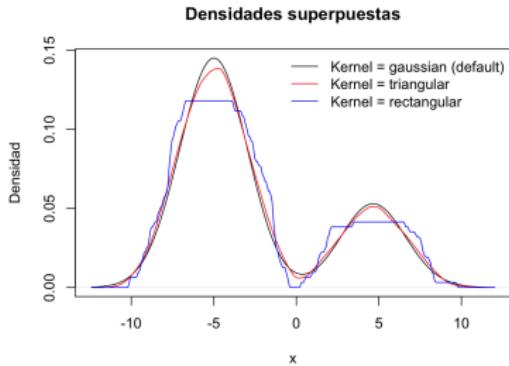
Si consideramos la densidad

$$f(x) = \frac{1}{4} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x+5)^2} + \frac{3}{4} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-5)^2}$$

# Estimador KDE

Si consideramos la densidad

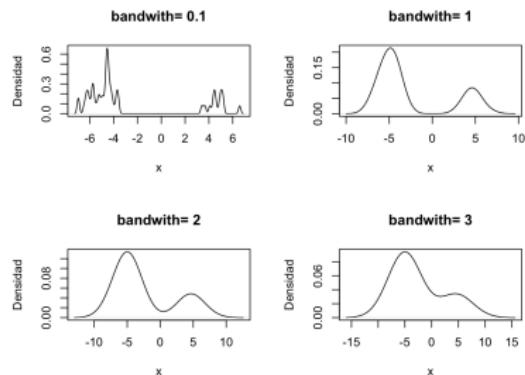
$$f(x) = \frac{1}{4} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x+5)^2} + \frac{3}{4} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-5)^2} \quad (\text{¡sí, es una densidad!})$$



```
w=rbinom(50,1,1/4)
datos = w*rnorm(50,5) +(1-w) * rnorm(50,-5)
```

```
plot(density(datos),
 main = "Densidades superpuestas",
 xlab = "x",
 ylab = "Densidad")

lines(density(datos, kernel="triangular",
 col="red")
lines(density(datos, kernel="rectangular",
 col="blue")
```



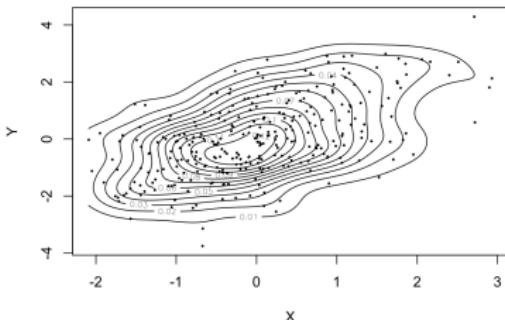
```
par(mfrow=c(2,2))
for(i in c(0.1,1,2,3)){plot(density(datos,bw=i),
 main = paste("bandwidth=",i),
 xlab = "x",
 ylab = "Densidad") }
```

# Estimador KDE 2D

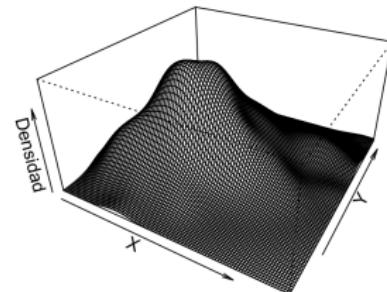
$$f_{KDE}(x, y) = \frac{1}{nbw_1 bw_2} \sum_{i=1}^n K\left(\frac{x-X_i}{bw_1}\right) K\left(\frac{y-Y_i}{bw_2}\right)$$

```
library(MASS)
set.seed(123)
1. Parámetros
n <- 300
mu <- c(0, 0)
Sigma <- matrix(c(1, 0.8,
 0.8, 2),
 nrow = 2, byrow = TRUE)
2. Generar datos
datos <- mvtnorm(n = n, mu = mu, Sigma = Sigma)
x <- datos[,1]
y <- datos[,2]
3. Gráfico de los datos
plot(x, y,
 pch = 16, cex = 0.5,
 xlab = "X", ylab = "Y",
 main = "Datos normales bivariados")
4. Estimación de densidad kernel normal (2D)
dens <- kde2d(x, y, n = 100)
5. Contornos de la densidad + puntos
contour(dens,
 xlab = "X",
 ylab = "Y",
 main = "Estimación de densidad kernel normal (2D)")
points(x, y, pch = 16, cex = 0.4)
6. Superficie 3D de la densidad
persp(dens,
 theta = 30, phi = 30,
 expand = 0.5,
 xlab = "X",
 ylab = "Y",
 zlab = "Densidad",
 main = "Densidad kernel normal (2D)")
```

Estimación de densidad kernel normal (2D)



Densidad kernel normal (2D)



# **Variables cualitativas**

## Serie estadística cualitativa

Disponemos de  $n = 96$  observaciones de una variable cualitativa  $X$  que representa el nivel de satisfacción de clientes en una escala del 1 al 5.

3	1	5	1	1	2	4	2	5	1	5	5	2	1	2	2
2	2	2	3	2	2	2	4	3	1	4	3	2	1	1	3
4	1	4	4	2	4	1	2	3	4	3	4	2	1	4	5
2	1	3	1	2	1	3	1	2	2	3	2	5	1	2	4
1	4	2	1	3	1	1	2	4	1	5	1	3	2	2	1
2	3	2	2	2	5	3	2	2	2	5	1	5	3	4	1

un resumen estadístico de esta serie:

	1	2	3	4	5
Cantidades	25	32	15	14	10
Proporciones	0.26	0.33	0.16	0.15	0.1
Cantidades cumuladas	25	57	72	86	96
Proporciones cumuladas	0.26	0.59	0.75	0.9	1

Las proporciones de cada modalidad constituye la distribución empírica de la variable observada.

## Otro ejemplo

```
muestra2=sample(x=c(1:6),size=100,
replace=T)
```

```
plot(table(muestra2))
barplot(table(muestra2))
summary(muestra2)
```

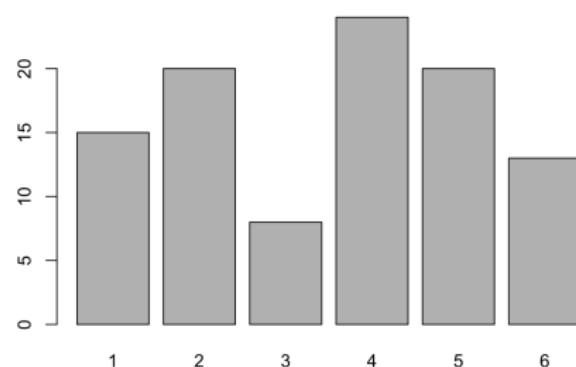
Min.	1st Qu.	Median	Mean	3rd Qu.
1.00	2.00	4.00	3.53	5.00

```
cumsum(table(muestra2))
```

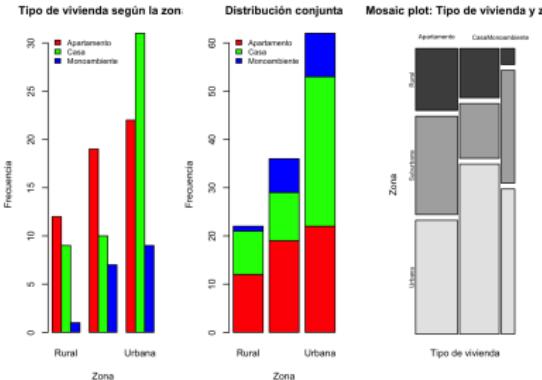
1	2	3	4	5	6
15	35	43	67	87	100

```
cumsum(proportions(table(muestra2)))
```

1	2	3	4	5	6
0.15	0.35	0.43	0.67	0.87	1.00



# Comparación de dos variables cualitativas

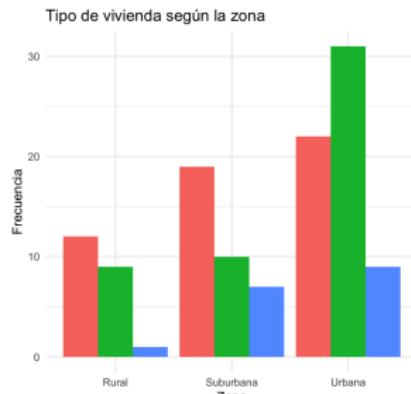


```
par(mfrow=c(1,3))
colores <- rainbow(nrow(tabla))

barplot(tabla,
 beside = TRUE,
 legend = FALSE,
 col=colores,
 xlab = "Zona",
 ylab = "Frecuencia",
 main = "Tipo de vivienda según la zona")
legend("topleft",
 legend = rownames(tabla),
 fill = colores, # clave para que aparezcan los colores
 cex = 0.8,
 bty = "n")

mosaicplot(tabla,
```

```
 main = "Mosaic plot: Tipo de vivienda y zona",
 xlab = "Tipo de vivienda",
 ylab = "Zona",
 color = TRUE)
```



```
library(ggplot2)
ggplot(datos, aes(x = Zona, fill = TipoVivienda)) +
 geom_bar(position = "dodge") +
 labs(title = "Tipo de vivienda según la zona",
 y = "Frecuencia") +
 theme_minimal()
```

# Independencia de dos variables cualitativas

Tabla de contingencia de plebiscito de una ley según preferencia política

	a favor	indiferente	en contra	total
democrato	138	83	64	285
republicano	64	67	84	215
total	202	150	148	500

# Independencia de dos variables cualitativas

Tabla de contingencia de plebiscito de una ley según preferencia política

	a favor	indiferente	en contra	total
democrato	138	83	64	285
republicano	64	67	84	215
total	202	150	148	500

EL test de independencia  $\chi^2$  es

$$\begin{cases} (H_0) : & X \text{ e } Y, \text{ son independientes} \\ (H_1) : & X \text{ e } Y \text{ no son independientes} \end{cases}$$

y tiene como estadística

$$\chi^2 = \sum_{i=1}^K \sum_{j=1}^L \frac{(n_{ij} - E_{ij})^2}{E_{ij}}$$

donde  $E_{ij} = \frac{n_i \cdot n_j}{n}$  es la cantidad teórica esperada en el casillero  $(i, j)$  si las variables fueran independientes.

Es positivo y tiene distribución  $\chi^2$  con  $(L - 1) \times (K - 1)$  grados de libertad. Más cerca está de 0, más las variables tienden a ser independientes.

## Comparación de dos variables cualitativas

	a favor	indiferente	en contra	total
democrato	$\frac{285(202)}{500} = 115,14$	$\frac{285(150)}{500} = 85,5$	$\frac{285(148)}{500} = 84,36$	285
republicano	$\frac{215(202)}{500} = 86,86$	$\frac{215(150)}{500} = 64,5$	$\frac{215(148)}{500} = 63,64$	215
total	202	150	148	500

	a favor	indiferente	en contra	total
democrato	138 (115.4)	83 (85.5)	64 (84,36)	285
republicano	64 (86.86)	67 (64.50)	84 (63,64)	215
total	202	150	148	500

El cálculo del estadístico del test es:

$$\begin{aligned}\chi^2 &= \frac{(138 - 115.4)^2}{115,14} + \frac{(83 - 85,5)^2}{85,5} + \frac{(64 - 86,86)^2}{86,86} + \frac{(64 - 86,86)^2}{86,86} \\ &\quad + \frac{(67 - 64,5)^2}{64,5} + \frac{(84 - 63,64)^2}{63,64} = 22,152\end{aligned}$$

con grados de libertad  $(2 - 1)(3 - 1) = 2$

Como  $P(\chi^2 > 22.152) = 0.001 < 0.05$  rechazamos la hipótesis nula de independencia.

# ¿Como hacemos esto en R?

Veamos como realizar un test de independencia entre preferencia política y posición frente al plebiscito, usando la tabla de contingencia:

```
Crear la tabla de contingencia
tabla <- matrix(c(138, 83, 64,64, 67, 84), nrow = 2,byrow = TRUE)

Nombrar filas y columnas
rownames(tabla) <- c("Democrata", "Republicano")
colnames(tabla) <- c("A favor", "Indiferente", "En contra")

Ver la tabla
tabla

Test chi-cuadrado de independencia
test_chi <- chisq.test(tabla)
test_chi
test_chi$expected
if(test_chi$p.value < 0.05){
 "Se rechaza H0: hay asociación entre preferencia política y postura frente al plebiscito."
} else {
 "No se rechaza H0: no se detecta asociación."
}

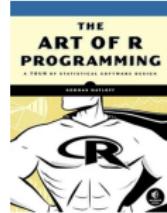
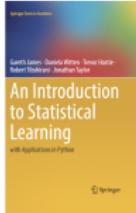
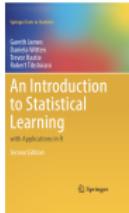
> tabla
 A_favor Indiferente En_contra
Democrata 138 83 64
Republicano 64 67 84
> # Test chi-cuadrado de independencia
> test_chi <- chisq.test(tabla)
> test_chi

Pearson's Chi-squared test

data: tabla
X-squared = 22.152, df = 2, p-value = 1.548e-05

> test_chi$expected
 A_favor Indiferente En_contra
Democrata 115.14 85.5 84.36
Republicano 86.86 64.5 63.36
```

# Referencias



- Venables, William N., et al. An introduction to R (2009)  
<https://cran.r-project.org/doc/manuals/R-intro.pdf>
- Wickham, Hadley; Cetinkaya-rundel, Mine; Grolemund, Garrett. R for data science: import, tidy, transform, visualize, and model data. O'Reilly Media, Inc., 2023.  
<https://r4ds.had.co.nz>
- Matloff, Norman. The art of R programming: A tour of statistical software design. No Starch Press, 2011.
- James, Gareth, Witten, Daniela, Hastie, Trevor, Tibshirani, Robert. An Introduction to Statistical Learning with Applications in R. Springer Texts in Statistics, 2021  
<https://www.statlearning.com>
- Philippe, Anne. "Notes de Cours sur le logiciel R." Université de Nantes, Laboratoire de Mathématiques Jean Leray (2018)