# *Data visualization and analysis in R*

Kristian Schultz

15.04.2025

www.sbi.uni-rostock.de

**SYSTEMS BIOLOGY BIOINFORMATICS ROSTOCK**

**Universität Rostock**

Traditio et Innovatio

gegründet 1419

# Data analysis with R



A software environment used to analyze **statistical information**, **graphical representation**, **reporting**, and **data modeling**
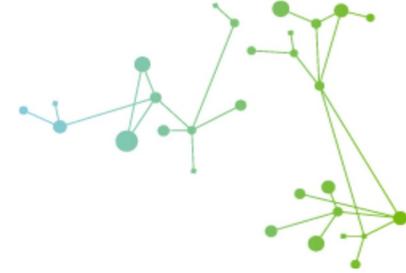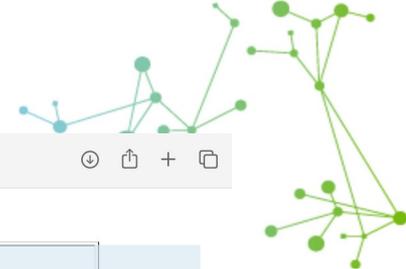
- *1992: Ross Ihaka and Robert Gentleman in the University of Auckland, New Zealand developed R.*

- *1995: The first version was released*

- *1997: CRAN (Comprehensive R Archive Network) was started*

- *2000: A stable beta version was released*

- *R has now thousands of packages, designed, maintained, and widely used by statisticians, biostatisticians, and geneticists*

# Installing R



**The Comprehensive R Archive Network**

**Download and Install R**

Precompiled binary distributions of the base system and contributed packages, **Windows and Mac** users most likely want one of these versions of R:

- Download R for Linux (Debian, Fedora/Redhat, Ubuntu)
- Download R for macOS
- Download R for Windows

R is part of many Linux distributions, you should check with your Linux package management system in addition to the link above.

**Source Code for all Platforms**

Windows and Mac users most likely want to download the precompiled binaries listed in the upper box, not the source code. The sources have to be compiled before you can use them. If you do not know what this means, you probably do not want to do it!

- The latest release (2022-10-31, Innocent and Trusting) R-4.2.2.tar.gz, read what's new in the latest version.

- Sources of R alpha and beta releases (daily snapshots, created only in time periods before a planned release).

- Daily snapshots of current patched and development versions are available here. Please read about new features and bug fixes before filing corresponding feature requests or bug reports.

- Source code of older versions of R is available here.

- Contributed extension packages

**Questions About R**

- If you have questions about R like how to download and install the software, or what the license terms are, please read our answers to frequently asked questions before you send an email.

**What are R and CRAN?**

R is 'GNU S', a freely available language and environment for statistical computing and graphics which provides a wide variety of statistical and graphical techniques: linear and nonlinear modelling, statistical tests, time series analysis, classification, clustering, etc. Please consult the R project homepage for further information.

https://cran.r-project.org

CRAN
Mirrors
What's new?
Search
CRAN Team

About R
R Homepage
The R Journal

Software
R Sources
R Binaries
Packages
Task Views
Other

Documentation
Manuals
FAQs
Contributed

# Installing R

# *Installing RStudio Desktop*



## https://posit.co

# Running R on Posit Cloud



https://posit.co

## *Packages in R*

- Collection of R functions, data and compiled code in a well-defined format to perform specific task.
- Allow to expend the functionality available to you in R programming.
- Packages are stored in a directory called the library.
- Installing packages:
    - using gui
    - using command line in console e.g., install.packages("ggplot2")
- Activate install package
    - library(<name of the package>) e.g., library(ggplot2)
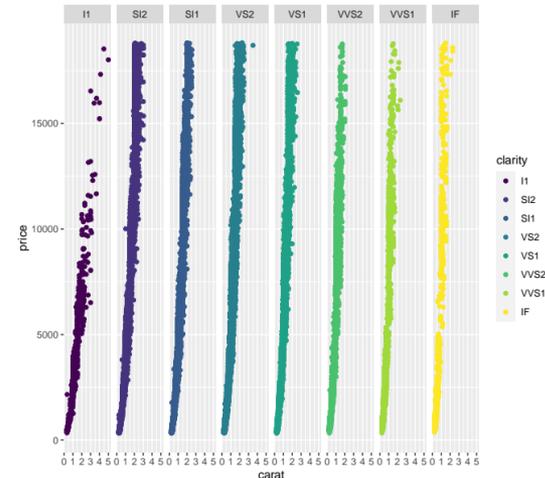


- Example from ggplot2 (ggplots comes with few example datasets, e.g., diamonds, iris)
    - qplot(data=diamonds, carat, price, colour=clarity, facets=.~clarity)

# The daily life of a data scientist..

# *Data visualization in R*

## Not this



*Twitter*
*@markabelan*

## But this

# Getting the data

# *Getting the data*



**kaggle**

+ Create

- Home
- Competitions
- Datasets
- Models
- Code
- Discussions
- Learn
- More

📋 Your Work

RECENTLY VIEWED

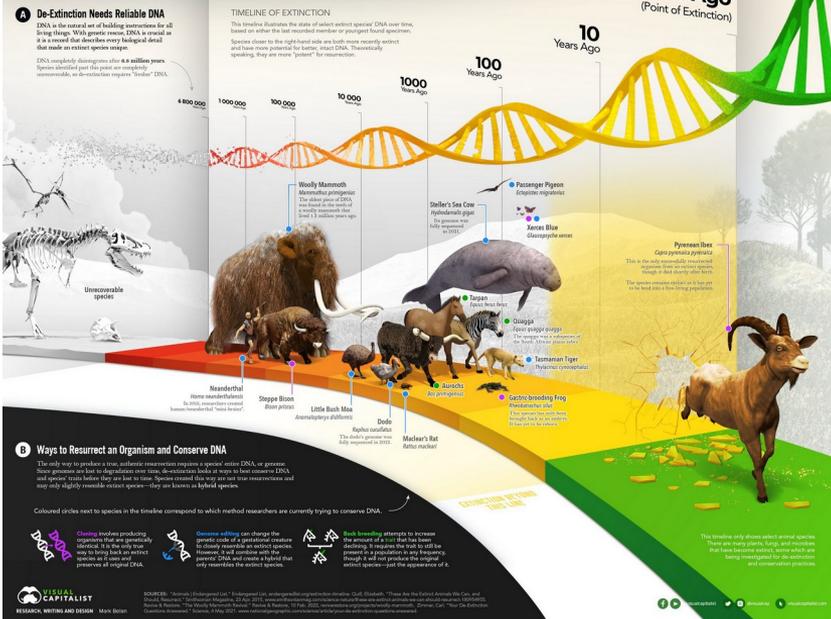- Diamond Analysis usin...
- Diamonds
- Iris Species
- Instagram fake accoun...
- Exploratory Data Analy...

🔍 Search

SHIVAM AGRAWAL · UPDATED 6 YEARS AGO

▲ 929    New Notebook    ⬇ Download (751 kB)  ⋮

# Diamonds

Analyze diamonds by their cut, color, clarity, price, and other attributes

---

Data Card    Code (415)    Discussion (8)

## About Dataset

### Context

This classic dataset contains the prices and other attributes of almost 54,000 diamonds. It's a great dataset for beginners learning to work with data analysis and visualization.

### Content

**price** price in US dollars (\$326--\$18,823)

**carat** weight of the diamond (0.2--5.01)

**cut** quality of the cut (Fair, Good, Very Good, Premium, Ideal)

**color** diamond colour, from J (worst) to D (best)

**clarity** a measurement of how clear the diamond is (I1 (worst), SI2, SI1, VS2, VS1, VVS2, VVS1, IF (best))

**x** length in mm (0--10.74)

**Usability** ⓘ
7.65

**License**
Unknown

**Expected update frequency**
Not specified

## Exploratory Data Analysis

- Exploratory Data Analysis (EDA) is an approach to analyzing datasets to summarize their main characteristics, often employing visual methods.

- The primary goal of EDA is to understand the underlying patterns, distributions, and relationships within the data.

- It involves techniques to identify outliers, detect patterns, test assumptions, and summarize the main features of the dataset.

### Exploratory vs Confirmatory Data Analysis

| EDA | CDA |
|---|---|
| • No hypothesis at first | • Start with hypothesis |
| • Generate hypothesis | • Test the null hypothesis |
| • Uses graphical methods (mostly) | • Uses statistical models |

# *Getting the data*



kaggle

+ Create

Home
Competitions
**Datasets**
Models
Code
Discussions
Learn
More

Your Work

RECENTLY VIEWED
Diamond Analysis usin...
Diamonds
Iris Species
Instagram fake accoun...
Exploratory Data Analy...

Search

SHIVAM AGRAWAL · UPDATED 6 YEARS AGO

▲ 929   New Notebook   ⬇ Download (751 kB)   ⬤   ⋮

# Diamonds

Analyze diamonds by their cut, color, clarity, price, and other attributes

Data Card   Code (415)   Discussion (8)

## About Dataset

### Context

This classic dataset contains the prices and other attributes of almost 54,000 diamonds. It's a great dataset for beginners learning to work with data analysis and visualization.

### Content

**price** price in US dollars (\$326--\$18,823)

**carat** weight of the diamond (0.2--5.01)

**cut** quality of the cut (Fair, Good, Very Good, Premium, Ideal)

**color** diamond colour, from J (worst) to D (best)

**clarity** a measurement of how clear the diamond is (I1 (worst), SI2, SI1, VS2, VS1, VVS2, VVS1, IF (best))
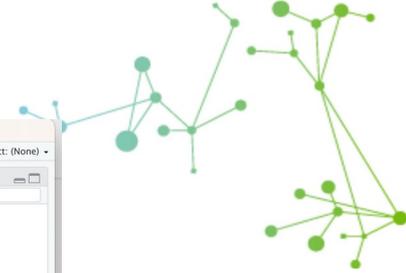
**x** length in mm (0--10.74)

**Usability** ⓘ
7.65

**License**
Unknown

**Expected update frequency**
Not specified

# Data visualization in R



RStudio

```
1  N <-10
2  counter<-0
3
4  for (x in rnorm(N))
5  {
6    if (x > -1 & x <1) {
7      counter<- counter +1
8    }
9  }
10 result <- counter/N
11
12 print(result)
13
14 seq
```

```
install.packages("tidyverse")
library(tidyverse)
library(ggplot)
install.packages("ggplot2")
library(ggplot2)
library(tidyverse)
mydata <-read.csv(file.choose())
install.packages("ggplot2")
install.packages("ggplot2")
library(ggplot2)
ggplot(data=mydata, aes(x=carat, y=price)) + geom_point()
```

```
> library(ggplot2)
```
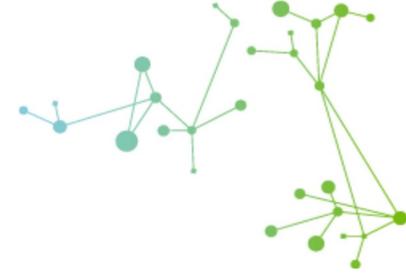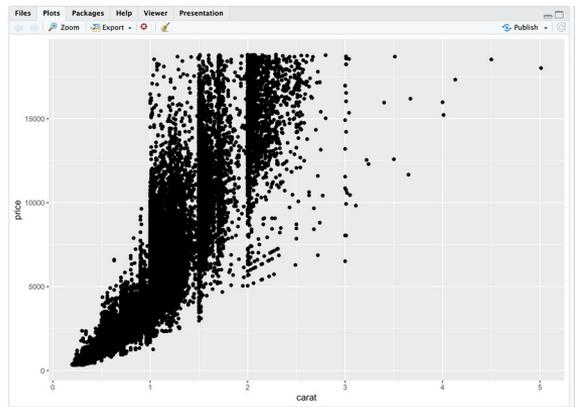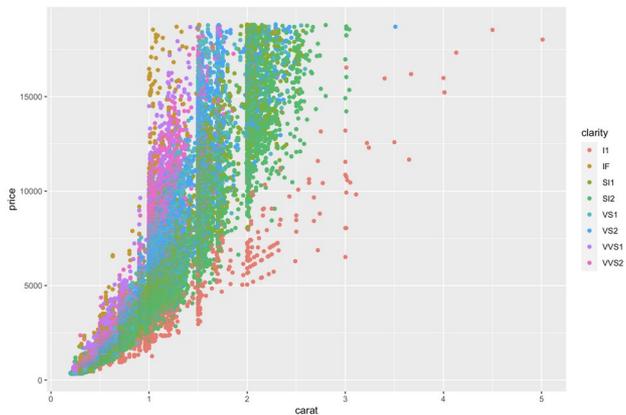
# Data visualization in R

# Data visualization in R
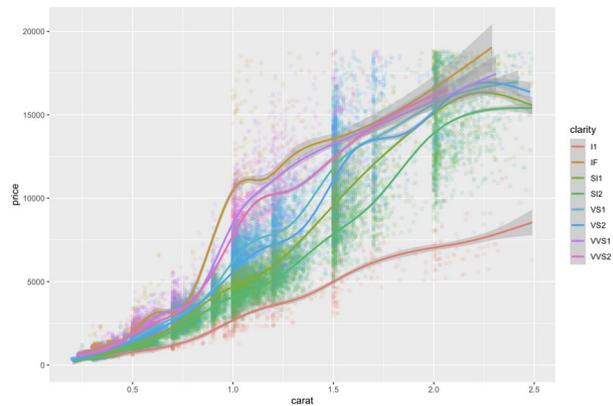
ggplot(data=mydata, aes(x=carat, y=price)) + geom_point()



ggplot(data=mydata, aes(x=carat, y=price, color=clarity)) + geom_point()



ggplot(data=mydata[mydata$carat<2.5,], aes(x=carat, y=price, color=clarity)) + geom_point(alpha=0.1) + geom_smooth()
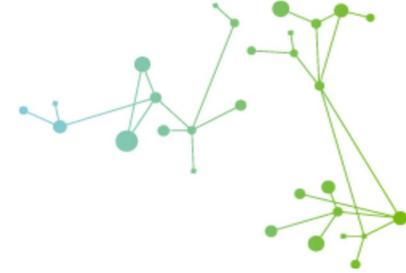
# Additional reading

https://carpentries-incubator.github.io/open-science-with-r/

https://marianattestad.com/blog

https://www.stats.ox.ac.uk/~evans/Rprog/LectureNotes.pdf

*Thank you for your attention!*