

Introduction to R (1)



Short course on modelling infectious disease dynamics in R

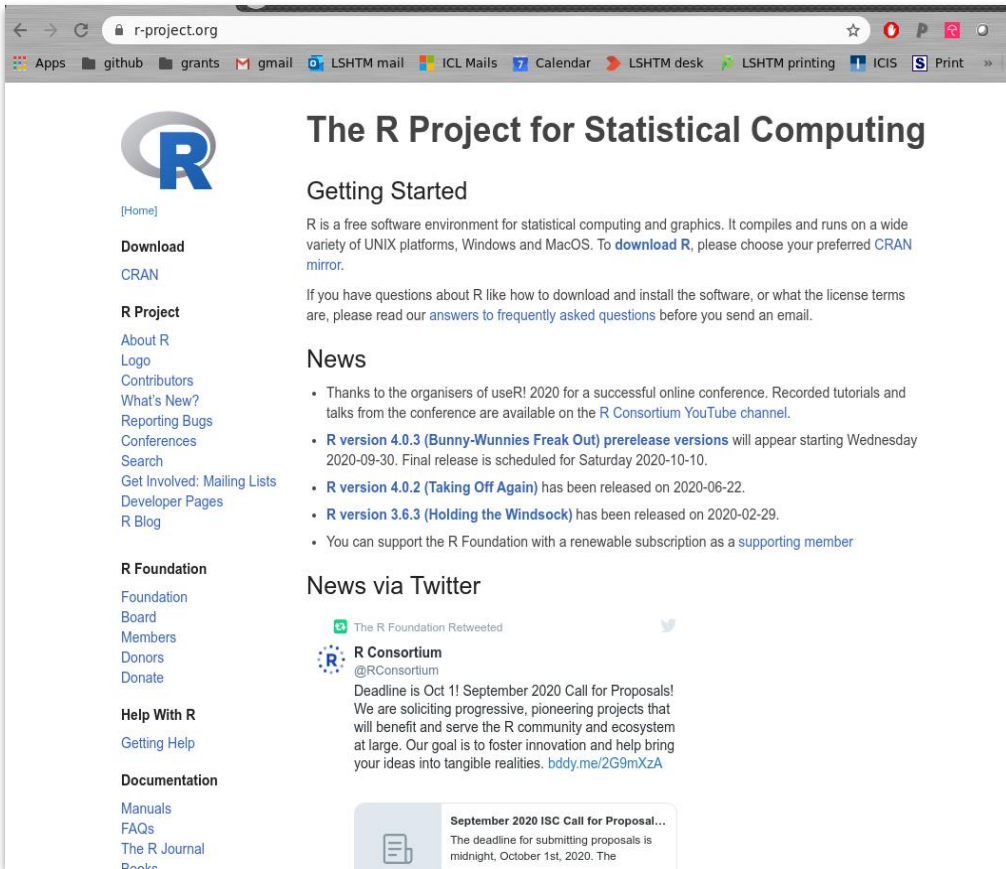
Ankara, Türkiye, September 2025

Dr Juan F Vesga

Aims of the session

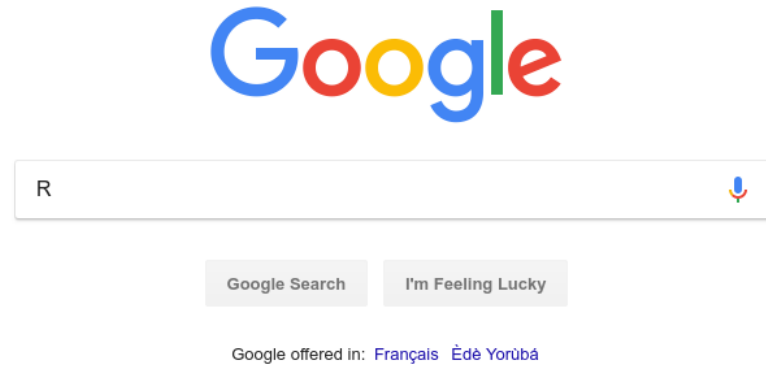
- Understand why R is a great option for epidemiology
- Get started with R and Rstudio
- Understand the value of R and Rstudio compared to other programming languages
- Familiarize with the R environment

What is ?



- a **free software** for **data analysis**
- an interpreted **programming language**, derived from 'S-plus'
- initially developed by **R. Ihaka** and **R. Gentleman** (1996)
- currently developed by the **R Core Team** (~20 people)
- **largest collection of tools for data analysis** (1,000s of contributors and specific packages)

Where can you get it?



- The **R** project: www.r-project.org
- archiving / distribution network CRAN: cran.r-project.org/mirrors.html
- available on Windows, MacOSX, Linux

What can you do with it?

- **basic statistics:** statistical tests, linear modelling, multivariate analysis
- **spatial statistics:** GIS, mapping, clustering
- **graph theory:** social sciences, network analysis, graph algorithm
- **genetics:** phylogenetic trees, genetic markers, genomics
- **epidemiology!**



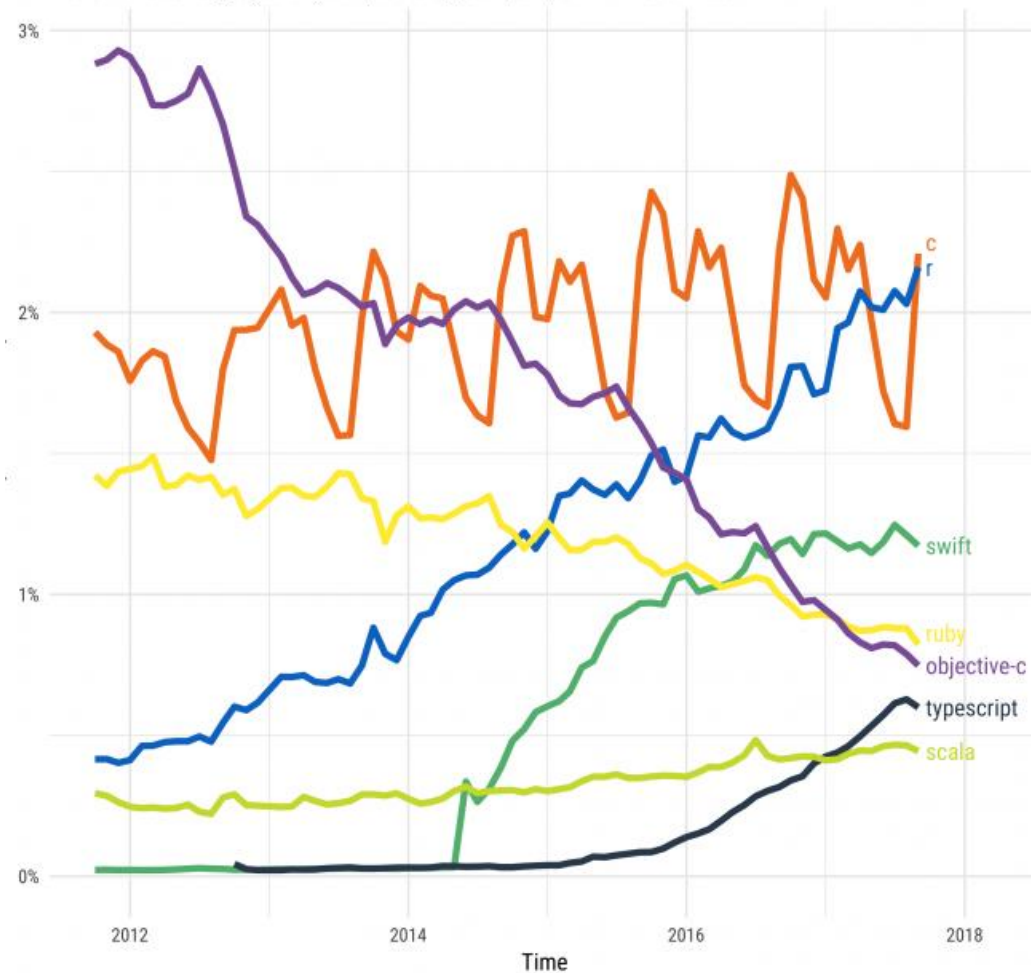
What does “free” mean?



- **Freedom** = *ability to make informed decisions*
- you don't pay for it
- the code is accessible by anyone
- anyone can use, modify and share the code

Stack Overflow Traffic to Programming Languages

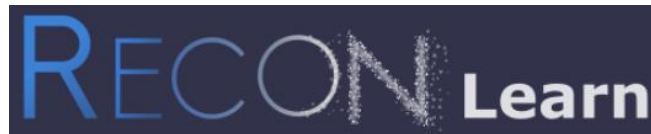
Based on visits to Stack Overflow questions from World Bank high-income countries.
The more-visited languages of Python, JavaScript, Java, C#, and PHP were omitted.



Growing use

and the 'epi' community

RECON



- R Epidemics Consortium (RECON)
 - NGO for the development of resources for health emergencies, with strong focus on outbreak response
 - 20+ packages for outbreak analytics
 - Share needs on [COVID-19 challenge](#)
 - www.repidemicsconsortium.org
- R4epis
 - Partnership between MSF + RECON
 - Analysis templates for field epidemiologists
 - R4epis.netlify.app
- RECON learn
 - Free / open training material by RECON
 - www.reconlearn.org

Alternatives to R

Python

- Both are frequently used in data science, and both are free
- Python is better for analyses that need scaling up
- R is better for investigative “one-time” analyses
- R more commonly used in medical research
- Python more commonly used in industry



Stata

- R and Stata most commonly used languages in medical research
- Commonly taught in medical statistics & epidemiology courses
- Requires a licence, but has official help
- Strong statistical analyses capabilities
- Less developed data visualization



Alternatives to R

SPSS

- Statistical Package for the Social Sciences
- Many researchers' first exposure to a statistical program
- Usually used as a “point and click”, although can write code
- Use is declining, and also requires a licence



Excel

- Very widely used
- Can be helpful for “having a look” at the data
- Generally not suitable for research-level inferential statistical analyses



Getting started



- get **R** for your system (download from CRAN)
- get a Graphical User Interface (GUI): **RStudio**, emacs + ESS, Tinn-R
- (or at least) get a text editor to write code: notepad++, emacs, vi, Tinn-R, ...

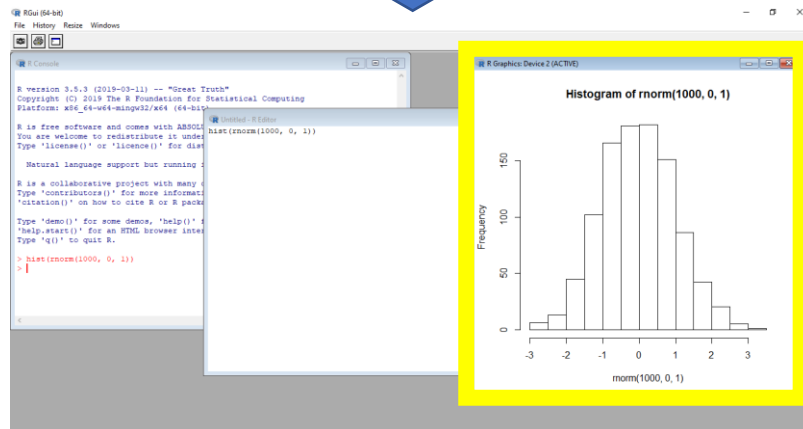
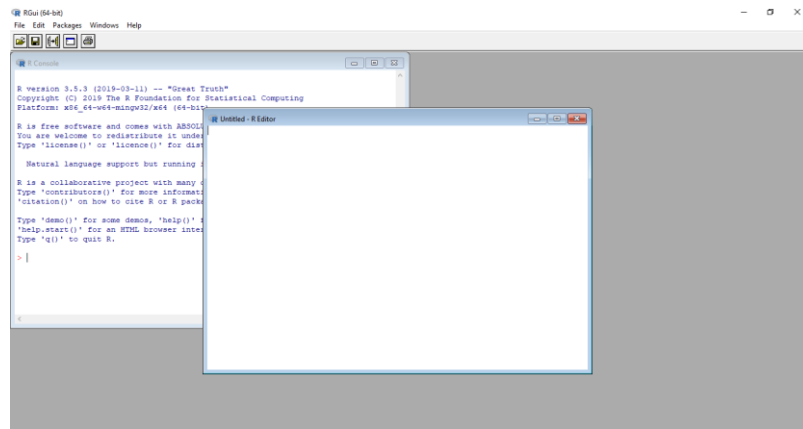
R Studio



- R Studio is an Integrated Development Environment (IDE) for R
 - i.e. it makes doing things in R easier and more structured
- Some advantages include:
 - An organised environment for your projects
 - Colour-coding, auto-completed brackets, auto-structuring
 - Keeping track of which packages have been installed, and which variables have been defined
 - Fixed quadrants for code and graphical outputs

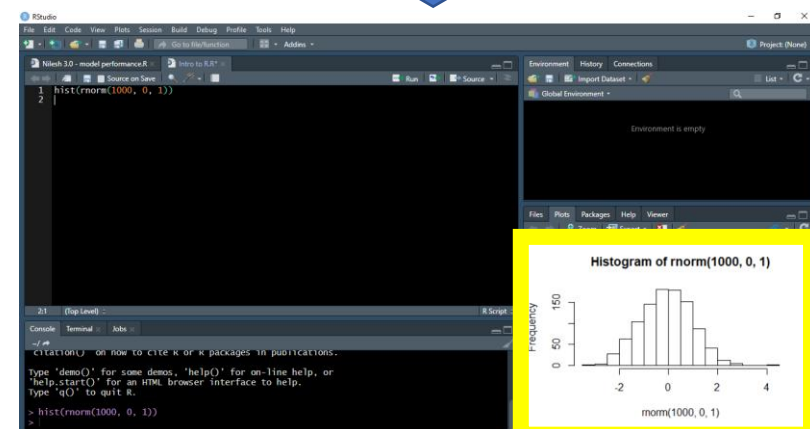
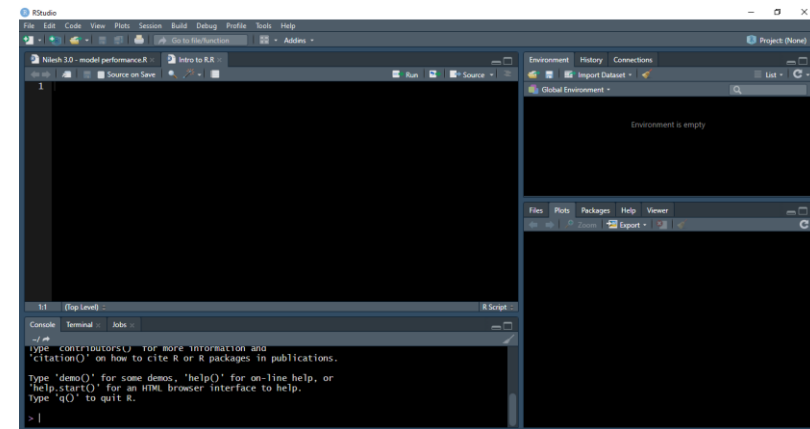
R (without R Studio)

graphs appear in a new window



R Studio

graphs appear in bottom-right quadrant



Getting familiar with RStudio

The screenshot displays the RStudio environment with the following components:

- Script Editor (Untitled1*):** Contains R code for loading packages, loading data, and plotting. The code is as follows:

```
1 library(incidence2)
2 library(outbreaks)
3 ?incidence
4 data(ebola_sim_clean, package = "outbreaks")
5 dat <- ebola_sim_clean$linelist
6
7
8 dat %>%
9   incidence(date_of_onset) %>%
10  plot()
11
12
```
- Environment Pane:** Shows the objects in the session. It lists 'dat' (5829 obs. of 11 variables) and 'ebola_sim_clean' (Large list (2 elements, 1.2 Mb)).
- Console:** Shows the execution of the code from the script. The output is as follows:

```
>
> library(incidence2)
> library(outbreaks)
> ?incidence
> data(ebola_sim_clean, package = "outbreaks")
> dat <- ebola_sim_clean$linelist
> dat %>%
+   incidence(date_of_onset) %>%
+   plot()
>
```
- Plots Pane:** Displays a line plot titled 'Daily Incidence'. The x-axis represents time from 2014-04-07 to 2015-04-07, and the y-axis represents 'Daily Incidence' from 0 to 40. The plot shows a sharp increase in incidence starting around 2014-07-07, peaking around 2014-10-07, and then gradually declining.

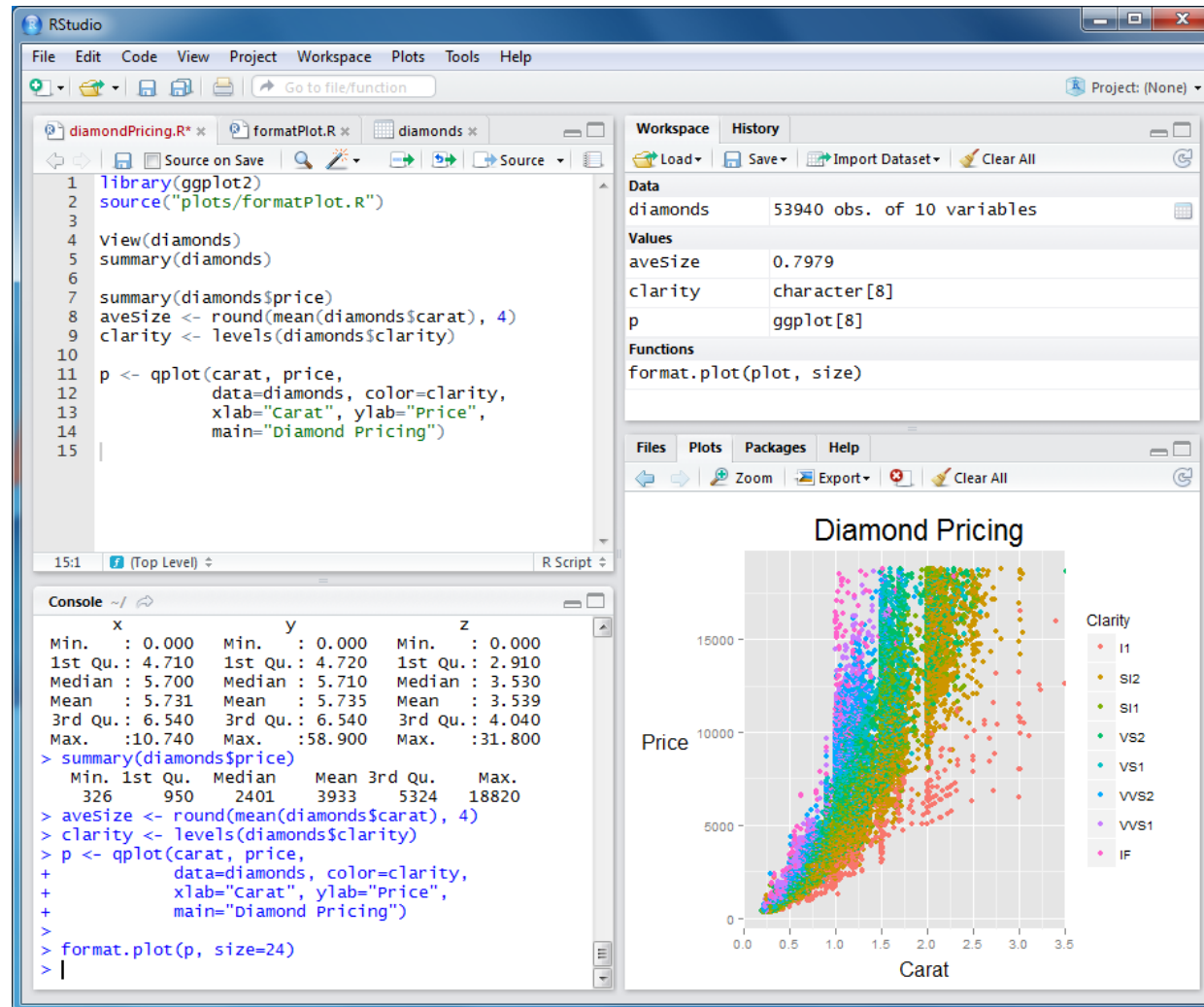
Script file:
write code here

R console:
Evaluate code
here

Overview of
objects in the
session

File browser,
Graphics, help,
packages loaded

And then...



Getting help



- `?foo` or `help("foo")`: access the help page of `foo`
- `??bar` or `help.search("foo")`: look for `foo` in help pages
- dedicated **mailing lists**: stat.ethz.ch/mailman/listinfo
- the **RECON forum**: <http://www.repidemicsconsortium.org/forum/>
- **google**

Someone's done it before!

- Stack **overflow** <https://stackoverflow.com/>
- Stack **exchange** (methods) <https://stats.stackexchange.com/>
- R-bloggers (chat, templates, examples) <https://www.r-bloggers.com/>

