

Iago Mosqueira

Introduction to R (and computing)

May 27, 2013

Why programming?

“Can one be a good data analyst without being a half-good programmer? The short answer to that is, ‘No’. The long answer to that is, ‘No!’.”

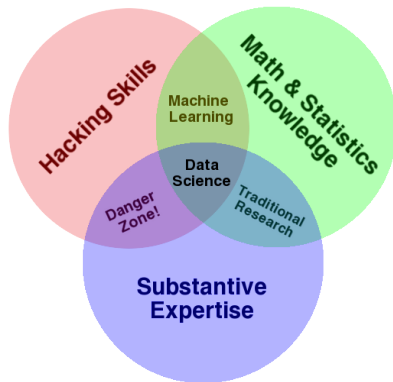
– Frank Harrell, 1999 S-PLUS User Conference, New Orleans (October 1999)

But this should be easy

“Managing fisheries is hard: it’s like managing a forest, in which the trees are invisible and keep moving around”

- Professor John Shepherd

Data analyst

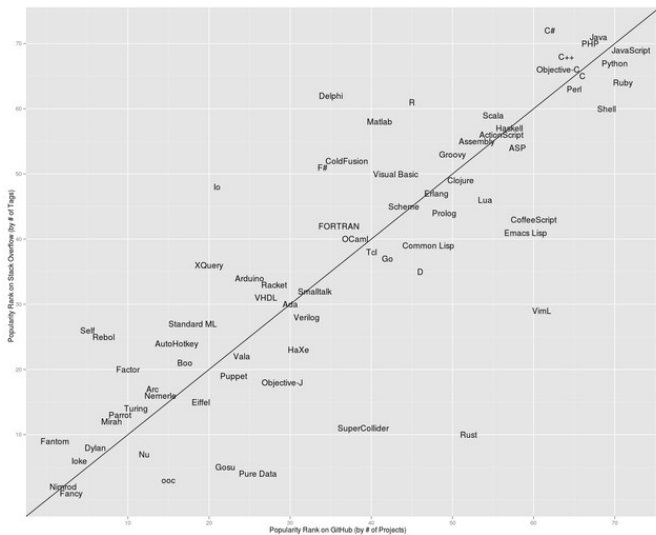


What is R



- Data analysis and statistics environment
- Interpreted computer language
- Open-source software project
- Active community of developers and practitioners
- Current version: 2.15.1, 2012-06-22, Roasted Marshmallows

Why R?



POLICYFORUM

COMPUTATIONAL SCIENCE

Troubling Trends in Scientific Software Use

Lucas N. Joppa,^{1*} Greg McNerny,¹² Richard Harper,¹ Lara Salido,³ Kenji Takeda,¹ Kenton O'Hara,¹ David Gavaghan,² Stephen Emmott¹

Software pervades every domain of science (1–3), perhaps nowhere more decisively than in modeling. In key scientific areas of great societal importance, models and the software that implement them define both how science is done and what science is done (4, 5). Across all science, this dependence has led to concerns around the need for open access to software (6, 7), centered on the reproducibility of research (1, 8–10). From fields such as high-performance computing, we learn key insights and best practices for how to develop, standardize, and implement software (11). Open and systematic approaches to the development of software are essential for all sciences. But for many scientists this is not sufficient. We

across all disciplines that are dependent upon a computational approach.

Surveying Species Distribution Modelers

We surveyed scientists across a single domain, species distribution modeling (SDM) (15) [see supplementary materials for details]. This strategic targeting separates our analysis from previous efforts in important ways, allowing an analysis spanning computational skill sets, while addressing the interplay between models and computation. Our ~400 respondents ranged from those who “find it difficult to use software” to those “very experienced and very technical.” Asking people to first identify with a scientific domain and addressing models and software through that

“Blind trust” is dangerous when choosing software to support research.

used “click-and-run” software with easy-to-manipulate user interfaces and dropped to 11% for those who used “syntax-driven” platforms. Further, 7, 9, and 18% of scientists cited “the developer is well-respected,” “personal recommendation,” and “recommendation from a close colleague,” respectively, as reasons for using software. Only 8% claimed they had validated software against other methods as a primary reason for choice; 79% expressed a desire to learn additional software and programming skills.

Many of these scientists rely on the fact that the software has appeared in a peer-reviewed article, recommendations, and personal opinion, as their reason for adopting software. This is scientifically misplaced, as the software code used to conduct the science

GPL v3.0

GNU General Public License

From Wikipedia, the free encyclopedia

"*GPL*" *redirects here*. For other uses, see *GPL (disambiguation)*.

The **GNU General Public License** (**GNU GPL** or simply **GPL**) is the most widely used^[5] [free software license](#). It was originally written by [Richard Stallman](#) for the [GNU Project](#).

The GPL is the first [copyleft](#) license for general use, which means that derived works can only be distributed under the same license terms. Under this philosophy, the GPL grants the recipients of a computer program the rights of the [free software definition](#) and uses copyleft to ensure the freedoms are preserved, even when the work is changed or added to. This is in distinction to [permissive free software licenses](#), of which the [BSD licenses](#) are the standard examples.

Contents [hide]

- 1 History
- 2 Versions
 - 2.1 Version 1
 - 2.2 Version 2
 - 2.3 Version 3
- 3 Terms and conditions
 - 3.1 Copyleft
- 4 Licensing and contractual issues
- 5 Copyright holders
- 6 Linking and derived works
 - 6.1 Libraries
 - 6.2 Communicating and bundling with non-GPL programs
- 7 The GPL in court
- 8 Compatibility and multi-licensing
 - 8.1 Multi-licensing
- 9 Adoption
- 10 Use for text and other media

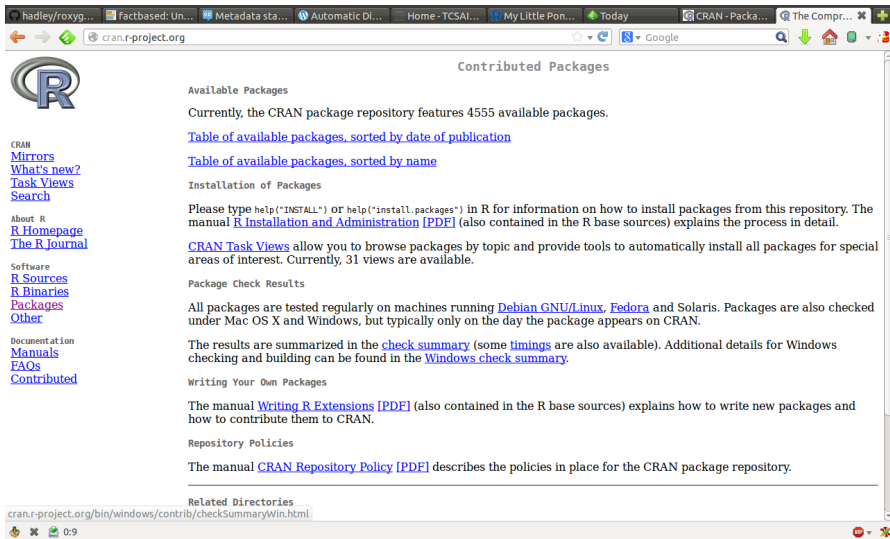
GNU General Public License



Free as in Freedom

GNU GPLv3 Logo

Author	Free Software Foundation
Version	3
Publisher	Free Software Foundation, Inc.
Published	29 June 2007
DFSG compatible	Yes ^[1]
FSF approved	Yes ^[2]
OSI approved	Yes ^[3]
Copyleft	Yes ^{[2][4]}
Linking from code with a different license	No (except for linking GNU AGPLv3 with GNU GPLv3 – see section)
Website	www.gnu.org/licenses



Contributed Packages

Available Packages

Currently, the CRAN package repository features 4555 available packages.

[Table of available packages, sorted by date of publication](#)

[Table of available packages, sorted by name](#)

Installation of Packages

Please type `help("INSTALL")` or `help("install.packages")` in R for information on how to install packages from this repository. The manual [R Installation and Administration \[PDF\]](#) (also contained in the R base sources) explains the process in detail.

[CRAN Task Views](#) allow you to browse packages by topic and provide tools to automatically install all packages for special areas of interest. Currently, 31 views are available.

Package Check Results

All packages are tested regularly on machines running [Debian GNU/Linux](#), [Fedora](#) and Solaris. Packages are also checked under Mac OS X and Windows, but typically only on the day the package appears on CRAN.

The results are summarized in the [check summary](#) (some [timings](#) are also available). Additional details for Windows checking and building can be found in the [Windows check summary](#).

Writing Your Own Packages

The manual [Writing R Extensions \[PDF\]](#) (also contained in the R base sources) explains how to write new packages and how to contribute them to CRAN.

Repository Policies

The manual [CRAN Repository Policy \[PDF\]](#) describes the policies in place for the CRAN package repository.

Related Directories

[cran.r-project.org/bin/windows/contrib/checkSummaryWin.html](#)

CRAN

[Mirrors](#)

[What's new?](#)

[Task Views](#)

[Search](#)

About R

[R Homepage](#)

[The R Journal](#)

Software

[R Sources](#)

[R Binaries](#)

[Packages](#)

[Other](#)

Documentation


[Manuals](#)

[FAQs](#)

[Contributed](#)

Task views


Blog with Knitr and Je... Welcome | Travelling... iagomosqueira/ICESSAI Home - TCSAI2012 Oráculo manual y art... The Comprehensive R...
cran.r-project.org ggplot2



CRAN Task Views

CRAN	Bayesian	Bayesian Inference
Mirrors	ChemPhys	Chemometrics and Computational Physics
What's new?	ClinicalTrials	Clinical Trial Design, Monitoring, and Analysis
Task Views	Cluster	Cluster Analysis & Finite Mixture Models
Search	DifferentialEquations	Differential Equations
	Distributions	Probability Distributions
	Econometrics	Computational Econometrics
	Environmetrics	Analysis of Ecological and Environmental Data
	ExperimentalDesign	Design of Experiments (DoE) & Analysis of Experimental Data
About R	Finance	Empirical Finance
R Homepage	Genetics	Statistical Genetics
The R Journal	Graphics	Graphic Displays & Dynamic Graphics & Graphic Devices & Visualization
	HighPerformanceComputing	High-Performance and Parallel Computing with R
Software	MachineLearning	Machine Learning & Statistical Learning
R Sources	MedicalImaging	Medical Image Analysis
R Binaries	Multivariate	Multivariate Statistics
Packages	NaturalLanguageProcessing	Natural Language Processing
Other	OfficialStatistics	Official Statistics & Survey Methodology
	Optimization	Optimization and Mathematical Programming
Documentation	Pharmacokinetics	Analysis of Pharmacokinetic Data
Manuals	Phylogenetics	Phylogenetics, Especially Comparative Methods
FAQs	Psychometrics	Psychometric Models and Methods
Contributed	ReproducibleResearch	Reproducible Research
	Robust	Robust Statistical Methods
	SocialSciences	Statistics for the Social Sciences
	Spatial	Analysis of Spatial Data
	Survival	Survival Analysis
	TimeSeries	Time Series Analysis


RStudio

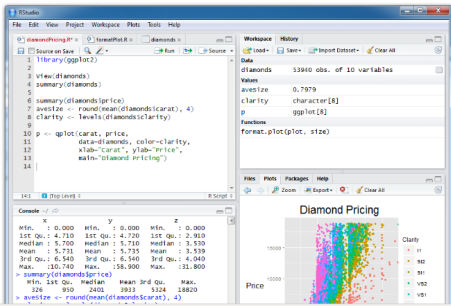


[Home](#) [Screenshots](#) [Download](#) [Docs](#) [Support](#) [Development](#) [Blog](#)

Welcome to RStudio

RStudio™ is a free and open source integrated development environment (IDE) for [R](#). You can run it on your desktop (Windows, Mac, or Linux) or even over the web using RStudio Server.

**Download RStudio**
for Windows, Mac or Linux




The screenshot shows the RStudio IDE interface. The top menu bar includes File, Edit, View, Project, Workspace, Plots, Tools, and Help. The main editor window displays a script for analyzing diamond data using ggplot2. The console window at the bottom shows the output of the summary and aveSize functions. The Plots pane on the right displays a scatter plot titled 'Diamond Pricing' with 'Price' on the y-axis and 'Clarity' on the x-axis, showing a positive correlation between the two variables.

```
1 library(ggplot2)
2
3 View(diamonds)
4 summary(diamonds)
5
6 summary(diamonds$price)
7 aveSize <- round(mean(diamonds$carat, 4)
8 clarity <- levels(diamonds$clarity)
9
10 p <- ggplot(carat, price,
11 data=diamonds, color=clarity,
12 xlab="carat", ylab="price",
13 main="diamond pricing")
14
```

Console output:

```
summary(diamonds$price)
  min. 1st Qu.  Median    mean 3rd Qu.  Max.
  326.   910.   2401.   3913.   5324.  18820.

aveSize <- round(mean(diamonds$carat, 4)
```



This is a smaller version of the RStudio IDE screenshot, showing the same script, console output, and scatter plot as the larger image.

Screencast
RStudio in 2 minutes

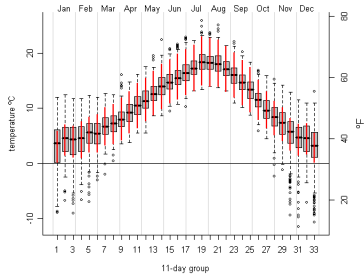
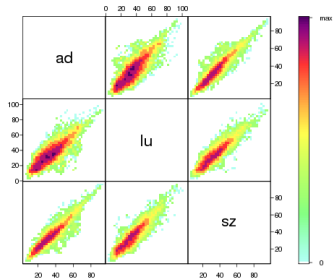
Basic features

- Numerous procedures (algebra, matrix, stats)
- Named storage (everything is an object)
- Functions
- Classes and methods (S3, S4)
- Special values (NA, NaN, Inf, NULL)
- Logical objects and boolean algebra
- `basic_features.R`

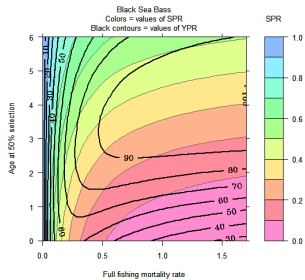
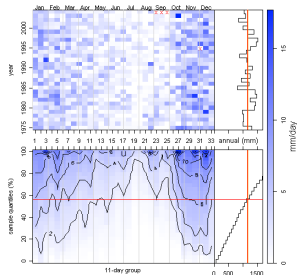
What else can it do?

- Data handling and storage
- Matrix algebra
- Regular expressions
- Statistics!
- OOP
- Programming
- Graphics

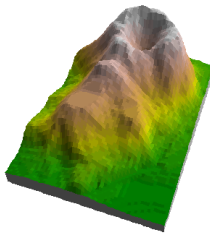
Eye candy



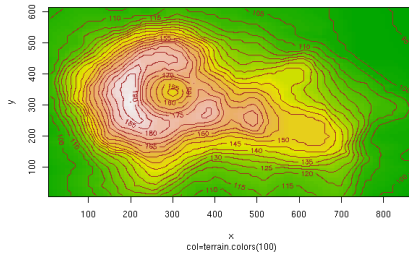
Eye candy



Eye candy



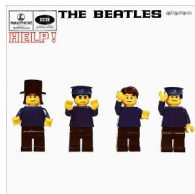
Maunga Whau Volcano



What doesn't it do

- No DB, but connections (SQL, NoSQL, Spreadsheets)
- No GUI, but IDE & GUI toolsets - CLI
- Slow, but C/C++, HPC
- No commercial support, but community
- Think for you

Help!



- Help for each function and data type
- ?mean
- ??mean
- ?help
- <http://rseek.org>
- stackoverflow, <http://stackoverflow.com/questions/tagged/r>
- Mailing lists

- Stock assessment and provision of management advice
 - ▶ Well tested, robust methods
 - ▶ Open to detailed inspection
- Data and model validation through simulation
- Risk analysis
- Capacity development & education
- Promote collaboration and openness in quantitative fisheries science
 - ▶ Open source
 - ▶ Community involvement
 - ▶ R as lingua franca
- Support the development of new models and methods
 - ▶ Extensible toolset
 - ▶ Links to other tools (ADMB, BUGS, ...)


start - FLR Project

flr-project.org/doku.php?id=start

Google

FLR PROJECT

TABLE OF CONTENTS




Training Course

Intro to R and FLR for Fisheries


3-5 JULY 2012 – Varese, Italy

Registration now open at [FishReg EC JRC](#)



The FLR library is a collection of tools in the [R](#) statistical language that facilitates the construction of bio-economic simulation models of fisheries and ecological systems. It is a generic toolbox, but is specifically suited for the construction of simulation models for evaluations of fisheries management strategies. The FLR library is under development by researchers across a number of laboratories and universities.

Love it !! Take me to the [Table of Contents](#) !!




You can install the **stable releases** for [R >= 2.13.0](#) with

```
install.packages(repos="http://flr-project.org/R")
```

or you can install the **development packages** for [R >= 2.14.0](#) with

```
install.packages(repos="http://flr-project.org/Rdevel")
```



You can get help by subscribing our [mailing list](#) and posting your doubt, or reading the **tutorials** on the ["Teach yourself FLR" wiki](#).

Tools of the trade

- Version Control Systems
 - ▶ CVS
 - ▶ SVN
 - ▶ git
- Editors & IDEs
- Literate Programming
 - ▶ Sweave
 - ▶ knitr
- Validation, Verification and Testing (VV&T)

Sexy data analysis



Setting up R & RStudio

- <http://cran.r-project.org>
- <http://rstudio.org>