

Barry Grant bjgrant@umich.edu http://thegrantlab.org

What is R?

R is a freely distributed and widely used programing **language** and **environment** for <u>statistical computing</u>, <u>data analysis</u> and <u>graphics</u>.



R provides an unparalleled interactive environment for data analysis.

It is script-based (*i.e.* driven by computer code) and not GUI-based (point and click with menus).

	4,	Salar	dbox	(79)	
R					

R version 3.2.2 (2015-08-14) -- "Fire Safety" Copyright (C) 2015 The R Foundation for Statistical Computing Platform: x86.64-apple-darwin13.4.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY. You are welcome to redistribute it under certain conditions. Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors. Type 'contributors()' for more information and 'citation()' on how to cite R or R pockages in publications.

iype 'demo()' for some demos, 'help()' for on-line help, or 'help.start()' for an HTML browser interface to help. iype 'q()' to quit R.



4. sandbox (R)

pico:sandbox> R

R version 3.2.2 (2015-08-14) -- "Fire Safety" Copyright (C) 2015 The R Foundation for Statistical Computing Platform: x86_64-apple-darwin13.4.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY. You are welcome to redistribute it under certain conditions. Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors. Type 'contributors()' for more information and 'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or 'help.start()' for an HTML browser interface to help. Type 'q()' to quit R.



pico:sandbox> R

Type "R" in your terminal

R version 3.2.2 (2015-08-14) -- "Fire Safety" Copyright (C) 2015 The R Foundation for Statistical Computing Platform: x86_64-apple-darwin13.4.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY. You are welcome to redistribute it under certain conditions. Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors. Type 'contributors()' for more information and 'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or 'help.start()' for an HTML browser interface to help. Type 'q()' to quit R.



pico:sandbox> R

Type "R" in your terminal

R version 3.2.2 (2015-08-14) -- "Fire Safety" Copyright (C) 2015 The R Foundation for Statistical Computing Platform: x86_64-apple-darwin13.4.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY. You are welcome to redistribute it under certain conditions. Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors. Type 'contributors()' for more information and 'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or 'help.start()' for an HTML browser interface to help. Type 'q()' to quit R.

This is the R prompt



pico:sandbox> R

Type "R" in your terminal

R version 3.2.2 (2015-08-14) -- "Fire Safety" Copyright (C) 2015 The R Foundation for Statistical Computing Platform: x86_64-apple-darwin13.4.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY. You are welcome to redistribute it under certain conditions. Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors. Type 'contributors()' for more information and 'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or 'help.start()' for an HTML browser interface to help. Type 'q()' to quit R.

This is the R prompt: Type q() to quit!

What R is NOT

A performance optimized software library for incorporation into your own C/C++ etc. programs.

A molecular graphics program with a slick GUI.

Backed by a commercial guarantee or license.

Microsoft Excel!

What about Excel?

- Data manipulation is easy
- Can see what is happening
- But: graphics are poor
- Looping is hard
- Limited statistical capabilities
- Inflexible and irreproducible



Use the right tool!

• There are many many things Excel just cannot do!



Christie Bahlai @cbahlai · 2h Weekly plug for scripted analyses:

Coauthor: "Can you change x,y,z about the analysis?" Me [not crying]: "Yes." [changes 2 lines of code]



Rule of thumb: Every analysis you do on a dataset will have to be redone 10–15 times before publication. Plan accordingly! Why use R? Productivity Flexibility Designed for data analysis

IEEE 2016 Top Programming Languages

Lan	guage Rank	Types	Spectrum Ranking
1.	С	🚺 🖵 🌲	100.0
2.	Java	⊕ 🗋 🖵	98.1
З.	Python	\bigoplus \Box	98.0
4.	C++	📮 🖵 🌲	95.9
5.	R	Ţ	87.9
6.	C#	🌐 🗋 🖵	86.7
7.	PHP	\oplus	82.8
8.	JavaScript	\oplus .	82.2
9.	Ruby		74.5
10.	Go	\bigoplus \Box	71.9

http://spectrum.ieee.org/computing/software/the-2016-top-programming-languages

R and Python: The Numbers

Popularity Rankings



Redmonk ranking, comparing the relative performance of programming languages on GitHub and Stack Overflow (September 2012 and January 2013, 2014, 2015)



- R is the "lingua franca" of data science in industry and academia.
- Large user and developer community.
 - As of Aug 1st 2016 there are 8811 add on R packages on <u>CRAN</u> and 1211 on <u>Bioconductor</u> - more on these later!
- Virtually every statistical technique is either already built into R, or available as a free package.
- Unparalleled exploratory data analysis environment.

Modularity	Core R functions are modular and work well with others
Interactivity	R offers an unparalleled exploratory data analysis environment
Infrastructure	Access to existing tools and cutting- edge statistical and graphical methods
Support	Extensive documentation and tutorials available online for R
R Philosophy	Encourages open standards and reproducibility

Modularity	Core R functions are modular and work well with others
Interactivity	R offers an unparalleled exploratory data analysis environment
Infrastructure	Access to existing tools and cutting- edge statistical and graphical methods
Support	Extensive documentation and tutorials available online for R
R Philosophy	Encourages open standards and reproducibility

Modularity

R was designed to allow users to interactively build complex workflows by interfacing smaller '**modular' functions** together.

get.seq() hmmer() pdbaln() pdbfit() pca() plot()

An alternative approach is to write a **single complex program** that takes raw data as input, and after hours of data processing, outputs publication figures and a final table of results.

All-in-one custom 'Monster' program

'Scripting' approach

Another common approach to bioinformatics data analysis is to write individual scripts in Perl/ Python/Awk/C etc. to carry out each subsequent step of an analysis



This can offer many advantages but can be challenging to make robustly modular and interactive.

Interactivity & exploratory data analysis

Learning R will give you the freedom to explore and experiment with your data.

"Data analysis, like experimentation, must be considered as a highly interactive, iterative process, whose actual steps are selected segments of a stubbily branching, tree-like pattern of possible actions". [J. W. Tukey]

Interactivity & exploratory data analysis

Learning R will give you the freedom to explore and experiment with your data.

"Data analysis, like experimentation, must be considered as a highly interactive, iterative process, whose actual steps are selected segments of a stubbily branching, tree-like pattern of possible actions". [J. W. Tukey]

Bioinformatics data is intrinsically **high dimensional** and frequently 'messy' requiring **exploratory data analysis** to find patterns - both those that indicate interesting biological signals or suggest potential problems.



R Features = functions()



How do we use R?

Two main ways to use R

_ 0 _ RStudio 4. sandbox (R) File Edit View Project Workspace Plots Tools Help pico:sandbox> R 0 damondfriding.R* x 9 format/Rot.R x Workspace History diamonds × 🔒 🔄 Source on Save 🛛 💁 🖉 -🚰 Load + 🗧 Save + 💽 Import Datacet + 🥑 Clear All -Run 😁 -Source -1 library(ggplot2) Data R version 3.2.2 (2015-08-14) -- "Fire Safety" diamonds \$3940 obs. of 10 variables 3 view(diamonds) Copyright (C) 2015 The R Foundation for Statistical Computing Values 4 summary(diamonds) 0.7979 avestize Platform: x86_64-apple-darwin13.4.0 (64-bit) 6 summary(diamondsSprice) avesize <- round(mean(diamonds)(carat) 3- Workspace and 8 cl. Code Edito R is free software and comes with ABSOLUTELY NO WARRANTY. 10 p 11 History You are welcome to redistribute it under certain conditions. 12 13 main-"Diamond Pricing" Type 'license()' or 'licence()' for distribution details. 14 Files Plots Fackages Help 💠 💿 🤰 Zoom 🛛 🗮 Export - 🔍 🅑 Cear All R Societ 1 14.1 C (Top Level) (Natural language support but running in an English locale Diamond Pricing : 0.000 : 0.000 Min. : 0.000 Min. E CARACTERIST R is a collaborative project with many contributors. 1st Qu.: 4.710 1st Qu.: 4.720 1st Qu.: 2.910 Median : 5.700 Median : 5.710 Median : 3.530 Type 'contributors()' for more information and Mean 3.510 Plots and file 3rd 4.040 'citation()' on how to cite R or R packages in publications. R Consol 1.800 Max. 148.7 Price 326 2401 1913 \$324 18820 VP Type 'demo()' for some demos, 'help()' for on-line help, or veSize <- round(mean(diamonds\$carat), 4) clarity <- levels(diamondsSclarity) V/82 'help.start()' for an HTML browser interface to help. p <- qplot(carat, price, Wist data-diamonds, color-clarity, xlab="carat", ylab="Price", main="Diamond Pricing") Type 'q()' to quit R. format.plot(plot-p, size-23) Carat 1. Terminal 2. RStudio

We will use RStudio today



Lets get started.

DO INTOLINGO, IT



Some simple R commands



Learning a new language is hard!

Error Messages

Sometimes the commands you enter will generate errors. Common beginner examples include:

- Incomplete brackets or quotes *e.g.* ((4+8)*20 <enter>
 - This eturns a + here, which means you need to enter the remaining bracket R is waiting for you to finish your input. Press <ESC> to abandon this line if you don't want to fix it.
- Not separating arguments by commas *e.g.*

plot(1:10 col="red")

• Typos including miss-spelling functions and using wrong type of brackets *e.g.*

exp{4}

Your turn! http://tinyurl.com/bioboot-R1

Doir Louisoir

Topics Covered:

Calling Functions Getting help in R Vectors and vectorization Workspace and working directory RStudio projects

Side-note: Use the code editor for R scripts



R scripts

- A simple text file with your R commands (*e.g.* day4.r) that contains your R code for one complete analysis
- Scientific method: complete record of your analysis
- Reproducible: rerunning your code is easy for you or someone else
- In RStudio, select code and type <ctrl+enter> to run the code in the R console
- Key point: <u>Save your R script!</u>

Side-note: RStudio shortcuts



Rscript: Third way to use R



Rscript --vanilla my_analysis.R

3. Rscript

From the command line! > Rscript --vanilla my_analysis.R # or within R: source(my_analysis.R)

Side-Note: R workspaces

- When you close RStudio, SAVE YOUR .R SCRIPT
- You can also save data and variables in an R workspace, but this is generally not recommended
- Exception: working with an enormous dataset
- Better to start with a clean, empty workspace so that past analyses don't interfere with current analyses
- rm(list = ls()) clears out your workspace
- You should be able to reproduce everything from your R script, so <u>save your R script, don't save your workspace!</u>

Help from within R

- Getting help for a function
- > help("log")
- > ?log
- Searching across packages
- > help.search("logarithm")
- Finding all functions of a particular type
- > apropos("log")

[7] "SSlogis" "as.data.frame.logical" "as.logical"
 "as.logical.factor" "dlogis" "is.logical"
[13] "log" "log10" "log1p" "log2" "logLik" "logb"
[19] "logical" "loglin" "plogis" "print.logLik" "qlogis"
 "rlogis"

R: Logarithms and Exponentials Find in Topic

log {base}

R Documentation

Logarithms and Exponentials

Description What the function does in general terms

log computes logarithms, by default natural logarithms, log10 computes common (i.e., base 10) logarithms, and log2 computes binary (i.e., base 2) logarithms. The general form log (x, base) computes logarithms with base base.

log1p (x) computes log(1+x) accurately also for |x| << 1 (and less accurately when x is approximately -1).

exp computes the exponential function.

expm1 (x) computes exp(x) - 1 accurately also for |x| << 1.

Usage How to use the function

```
log(x, base = exp(1))
logb(x, base = exp(1))
log10(x)
log2(x)
```

log1p(x)

exp(x) expml(x)

Arguments What does the function need

a numeric or complex vector.

base a positive or complex number: the base with respect to which logarithms are computed. Defaults to e=exp(1).

Details

All except logb are generic functions: methods can be defined for them individually or via the <u>Math</u> group generic.

log10 and log2 are only convenience wrappers, but logs to bases 10 and 2 (whether computed via log or the wrappers) will be computed more efficiently and accurately where supported by the OS. Methods can be set for them individually (and otherwise methods for log will be used).

logb is a wrapper for log for compatibility with S. If (S3 or S4) methods are set for log they will be dispatched. Do not set S4 methods on logb itself.

All except log are primitive functions.

R: Logarithms and Exponentials Find in Topic

Value What does the function return

A vector of the same length as x containing the transformed values. log(0) gives -Inf, and log(x) for negative values of x is NaN. exp(-Inf) is 0.

?log

For complex inputs to the log functions, the value is a complex number with imaginary part in the range [-pi, pi]: which end of the range is used might be platform-specific.

S4 methods

exp, expm1, log, log10, log2 and log1p are S4 generic and are members of the <u>Math</u> group generic.

Note that this means that the S4 generic for log has a signature with only one argument, x, but that base can be passed to methods (but will not be used for method selection). On the other hand, if you only set a method for the Math group generic then base argument of log will be ignored for your class.

Source

log1p and expm1 may be taken from the operating system, but if not available there are based on the Fortran subroutine dlnrel by W. Fullerton of Los Alamos Scientific Laboratory (see <u>http://www.netlib.org/slatec/fnlib/dlnrel.f</u> and (for small x) a single Newton step for the solution of log1p(y) = x respectively.

References

Becker, R. A., Chambers, J. M. and Wilks, A. R. (1988) The New S Language. Wadsworth & Brooks/Cole. (for log, log10 and exp.)

Chambers, J. M. (1998) Programming with Data. A Guide to the S Language. Springer. (for logb.)

See Also Discover other related functions

Trig, sqrt, Arithmetic.

Examples Sample code showing how it works

log(exp(3)) log10(1e7) # = 7

```
x <- 10^-(1+2*1:9)
cbind(x, log(1+x), log1p(x), exp(x)-1, expm1(x))
```

[Package base version 3.0.1 Index]

Optional Exercise

Use R to do the following. Create a new script to save your work and code up the following four equations:

$$1 + 2(3 + 4)$$
$$\ln(4^{3} + 3^{2+1})$$
$$\sqrt{(4+3)(2+1)}$$
$$\left(\frac{1+2}{3+4}\right)^{2}$$

Learning Resources

- TryR. An excellent interactive online R tutorial for beginners.
 < <u>http://tryr.codeschool.com/</u> >
- RStudio. A well designed reference card for RStudio.
 < <u>https://help.github.com/categories/bootcamp/</u> >
- DataCamp. Online tutorials using R in your browser.
 < <u>https://www.datacamp.com/</u> >
- R for Data Science. A new O'Reilly book that will teach you how to do data science with R, by Garrett Grolemund and Hadley Wickham.

< <u>http://r4ds.had.co.nz/</u> >