

# Data Analysis Tools in an Age of Computer Technology

John Maindonald

Centre for Mathematics and Its Applications,  
Australian National University

May 28, 2013



# The Two R's of R

Ross Ihaka and Robert Gentleman developed the initial version of R at Auckland University in NZ. They set out to create an open source version of the earlier AT&T S language.

# From Auckland to Canberra



# Auckland University and Mt Albert Research Centre



# Computing for Statistics — Then & Now

---

1970s

---

As pencil & paper, but faster

Rely on statistical theory

(eg, Central Limit Theorem)

Data was similar to 1930s

Many incompatible programs

---

2013

---

Improve on pencil & paper

Theory, plus computation

(simulation, bootstrap, ...)

New types of data; 'big' data

Incompatible programs + R

---

# Why R? — Data Analysis

- ▶ Designed to encourage good statistical practice
  - \_ analyses and graphics go hand in hand
- ▶ User can do most things without writing new code
  - \_ but can use new code where necessary
- ▶ New methods appear first in R, later elsewhere
- ▶ R is much more than a language and a system
  - \_ It is a community of users and of expertise
  - \_ Working co-operatively makes great things possible
- ▶ More reasons below (and some negatives!)

## Why R? – Two Opinons

- ▶ I use R exclusively to perform data analysis, and Python for more generic programming tasks (e.g. workflow control of a computer model).  
[Paul Hemstra: stackexchange]
- ▶ Since the language has been around for ever, lots of people have done things that you're likely to want to do. . . . , when faced with a hard problem, . . . download the package and get to work. And R “just works”: you give it a dataset, and it knows what summary statistics are useful. You give it some results, and it knows what plots . . . . [Ben Dundee: stackexchange]

<http://programmers.stackexchange.com/questions/181342/r-vs-python-for-data-analysis>

## R: Language++ & Library++

- ▶ The R System is a Library of packages
- ▶ Each package has access, in principle, to all others
- ▶ New packages leverage off existing code –
  - \_ I/O, data manipulation, graphics, statistics,
  - \_ spatial interpolation, . . .
- ▶ A package management system handles dependencies
- ▶ Minimal consistency is enforced across packages
  - \_ Documentation, Coding standards, Directory tree

# R's Role in the CoheRent Whole

The R project is

- ▶ an international community
  - \_ statistical experts, expert users,
  - \_ software development experts
- ▶ a repository for software and documentation
  - \_ statistics, graphics, . . . , and much else
- ▶ an information resource
  - \_ again, for statistics and much else
- ▶ a marketplace for ideas
  - \_ It brings together diverse skills & insights
  - \_ Ideas reflect diverse scientific backgrounds.



## UmbRella R – Examples

- ▶ Spatial Analysis
  - \_ Initially, interfaces to GRASS, SAGA, etc.
  - \_ Steadily, abilities have moved to R
- ▶ Motion Charts
  - \_ For now, *googleVis* interfaces to Google apps
  - \_ Work at Auckland Univ will make these part of R.
- ▶ Commentary
  - \_ Interfacing is a good way to start
  - \_ Gradually, transfer functionality to R

# What is Wrong with R?

- ▶ The language model is dated
- ▶ R is powerful, but not powerful enough
  - \_ Users now want even more powerful tools
- ▶ R code is slow
  - \_ hence call C code if faster execution is needed
- ▶ R has been a first draft of a rewrite of S
  - \_ Much has been learned as the draft has evolved
- ▶ Some specific areas require attention
  - \_ e.g., the graphics for `lm` (regression) models.

# On Balance, Virtues vs Defects

- ▶ R is a Mature System with a Dated Language Model
  - \_ Maturity wins!
- ▶ Slow R code is mainly a problem for power users
- ▶ The Package Management System is a big plus
- ▶ An evolving R will be with us for a while yet!
- ▶ Challenge: Use what we now know to improve R.
  - \_ Easiest if changes to base R can be avoided
  - \_ *knitr* is a good example of what is possible

# Rethinking Statistical Analysis

- ▶ Simulation is the basic tool for estimating sampling distributions
- ▶ Sometimes we are lucky, and mathematical results make the recourse to simulation unnecessary
- ▶ Simulation does make theoretical assumptions
  - \_ Theoretical assumptions are inescapable!
- ▶ Alternatives – the bootstrap
  - \_ Here the data determines the distribution
  - \_ This introduces data sampling error

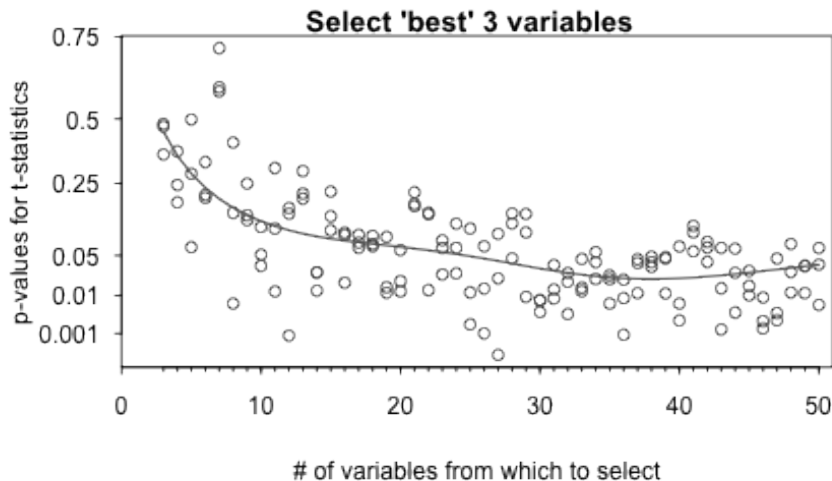
## Simulation – Variable Selection

y variable	x1	x2	x2	. . .
Noise	Noise	Noise	Noise	. . .
rnorm(100)	rnorm(100)	rnorm(100)	rnorm(100)	. . .

- ▶ With no variable selection (3 explanatory variables  $x_1$ ,  $x_2$  and  $x_3$ ),  $p$ -values have an expected value of 0.5
- ▶ As the number of variables available for selection increases, selecting the best 3 gives ever smaller notional  $p$ -values.

## Simulation – Variable Selection

Here variable values are random normal noise



Message: Stepwise regression has serious traps!

# Rethinking Regression

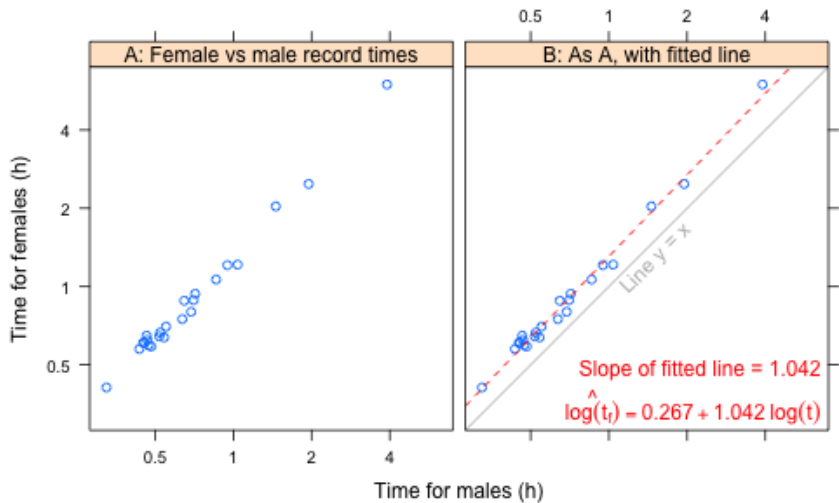
- ▶ Most R model objects can both be printed and plotted
  - \_ *plot()*, with an *lm* object, gives diagnostic plots
- ▶ In small datasets, some features are just noise
  - \_ Answer – do such plots for simulated data
  - \_ Hence *plotSimDiags()* in DAAG ( $\geq 1.16$ )

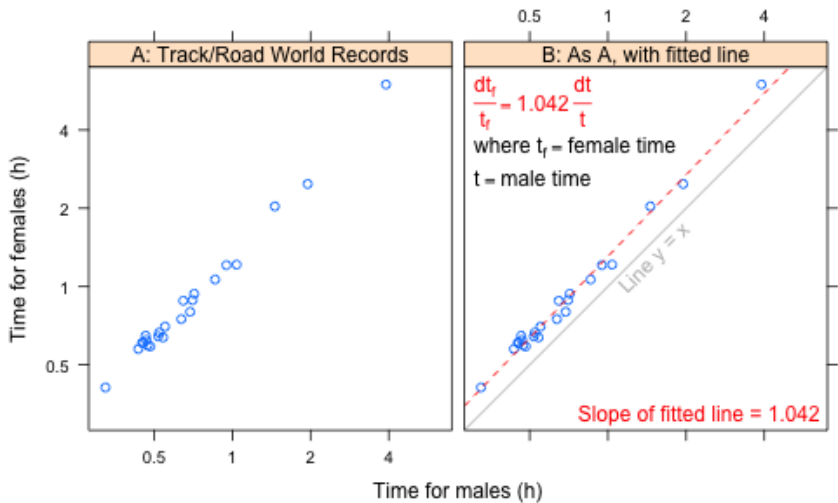
## Rethinking Regression – Data

The following are male and female record times for hill races in Northern Ireland. Race distances are in miles, extent of climb is in feet, and times are in hours.

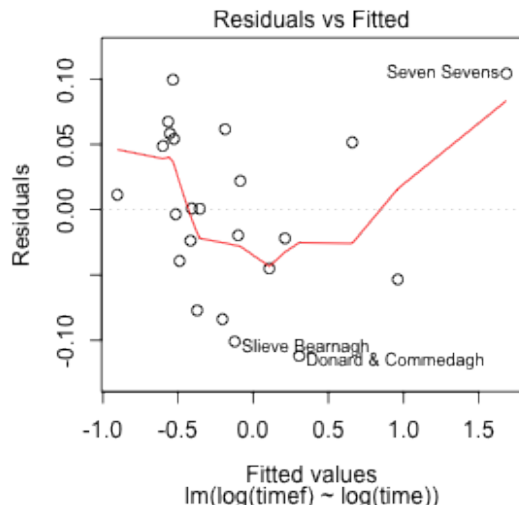
	dist	climb	time	timef
Binevenagh	7.5	1740	0.8583	1.0644
Slieve Gullion	4.2	1110	0.4667	0.6231
Glenariff Mountain	5.9	1210	0.7031	0.8869

... 22 further rows of data

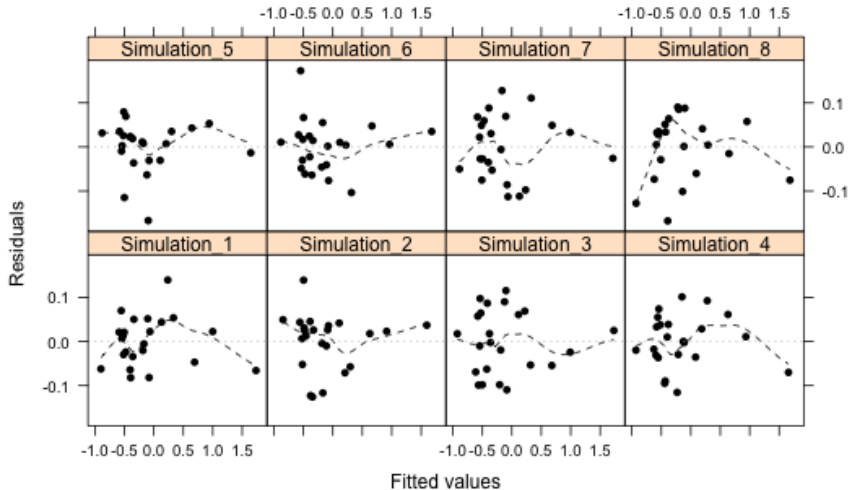




## Diagnostic plot 1



### Residuals vs fitted (Regress female on male times)



These simulate the sampling distribution of diagnostic plot 1.

## Fixing termplot()

- ▶ Currently, useful for a limited range of models
  - \_ For these, it gives highly helpful insight
- ▶ Output is not ideal for transformation terms
  - [e.g.,  $\log(x)$  terms]
- ▶ Shows partial regression effects
  - \_ currently, for main effects terms only

# Rethinking the Textbook

Make the web an adjunct to the printed page

The same issues arise for course notes.

- ▶ Texts commonly have their own R package
- ▶ Provide data sets (obviously) & functions
- ▶ Use functions to package any major code chunks
  - \_ For ephemeral fn's, create by sourcing scripts \_ e.g., *gamclass* has functions like `fig2.1()`
- ▶ Vignettes, R markdown, ...
- ▶ There is much more to be done and said ...
  - \_ This is, after all, a time of rapid computing related technological change!

# Acknowledgements

- ▶ Friedrich Leisch and R-core, for *Sweave()*
- ▶ Yihui Xie for *knitr*, which improves on *Sweave()*
  - \_ Slides were created using *knitr*.