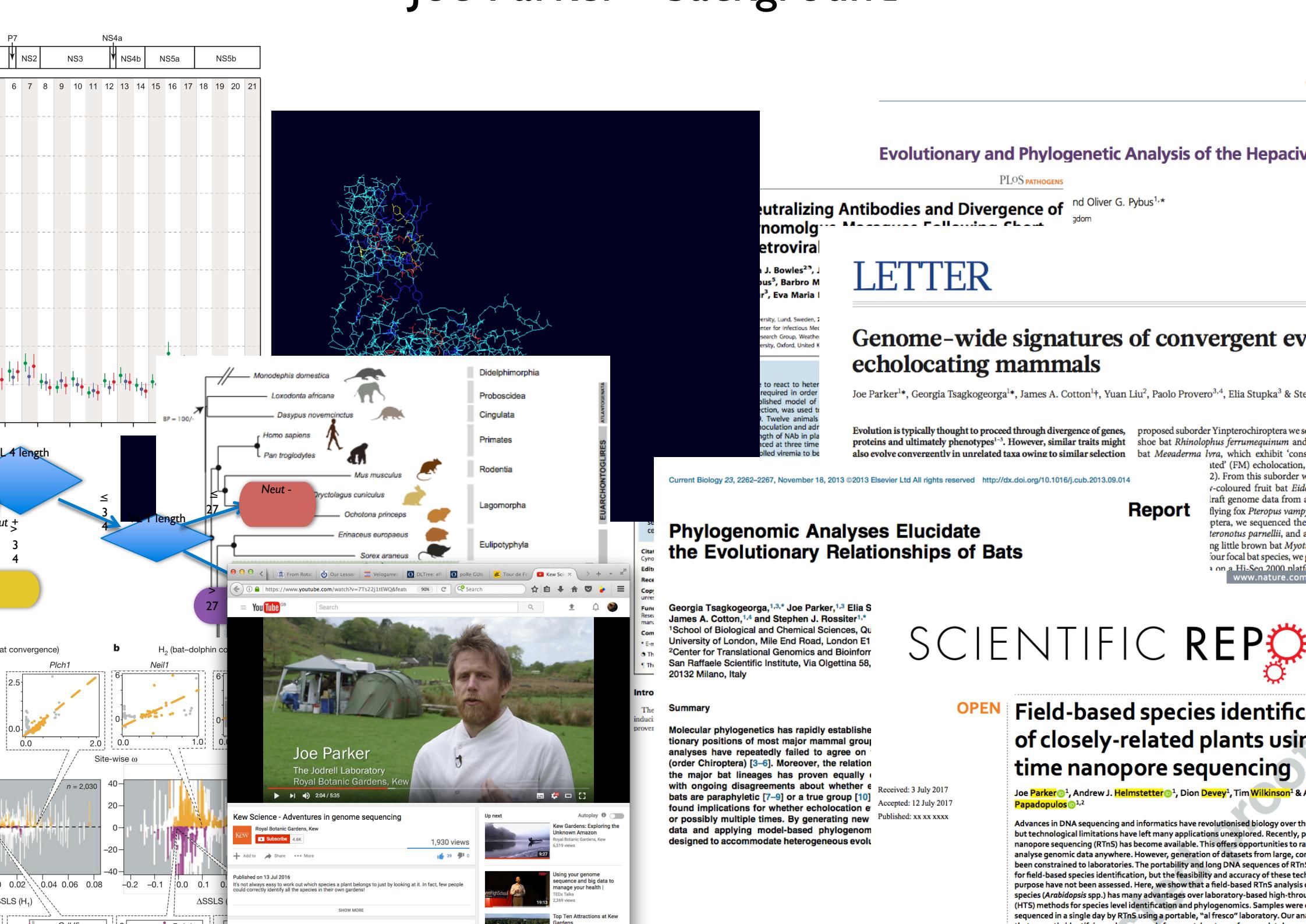


Statistics & Bioinformatics

Joe Parker



Goals

- Equip you to be a bioinformatician
- (in a statistical sense)
- That's it

Assessment

What do we do?

- Collect empirical observations of phenomena
- Spot patterns
- Propose models that explain them
- Simulate guesses according to the models
- Compare them, refine, repeat
- ... SCIENCE

Bioinformatics and stats

- **Toolbox:** big data, replicates, moar
- **Concepts:** statistics, probability, populations
- **Techniques:** multivariate analyses, programming, genomics

NOT

- Stats 101
- Programming course
- Mol. Biol.
- Genomics
- Machine-learning / AI / any advanced stuff

Roadmap

Populations, probability and statistics

- Almost all biological phenomena vary continuously
- Probability: assess outcome given model
- Equally, for a large N how many X?
- How likely is X?
- Statistics:
 - How likely is it that X occurs by chance?
 - By some mode?
 - From same distribution (process) as Y?

The P-value

- It isn't: magic
- Working definition: probability of seeing a value as or more extreme than the observed, assuming model, parameters
- “Is this result random?” (ahem)

Descriptive stats

- Count
- Mean, mode, median
- Range
- Confidence intervals

Applications

- Trees on $h \{R\}$. Two species or one?
- k nonsynonymous changes of N aa replacements. Convergent?
- 913 patients with cancer, 50,000 SNPs, 2031 haplotypes. Will patient X get cancer?

Introduction to R

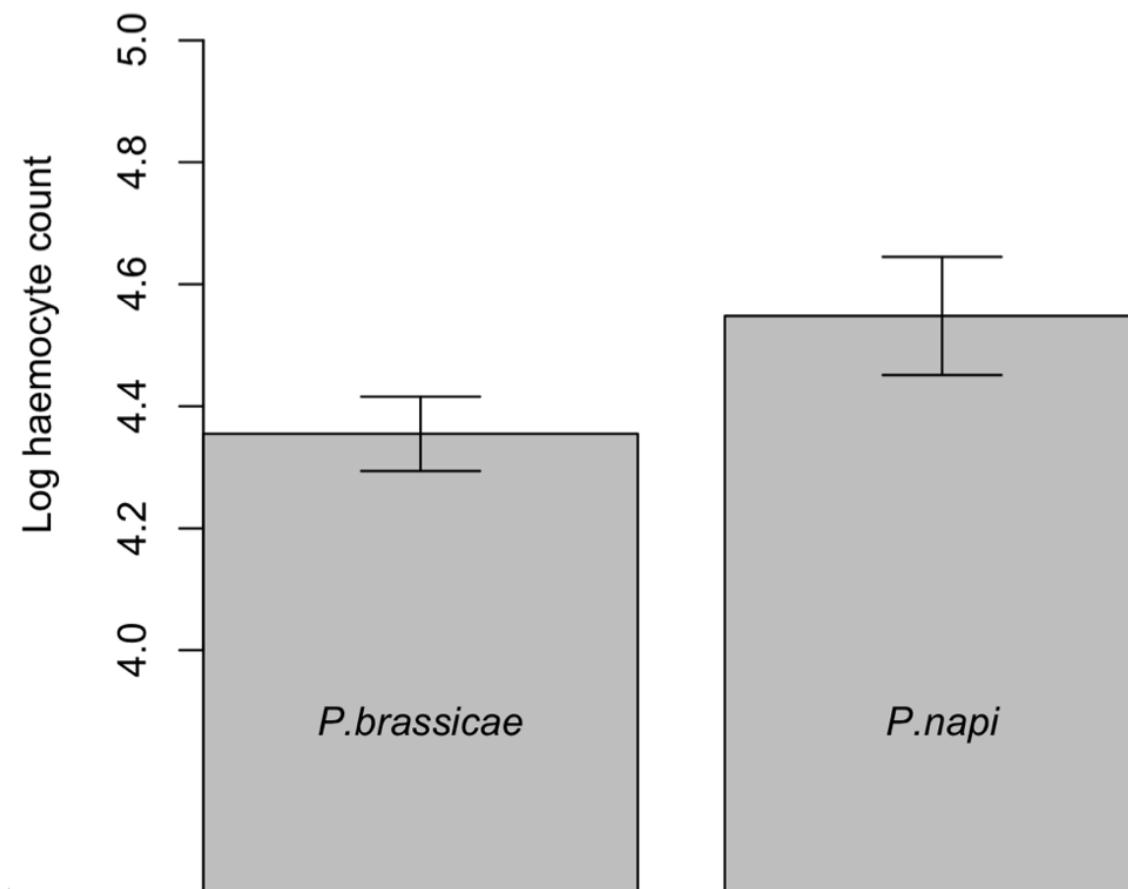
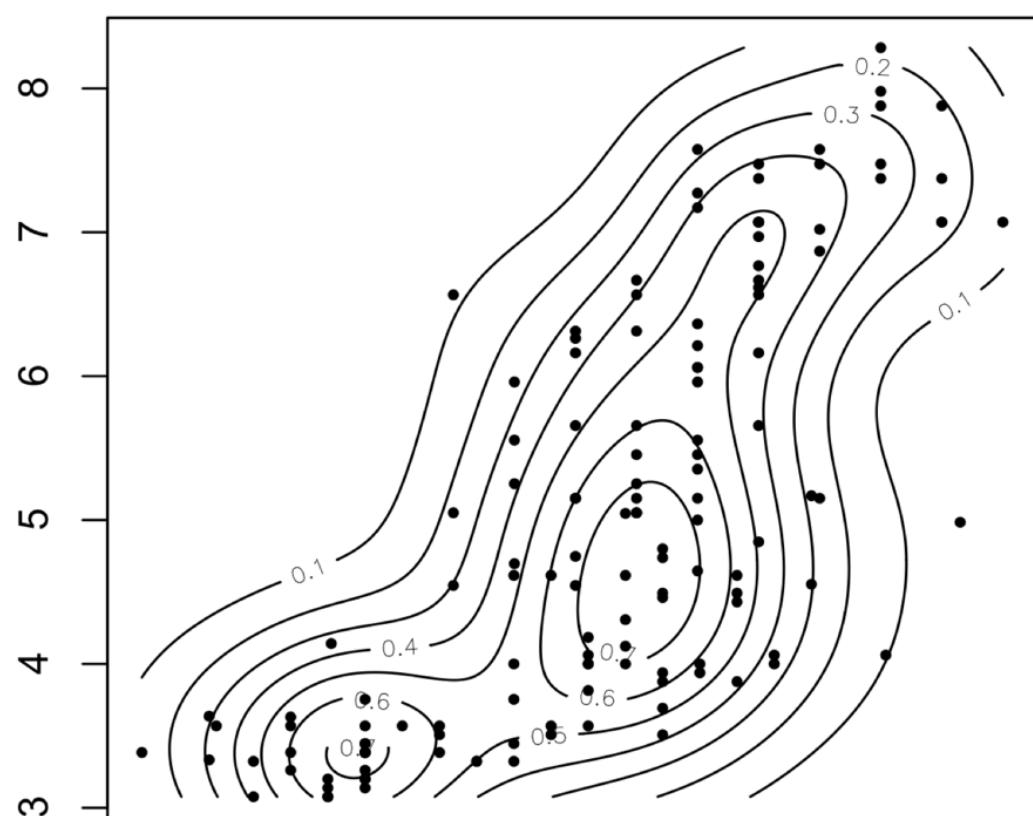
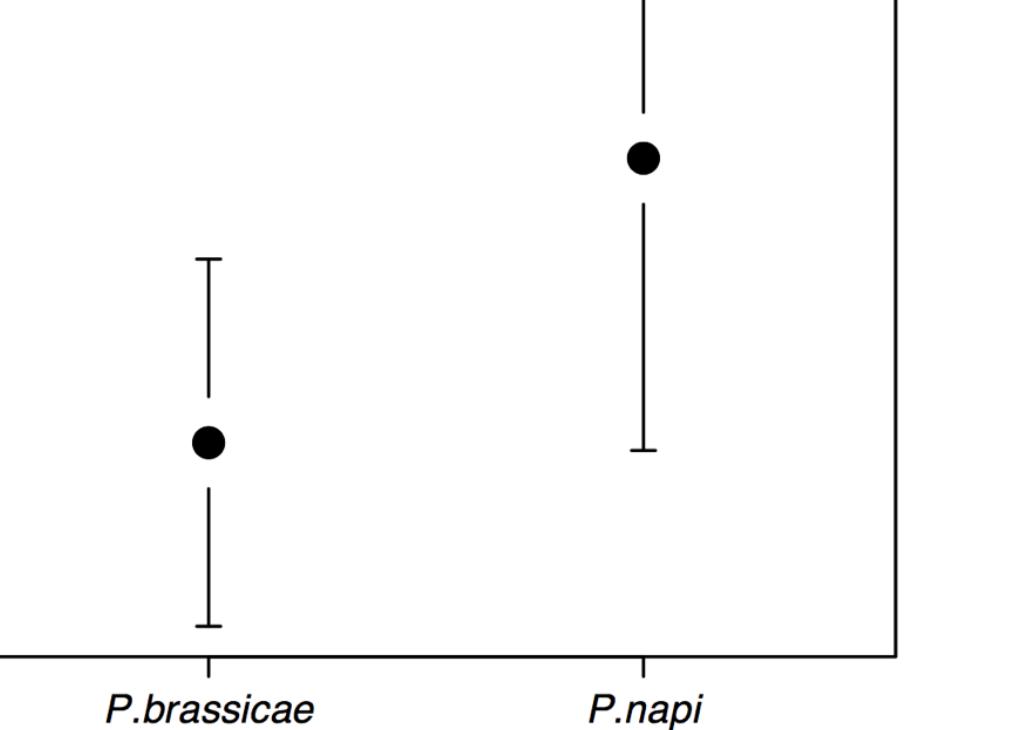
Joe Parker (Rob Knell)

What is R?

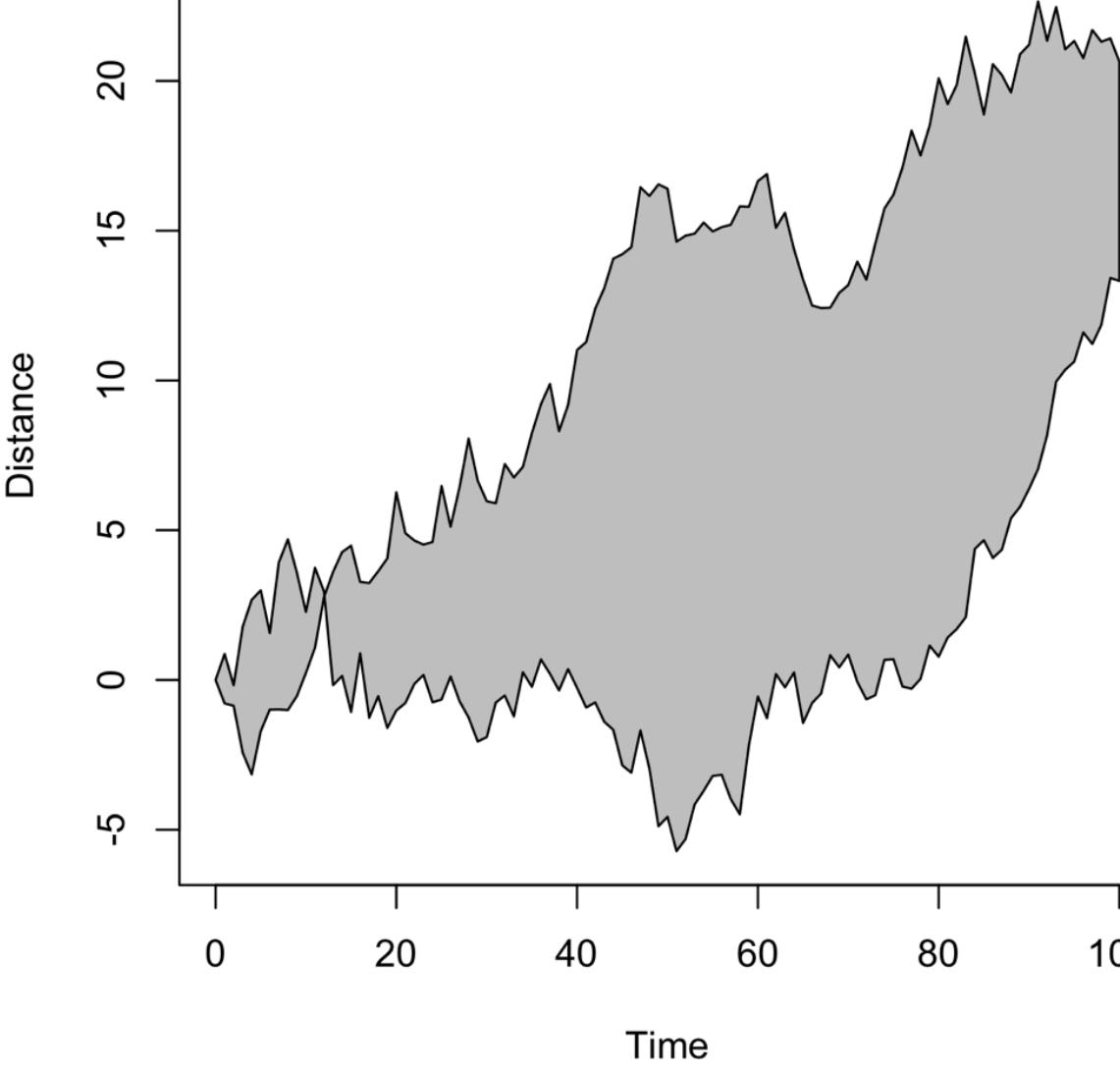
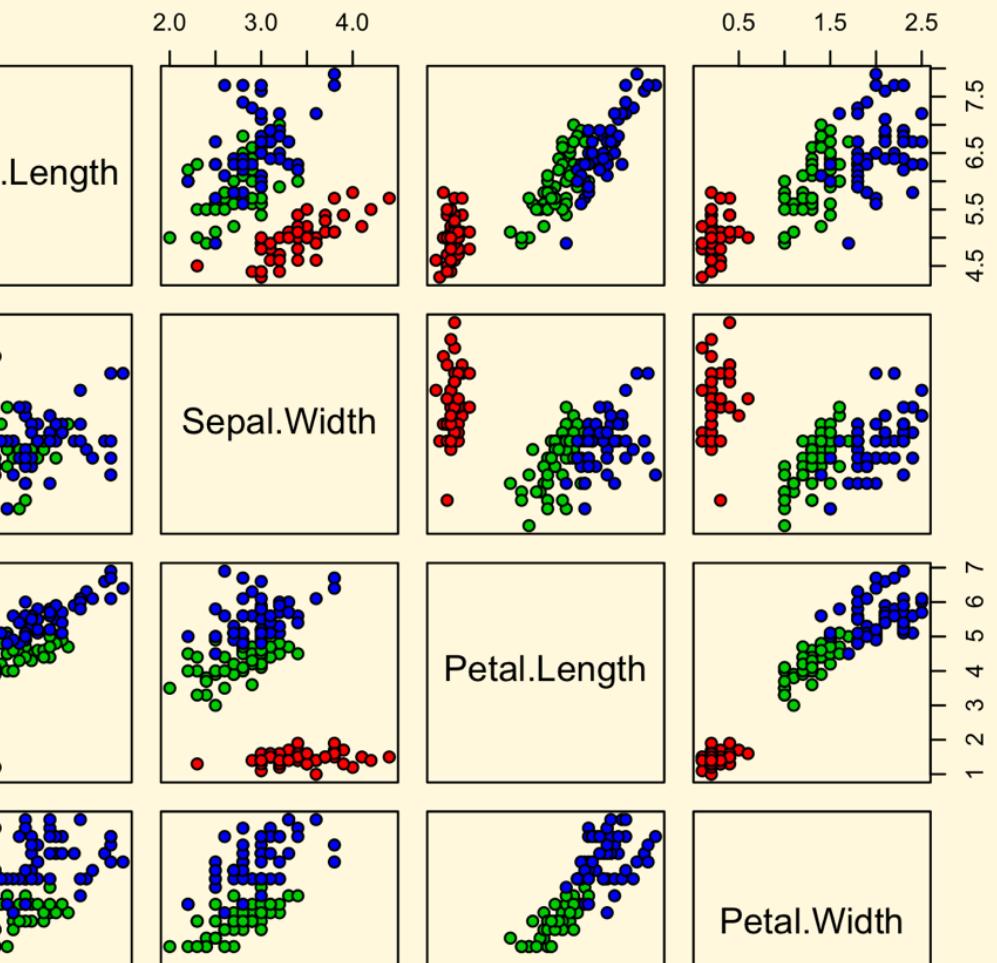
- Open source statistical programming language
- Based on a commercial language called “S”
- Powerful and flexible software that can do any statistical analysis you can think of and much more

What can I do with R?

- With the base R installation you can carry out a tremendous variety of statistical analyses
- You can draw publication-quality graphs
- With add-on packages you can do just about anything

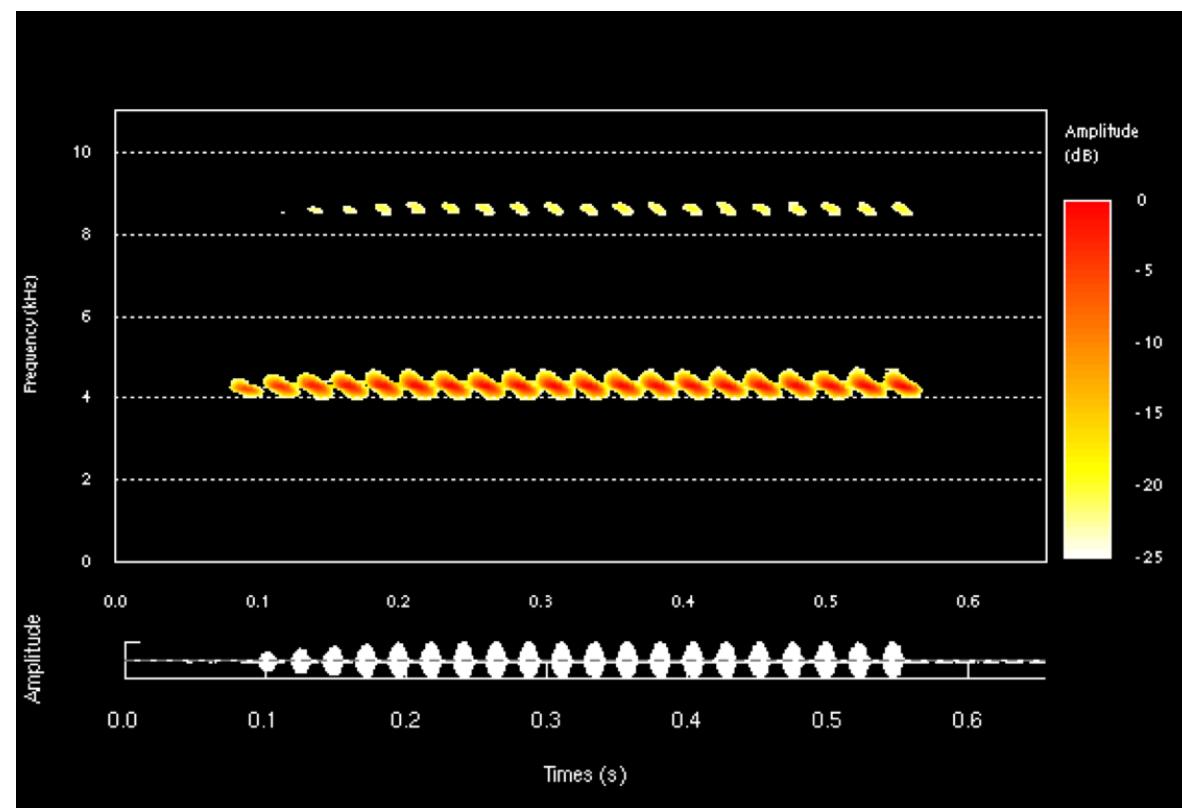
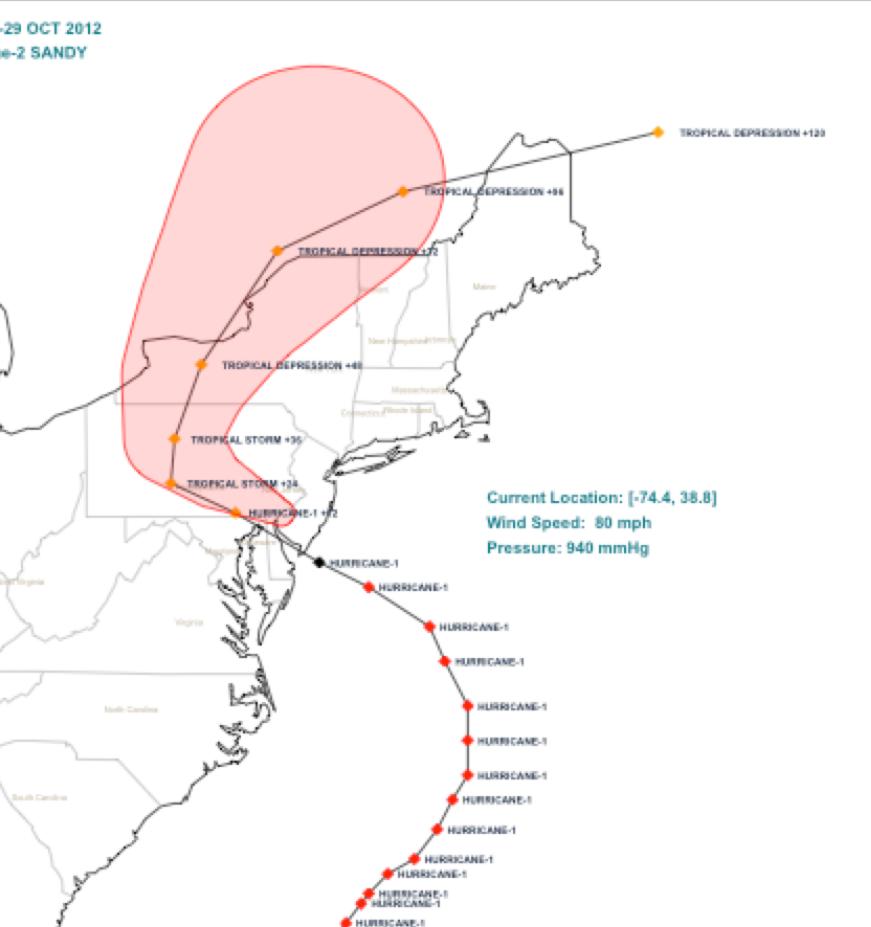


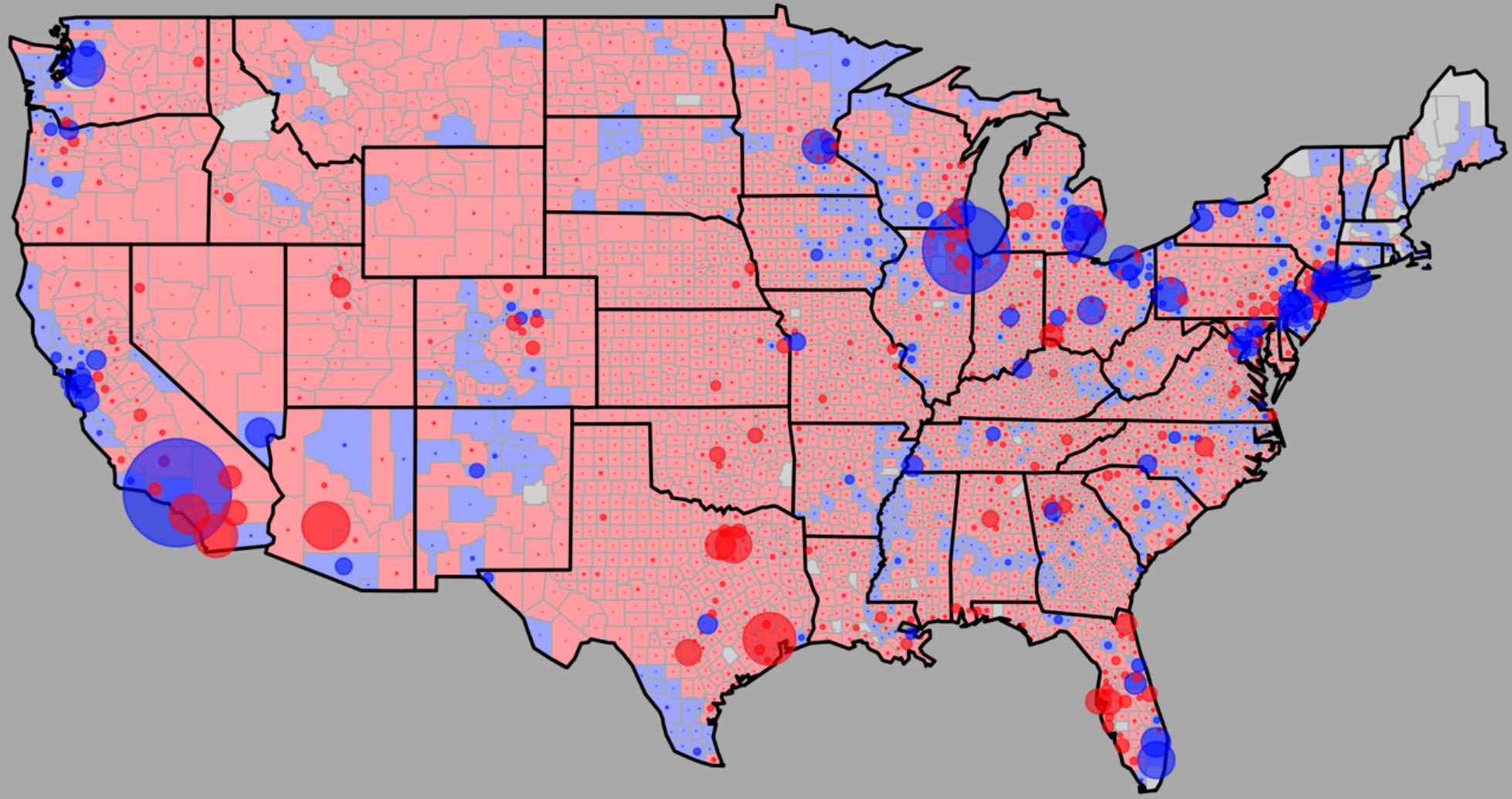
Edgar Anderson's Iris Data



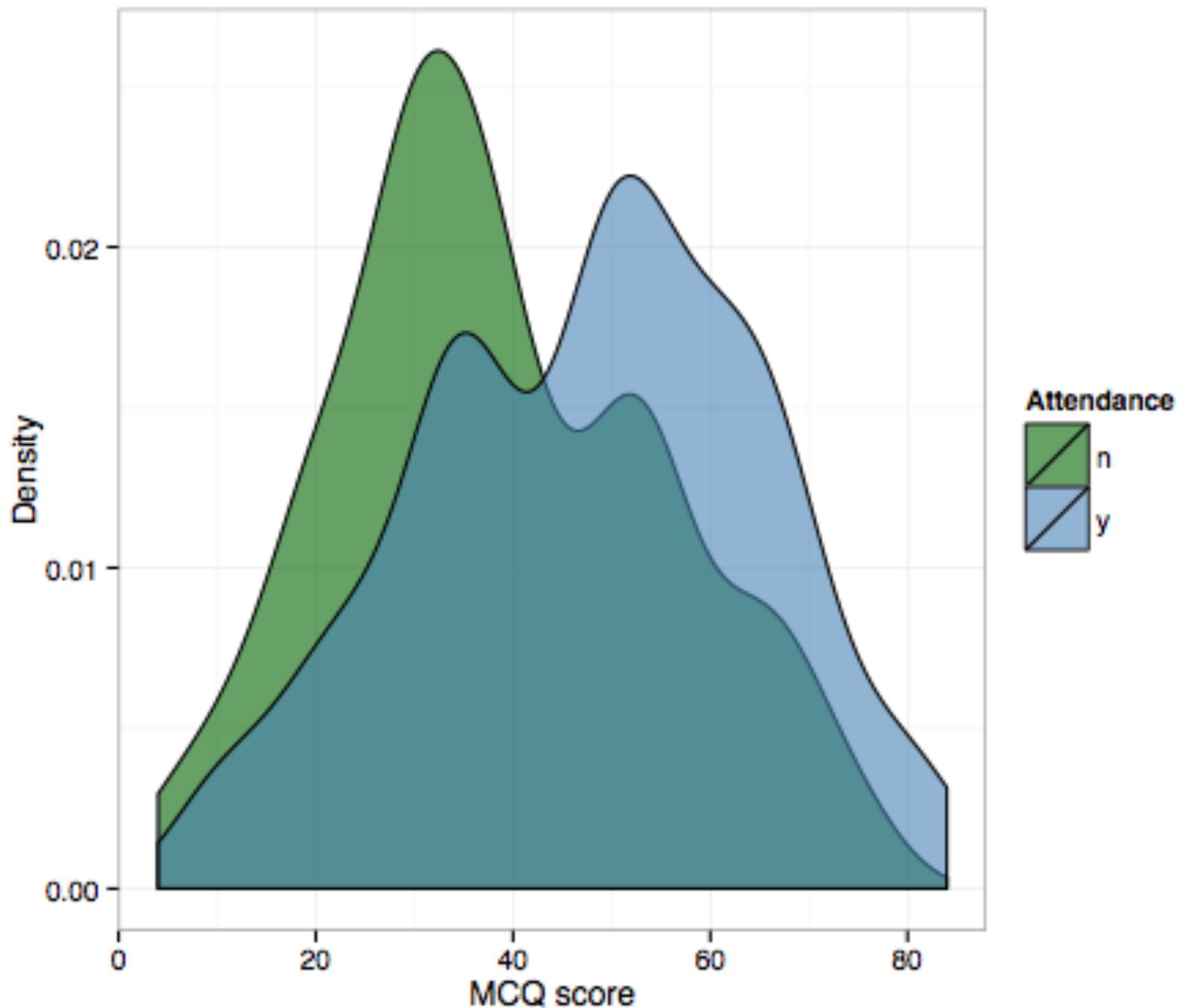
Sonogram of cricket song

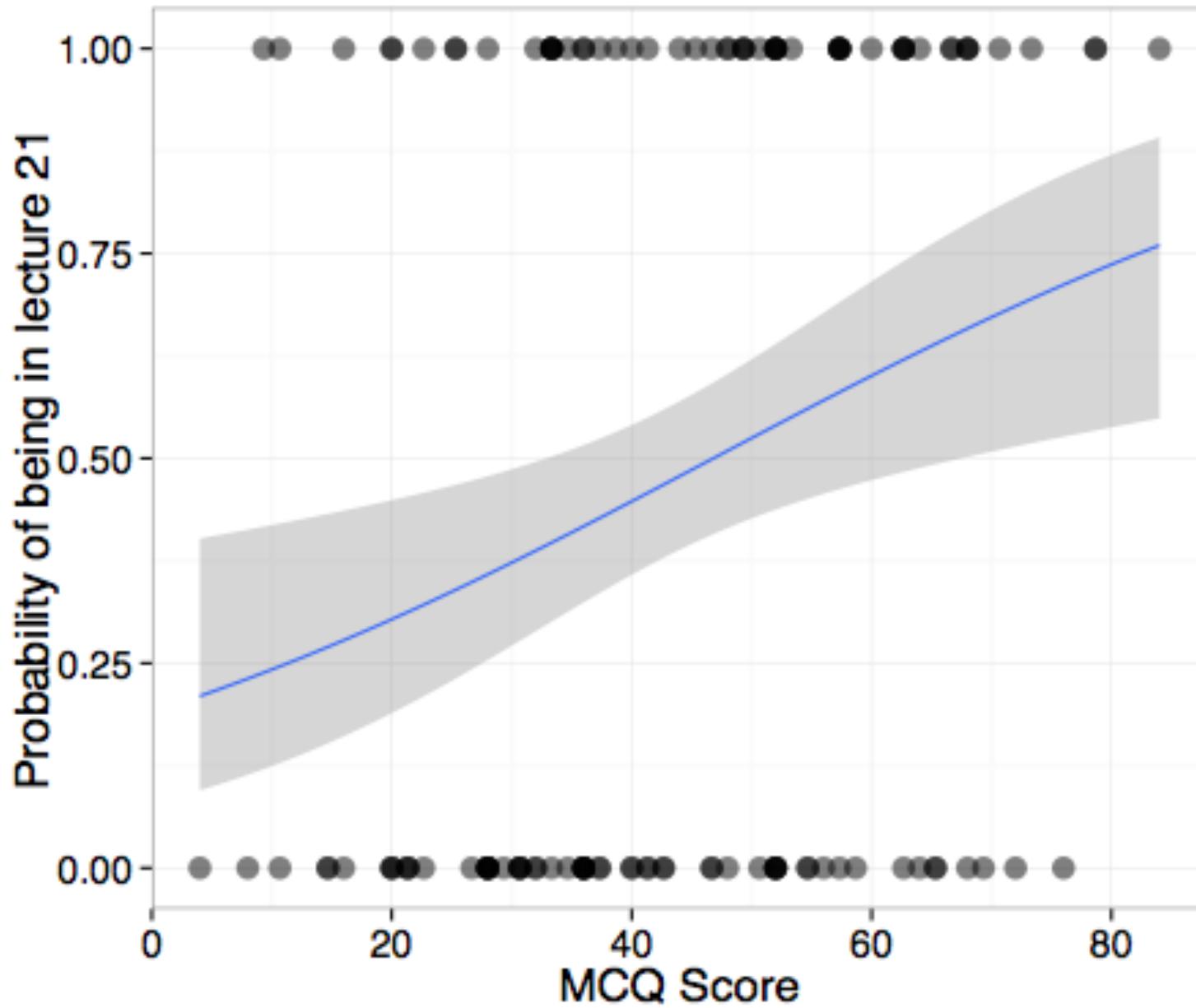
Hurricane Sandy predicted path

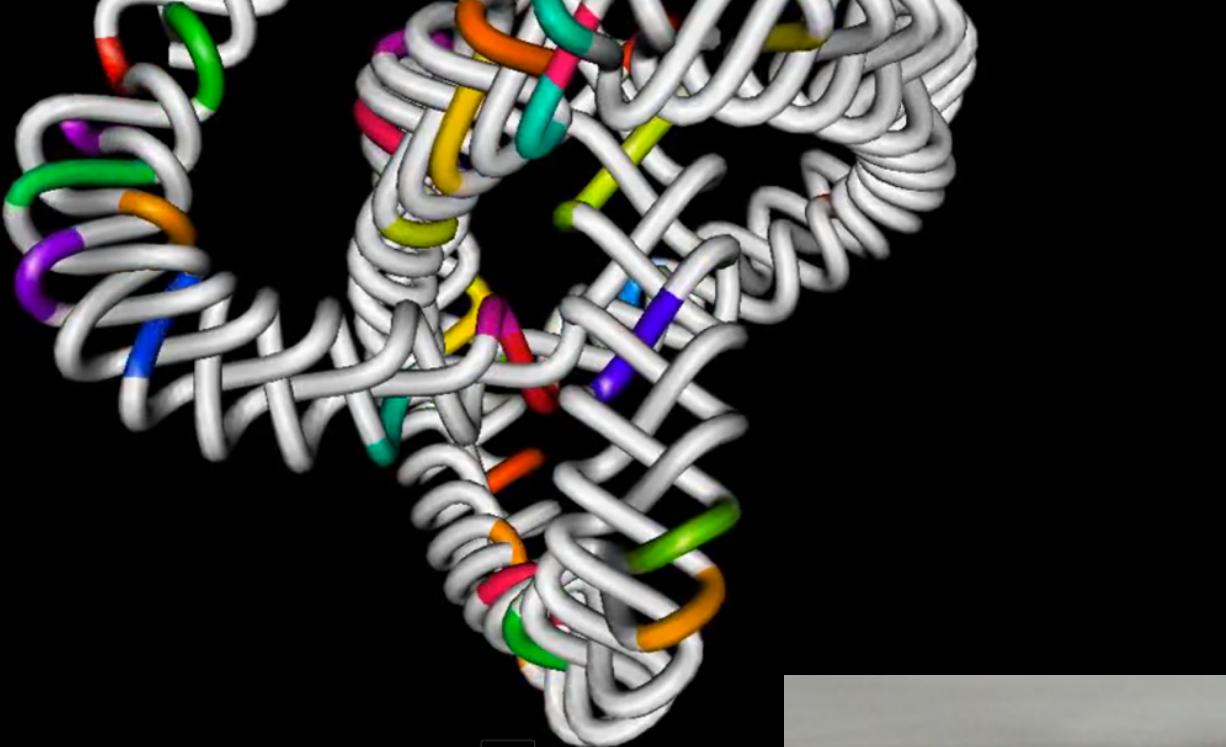




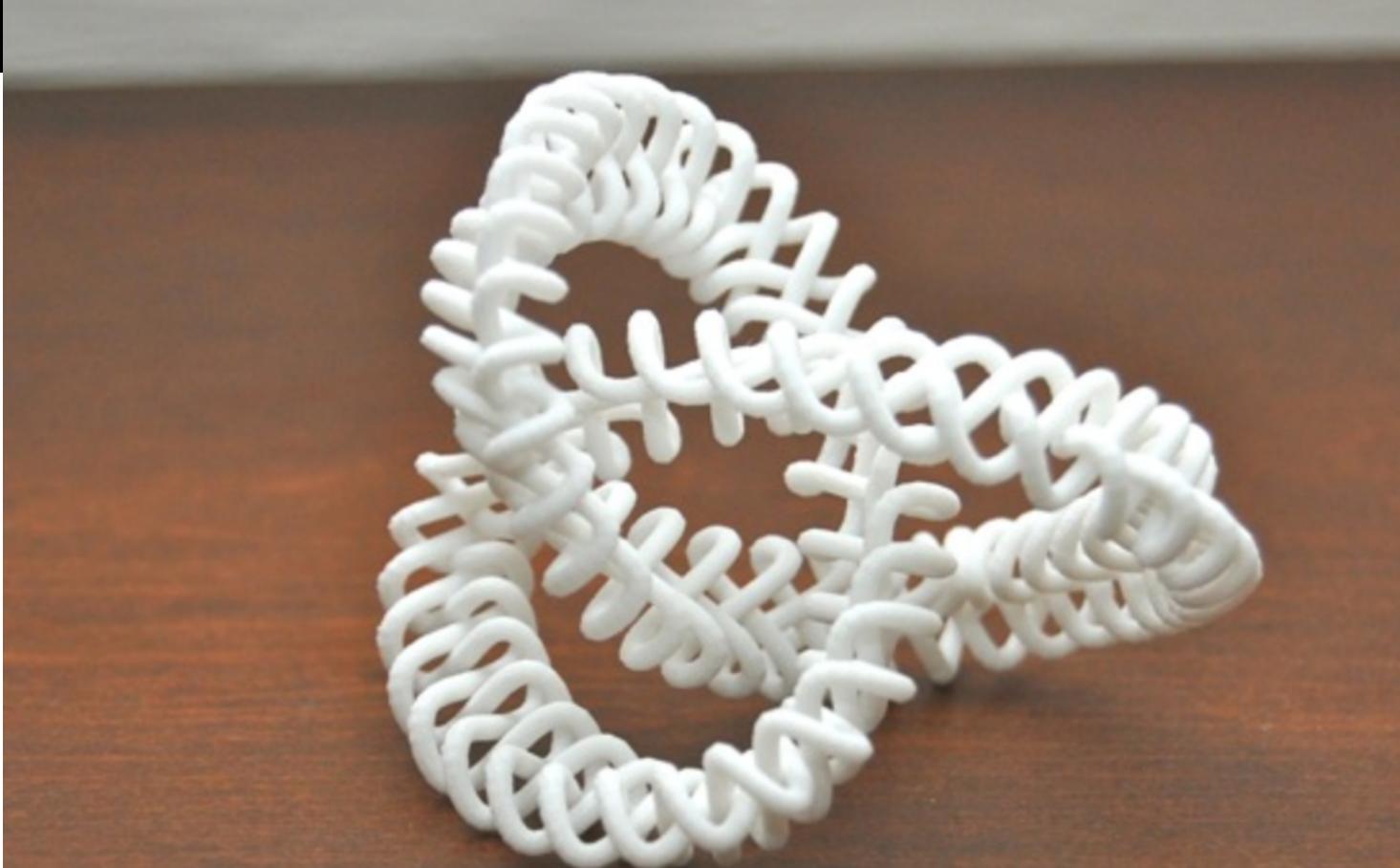
US Presidential election map







rgl package



How to get R

cran.r-project.org

The Comprehensive R Archive Network

Download and Install R

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

- [Download R for Linux](#)
- [Download R for MacOS X](#)
- [Download R for Windows](#)

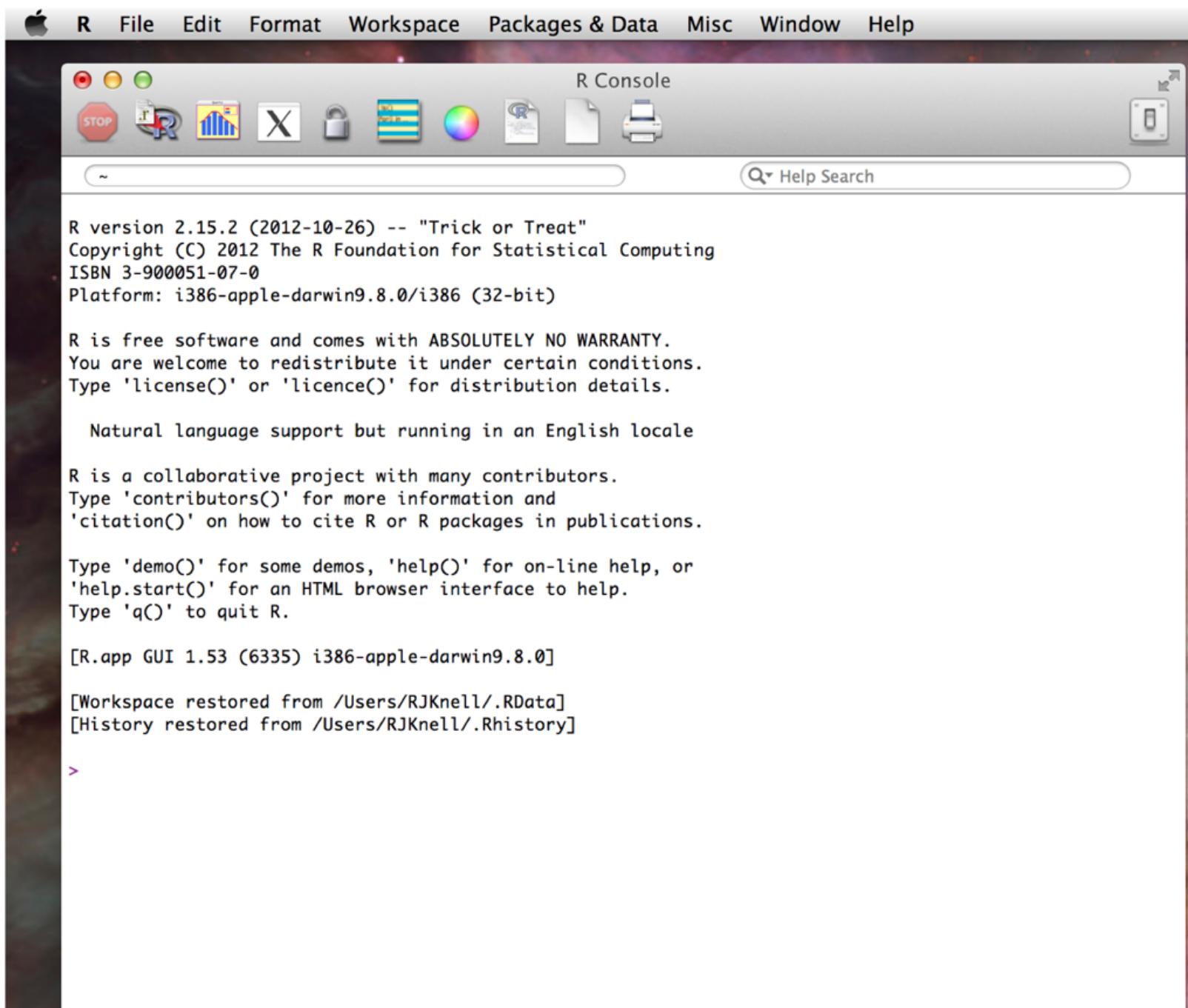
R is part of many Linux distributions, you should check with your Linux package management system to the link above.

Source Code for all Platforms

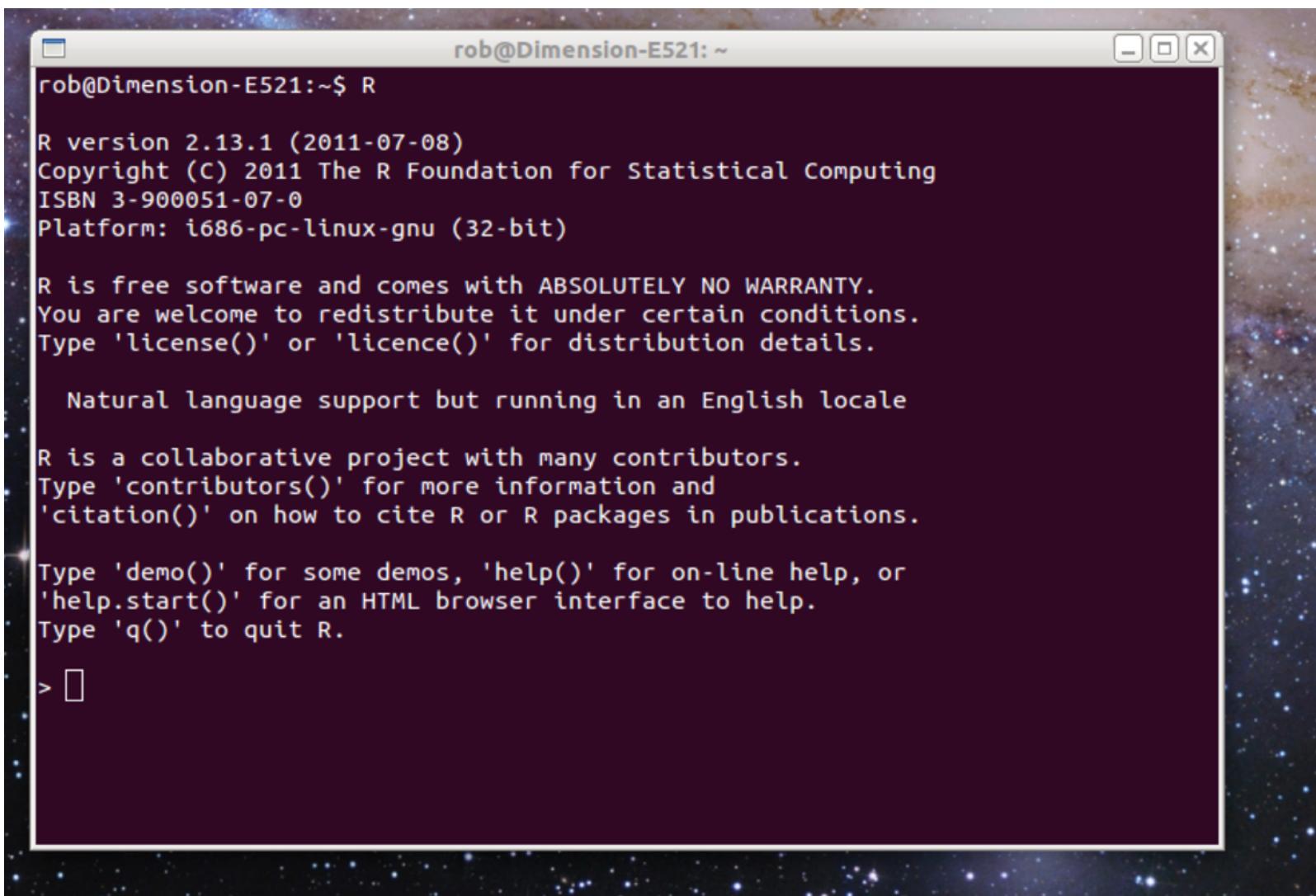
Windows and Mac users most likely want to download the precompiled binaries listed in the upper box source code. The sources have to be compiled before you can use them. If you do not know what this means probably do not want to do it!

- The latest release (2012-10-26, Trick or Treat): [R-2.15.2.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 [features and bug fixes](#) before filing corresponding feature requests or bug reports.

R basics



R basics



A screenshot of a terminal window titled "rob@Dimension-E521: ~". The window displays the R startup message. The text is white on a dark background. It includes the R version (2.13.1), copyright information (The R Foundation for Statistical Computing), ISBN, platform (i686-pc-linux-gnu (32-bit)), a disclaimer about warranty, instructions for redistribution, natural language support, collaborative project information, and help resources. At the bottom, it shows a prompt starting with '>'. The background of the slide features a starry space image.

```
rob@Dimension-E521:~$ R

R version 2.13.1 (2011-07-08)
Copyright (C) 2011 The R Foundation for Statistical Computing
ISBN 3-900051-07-0
Platform: i686-pc-linux-gnu (32-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.

> □
```

software

- You type an expression or instruction at the command prompt and press “enter”
- If R can understand what you’ve typed, it’ll do what you’ve asked...
- ...or what it thinks you’ve asked: unlike Python, R contains some fairly high-level functions

Data exploration and description

Exploratory data analysis

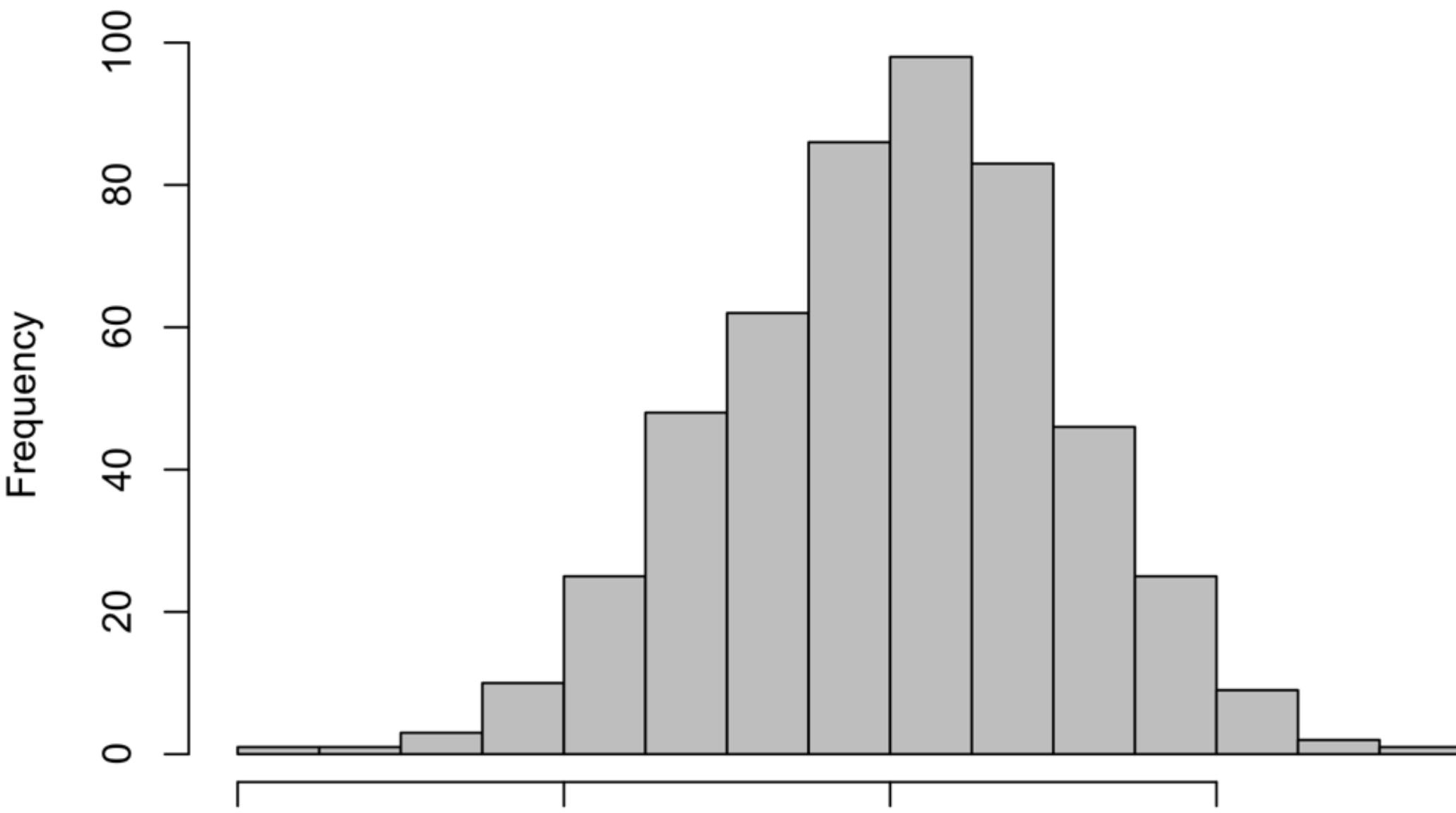
- Gain an understanding of the nature of the data
- Get clues regarding the relationships between variables
- Detect errors
- Know which analyses might be necessary
- ...

Exploratory data analysis

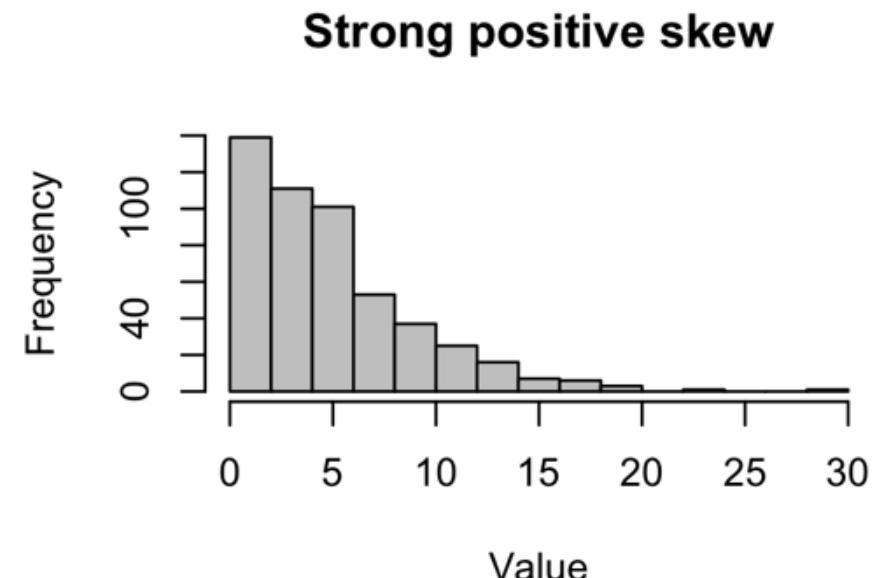
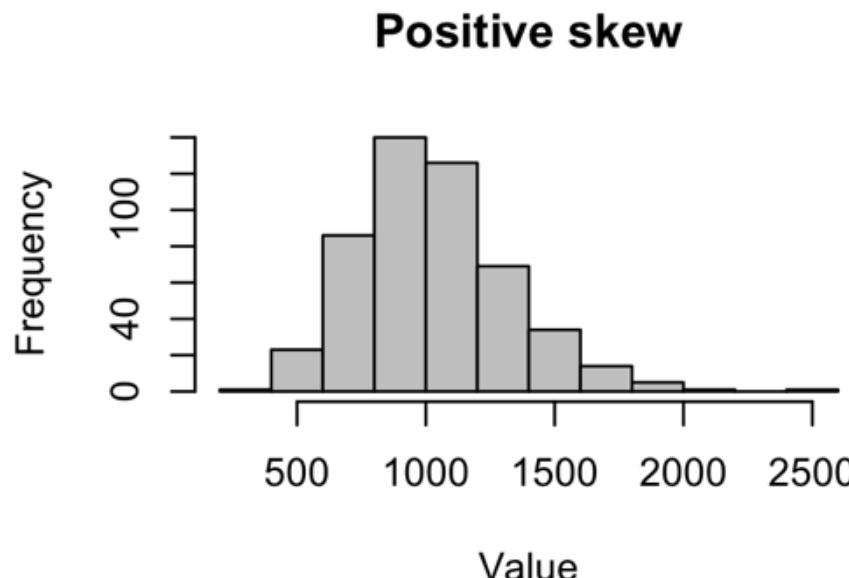
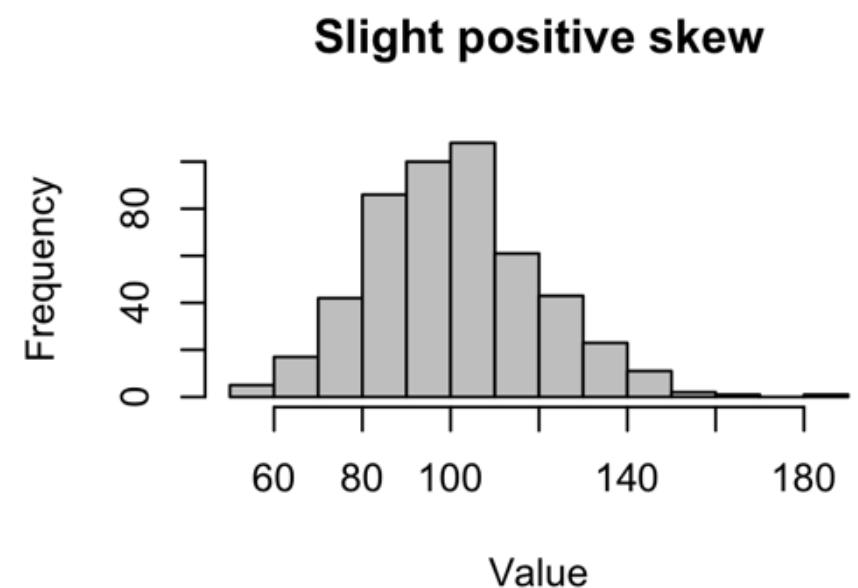
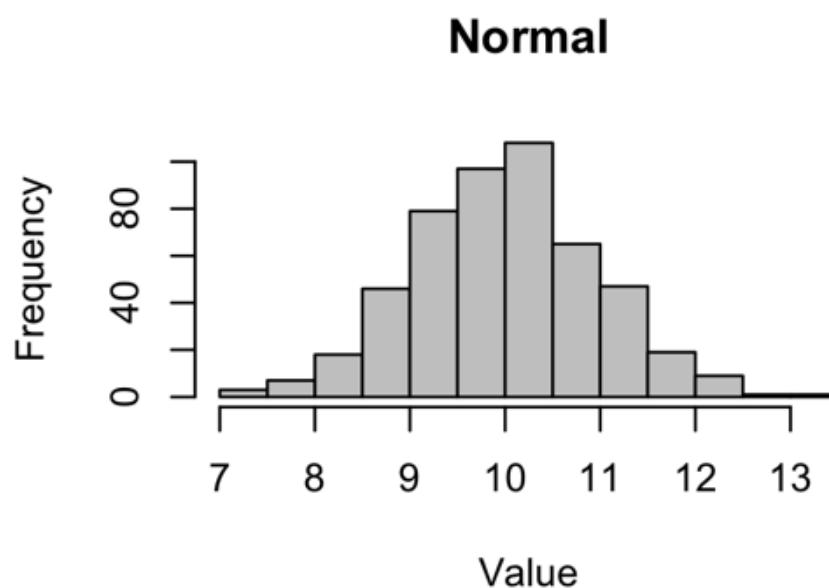
- Numerical and graphical descriptions of data
- Need to revise some basic numerical descriptives before moving on to graphical techniques

Frequency distributions

500 normally distributed random numbers



Frequency distributions



Descriptive statistics

- Measures of central tendency
 - Mean, median, mode
- Measures of dispersion
 - Standard deviation, variance, IQR

Central tendency

Mean $\bar{x} = \frac{1}{n} \sum x$

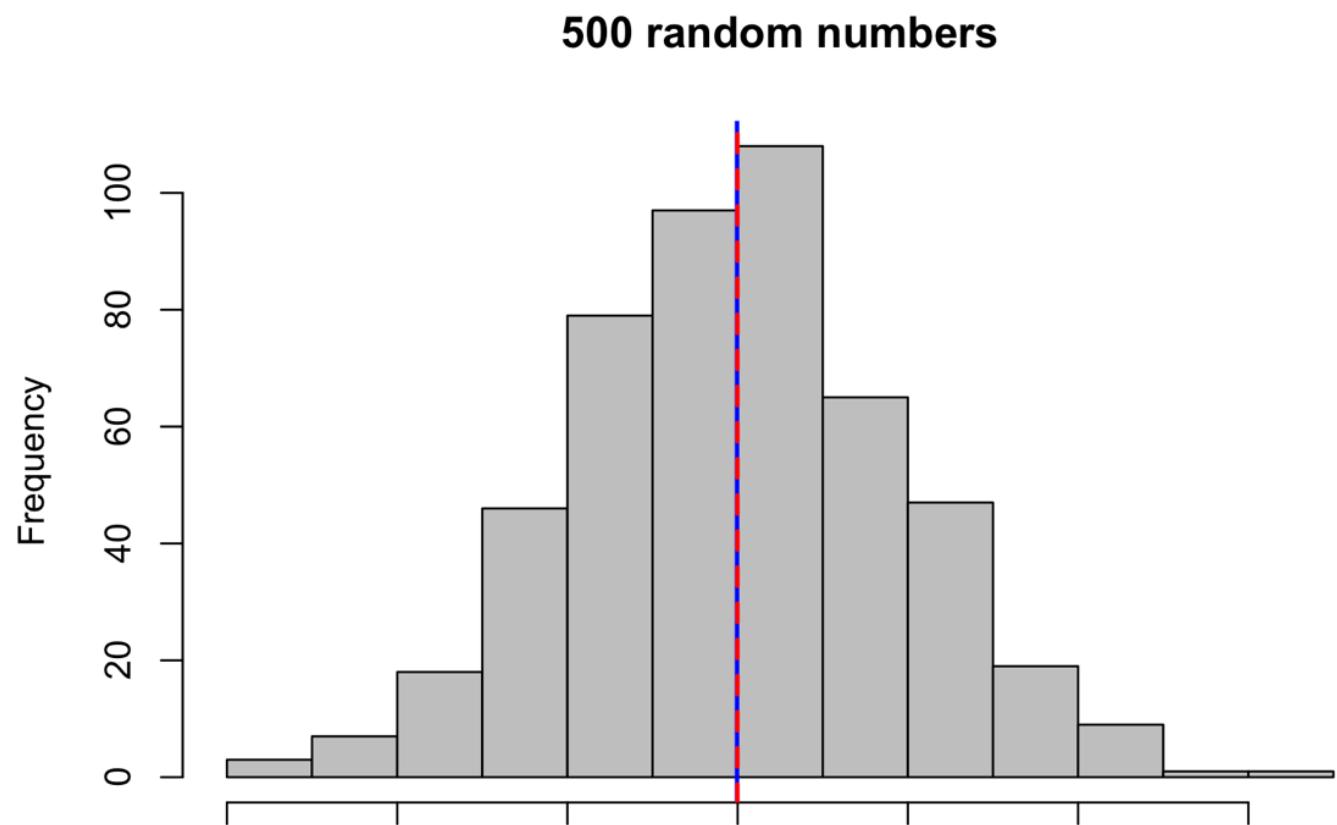
Median: rank data, take middle value.

If even number of data take the average of the two middle values

Mode: most common value

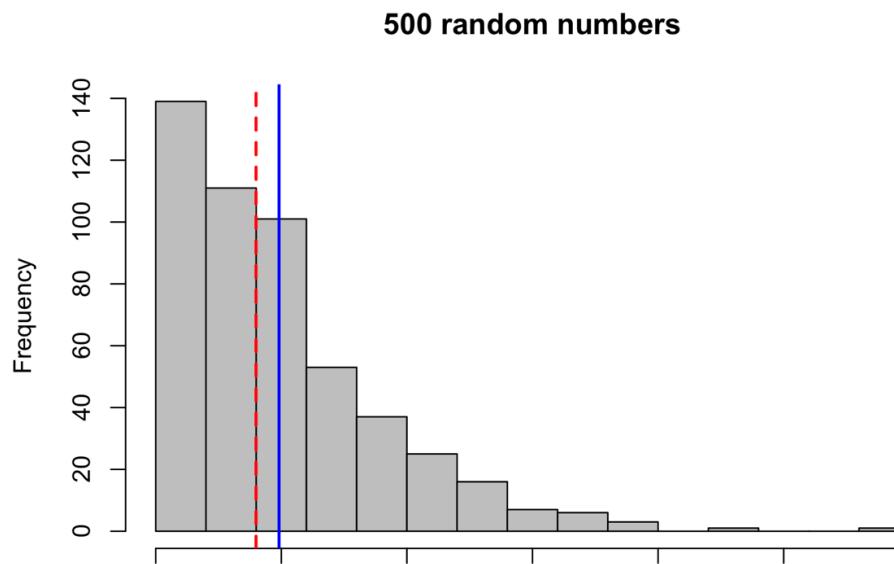
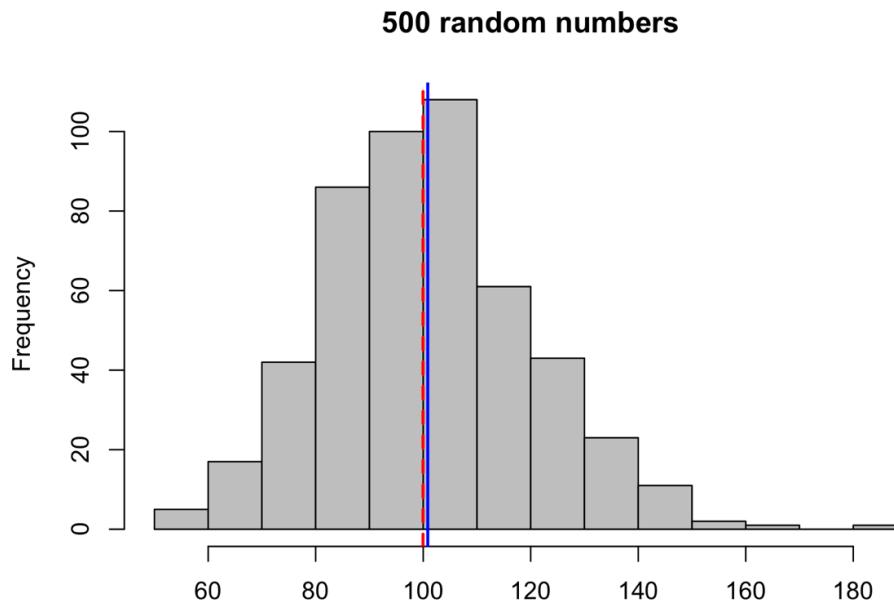
Central tendency

For normally distributed
data, mean (blue) =
median (red) = mode



Central tendency

As the distribution is more skewed, the mean, mode and median become less similar



Measures of dispersion

Variance

$$s^2 = \frac{\sum(x - \bar{x})^2}{n - 1}$$

Standard deviation

$$s = \sqrt{s^2}$$

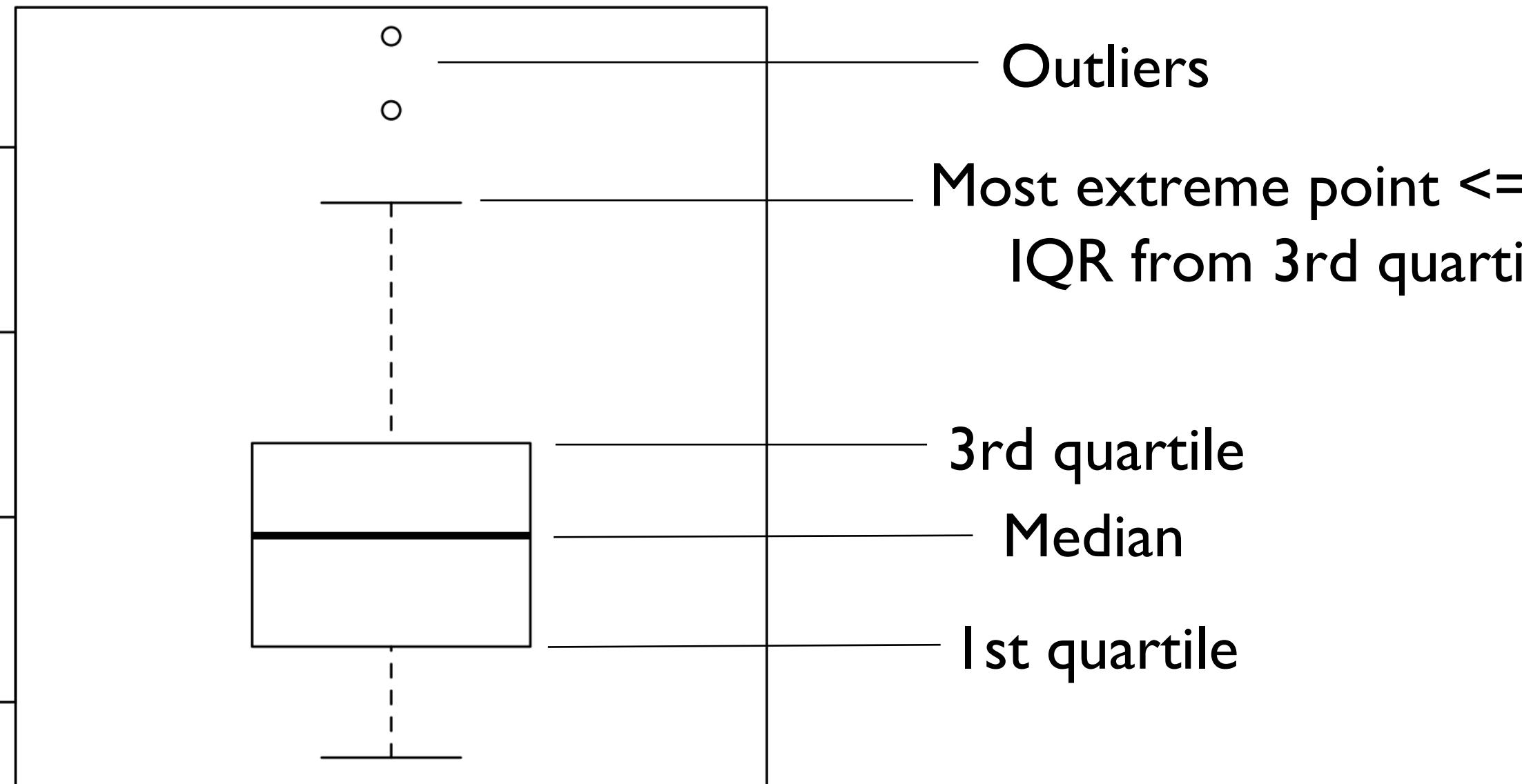
Inter-quartile range (IQR)

Range from 1st quartile (top of lower 25%) to third quartile

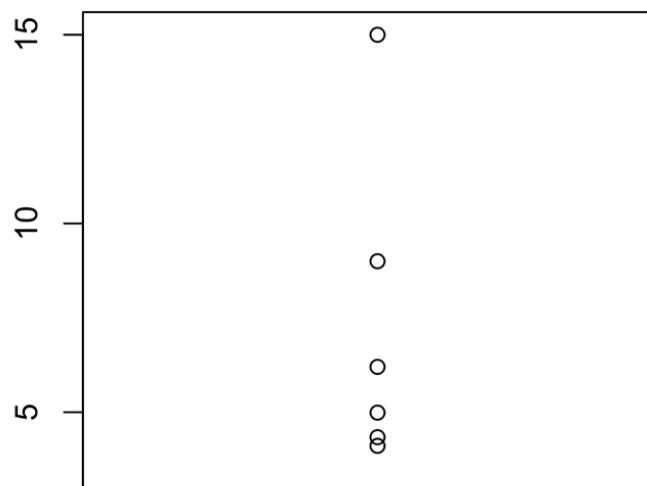
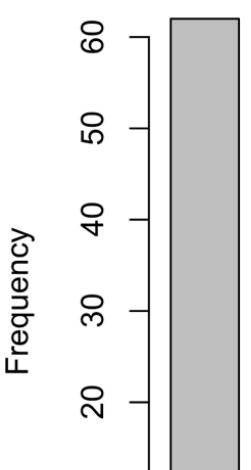
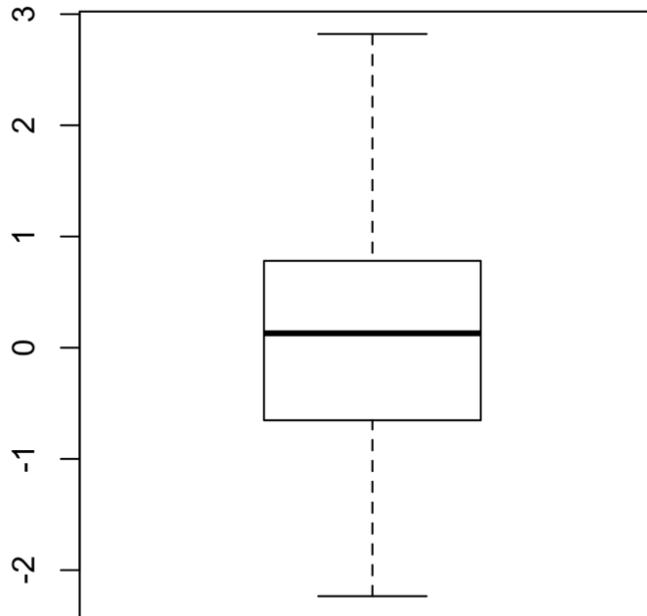
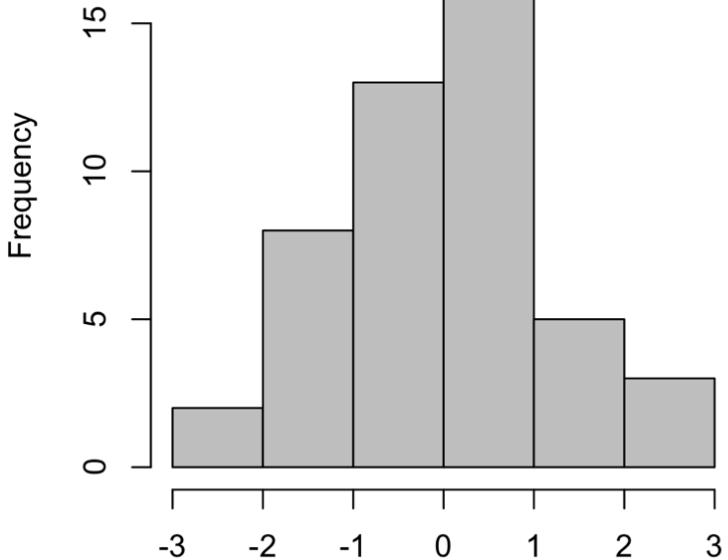
Graphics for exploratory analysis: univariate data

- Histograms
- Boxplots

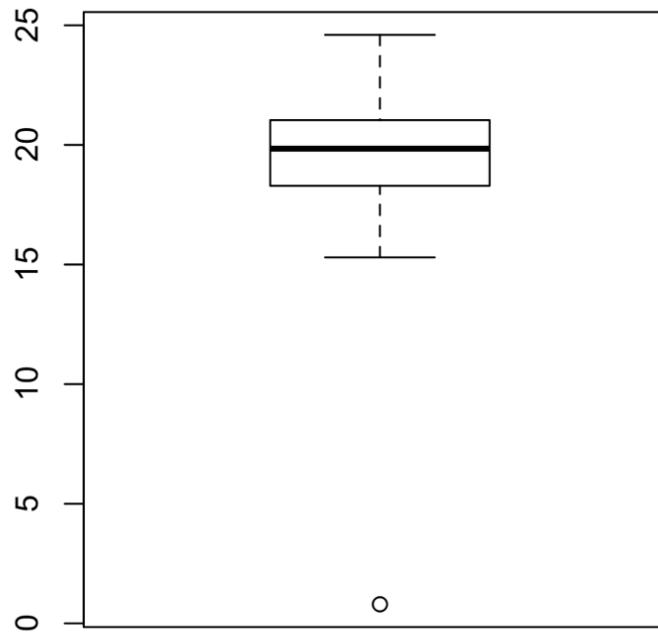
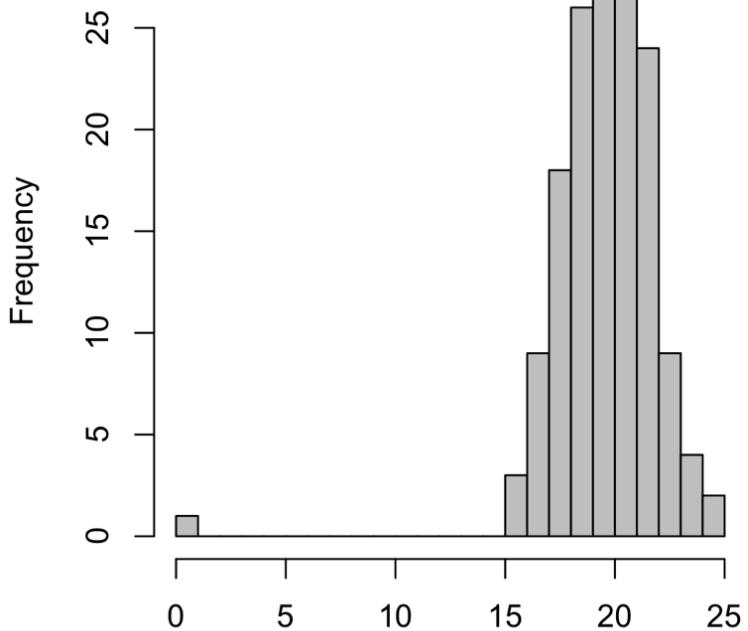
Boxplot



boxplots



boxplots

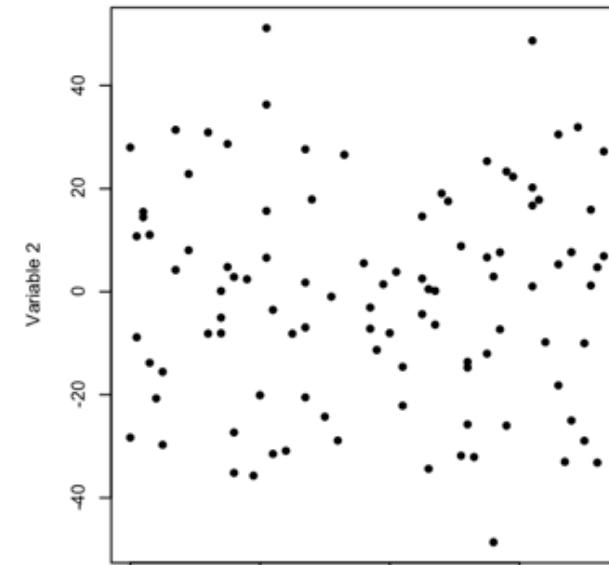
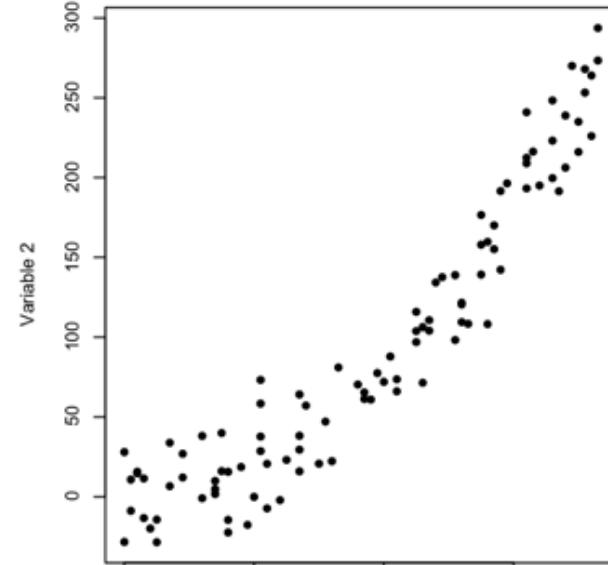
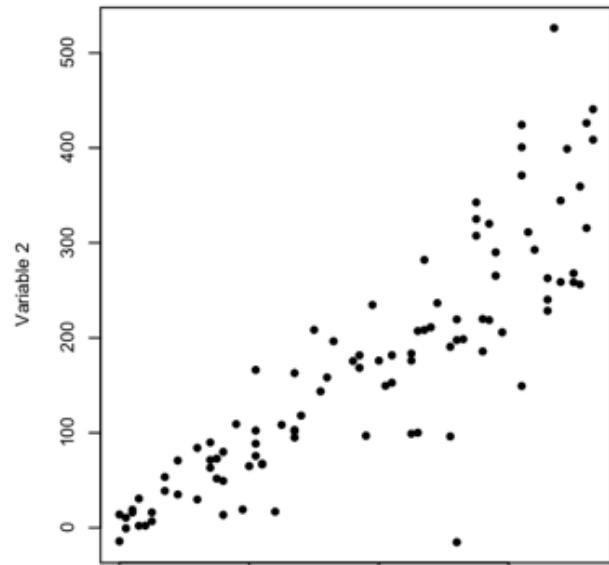
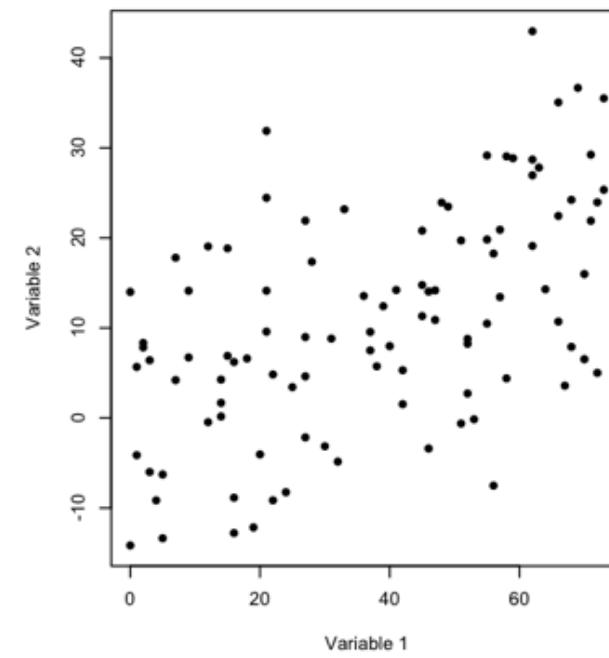
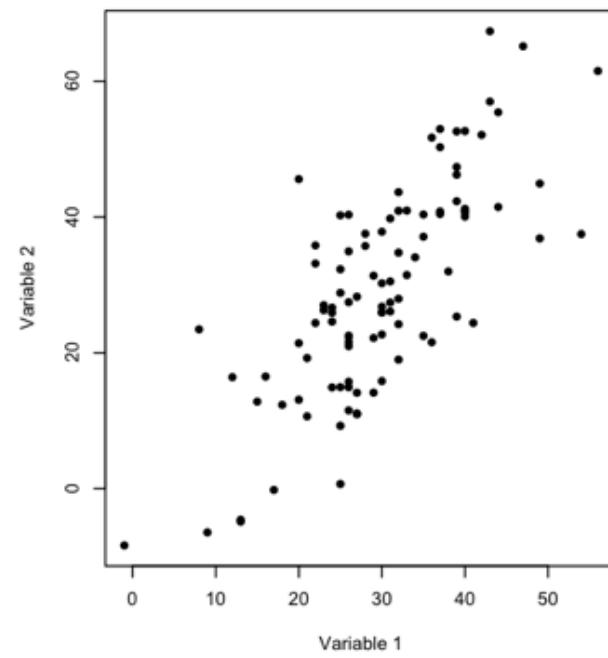
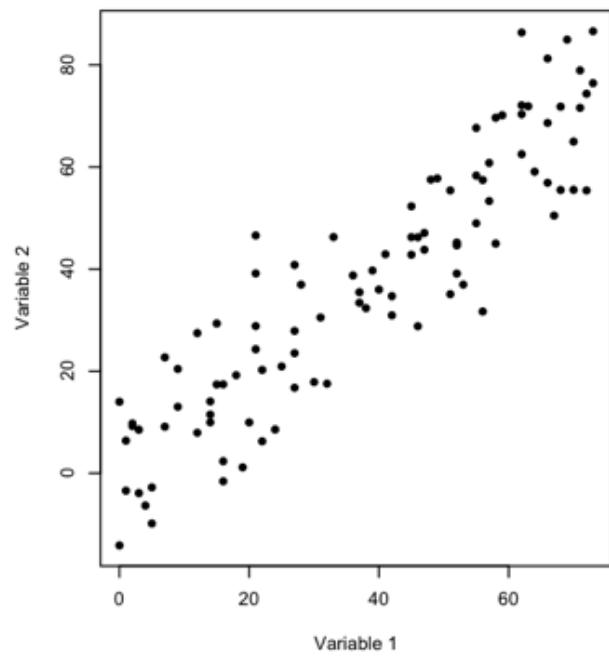


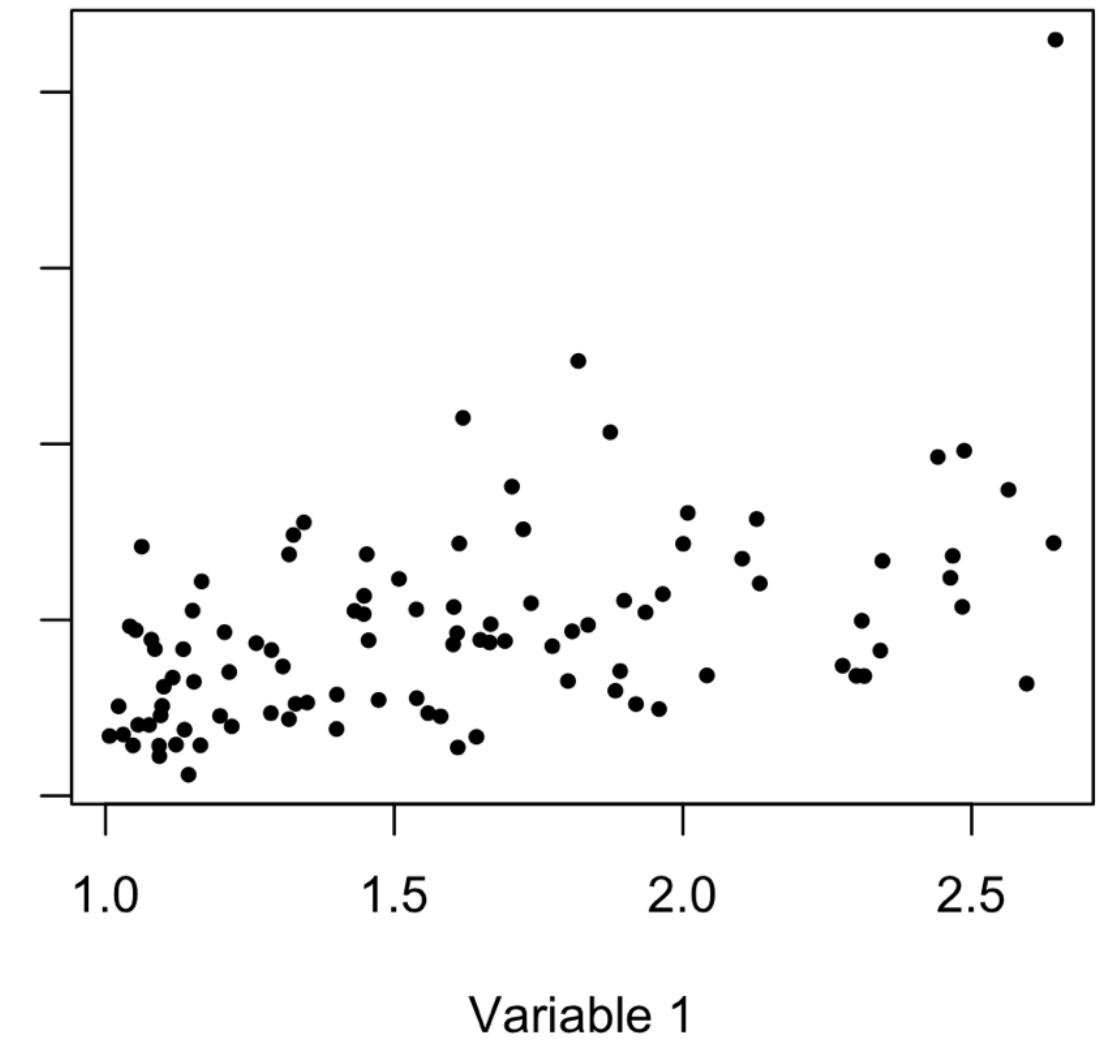
Univariate data

- Using exploratory stem-and-leaf plots, histograms and boxplots can:
- Tell us about the shape of the frequency distribution
- Help us identify outliers
- Help us identify possible errors

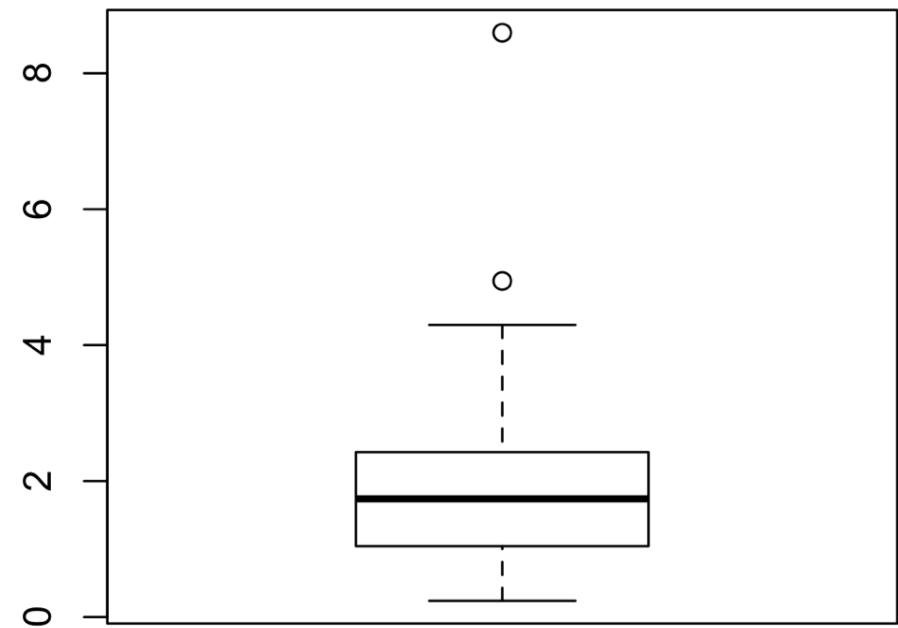
Graphics for exploratory analysis: bivariate data

- Scatter plots: use with two continuous variables
- Stripplots and boxplots: use with one continuous variable and one categorical variable
- Barplots: use with frequency data

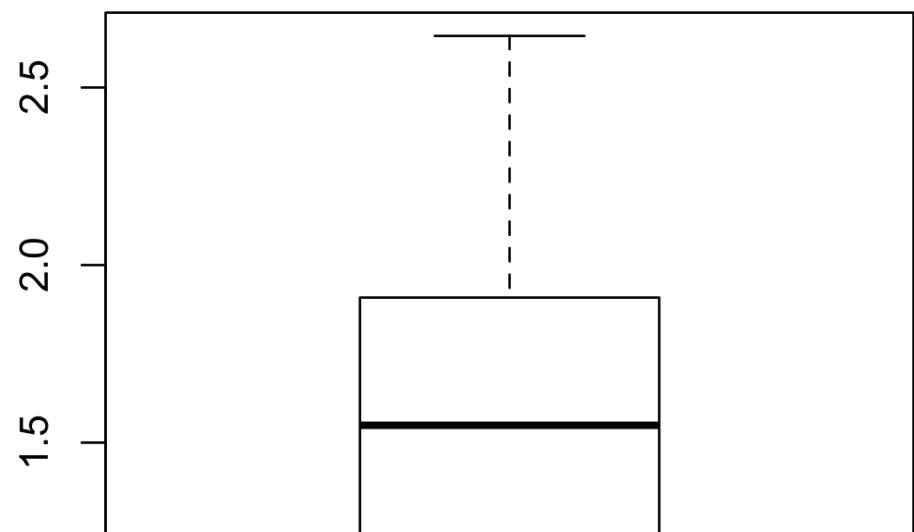


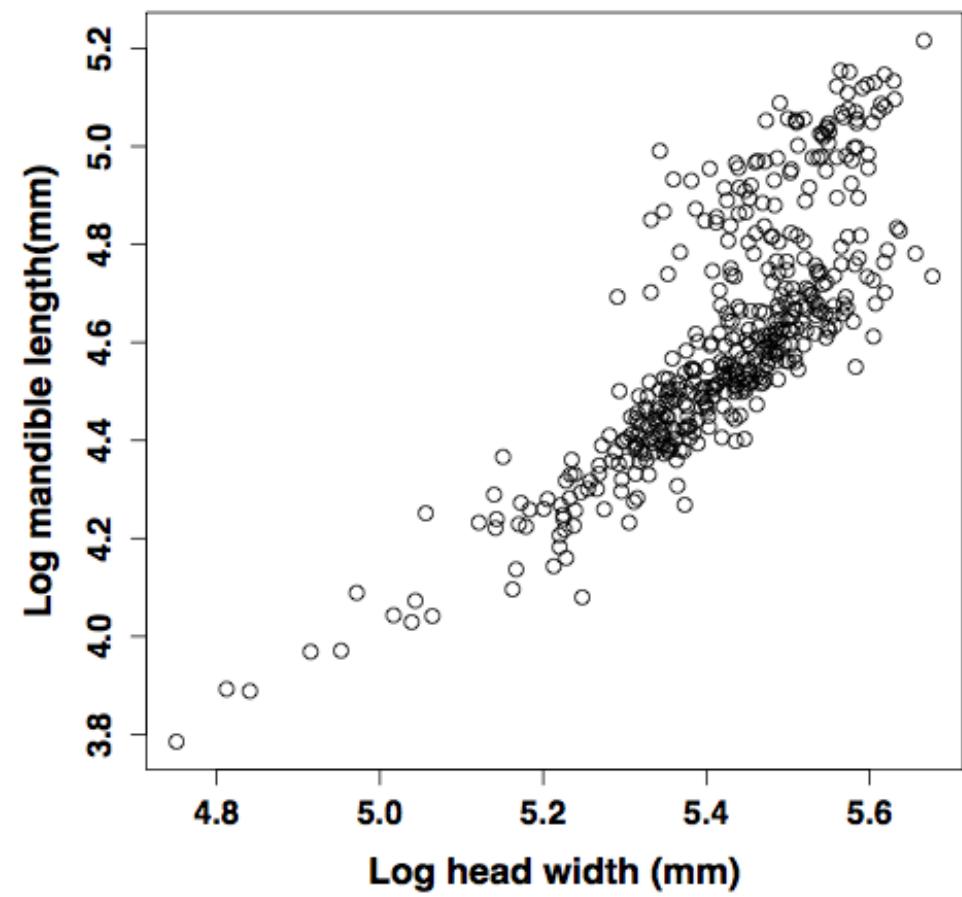
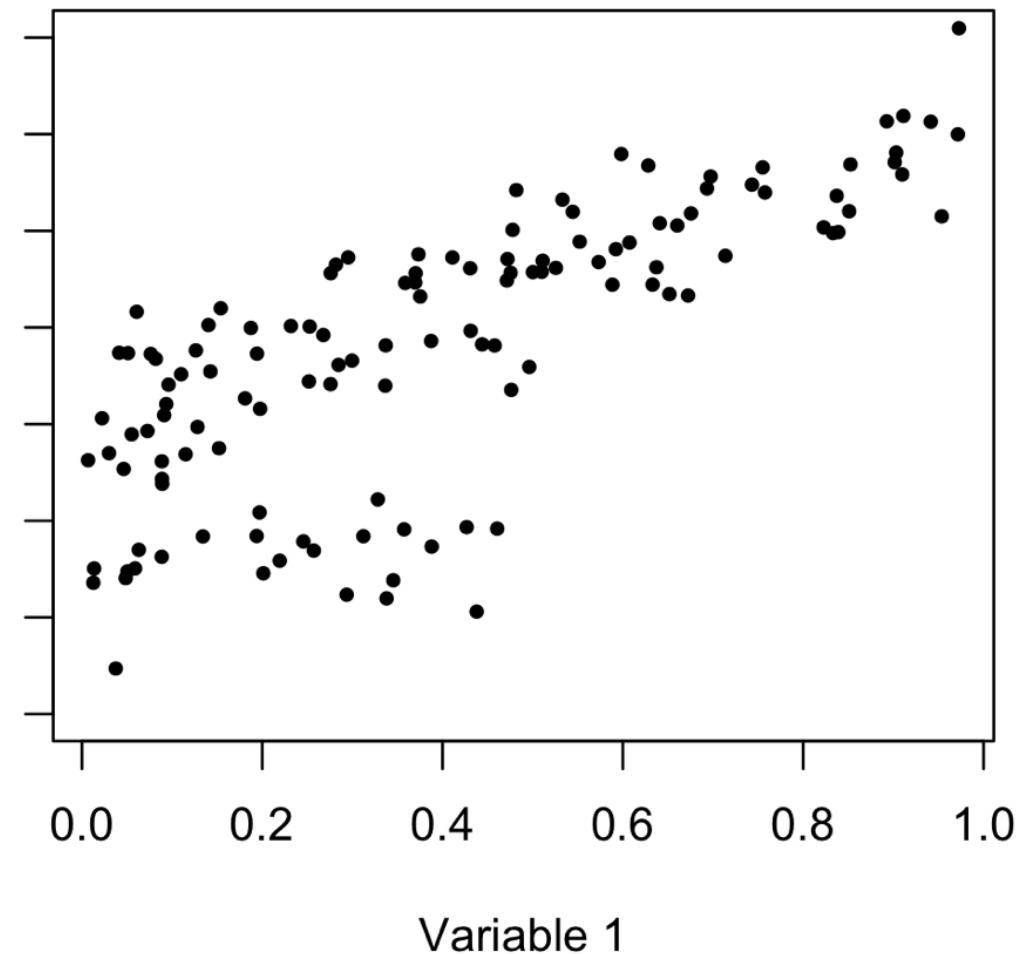


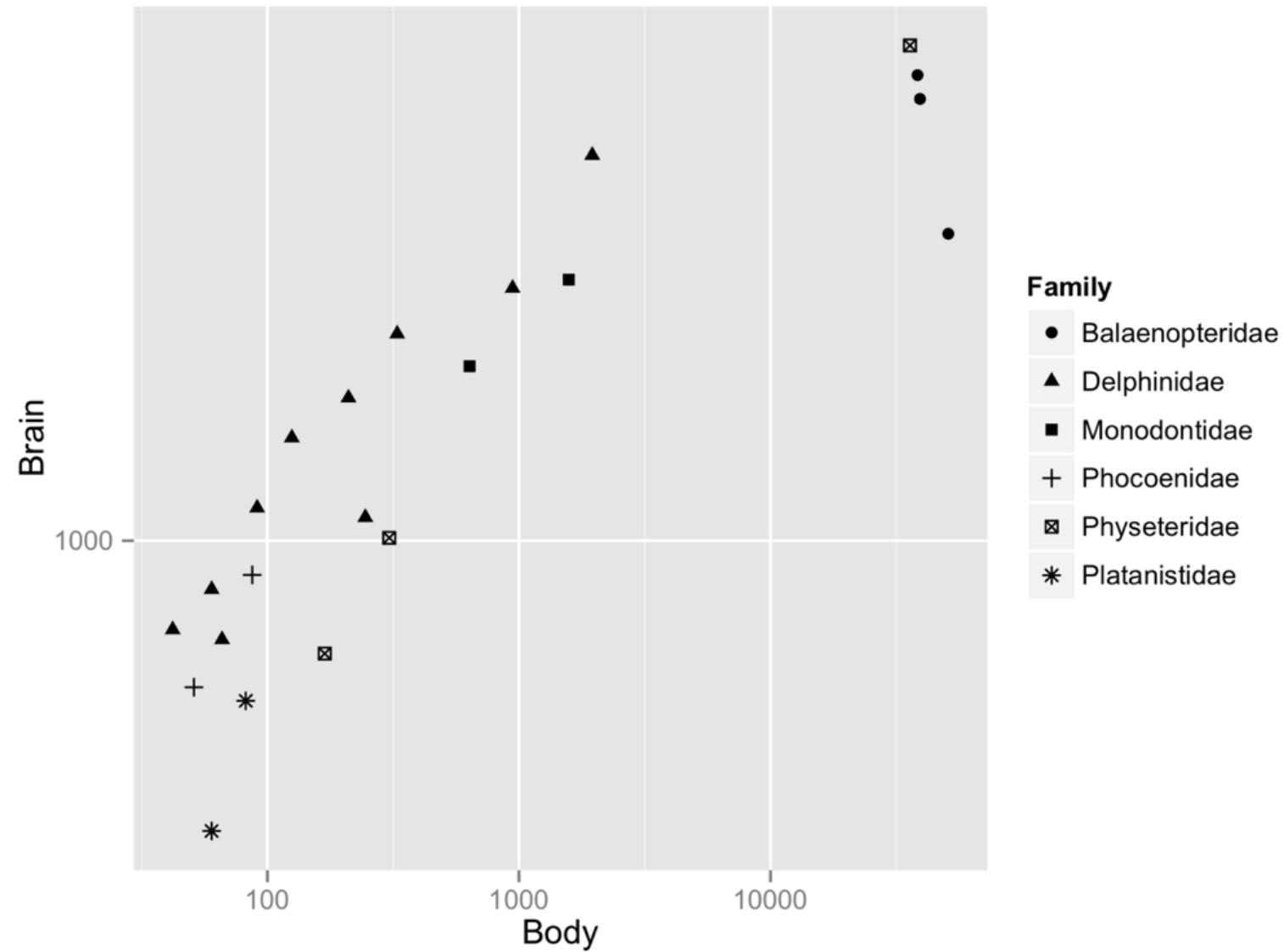
Variable 2



Variable 1



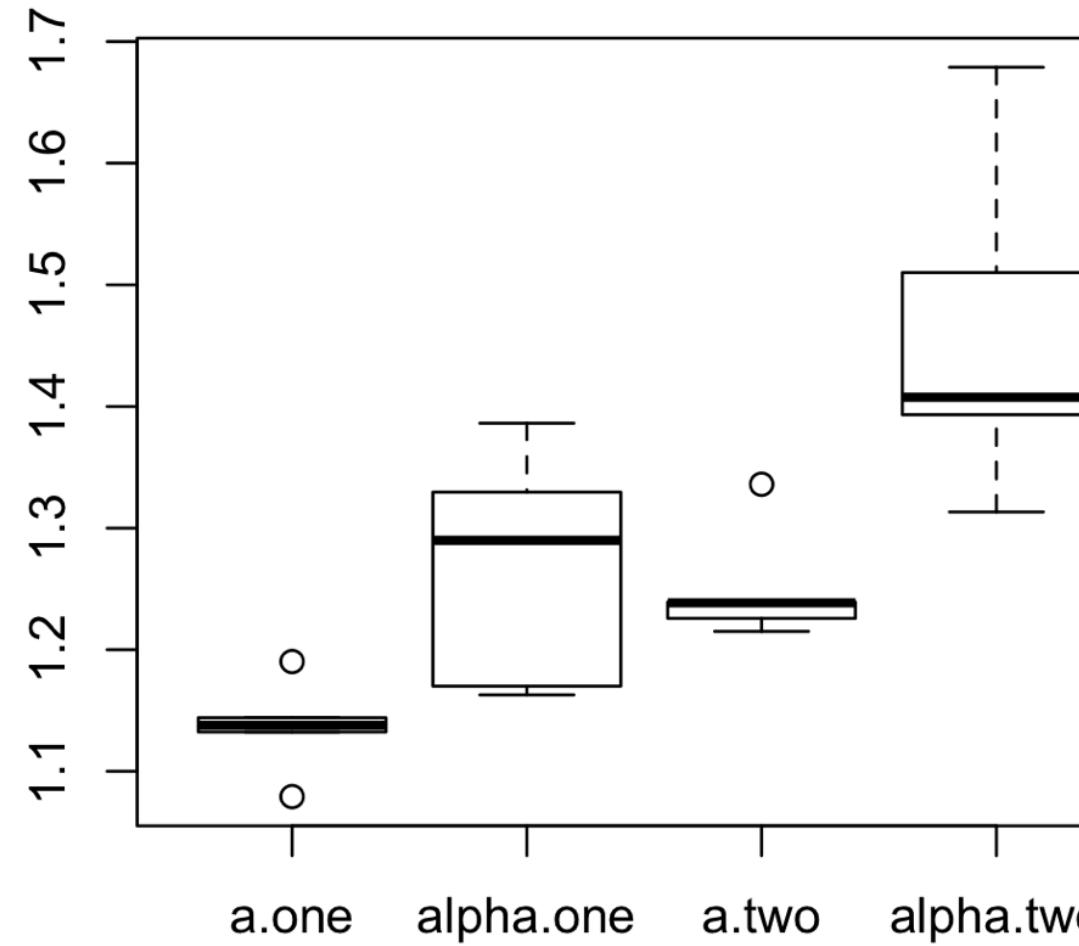
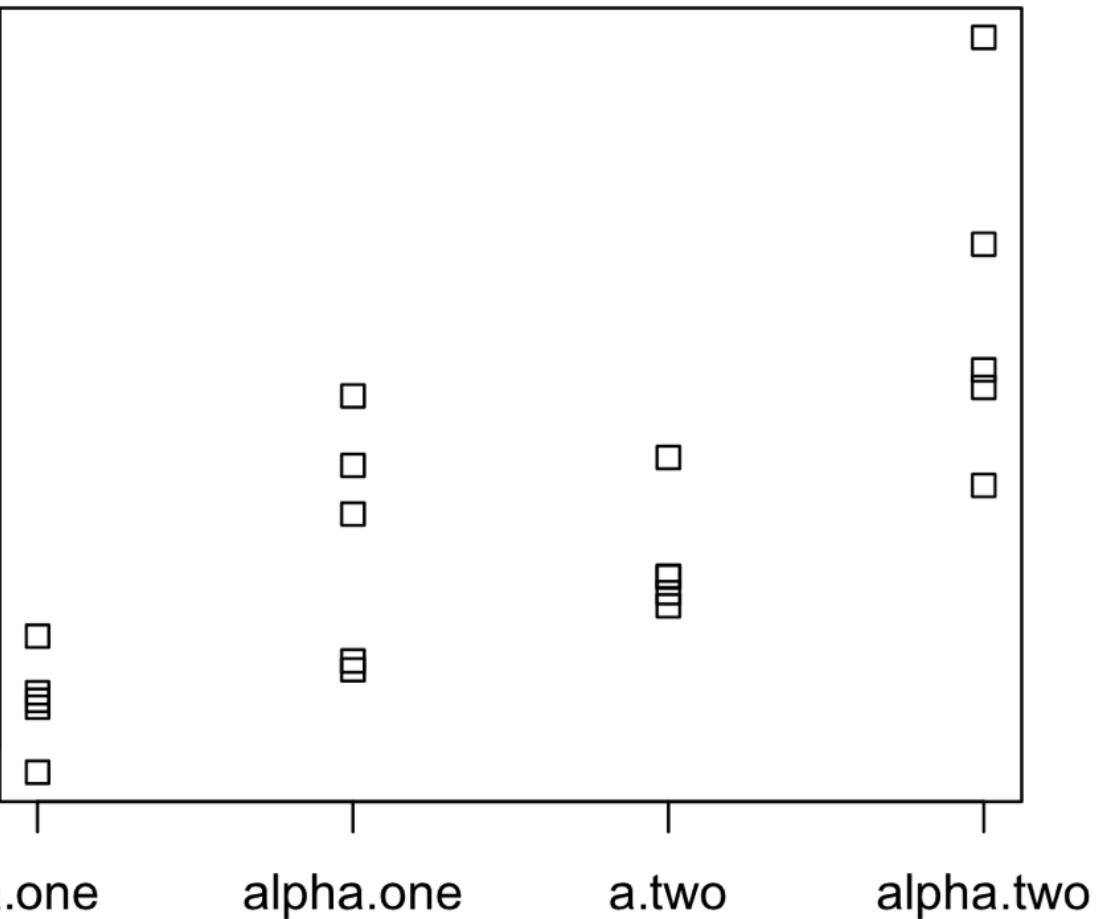




Scatterplots

- We can use scatterplots to:
 - See relationships between two variables
 - Check for non-linearity
 - Check for outliers and errors
 - Check for changes in variance
 - Check for structure in the data

Stripplots and boxplots



Stripplots and boxplots

- Use when you have one continuous variable and one categorical variable
- Stripplots are less useful when n is high
- Give an impression of how the continuous variable depends on the categorical
- Allow you to identify outliers, errors and patterns in variance

Recap