

Introduction to An Introduction to Computational Data Analysis for Biology

<http://jarrettbyrnes.info/biol697>

Jarrett Byrnes

UMass Boston

Sept 4, 2012

Outline for Today

1. Why this course?
2. Who are we?
3. How will we approach the work?
4. How will this course work?
5. R!

What is this Course About?

- ▶ Introduction to
- ▶ Computational
- ▶ Data
- ▶ Analysis
- ▶ for Biology

What is this Course About?

- ▶ Introduction to - starting with the basics
- ▶ Computational
- ▶ Data
- ▶ Analysis
- ▶ for Biology

What is this Course About?

- ▶ Introduction to - starting with the basics
- ▶ Computational - programming & other computational tools
- ▶ Data
- ▶ Analysis
- ▶ for Biology

What is this Course About?

- ▶ Introduction to - starting with the basics
- ▶ Computational - programming & other computational tools
- ▶ Data - collection, curation, maintenance of information
- ▶ Analysis
- ▶ for Biology

What is this Course About?

- ▶ Introduction to - starting with the basics
- ▶ Computational - programming & other computational tools
- ▶ Data - collection, curation, maintenance of information
- ▶ Analysis - statistics
- ▶ for Biology

What is this Course About?

- ▶ Introduction to - starting with the basics
- ▶ Computational - programming & other computational tools
- ▶ Data - collection, curation, maintenance of information
- ▶ Analysis - statistics
- ▶ for Biology - SCIENCE FIRST

What is this Course About?

- ▶ Introduction to - starting with the basics
- ▶ Computational - programming & other computational tools
- ▶ Data - collection, curation, maintenance of information
- ▶ Analysis - statistics
- ▶ for Biology - SCIENCE FIRST

What I want for you:

To be able to go from your ideas about a system to a model, fit and evaluated with the appropriate data.

Course Goals

Course Goals

1. Learn how to think about your research in a systematic way to design efficient observational & experimental studies.

Course Goals

1. Learn how to think about your research in a systematic way to design efficient observational & experimental studies.
2. Understand how to get the most bang for your buck from your data.

Course Goals

1. Learn how to think about your research in a systematic way to design efficient observational & experimental studies.
2. Understand how to get the most bang for your buck from your data.
3. Make you effective collaborators with statisticians.

Course Goals

1. Learn how to think about your research in a systematic way to design efficient observational & experimental studies.
2. Understand how to get the most bang for your buck from your data.
3. Make you effective collaborators with statisticians.
4. Make you comfortable enough to learn and grow beyond this class.

Why a Computational Focus?

```
library(plyr)

d_ply(eelgrass, .genotypes, function(x) {
  print(summary(lm(shoots ~ geese, data = x)))
})
```

Why a Computational Focus?

```
library(plyr)

d_ply(eelgrass, .genotypes, function(x) {
  print(summary(lm(shoots ~ geese, data = x)))
})
```

- ▶ Programming is a necessary skill for everything

Why a Computational Focus?

```
library(plyr)

d_ply(eelgrass, .genotypes, function(x) {
  print(summary(lm(shoots ~ geese, data = x)))
})
```

- ▶ Programming is a necessary skill for everything
- ▶ We live in the era of big data

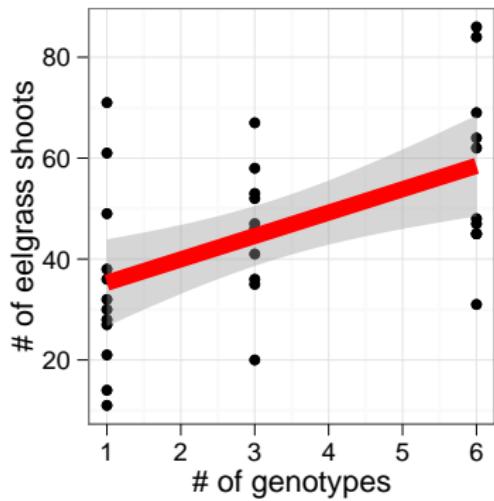
Why a Computational Focus?

```
library(plyr)

d_ply(eelgrass, .genotypes, function(x) {
  print(summary(lm(shoots ~ geese, data = x)))
})
```

- ▶ Programming is a necessary skill for everything
- ▶ We live in the era of big data
- ▶ Comfort with algorithmic thinking helps your science

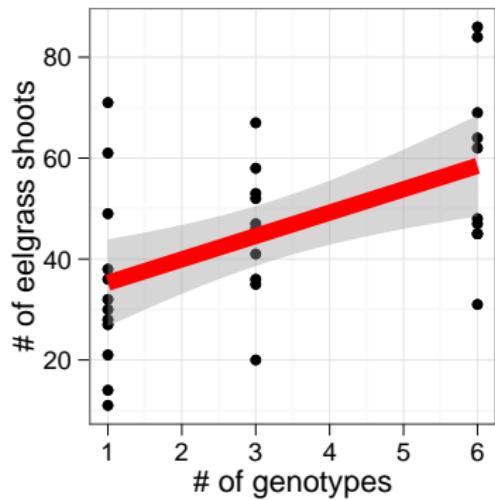
How will we use statistics?



How will we use statistics?

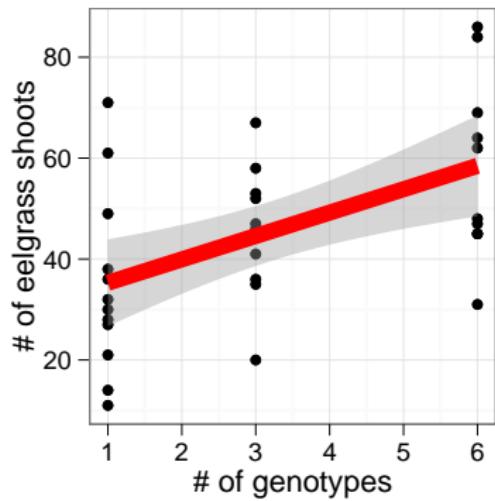
- ▶ Estimation

- ▶ Parameter in model
- ▶ Variance in parameter estimation



How will we use statistics?

- ▶ Estimation
 - ▶ Parameter in model
 - ▶ Variance in parameter estimation
- ▶ Model Evaluation
 - ▶ What parameters should be included in a model?
 - ▶ Does a model fit the data?
 - ▶ Comparison of competing hypotheses



Two Different Skillsets

- ▶ Statistics

- ▶ Programming

Questions?

Who Am I?

A wide-angle photograph of a coastal landscape. The foreground features a dark, sandy beach with white-capped waves crashing against a rocky shoreline. To the right, a steep, green-covered hillside slopes down to the water. In the middle ground, a large, dark rock formation juts out from the sea. The background shows a vast, calm ocean meeting a range of green hills under a heavy, overcast sky.

1. Marine Ecologist



2. Climate Impacts in Kelp Forests



Jackie Sones



3. Biodiversity's Effect on Ecosystem Function

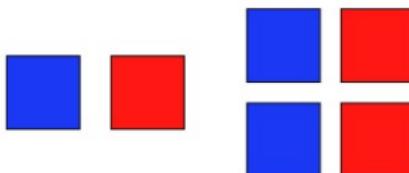
Kris Aquilino

Monocultures



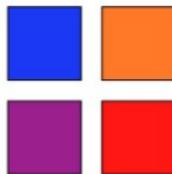
N=6 for each species and density

2 Species Combinations



N=3 for each combination and density

Polyculture

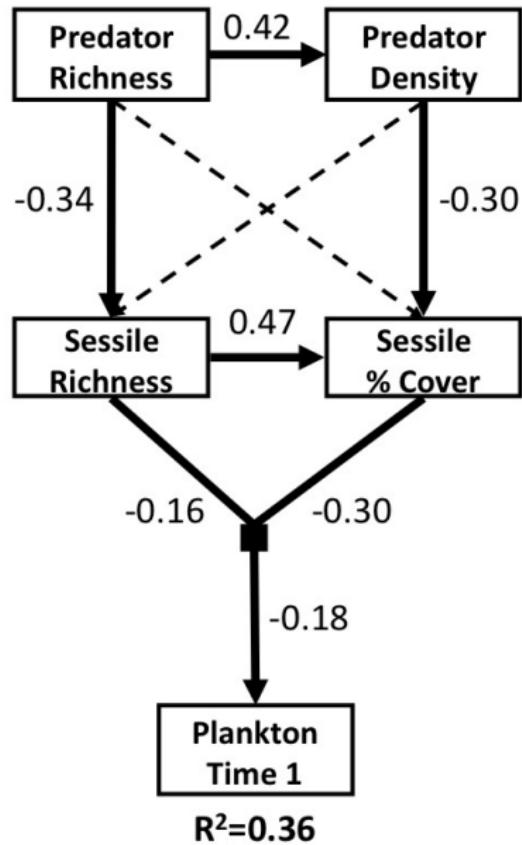


N=8

4. Experimentalist



5. Big Data from Large Networks



6. Link Scientific Understanding of a Complex World to Statistical Models

Who Are You?

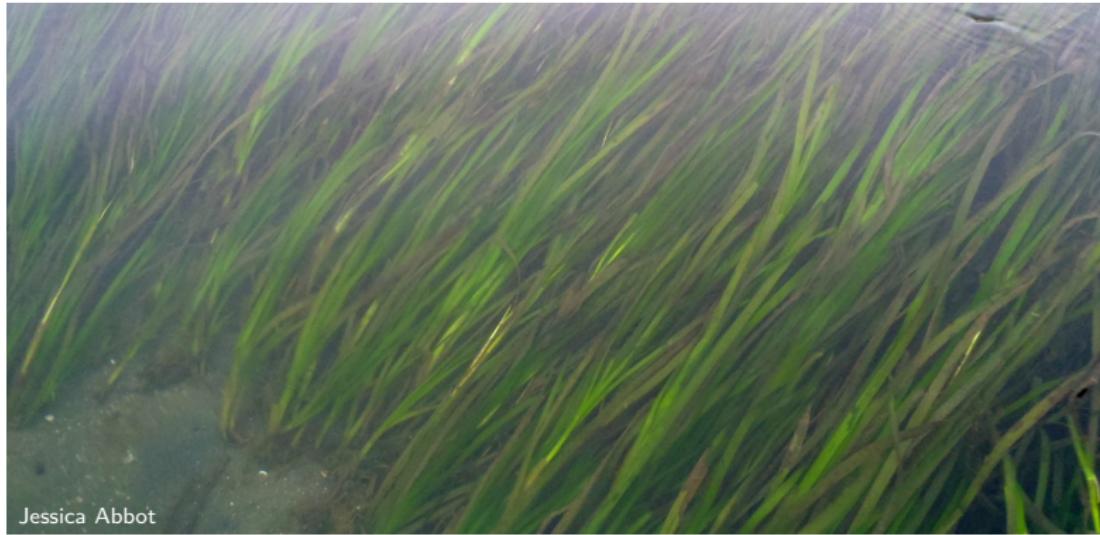
Who Are You?

1. Name
2. Lab
3. Brief research description
4. Why are you here?

Our Approach to Data Analysis

Data from Reusch et al. 2005 PNAS

Start with a Question



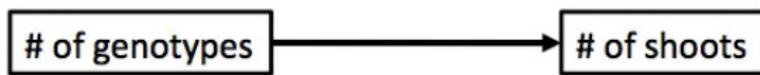
Jessica Abbot

Does seagrass genetic diversity increase productivity?

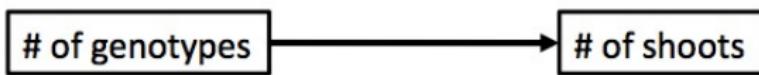
Build an Understanding of the System

1. Literature
2. Observation
3. Natural History

Construct a Model of the System

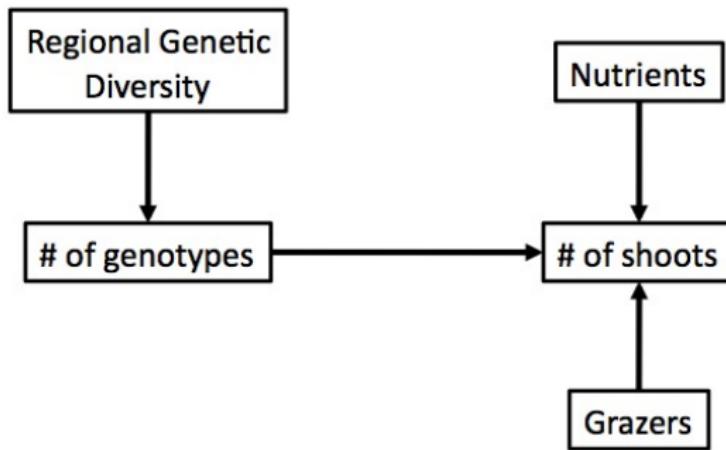


Construct a Model of the System



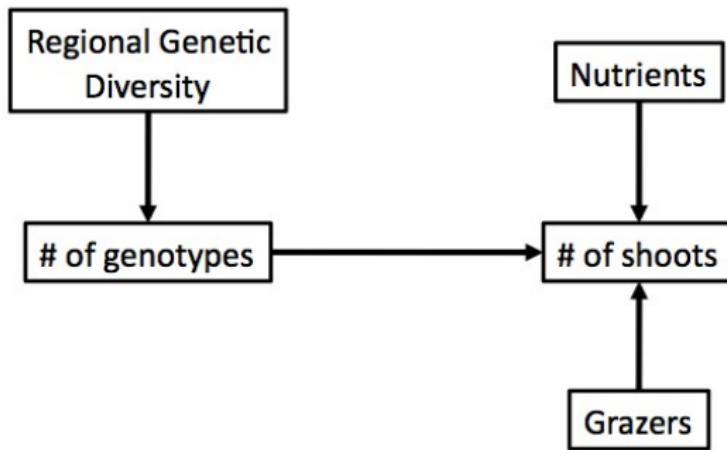
Causal Graph

Construct a Model of the System



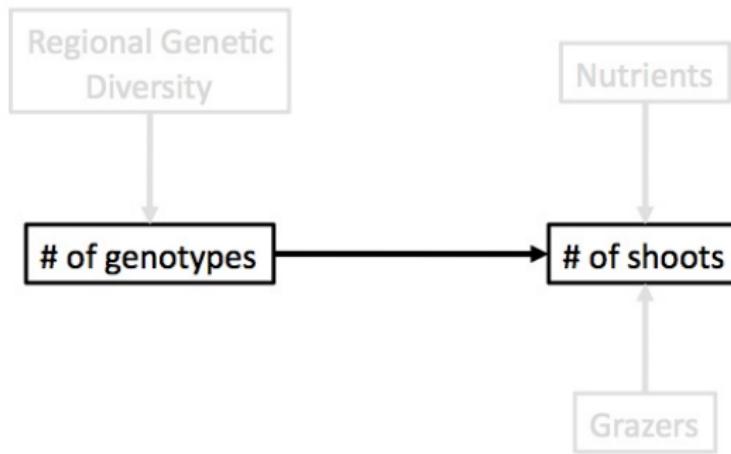
Causal Graph

Construct a Model of the System



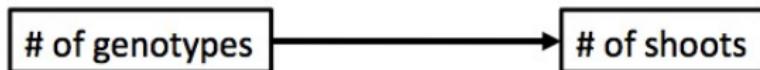
Causal Graph

Construct a Model of the System



Causal Graph

Collect the Data to Best Estimate & Test the Model

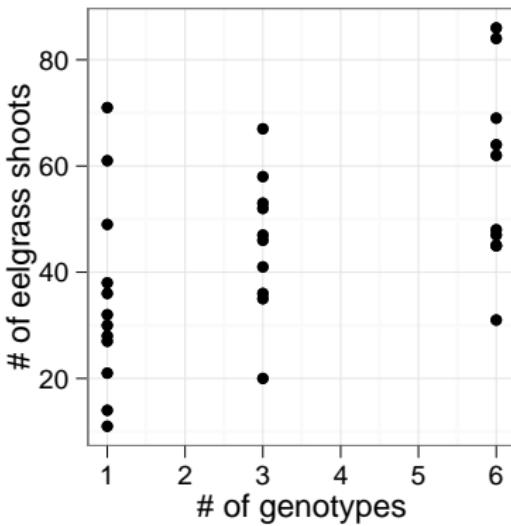


1
genotype

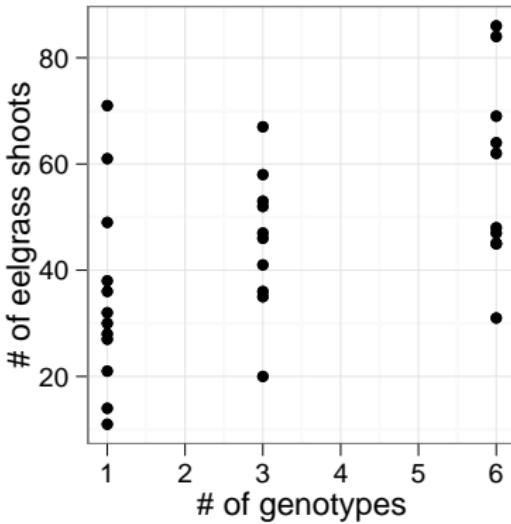
3
genotypes

6
genotypes

Look at Your Data

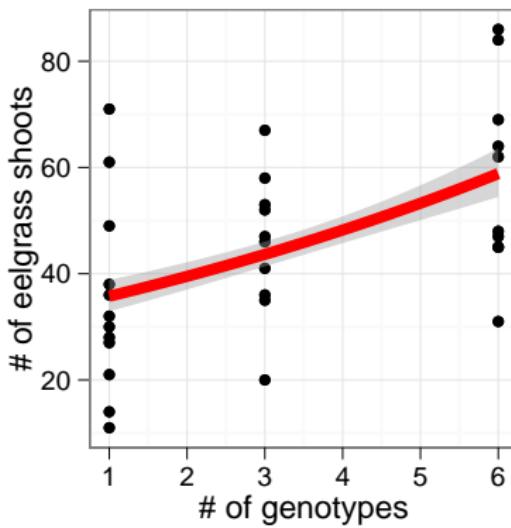


Look at Your Data



Fit a model(s), chosen to suit this data

Analysis!



Build Open Reproducible Research

Many Methods of Sharing Data, Methods, and Results Beyond Publication

1. GitHub - public code repository
2. FigShare - share key figures, get a doi
3. Blog - open 'notebook'
4. Dryad or Other Repository - post-publication data sharing

Lecture/Lab/Labinar?

- ▶ I will yammer on

Lecture/Lab/Labinar?

- ▶ I will yammer on
- ▶ R lab will be part of class

Lecture/Lab/Labinar?

- ▶ I will yammer on
- ▶ R lab will be part of class
- ▶ Notes available at <http://jarrettbyrnes.info/biol697>

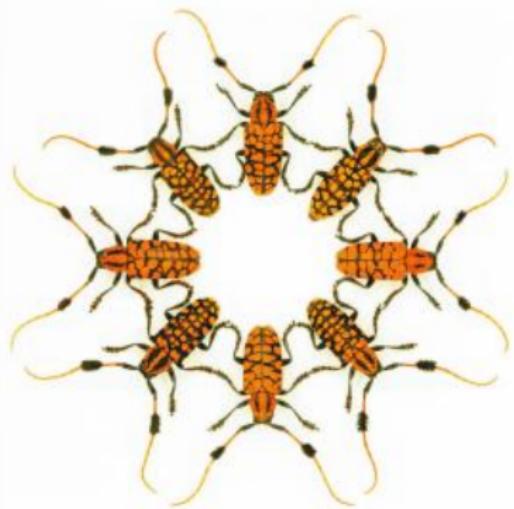
Lecture/Lab/Labinar?

- ▶ I will yammer on
- ▶ R lab will be part of class
- ▶ Notes available at <http://jarrettbyrnes.info/biol697>
- ▶ Slide source available at
<http://github.com/jbyrnes/biol697>

Special Topics

Additional special topics mini-labinars, e.g. knitr & LaTeX

Readings for Class



The Analysis of Biological Data
WHITLOCK · SCHLUTER

Whitlock, W.C. and Schluter, D.
(2008) The Analysis of Biological
Data. Roberts and Company
Publishers.

<http://www.zoology.ubc.ca/~whitlock/ABD/teaching/index.html>

Readings for Class



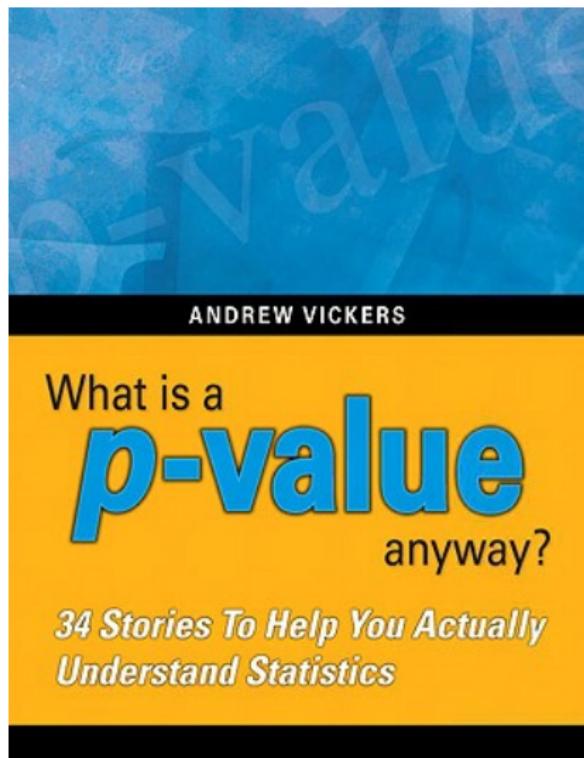
R

IN A NUTSHELL

A Desktop Quick Reference

Adler, J. (2009) R in a Nutshell:
A Desktop Quick Reference.
O'Reilly.

Reflections



Media. Vickers, A. (2009) What is a p-value anyway? 34 Stories to Help You Actually Understand Statistics. Addison Wesley.

Write a weekly reflection. 1 page.
Graded for participation (10%).
1 entry posted per week for discussion.

<http://learningdata.wordpress.com/>

Problem Sets

- ▶ 40% of your grade
- ▶ Adapted from Whitlock and Schluter
- ▶ Will often require R
- ▶ Turn in all code, and it must be understandable

Practical Exams

- ▶ 20% Midterm, 30% final
- ▶ Real world data analysis problems
- ▶ Will require R
- ▶ Turn in all code, and it must be understandable

Extra Credit: Your Work

- ▶ 10% Extra
- ▶ Report on your own data
- ▶ Cogently present what you did, why you did it, and the results & interpretation
- ▶ Data & Code must be accessible & understandable
- ▶ Extra points for putting work online so others can use & view your work

Topics

1. Data & Data Management

Topics

1. Data & Data Management
2. Biological Processes & Statistical Distributions

Topics

1. Data & Data Management
2. Biological Processes & Statistical Distributions
3. Data Visualization

Topics

1. Data & Data Management
2. Biological Processes & Statistical Distributions
3. Data Visualization
4. Simulation & Basic Hypothesis Testing

Topics

1. Data & Data Management
2. Biological Processes & Statistical Distributions
3. Data Visualization
4. Simulation & Basic Hypothesis Testing
5. Sampling Design

Topics

1. Data & Data Management
2. Biological Processes & Statistical Distributions
3. Data Visualization
4. Simulation & Basic Hypothesis Testing
5. Sampling Design
6. Fitting Linear Models: Least Squares

Topics

1. Data & Data Management
2. Biological Processes & Statistical Distributions
3. Data Visualization
4. Simulation & Basic Hypothesis Testing
5. Sampling Design
6. Fitting Linear Models: Least Squares
7. Fitting Linear Models: Likelihood

Topics

1. Data & Data Management
2. Biological Processes & Statistical Distributions
3. Data Visualization
4. Simulation & Basic Hypothesis Testing
5. Sampling Design
6. Fitting Linear Models: Least Squares
7. Fitting Linear Models: Likelihood
8. Generalized Linear Models

Topics

1. Data & Data Management
2. Biological Processes & Statistical Distributions
3. Data Visualization
4. Simulation & Basic Hypothesis Testing
5. Sampling Design
6. Fitting Linear Models: Least Squares
7. Fitting Linear Models: Likelihood
8. Generalized Linear Models
9. Experiments & the Linear Model (ANOVA)

Topics

1. Data & Data Management
2. Biological Processes & Statistical Distributions
3. Data Visualization
4. Simulation & Basic Hypothesis Testing
5. Sampling Design
6. Fitting Linear Models: Least Squares
7. Fitting Linear Models: Likelihood
8. Generalized Linear Models
9. Experiments & the Linear Model (ANOVA)
10. Multiple Continuous Predictors

Topics

1. Data & Data Management
2. Biological Processes & Statistical Distributions
3. Data Visualization
4. Simulation & Basic Hypothesis Testing
5. Sampling Design
6. Fitting Linear Models: Least Squares
7. Fitting Linear Models: Likelihood
8. Generalized Linear Models
9. Experiments & the Linear Model (ANOVA)
10. Multiple Continuous Predictors
11. What should I sample? Simpson's Paradox

Topics

1. Data & Data Management
2. Biological Processes & Statistical Distributions
3. Data Visualization
4. Simulation & Basic Hypothesis Testing
5. Sampling Design
6. Fitting Linear Models: Least Squares
7. Fitting Linear Models: Likelihood
8. Generalized Linear Models
9. Experiments & the Linear Model (ANOVA)
10. Multiple Continuous Predictors
11. What should I sample? Simpson's Paradox
12. Interactions & Nonlinearities

Topics

1. Data & Data Management
2. Biological Processes & Statistical Distributions
3. Data Visualization
4. Simulation & Basic Hypothesis Testing
5. Sampling Design
6. Fitting Linear Models: Least Squares
7. Fitting Linear Models: Likelihood
8. Generalized Linear Models
9. Experiments & the Linear Model (ANOVA)
10. Multiple Continuous Predictors
11. What should I sample? Simpson's Paradox
12. Interactions & Nonlinearities
13. Bootstrapping

Topics

1. Data & Data Management
2. Biological Processes & Statistical Distributions
3. Data Visualization
4. Simulation & Basic Hypothesis Testing
5. Sampling Design
6. Fitting Linear Models: Least Squares
7. Fitting Linear Models: Likelihood
8. Generalized Linear Models
9. Experiments & the Linear Model (ANOVA)
10. Multiple Continuous Predictors
11. What should I sample? Simpson's Paradox
12. Interactions & Nonlinearities
13. Bootstrapping
14. Model Comparison

Questions?

What is R?

R is an integrated suite of software facilities for data manipulation, calculation and graphical display. It includes

- ▶ *an effective data handling and storage facility,*
- ▶ *a suite of operators for calculations on arrays, in particular matrices,*
- ▶ *a large, coherent, integrated collection of intermediate tools for data analysis, graphical facilities for data analysis and display either on-screen or on hardcopy, and*
- ▶ *a well-developed, simple and effective programming language which includes conditionals, loops, user-defined recursive functions and input and output facilities.*

From <http://r-project.org>

What is R?

- ▶ A programming language uniquely developed for statistical analysis

Why R?

1. Free

Why R?

1. Free
2. Huge growing community

Why R?

1. Free
2. Huge growing community
3. Packages to do almost anything

Why R?

1. Free
2. Huge growing community
3. Packages to do almost anything
4. Makes reusable research easy

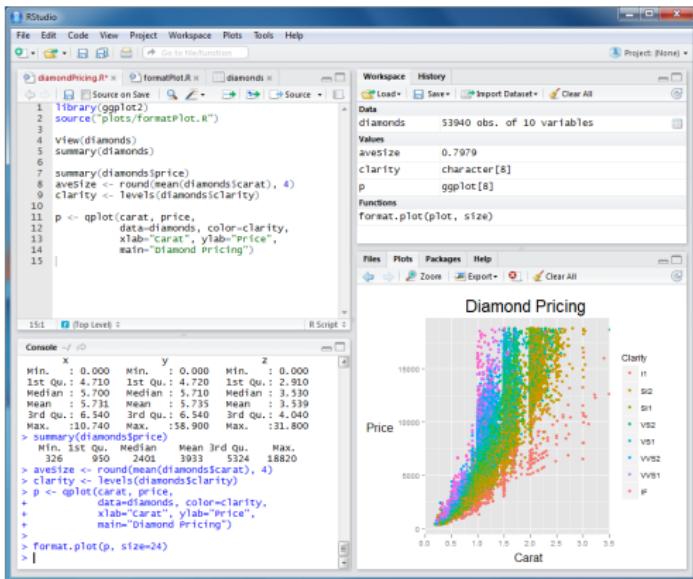
Why R?

1. Free
2. Huge growing community
3. Packages to do almost anything
4. Makes reusable research easy
5. C-based language

Why R?

1. Free
2. Huge growing community
3. Packages to do almost anything
4. Makes reusable research easy
5. C-based language
6. Syntax naturally matches analytical thinking

What is R Studio?



- ▶ Cross-Platform Graphical User Interface for R
- ▶ It is not R

Let's Fire It Up!

Open R-Studio.

Don't have it? Download it from <http://rstudio.org>

What do you see?

```
R version 2.14.2 (2012-02-09)
Copyright (C) 2012 The R Foundation for Statistical Computing
ISBN 3-900051-07-0
Platform: i386-apple-darwin9.8.1 (32-bit)

It 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

It 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,
'help.start()' for an HTML browser interface to help,
'type()' to quit R.
```

[Workspace loaded from ./RData]

24 obs. of 4 variables

Values

- gcl.working.env
- l
- l.vec
- log

Functions

- all
- apply
- isolate
- lapply
- lvec
- log
- mean
- plot
- print
- sum

PlotMeans(response, factor1, factor2, error.bars = c("se", "sd", "conf.int", "none"), level = 0.95, slab = deparse(substitute(factor1))), ylab = deparse(substitute(factor2)), legend.lab = deparse(substitute(factor2)), main = "Plot of Means", pch = 1:nlevels.l, lty = 1:nlevels.l, col = palette(c), ylab = NA)

PlotMeans(response, factor1, factor2, error.bars = c("se", "sd", "conf.int", "none"), level = 0.95, plot.type = "points", width.css = 2, bar.col = "Lightblue", slab = deparse(substitute(factor1)), pch = palette(mean.c), main = "Plot of Means", pch = 1:nlevels.l, lty = 1:nlevels.l, col = palette(c), ylab = NA)

PlotMeans(response, factor1, factor2, error.bars = c("se", "sd", "conf.int", "none"), level = 0.95, plot.type = "jitter", main = "Plot of Means", pch = NA, lty = 1:nlevels.l, col = palette(c), error.bars.lty = 1, dividers = TRUE, just.legend = T, cex.lab = 3, cex.axis = 1.4, cex.legends = 1, pch.col = "black", type = "o", pch = NA, just.outside = 1, has.legend = T, cex.lab = 3, cex.axis = 1, mar = c(5, 4, 4, 4), xpd = 0.1)

What do you see?

code editor

The screenshot shows the RStudio interface. The top navigation bar includes tabs for 'Untitled 1', 'Source on Save', 'Run', 'File', 'Edit', 'Load', 'Tools', 'Import Dataset', and 'Clear All'. Below the navigation bar is the 'Workspace' pane, which lists variables and functions:

- Variables:**
 - g1: working.env
 - a: num
 - b: vec
 - msg: "JGR"
- Functions:**
 - all: Fortify, tapply
 - allclose
 - boxplot
 - boxplot(response, factor1, factor2, error.bars = c("se", "sd", "conf.int", "none"), level = 0.95, slab = deparse(substitute(factor1))), ylab = paste("mean", deparse(substitute(factor2))), legend.lab = deparse(substitute(factor2)), main = "Plot of Mean", pch = 1:nlevels.J, lty = 1:nlevels.J, col = palette(J), ylab = NA
 - boxplot(response, factor1, factor2, error.bars = c("se", "sd", "conf.int", "none"), level = 0.95, plot.type = "points", width.css = 2, bar.col = "lightblue", slab = deparse(substitute(factor1)), pch = palette(mean\$ef), main = deparse(substitute(factor2))), legend.lab = deparse(substitute(factor2)), main = "Plot of Mean", pch = 1:nlevels.J, lty = 1:nlevels.J, col = palette(J), ylab = NA, xlab = "Mean", ylab = NA, xaxt = "n", yaxt = "n", xaxs = "r", yaxs = "r", xpd = 0, xaxp = c(0, 1, 0.5, 1), yaxp = c(0, 1, 0.5, 1), xaxs = "r", yaxs = "r", xpd = 0, xaxp = c(0, 1, 0.5, 1), yaxp = c(0, 1, 0.5, 1), pch = NA, lty = 1:nlevels.J, col = palette(J), error.bars.lty = 1, dividers = TRUE, just = "left", t.cex = 1, cex.main = 1, mar = c(5, 4, 4, 4), b3d = 0.1)

The 'Console' pane at the bottom left shows the R version information and a warning message:

```
R version 2.14.2 (2012-02-29)
Copyright (C) 2012 The R Foundation for Statistical Computing
ISBN 3-900051-07-0
Platform: i386-apple-darwin9.8.1 (32-bit)
```

A warning message follows:

```
Warning message:
In file.rename("exercises/lec01.R", "exercises/lec01.R~") :
  cannot open file 'exercises/lec01.R': No such file or directory
```

The 'File' menu at the top right includes options for New, Open, Save, Projects, Classes, and Help.

What do you see?

code editor

The screenshot shows the RStudio interface. The top navigation bar includes tabs for 'Untitled 1', 'Source on Save', 'Run', 'File', 'Edit', 'Load', 'Save', 'Import Dataset', and 'Clear All'. Below the navigation bar are two panes: 'Workspace' and 'History'. The 'Workspace' pane lists variables and functions: 'PQ_RU' (24 obs. of 4 variables), 'gci.working.env' (environment[4]), 'i' (integer[3]), 'j' (vec), 'n' (vec), and 'pch' (character[2]). The 'History' pane contains R code for plotting data. The bottom left is the 'Console' pane, which displays the R command line and its output. The bottom right is the 'File Explorer' pane, showing a directory structure for 'R sem 2012' containing 'exercises', 'introduction', 'lectures', 'Practical_Javaan_exercise', 'readings', and 'soc_analisis_2012.docx'. The bottom status bar shows file size (117.7 KB), date (Feb 9, 2012, 2:26 PM), and zoom level (100%).

console

```
R version 2.14.2 (2012-02-09)
Copyright (C) 2012 The R Foundation for Statistical Computing
ISBN 3-900051-07-0
Platform: i386-apple-darwin9.8.1/darwin9.8.1 (34-bit)

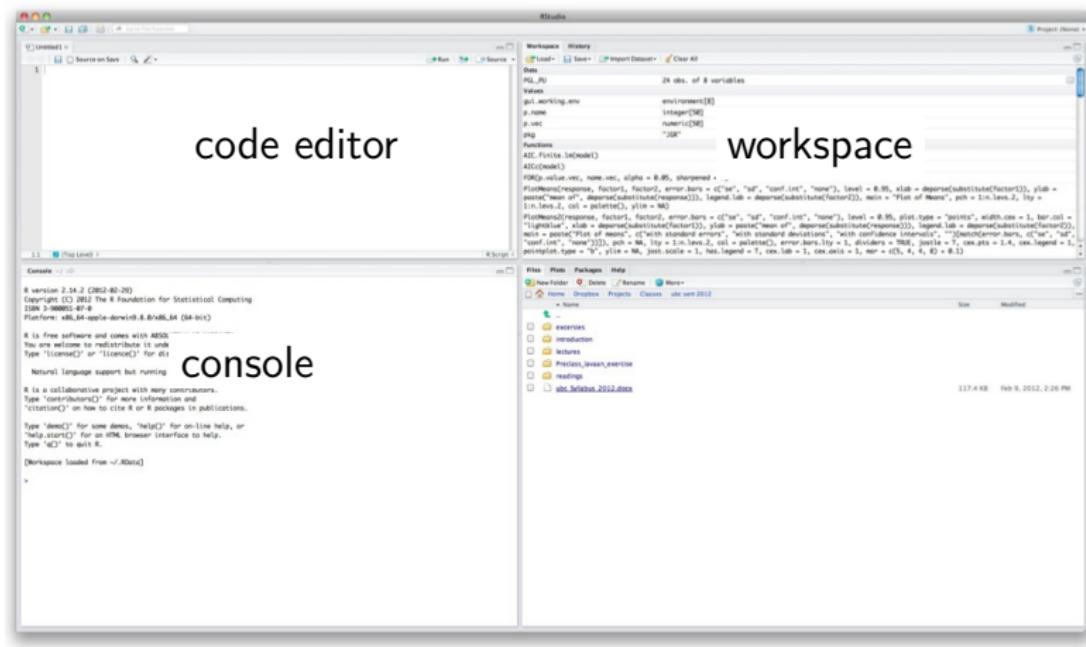
It is Free Software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under
Type 'license()' or 'licence()' for dclt
Natural language support but running

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,
'help.start()' for an HTML browser interface to help,
'type()' for writing R.
```

[Workspace loaded from ~/Rsem]

What do you see?



What do you see?

The image shows a screenshot of the RStudio interface. It features four main tabs arranged horizontally:

- code editor**: The leftmost tab, showing R code for plotting data.
- workspace**: The second tab, displaying the current workspace with objects like `pd_01`, `values`, `gci.working.env`, `a`, `b`, `c`, `vec`, and `plot`.
- console**: The third tab, showing the R command-line interface with various commands entered and their outputs.
- misc tabs**: The rightmost tab, showing a file browser with a directory structure for exercises, introduction, lectures, practice, readings, and a specific file named `sec_analyses_2012.dock`.

The Console and Math

```
1 + 1
```

```
## [1] 2
```

Everything is an Object

```
a.number <- 1 + 1
```

Everything is an Object

```
a.number <- 1 + 1
```

```
a.number
```

```
## [1] 2
```

Everything is an Object

```
a.number <- 1 + 1
```

```
a.number
```

```
## [1] 2
```

Note: Comment Your Code as You Write with

The text after # is not evaluated.

```
# This is going to be the number two  
a.number <- 1 + 1
```

Note: Comment Your Code as You Write with

The text after # is not evaluated.

```
# This is going to be the number two  
a.number <- 1 + 1
```

```
##### -----  
# You can get creative with comments to separate code  
# blocks And write a lot, which is good practice  
##### -----
```

Functions Work on Objects

```
sin(a.number)
```

```
## [1] 0.9093
```

Functions Work on Objects

```
sin(a.number)  
## [1] 0.9093
```

How to get help for a function

```
`?`(cos)  
  
help(cos)  
  
`?`(`?`("cosine function"))
```

Lots of Object Types - like Data!

```
head(cars, n = 3) #note the n= argument!  
  
##   speed dist  
## 1     4     2  
## 2     4    10  
## 3     7     4
```

Lots of Object Types - like Data!

```
head(cars, n = 3) #note the n= argument!  
  
##   speed dist  
## 1     4     2  
## 2     4    10  
## 3     7     4
```

Try looking at all of cars

Lots of Object Types - like Data!

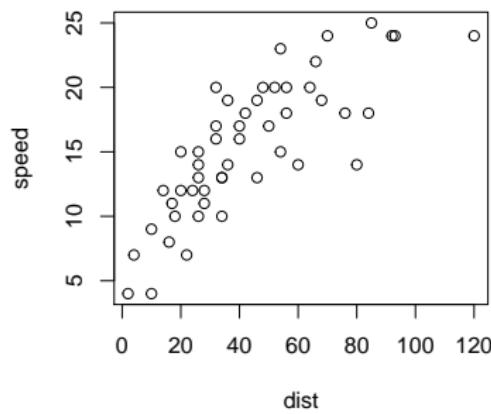
```
head(cars, n = 3) #note the n= argument!  
  
##   speed dist  
## 1     4     2  
## 2     4    10  
## 3     7     4
```

Try looking at all of cars
Can be lots of information stored in an object

```
names(cars)  
  
## [1] "speed" "dist"
```

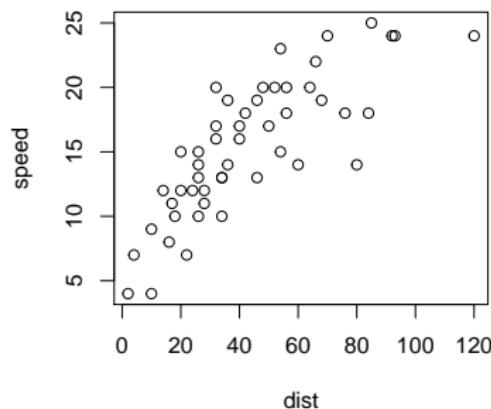
Graphics are a Snap

```
plot(speed ~ dist, data = cars)
```



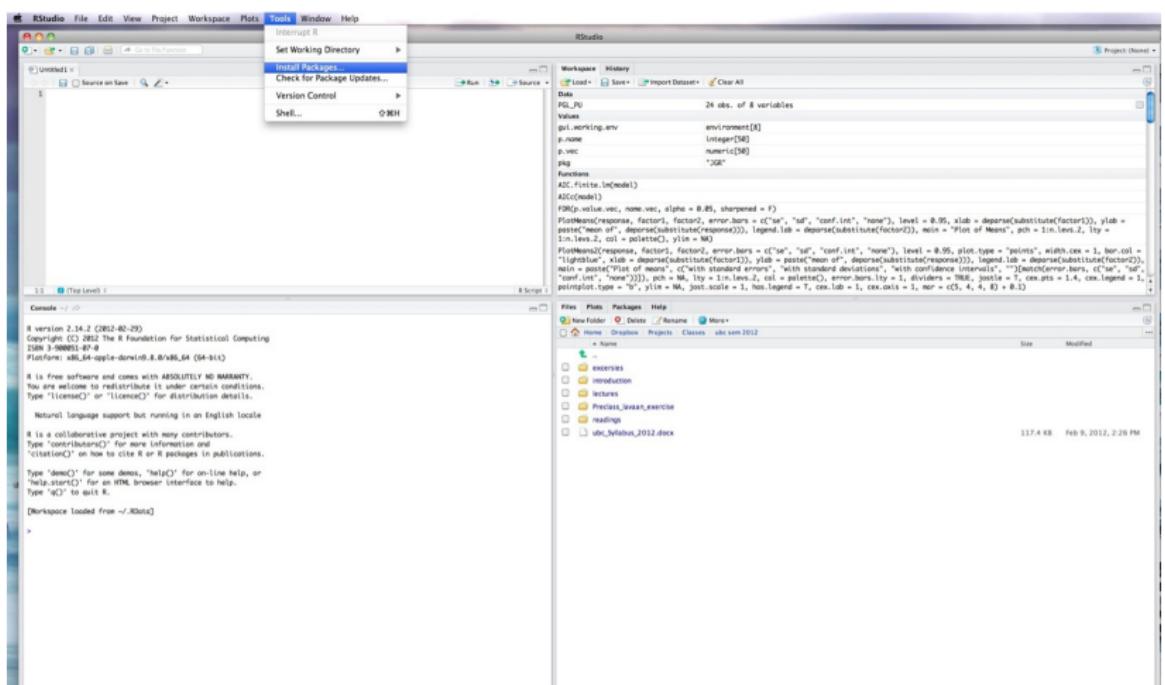
Graphics are a Snap

```
plot(speed ~ dist, data = cars)
```

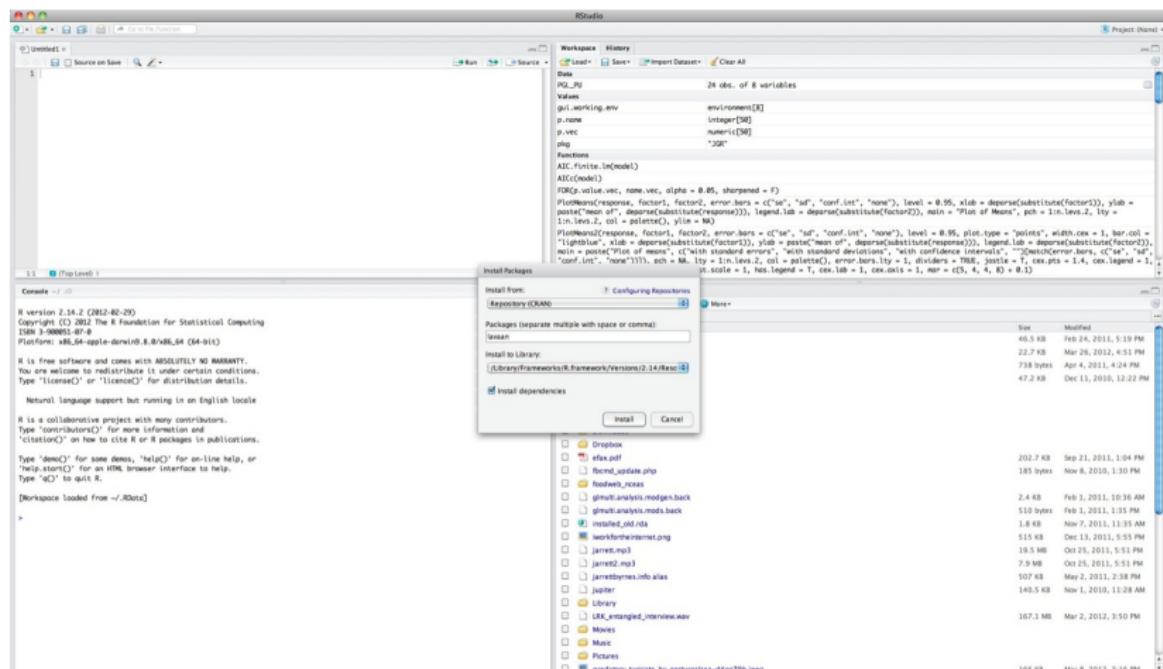


Look at ?plot to see other arguments to change appearance

Installing Packages



Installing Packages



Installing Packages

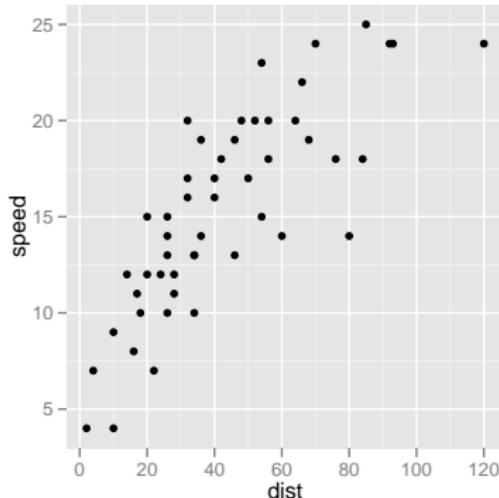
You can also install packages from the command line.

```
install.packages("ggplot2", repos = "http://cran.case.edu/",  
dependencies = TRUE)
```

Using one of the above methods, install the package ggplot2 and its dependencies now.

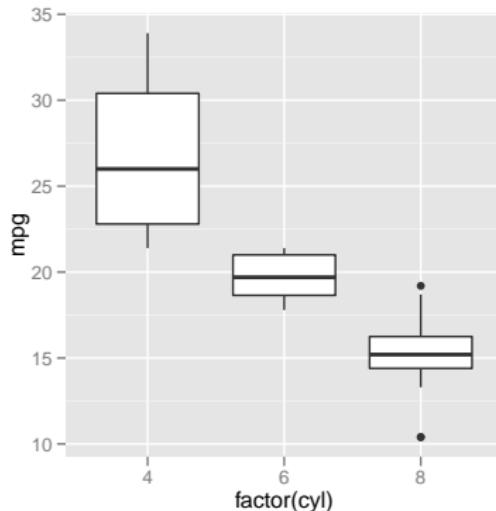
Using a Package

```
library(ggplot2)  
  
qplot(dist, speed, data = cars)
```



You Try It

- ▶ Load ggplot2 and look at the mtcars data set
- ▶ Look at the qplot help file & demos
- ▶ Make two plots



Next time

- ▶ Data Management!
- ▶ Contact me if you are not enrolled
- ▶ Read chapter 1 of the Nutshell
- ▶ Read P-Values chapters 1, 32-34 & ponder