



Demystifying statistics & programming in R

Dr Neil Burns – neil.burns@sruc.ac.uk

How the session will work

Mix of presentation are R coding

Flit between PowerPoint and R:

- You can follow along either as I run code or
- Run the code yourself as we proceed

This is now a github slide

Section headings in R will be highlighted like:

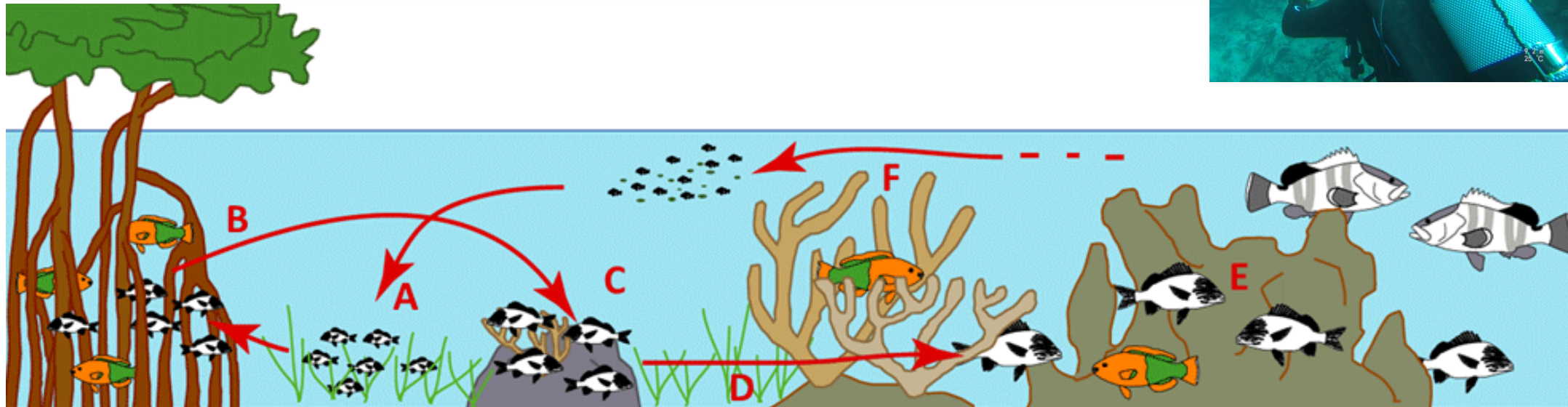
```
#####  
#####      1. some sort of title      #####  
#####
```



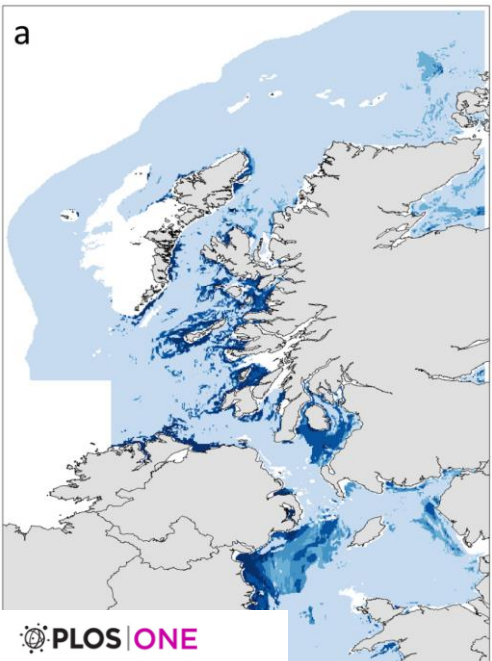
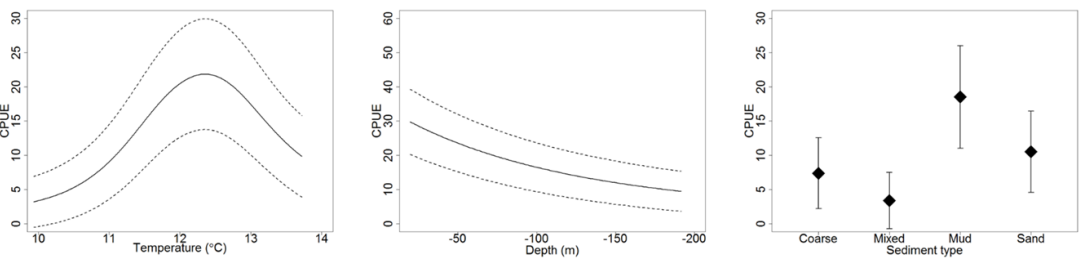
Research Interests

Population ecology and ecosystem health

- Develop methods to understand:
 - Demographic distributions of marine species
 - Spatial and temporal life stage connectivity
 - Relationship between connectivity and ecosystem health



Research Interests

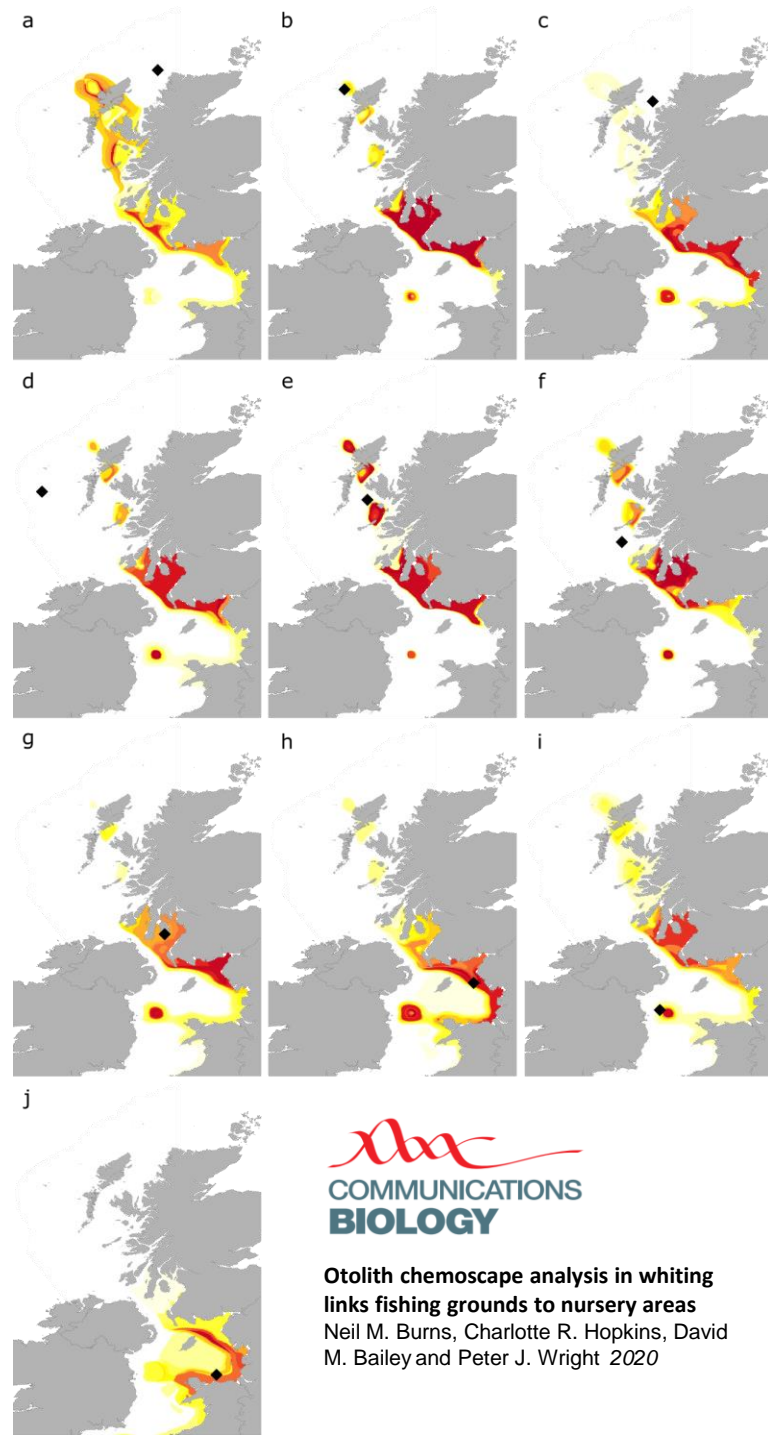


PLOS ONE

RESEARCH ARTICLE

A method to improve fishing selectivity through age targeted fishing using life stage distribution modelling

Neil M. Burns^{1,2*}, David M. Bailey¹, Peter J. Wright²



COMMUNICATIONS BIOLOGY

Otolith chemoscape analysis in whiting links fishing grounds to nursery areas
Neil M. Burns, Charlotte R. Hopkins, David M. Bailey and Peter J. Wright 2020



CRAN
[Mirrors](#)
[What's new?](#)
[Task Views](#)
[Search](#)

[About R](#)
~ ~ ~

The Comprehensive R Archive Network

Download and Install R

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

- [Download R for Linux](#)
- [Download R for \(Mac\) OS X](#)
- [Download R for Windows](#)

R is part of many Linux distributions, you should check with your Linux package management system in addition 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, not the source code. The sources have to be compiled before you can use them.



CRAN
[Mirrors](#)
[What's new?](#)
[Task Views](#)
[Search](#)

[About R](#)
[R Homepage](#)
[The R Journal](#)

R for Windows

Subdirectories:

[base](#)

Binaries for base distribution. This is what you want to **install R for the first time**.

[contrib](#)

Binaries of contributed CRAN packages (for R \geq 2.13.x; managed by Uwe Ligges). There is also information on [third party software](#) available for CRAN Windows services and corresponding environment and make variables.

[old contrib](#)

Binaries of contributed CRAN packages for outdated versions of R (for R $<$ 2.13.x; managed by Uwe Ligges).

[Rtools](#)

Tools to build R and R packages. This is what you want to build your own packages on Windows, or to build R itself.

Please do not submit binaries to CRAN. Package developers might want to contact Uwe Ligges directly in case of questions / suggestions related to Windows binaries.

You may also want to read the [R FAQ](#) and [R for Windows FAQ](#).

Note: CRAN does some checks on these binaries for viruses, but cannot give guarantees. Use the normal precautions with downloaded executables.



CRAN
[Mirrors](#)
[What's new?](#)
[Task Views](#)
[Search](#)

[About R](#)
[R Homepage](#)
[The R Journal](#)

[Software](#)
[R Sources](#)

R-4.0.3 for Windows (32/64 bit)

[Download R 4.0.3 for Windows](#) (85 megabytes, 32/64 bit)

[Installation and other instructions](#)

[New features in this version](#)

If you want to double-check that the package you have downloaded matches the package distributed by CRAN, you can compare the [md5sum](#) of the .exe to the [fingerprint](#) on the master server. You will need a version of md5sum for windows: both [graphical](#) and [command line versions](#) are available.

Frequently asked questions

- [Does R run under my version of Windows?](#)
- [How do I update packages in my previous version of R?](#)
- [Should I run 32-bit or 64-bit R?](#)

Please see the [R FAQ](#) for general information about R and the [R Windows FAQ](#) for Windows-specific information.

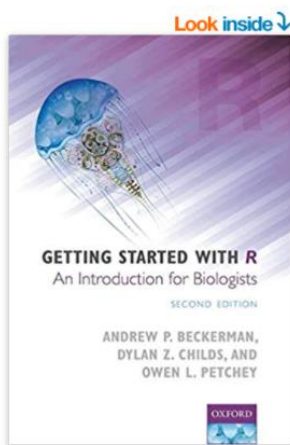
RStudio

Take control of your R code

RStudio is an integrated development environment (IDE) for R. It includes a console, syntax-highlighting editor that supports direct code execution, as well as tools for plotting, history, debugging and workspace management. [Click here to see more RStudio features.](#)

RStudio is available in **open source** and **commercial** editions and runs on the desktop (Windows, Mac, and Linux) or in a browser connected to RStudio Server or RStudio Server Pro (Debian/Ubuntu, Red Hat/CentOS, and SUSE Linux).

<https://rstudio.com/products/rstudio/>



Getting Started with R: An Introduction for Biologists Paperback – 26 Mar. 2017

by Andrew P. Beckerman ▾ (Author), Dylan Z. Childs ▾ (Contributor), Owen L. Petchey (Contributor)

★★★★☆ ▾ 48 ratings

> See all formats and editions

Kindle Edition
£14.62

Hardcover
£50.99 ✓prime

Paperback
£22.99 ✓prime

Read with Our **Free App**

5 New from £49.99

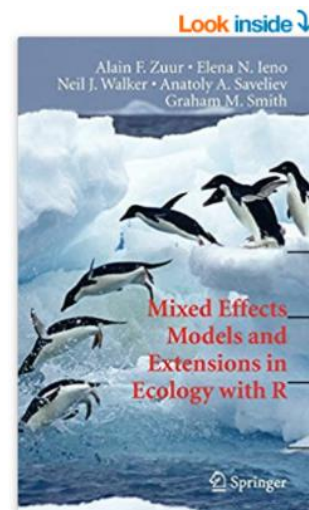
8 Used from £18.17
23 New from £18.79

FREE delivery: **Saturday, Dec 5**

Order within 10 hrs [Details](#)

Note: This item is eligible for **click and collect**. [Details](#)

R is rapidly becoming the standard software for statistical analyses, graphical presentation of data, and programming in the natural, physical, social, and engineering sciences. Getting Started with R is now the go-to introductory guide for biologists wanting to learn how to use R in their research. It teaches readers



Mixed Effects Models and Extensions in Ecology with R

Hardcover – 12 Mar. 2009

by Alain Zuur (Author), Elena N. Ieno (Author), Neil Walker (Author)

★★★★☆ ▾ 37 ratings

> See all formats and editions

Kindle Edition
£69.98

Hardcover
£73.66 ✓prime

Paperback
£109.99 ✓prime

Read with Our **Free App**

3 Used from £90.51
11 New from £69.19

2 Used from £123.17
11 New from £91.14

FREE delivery: **Thursday, Dec 3** [Details](#)

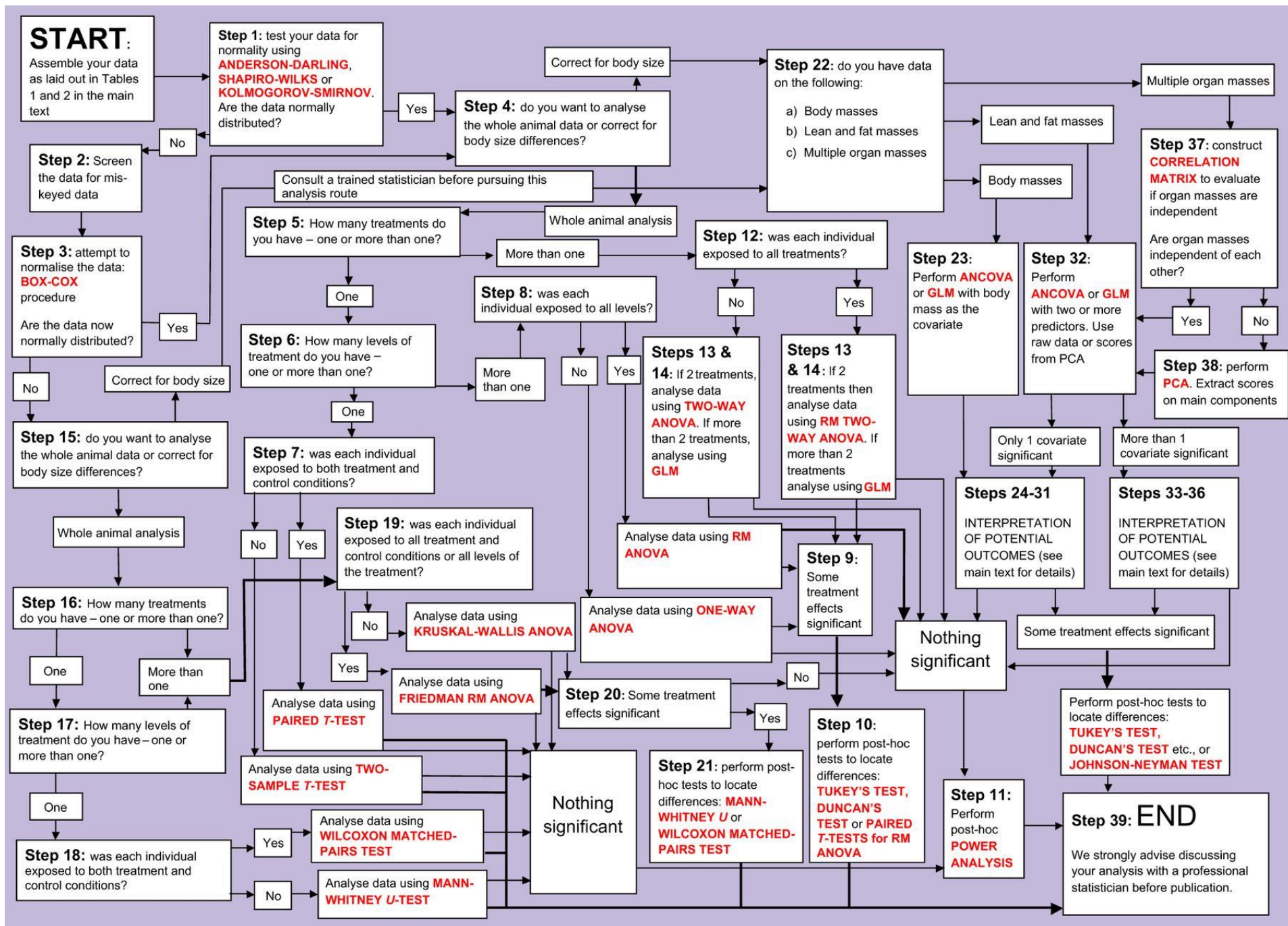
Note: This item is eligible for **click and collect**. [Details](#)

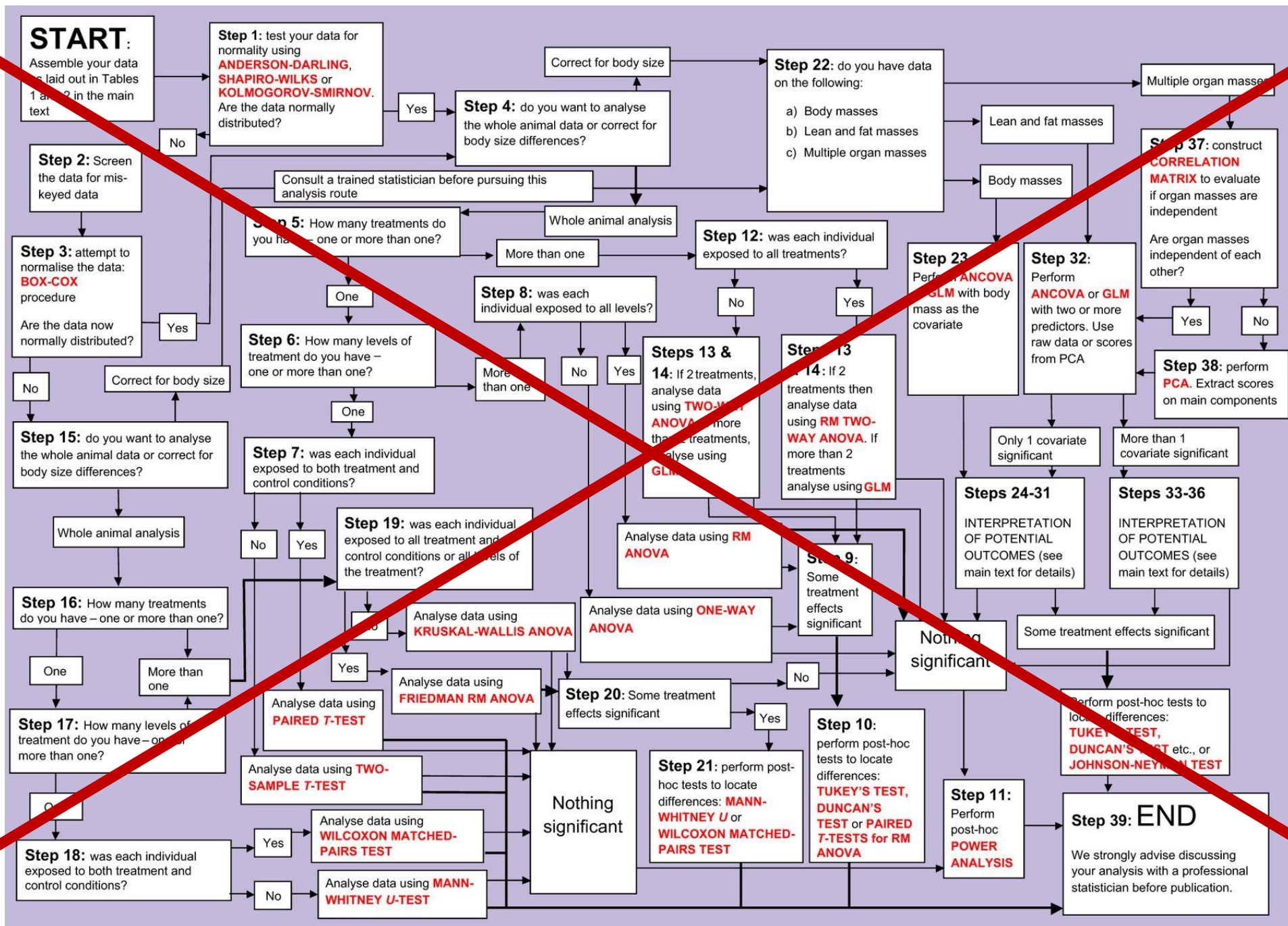
This book discusses advanced statistical methods that can be used to analyse environmental collected data are measured repeatedly over time, or space a GLMM or GAMM methods. The book starts by revising regression, additive

R script set up...oh yeah and making it look cool!

```
#####  
#####      1. Setting up your script      #####  
#####
```







They are just robots

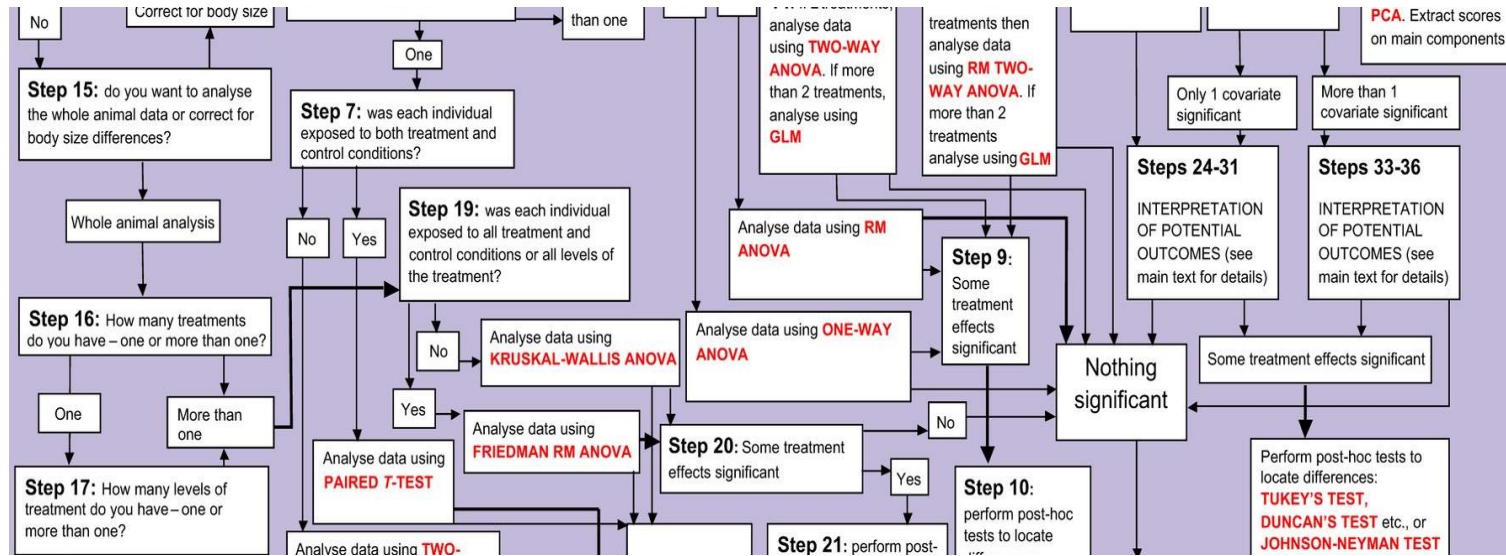


Add Thomas
Bayes



Stained glass window Latin square, honouring [Ronald Fisher](#)

Why use someone else's robot?



Linear models (lm)
Generalised Linear Models (glm)
Generalised Additive Models (gam)
Mixed effects versions of (lmm, glmm, gamm)

Statistical modelling
using glms

Why R?



But...we are getting ahead of ourselves.

	A	B	C	D
1	Site	Coral_colour	Perc_cover	S_abund
2	1A	Blue	20.5	23
3	1A	Green	9.5	2
4	1A	Red	70	19
5	1B	Blue	50	18
6	1B	Green	25.5	22
7	1B	Red	25.5	6

- Data into R from a .csv file
- Each column is a variable (explanatory (x) and response (y))
- Keep the column names short and simple but meaningful (to you)
- No spaces in column names_use_as_separator
- Be consistent with capitals

The very 1st thing you need to do = UNDERSTAND YOUR DATA!!!!

Our toy example



Abundance of sharks (y_1) ~

It's a whole number

They are big and there are not loads of them so probably less than 60



Coral percentage cover (x_1)

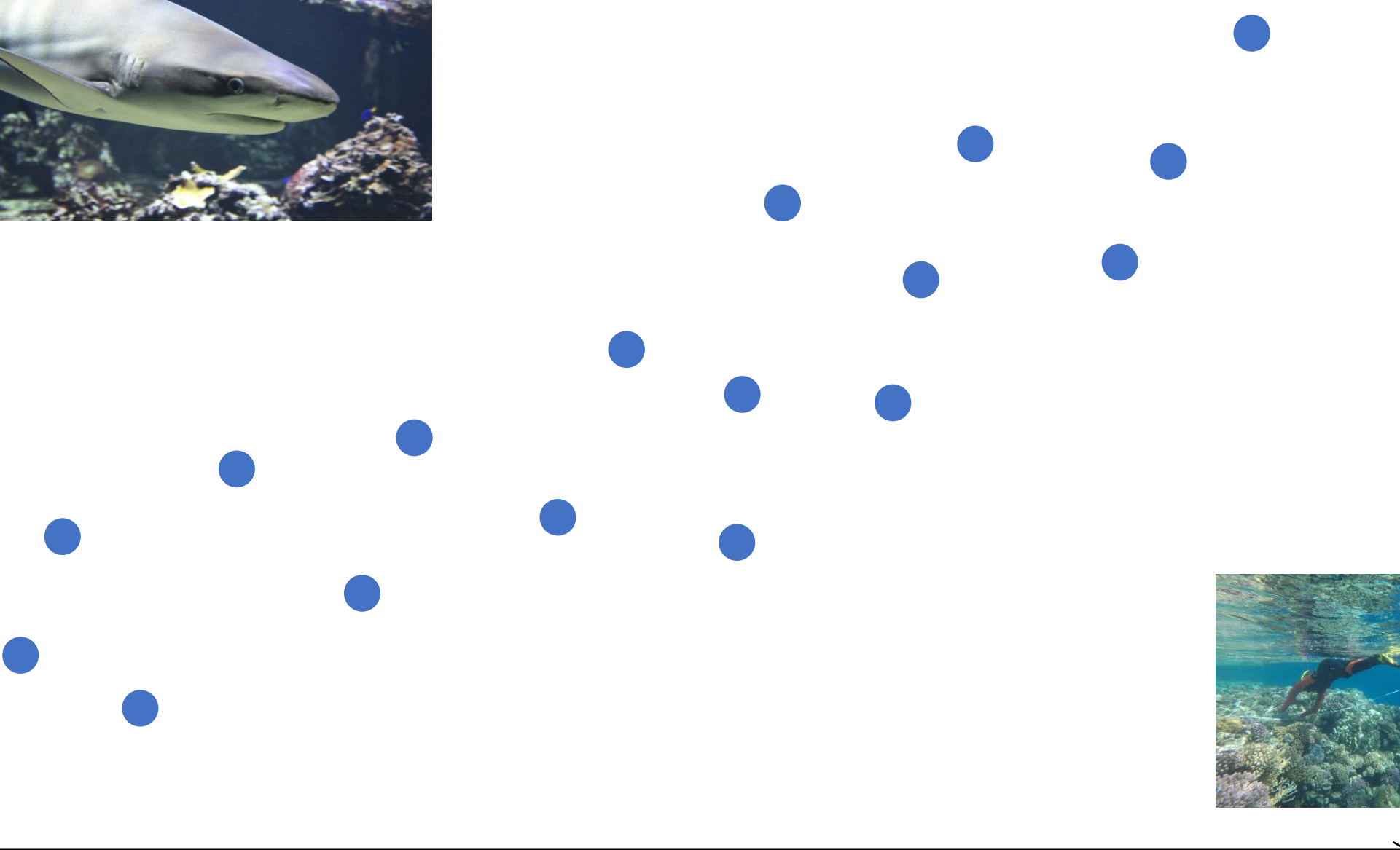
0 to 100%

Coral colour (x_2)

Blue, Brown & green



Shark
abundance



Coral percentage cover





Shark
abundance

Coral percentage cover



Back to R

```
#####  
##### 2. Understanding our data #####  
#####
```



Building useful robots and making inferences

Build the model from the knowledge of your data then use **information theory** to select models (and make inferences).

- Model selection
 - Picking the best of a bad bunch.
- Model validation
 - So how rubbish is it?



By using AIC we are selecting the most “cost” effective model. With cost being measured in degrees of freedom.

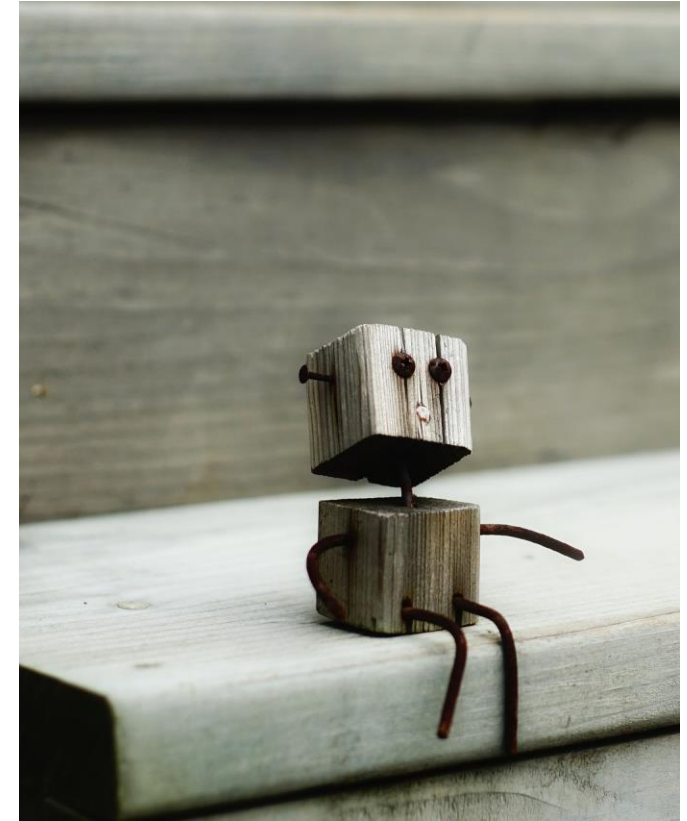
AIC balance between over and underfitting

Model selection

Using AIC (log-likelihood ratio tests are also useful for those who like a p-value)

Selection recipe

- Start with the most complex model and work “back” towards the most simple
- Use AIC to choose (3 rules)
 - Simple models are best
 - Small AIC is best
 - If these rules contradict (ie the more complex model has smaller AIC) then AIC should be different by more than 2



“I used backwards stepwise model selection to ...”

Weirdly it's not an exact science...

General aim of model validation is:

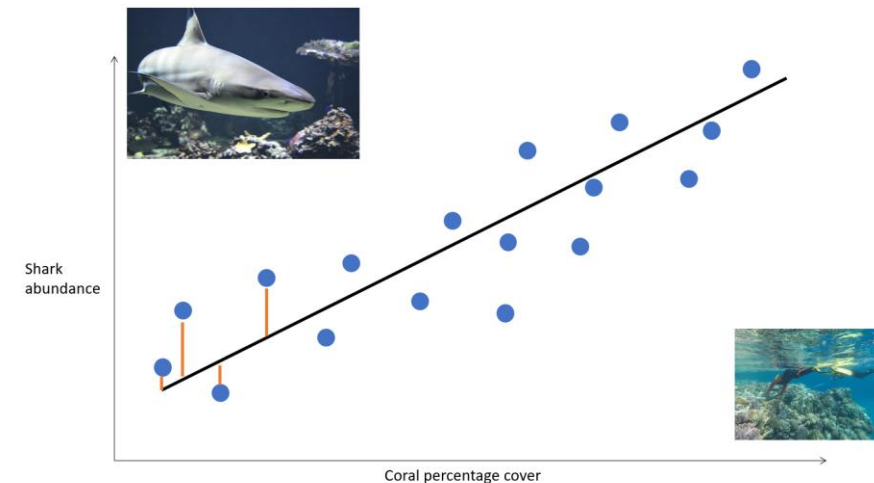
to have normally distributed residuals

and

no obvious patterns in the residuals



WTF is a residual?



Final interpretation

Simulation allows us to
interpret our complex models

A word about coefficients and factors >>>>>>>

$$Y = mX + C \quad (+ er)$$

↑ ↑
gradient intercept

$$Y = mX_1 + mX_2 + C \quad (+ er)$$

Our mod3

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.929914   0.417166   4.626 1.14e-05 ***
perc_cov     0.116220   0.007622  15.248 < 2e-16 ***
---
```

Our mod2

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  5.866607   0.513668  11.421 < 2e-16 ***
perc_cov     0.094528   0.008575  11.024 < 2e-16 ***
cor_colBrown -0.986790   0.143164  -6.893 5.75e-10 ***
cor_colGreen -0.580381   0.111171  -5.221 1.03e-06 ***
---
```

Our mod1

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  5.49784   0.60079   9.151 1.17e-14 ***
perc_cov     0.10071   0.01005  10.026 < 2e-16 ***
cor_colBrown -1.61468   1.00986  -1.599 0.11320
cor_colGreen  3.69075   1.38125   2.672 0.00889 **
perc_cov:cor_colBrown 0.01524   0.02008   0.759 0.44973
perc_cov:cor_colGreen -0.08092   0.02581  -3.135 0.00229 **
---
```


Add some hellish staring into picture here

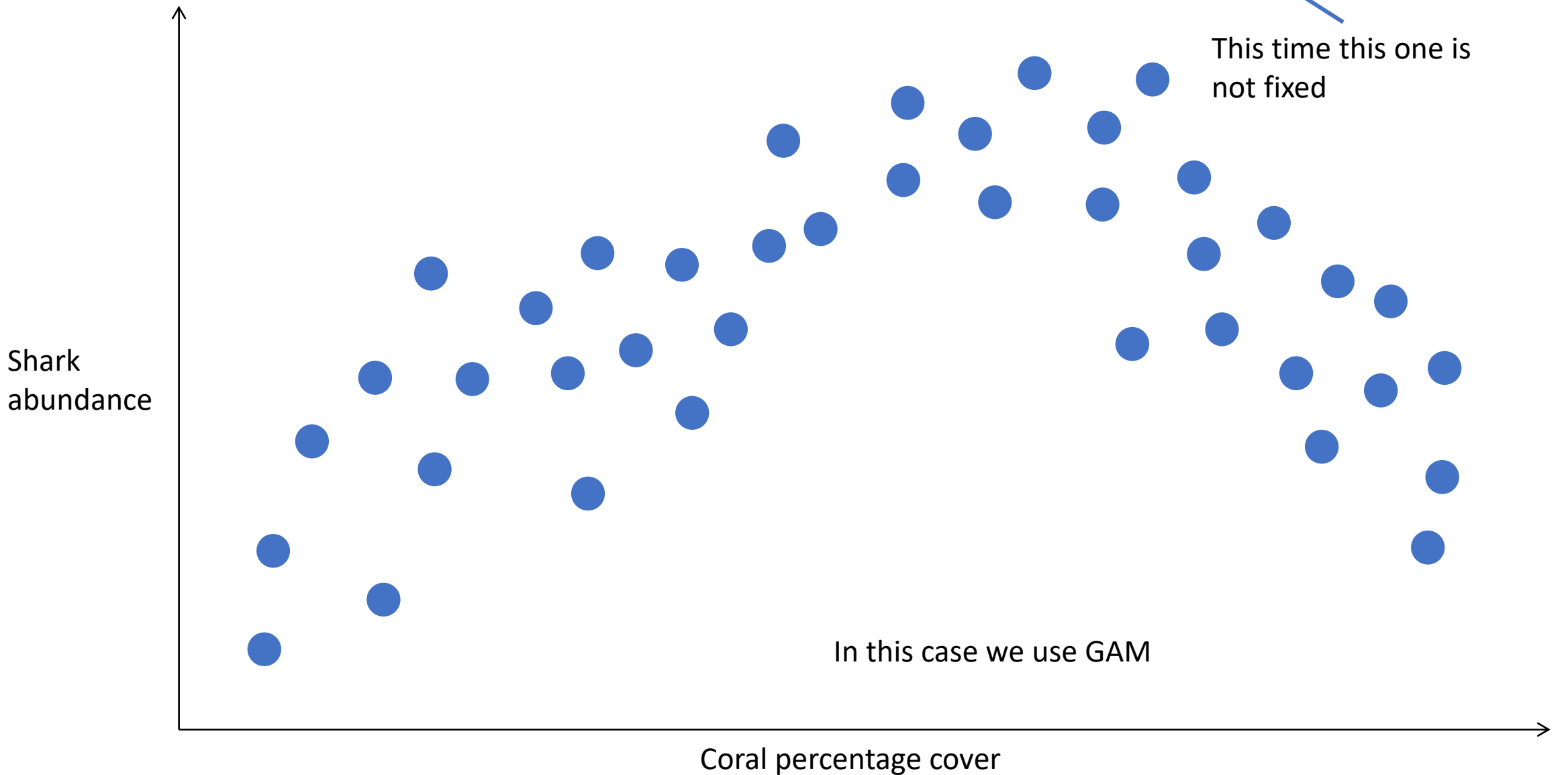
Back to R

```
#####  
##### 3. Models and inference #####  
#####
```

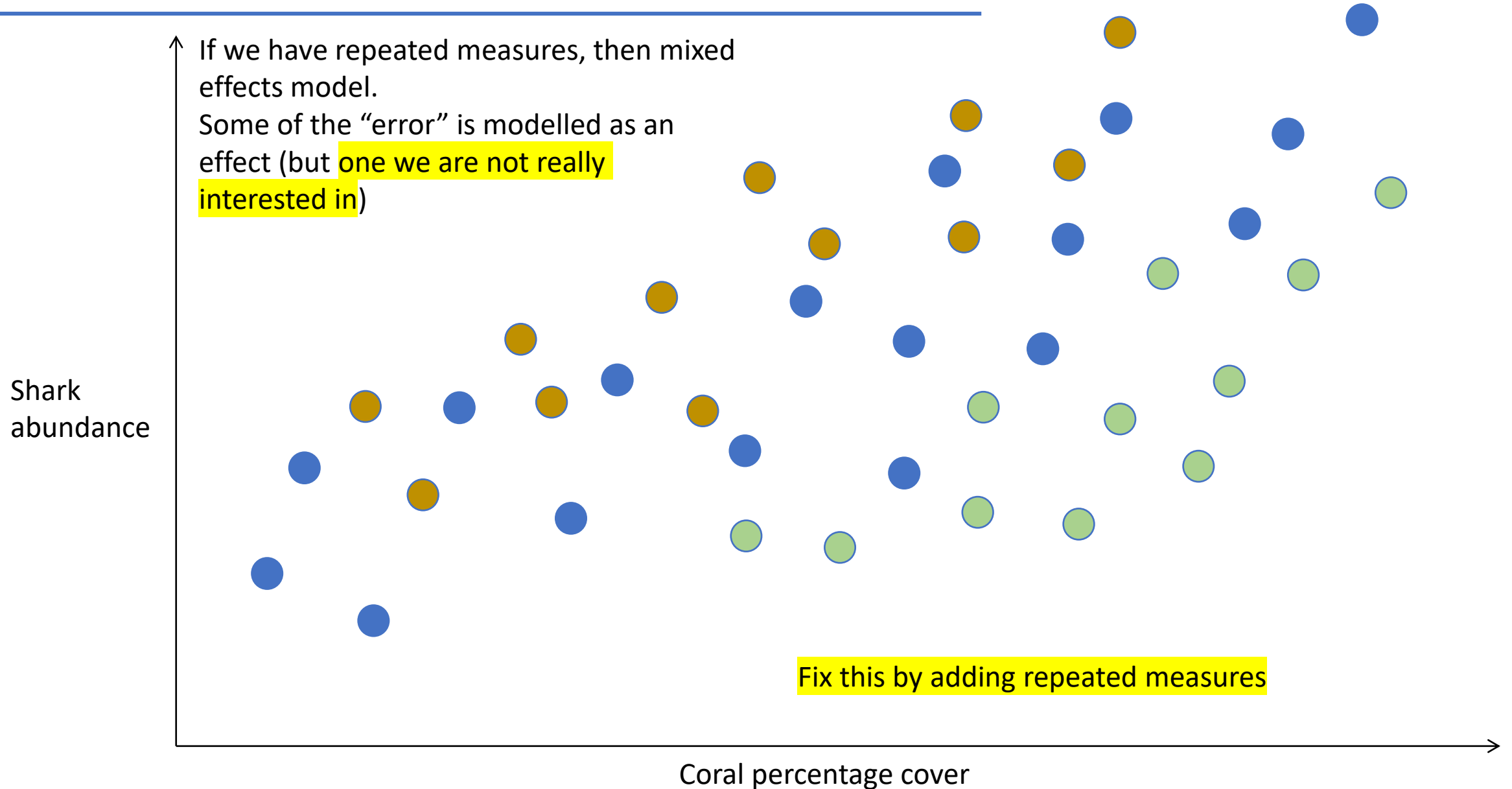


Ahh wait, but what if...?

$$Y = mX + C (+ er)$$



Ahh wait, but what if...?



Summary and questions

- Graph and understand your data
- Avoid spending too long down rabbit holes
- Use the tools we have discussed to build the robot
 - and the “recipe” to perform model selection
- We can use the [*PGR Stats in R team*](#) to discuss and resolve problems

