

# Course introduction

---

Michael Noonan

Biol 520C: Statistical modelling for biological data

1. Course Overview
2. Design- vs. Model-based Inference
3. What is modelling?
4. Components of a model

# Course Overview

---

Name: Michael Noonan

Office: SCI 379

Email: michael.noonan [at] ubc.ca (use subject heading BIOL520C in all email communication)

Office Hours: Mon 10-12h; Thurs 10h-11h, or by appointment arranged via email.

Course Website: <https://noonanm.github.io/Biol520C/index.html>

- Focus is on model-based inference (i.e., combining data with models to generate mechanistic descriptions of biological patterns).
- Building from simple linear regression, you'll learn regression methods for handling the most routinely encountered features in biological data (hierarchical data structures, non-Gaussian error distributions, non-linearity, autocorrelation, etc...).
- Emphasis on statistical best practices.
- How to use open source software (R) to apply these analyses.

- Basic statistics (concepts like means, medians, variances, probability distributions, regression should be familiar to you).
- Math you should be familiar with basic calculus (derivatives and integrals) and linear algebra (operations on vectors and matrices).
- Computer programming (we will be using R, but the course is not focused on 'how to code').
- Methods for handling *ad hoc*, corner cases.

- For each topic, there will be a core lecture and an associated practical assignment.
- Lectures will cover the core concepts of the course. Lecture slides will be posted on the course website the evening prior to the lecture. You are encouraged to take notes, and to ask questions in the lectures. All lectures will be recorded and made available to you.
- The practicals use structured tutorials to guide you on the use of the open-source software program R for applying the methods learned in the lectures to data. The lectures and practicals are designed to be *complementary* and not all the material in the practicals will be covered in the lectures and vice versa.

Practicals (10)	40%	Due on ~weekly basis
Participation from practicals	10%	Exact sched. in course outline
Hypothesis and expected outcome(s)	5%	Week 4
Paper	35%	Week 14
Presentation	10%	Weeks 13 & 14
Total	100%	



Beginning this week, you will be asked to complete practical assignments on an ~weekly basis. There will be a total of 10 practicals to be completed throughout the course.

The course web page on github will host the practicals, and the various datasets associated with each practical. Lectures will be given twice per week. After the second lecture, we will have covered all of the material that is needed to complete the week's practical assignment material which is due before the **start of the following Thursday** lecture (to be submitted via canvas).

**Grading:** Each practical assignment is worth a total of 5% of your total grade. Of this, 1% is given for submitting the tutorial on time, irrespective of whether or not the answers are correct (participation). The remaining 4% comes from the answers provided. Late practicals will be accepted, but will only be worth a maximum of 4%.

You will be required to apply the modelling tools covered in the lectures on a dataset of your choosing and write a short paper comprised of 6 sections:

- **Introduction:** Provide a brief description of the study system from which the data come and an outline of what questions you intend on exploring with the data. **(12.5%)**
- **Methods:** Describe how the data were collected, what variables are included, and what analyses were applied. **(20%)**
- **Results:** Length: Describe your statistical findings. Tables and figures should be used throughout. **(20%)**
- **Discussion:** Provide a brief summary of your findings and place them in a biological context. **(12.5%)**
- **References:** Include references to all necessary literature and statistical packages employed. **(5%)**
- **Appendix:** The appendix material should include an R markdown document that details every step of the analyses. **(30%)**

**Datasets:** To complete these assignments, you will have access to a number of pre-selected datasets. You can opt to use your own data to complete these assignments if you prefer, and are encouraged to do so, but you must seek instructor approval. If you intend on using your own data, it is recommended that you discuss this with me as early as possible.

**Late Assignments:** You are to submit your paper by the end of the day on Dec 11th. Late papers will have 10% deducted per day that they are overdue, and will receive a grade of zero if more than 10 days late without a valid excuse.

Prior to submitting your papers to the instructor, students will be required to give a 10-minute presentation to the class.

The presentation should include all of the sections that are included in the paper, however the appendix detailing the R code that was used should be integrated into the methods section of the presentation.

**Grading:** Grading will be based on the rubric available on Canvas.

At the **end of week four**, you will be expected to submit a one page proposal describing the study system you will be working on for your course project, the initial hypothesis, the expected outcome, and the data that you will be using to address this question.

**Grading:** All submitted proposals will receive a grade of 5/5. Late assignments will receive a grade of 0. In addition, no other assignment related to the core project will be accepted until the study proposal has been submitted.

There is no textbook for this course, but if you are interested in expanding your knowledge beyond what is covered, the following are recommended:

- Hilborn R, Mangel M. The ecological detective: confronting models with data. 1997. Princeton University Press. ~ \$90
- Zuur, A et al. (2009). Mixed effects models and extensions in ecology with R. Springer. ~ \$145
- Bolker, B. M. (2008). Ecological models and data in R. Princeton University Press. ~ \$90



Week	Lecture Topics
1	Course introduction; Regression refresher
2	Probability theory; Likelihood; Maximum likelihood
3	Mult. linear regression; Param inter.; Interpreting residuals
4	Mixed effects models; Model Selection; Information criterion
5	Model Selection; Model averaging; Heteroskedasticity;
6	Temporal autocorrelation; Spatial Autocorrelation
7	Phylogenetic inertia & Phylogenetically controlled regression
8	Logistic and Poisson regression
9	Zero-inflated data; Non-linear modelling; Deterministic functions
10	Mid-term break no lectures
11	Stochastic simulation and power analysis; Course Overview
12	Independent Project Work
13	Student presentations (10%)
14	Student presentations (10%) & Term paper due (35%)

# **Design- vs. Model-based Inference**

---

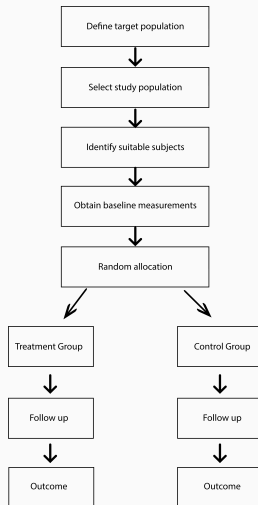


Science is a process of learning about nature. As scientists, we weigh competing ideas about how the world works (hypotheses) against observations (data). But our descriptions of the world are almost always incomplete, our observations have error, and important data is often missing... So, how do we accurately compare what we observe with what we hypothesize without bias?

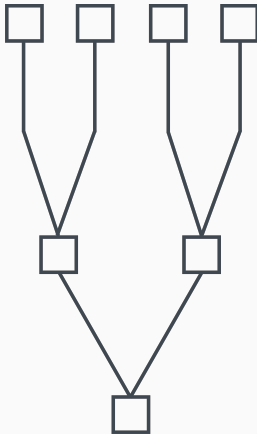
## Statistics

The process of making scientific inference can be split into two broad categories: design-based and model-based

- In design-based inference, most of the focus is on experimental design
- You assume your sample is a random sample of a target population
- Inference is focused towards a target population
- Goal is to be able to demonstrate that 'x causes y'
- Data are typically analysed by comparing means and variances across groups (e.g., ANOVAs, *t*-tests, etc...)



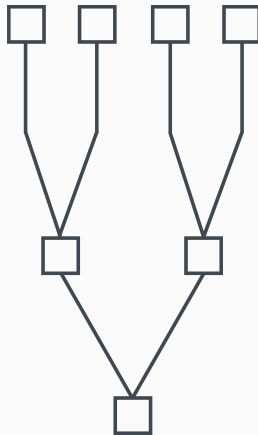
- Devise a hypotheses
- Devise an experiment with outcomes that will clearly accept or reject the hypothesis
- Carry out the experiment so as to get a clean result
- Recycle the procedure to refine the remaining possibilities (Platt, 1964)



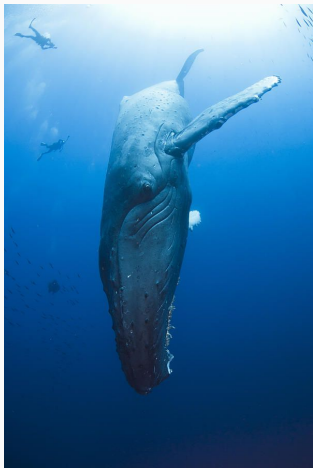
Core of design-based inference is confronting single hypotheses with data

Platt's decision tree is based on:

- i) Clear, distinct hypotheses
- ii) Unambiguous outcomes
- iii) A relationship between statistical significance and biological relevance



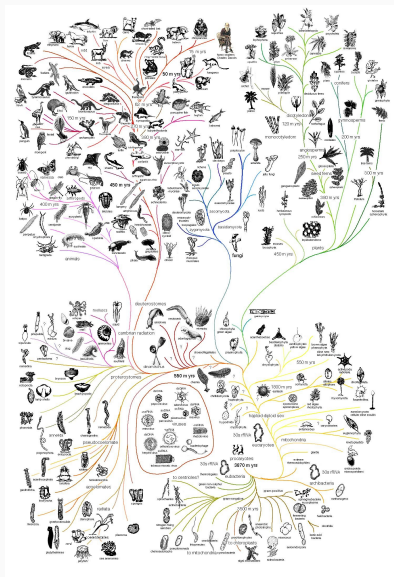
Many biological processes have long time-scales.



Humpback whales (*Megaptera novaeangliae*) can live for 50+ years. Source: David Valencia



Bristlecone pines (*Pinus longaeva*) live for thousands of years. Source: wired.com



Source: Chris King

Many biological processes have long time-scales.

Many biological systems have very poor reproducibility.





Source: Tom and Pat Leeson

Many biological processes have long time-scales.

Many biological systems have very poor reproducibility.

How can you design a controlled experiment in a wild population?

What if you're interested in species conservation?

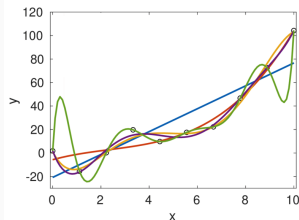
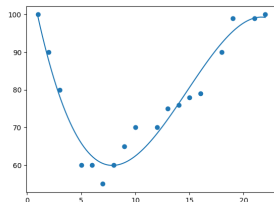


Source: Wikipedia

In 1988 the wild pop. of black footed ferrets (*Mustela nigripes*) was down to 18 ind.

What do you do if a power analysis says you need 20 animals?

- In model-based inference, most of the focus is on identifying an unknown deterministic model
- Inference can be extrapolated beyond the target population
- Goal is to provide a theoretical framework for why 'x causes y'
- In model-based inference you make distributional assumptions to make your response a random variable
- Data are typically analysed by fitting a model to data and interpreting the parameter estimates



# What is modelling?

---

**Hypothesis:** An idea, supposition, or otherwise unproven theory used as the basis for further investigation.

**Model:** A generalised description of some phenomenon.

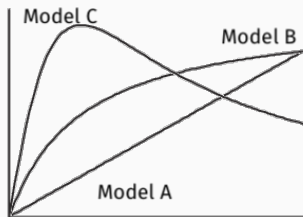
Model  $\neq$  Hypothesis

A single hypothesis can be represented by multiple models.

**Hypothesis:** Body mass  $M$  increases with age  $L$

**Models:**

- $M = aL$  Model A: Body mass is proportional to age
- $M = \frac{AL}{1+bL}$  Model B: Body mass saturates as age increases
- $M = aLe^{-bL}$  Model C: Body mass increases and then decreases as age increases



Source: Hillborn and Mangel 1997

The equation of a model is a very specific expression of the hypothesis. In other words, models help us clarify verbal descriptions of nature and mechanisms.

Models help us understand which parameters and processes are important, and which ones are not.

No model is completely correct.



The goal of modelling is not to provide a perfect description of the world, but to distill a process down to the most important components

Complicated models with lots of parameters usually provide better fits, but if the model is as complicated as nature itself why bother with modelling? Just go for a walk in the woods and be happy.

Don't fall in love with a model, the important thing is the system

Complex models provide more numerical precision, but simple models are more interpretable.

Simple models risk leaving out important parameters, complex models need lots of data for good parameter estimation

How complex should a model be?

Short answer: Let the data tell you.

Long answer: There are methods for this that we'll cover in later lectures.

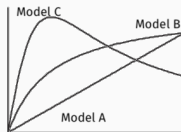
## Components of a model

---

Models are comprised of two main components:

**Deterministic part:** Describes the shape of the relationship (i.e., your hypothesis).

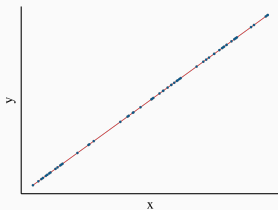
- Model A:  $M = aL$
- Model B:  $M = \frac{AL}{1+bL}$
- Model C:  $M = aLe^{-bL}$



**Stochastic part:** Describes the randomness of the process (i.e., captures the noise in a system).

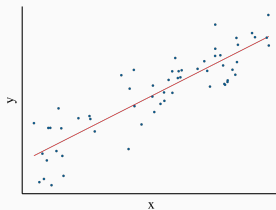
## Deterministic models

- No components are uncertain
- Outcome is always the same
- $y_i = \beta_0 + \beta_1 x_i$



## Stochastic models

- Some components are uncertain and characterised by probability distributions
- Outcome is variable
- $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$



## References

---

Platt, J.R. (1964). Strong inference. *Science*, 146, 347–353.

Hilborn R, Mangel M. The ecological detective: confronting models with data. 1997. Princeton University Press. Chapter 1