# Course introduction

Michael Noonan

August 11, 2020

# Table of contents

# Housekeeping

# Slide title

Stuff

- Stuff

- Focus is on model-based inference

- Focus is on model-based inference (i.e., combining data with models to generate mechanistic descriptions of biological patterns)

- 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

- 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...)

- 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

- 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 anaylses

# What this course is not about

- Basic statistics (concepts like means, medians, variances, probability distributions, regression should be familiar to you)

- Basic statistics (concepts like means, medians, variances, probability distributions, regression should be familiar to you)

- Math

- Basic statistics (concepts like means, medians, variances, probability distributions, regression should be familiar to you)

- Math... but you should be familiar with basic calculus (derivatives and integrals) and linear algebra (operations on vectors and matrices)

- Basic statistics (concepts like means, medians, variances, probability distributions, regression should be familiar to you)

- Math... but 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 high level coding)

- Basic statistics (concepts like means, medians, variances, probability distributions, regression should be familiar to you)

- Math... but 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 high level coding)

- Methods for handling *ad hoc*, corner cases

# Design- vs. Model-based Inference

# The Scientific Method

Science is a process of learning about nature. As scientists, we weigh competing ideas about how the world works (hypotheses) against observations (data).

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...

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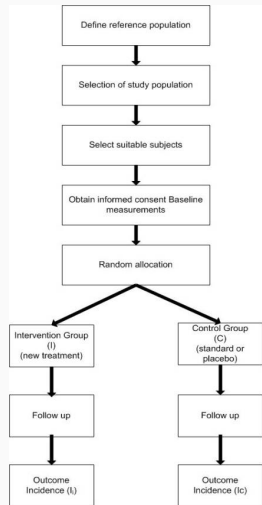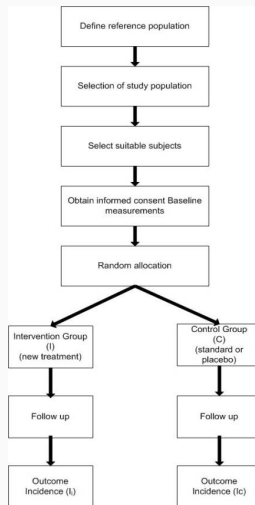
Science is a process of learning about nature. As scientists, we weigh competing ideas about how the world works (hypotheses) against observations (data). But our descriptions of the world are almost always incomplete, our observations have error, and important data is often missing... So, how do we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

## Statistics

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?
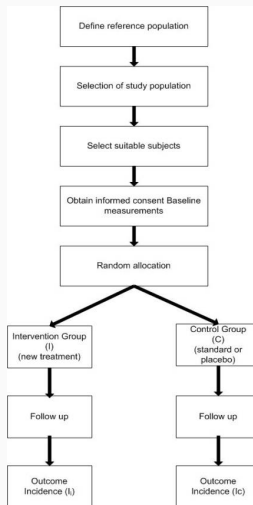
## Statistics

The process of making scientific inference can be split into two broad categories:

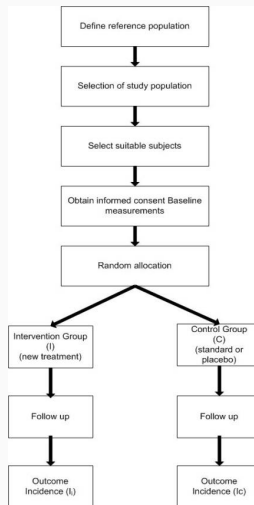Science is a process of learning about nature. As scientists, we weigh competing ideas about how the world works (hypotheses) against observations (data). But our descriptions of the world are almost always incomplete, our observations have error, and important data is often missing... So, how do we 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

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

# Design-based inference

- In design-based inference, most of the focus is on experimental design

# Design-based inference

- In design-based inference, most of the focus is on experimental design

- You assume your sample is a random sample of a target population
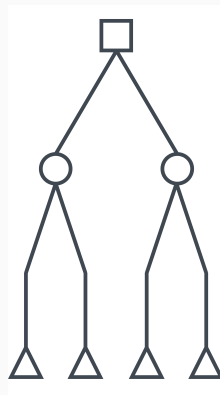
- In design-based inference, most of the focus is on experimental design

- You assume your sample is a random sample of a target population

- Inference is focused towards a target population

- 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'

- 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 a hypotheses

- Devise an experiment with outcomes
  that will clearly accept or reject the
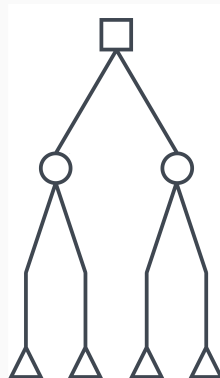  hypothesis

- 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

- 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

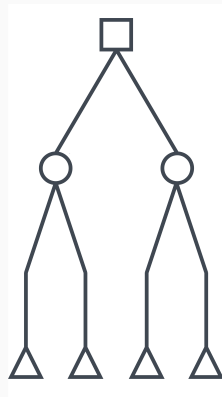- Recycle the procedure to refine the remaining possibilities (Platt, 1964)

- 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:
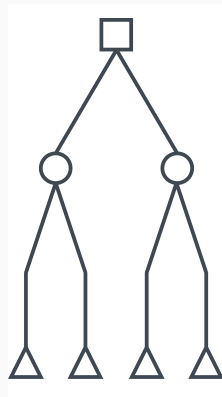
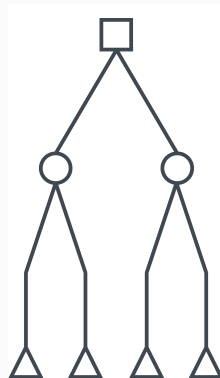Platt's decision tree is based on:

i) Clear, distinct hypotheses

Platt's decision tree is based on:

  i) Clear, distinct hypotheses

  ii) Unambiguous outcomes

Platt's decision tree is based on:

  i) Clear, distinct hypotheses

 ii) Unambiguous outcomes

iii) A relationship between statistical
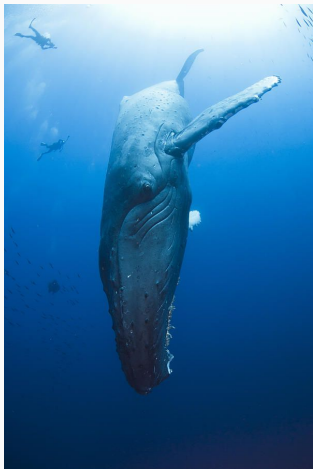      significance and biological relevance

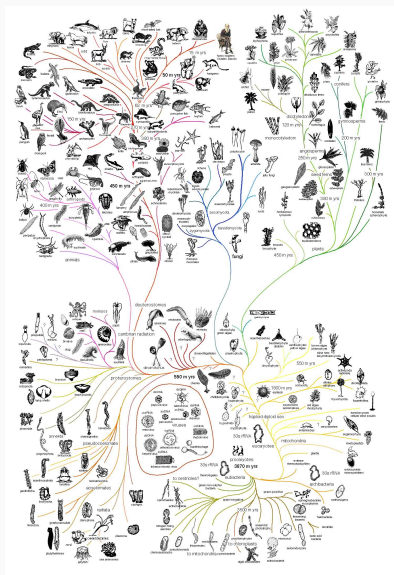Many biological processes have long time-scales.

Sperm whales (*Physeter macrocephalus*) can live for 70+ years

Sperm whales (*Physeter macrocephalus*) can live for 70+ years



Bristlecone pines (*Pinus longaeva*) live for thousands of years

Many biological processes have long time-scales.

Many biological processes have long time-scales.

Many biological systems have very poor reproducibility.

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?

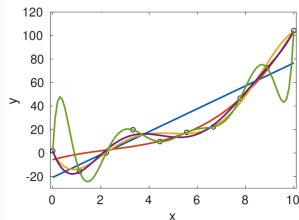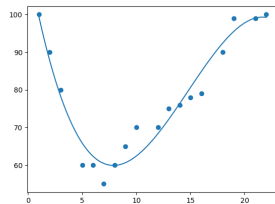What if you're interested in species conservation?



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

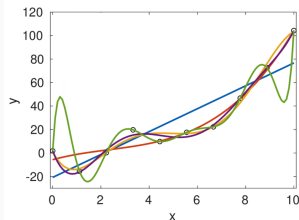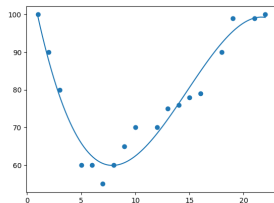What if you're interested in species conservation?



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

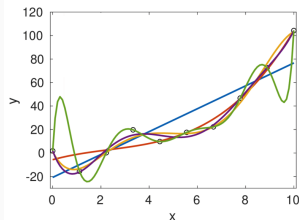What do you do if a power analysis says you need 20 animals?

- In model-based inference, most of the focus is on identifying an unknown deterministic model

- In model-based inference, most of the focus is on identifying an unknown deterministic model

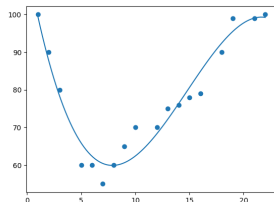- Inference can be extrapolated beyond the target population

- In model-based inference, most of the focus is on identifying an unknown deterministic model

- Inference can be extrapolated beyond the target population

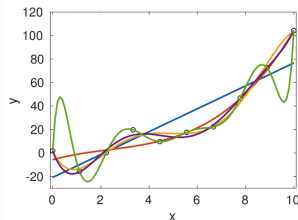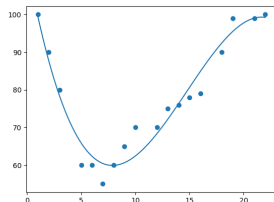- Goal is to provide a theoretical framework for why 'x causes y'

- In model-based inference, 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
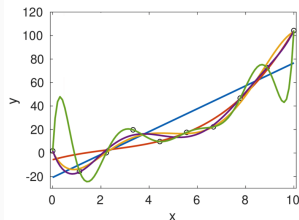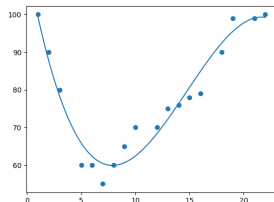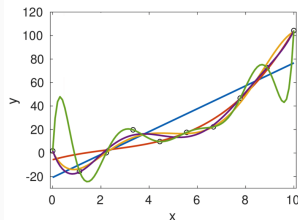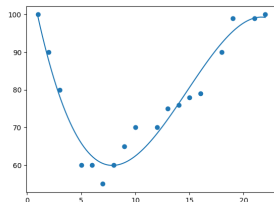
# Model-based inference



- 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:**

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

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

**Model:**

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

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

**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.

A single hypothesis can be represented by multiple models.

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

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

The equation of a model is a very specific expression of the hypothesis.

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.

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.

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

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

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

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.

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

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

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

Simple models risk leaving out important parameters

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

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?

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.

**Model complexity**

THE UNIVERSITY OF BRITISH COLUMBIA
Okanagan Campus

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.

Biol 520C: Statistical modelling for biological data 90

# Components of a model
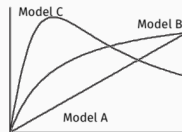
Models are comprised of two main components:

Models are comprised of two main components:

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

Models are comprised of two main components:

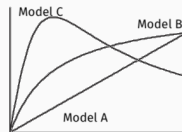**Deterministic part**: Describes the shape of the relationship (i.e., your hypothesis).

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

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).
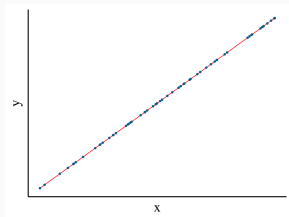
**Deterministic models**

**Deterministic models**

- No components are uncertain

**Deterministic models**

- No components are uncertain

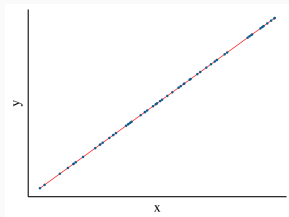- Outcome is always the same

**Deterministic models**

- No components are uncertain

- Outcome is always the same

- $y_i = \beta_0 + \beta_1 x_i$
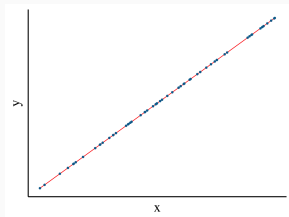
**Deterministic models**

- No components are uncertain

- Outcome is always the same

- $y_i = \beta_0 + \beta_1 x_i$



**Stochastic models**

### Deterministic models

- No components are uncertain

- Outcome is always the same
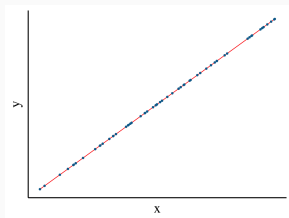
- $y_i = \beta_0 + \beta_1 x_i$



### Stochastic models

- Some components are uncertain and characterised by probability distributions

## Deterministic models

- No components are uncertain

- Outcome is always the same
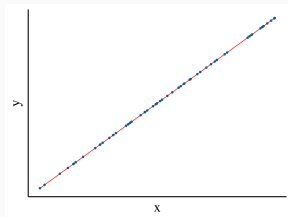
- $y_i = \beta_0 + \beta_1 x_i$



## Stochastic models

- Some components are uncertain and characterised by probability distributions

- Outcome is variable

**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
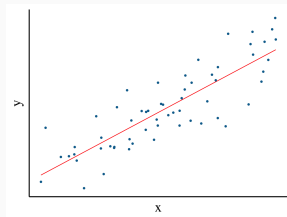
- Outcome is variable

- $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$

# References

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