

Introduction to Epidemiological and Biostatistical Thinking

UW Neurology Fellowship

Marlena Bannick

6/18/2020

PhD Student, University of Washington Dept. of Biostatistics
Researcher, Institute for Health Metrics and Evaluation

Goal

Introduce you to epidemiological thinking and key (bio)statistical concepts that you can use to critically interpret scientific studies in health and medicine.

Learning Objectives I.

1. **Basics.** Identify key elements of an epidemiological study and how they relate to the scientific question
2. **Study Design.** Recognize the basic types of epidemiological study design and identify when each design is appropriate for the scientific question
3. **Bias.** Recognize sources of bias in study designs or measurements and understand how they might affect your ability to answer the scientific question

Learning Objectives II.

4. **Modeling.** Understand how you can formulate your understanding about a data generating process, assumptions, and a hypothesis to test in a statistical model
5. **Inference.** Recognize the distinction between an effect size, a confidence interval, and a p-value as they relate to parameters that are estimated in a statistical model

A epidemiological study should be generated by a *scientific question of interest*. Broadly, you can think of these scientific questions falling into two main categories:

- **Descriptive:** What is the incidence rate of ischemic stroke (IS) in women aged 45 - 60 years old?
- **Inferential:** What is the effect of an experimental treatment on mortality following ischemic stroke in women aged 45 - 60?

From a statistical point of view it is not a clean distinction because you still use statistical tools to *infer* the incidence rate for a descriptive study.

The questions *who, what, where, when* have never been more important than in the context of epidemiology!

Having a well-defined scientific question means having clear answers for the following components:

- **Exposure:** What is the group in study exposed to that you want to measure the effect of, and over what period of time?
- **Population:** Who is the group being studied?
- **Outcome:** What outcome is being studied (either in relation to the exposure or on its own) and over what period of time?

The *why* is also important! Epidemiological studies should serve some purpose.

Once you've defined your target exposure, outcome, and population that makes up your scientific questions, understanding **measurement** of the outcomes is of utmost importance.

Some common outcome measurements in the context of health sciences are

- **prevalence**: proportion of a population with an outcome
- **incidence**: rate of getting the outcome among individuals in a population that did not already have the outcome (“risk”)
- **remission**: rate of returning to be outcome-free among those that had the outcome

Think about denominators!

What are the exposure, outcome, and population for each of these scientific questions?

- **Descriptive:** What is the incidence rate of ischemic stroke (IS) in women age 45 - 60 years old?
- **Inferential:** What is the effect of an experimental treatment on mortality following ischemic stroke in women age 45 - 60?

Table 1: Basic Elements of Study Design

	Descriptive	Inferential
Exposure		
Outcome		
Population		

What are the exposure, outcome, and population for each of these scientific questions?

- **Descriptive:** What is the incidence rate of ischemic stroke (IS) in women age 45 - 60 years old?
- **Inferential:** What is the effect of an experimental treatment on mortality following ischemic stroke in women age 45 - 60?

	Descriptive	Inferential
Exposure		experimental treatment
Outcome	ischemic stroke (IS)	death from IS
Population	women age 45-60 without IS	women age 45-60 with IS

How would you make these questions more precise?

With a binary exposure and a binary outcome, the results of a study will look something like this 2x2 table:

Table 2: Example 2x2 Table

	Outcome	No Outcome
Exposed	a	c
Unexposed	b	d

But there are *so many ways* to obtain that 2x2 table, so it is imperative to understand the study design behind the data!

Understanding study design will make it clear **what are the valid analyses** that can be performed on the data in that table.

Study Design

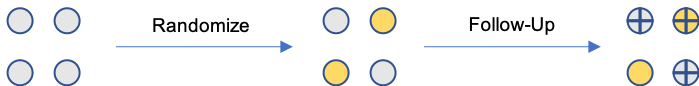
Starting with what is typically considered the studies that will provide the “strongest” evidence of a *causal* relationship between an exposure and an outcome:

- **Randomized controlled trials:** participants are *randomly* assigned to an exposure treatment or a control and followed up over time to record outcomes
- **Cohort studies:** participants are selected based on their exposure status and followed up to record outcomes
- **Case control studies:** participants are selected based on their outcome status and we inquire about exposure in the past
- **Cross-sectional studies:** measure exposure and outcome of participants at the same point in time (no temporal element)
- **Case reports:** report on the outcome status of one or a handful of interesting cases

Study Design

● Exposed ● Unexposed + Has outcome

Randomized Controlled Trial



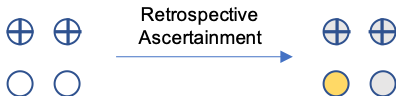
Cohort Study



Cross-Sectional



Case Control Study



Case Report



Biases in the epidemiological context are any factors in your study that *prevent* you from being able to answer your precise scientific question.

Biases may result from systematically incorrect measurements of the outcome, the exposure, or the population.

Some study designs may avoid certain types of bias, but it is crucial to always be on the lookout for sneaky biases when designing, analyzing, or reading a study.

Examples of biases include:

- **Selection bias**: the population that you want to study is not the population that is actually in your study
- **Confounding bias**: the relationship between exposure and outcome among those in your study is *confounded* by other variables (more later)
- **Recall bias**: individuals are being asked about exposures or outcomes that they do not remember correctly
- **Social desirability bias**: individuals are not comfortable disclosing their true exposure or outcome status for fear of judgement by others

This is by no means an exhaustive list. See [a catalogue of bias](#) for a taxonomy and more examples.

Biases & Study Design

- Randomized controlled trials are designed to *eliminate bias*: statistically speaking, we do not expect there to be significant differences in the characteristics of the treatment groups
- Observational study designs like cohort studies and case control studies *observe* what's already happening – what if those that are exposed also have characteristics that make it more likely that they will have the outcome (**confounding**: more later)?
- Studies that rely on participants to self-asertain, or to recall things from the past (e.g. case control studies) may result in systematic measurement error of exposure or outcome

Selection Bias

Recall our example inferential question: **What is the effect of an experimental treatment on mortality following ischemic stroke in women age 45 - 60?**

Consider the following sampling strategies:

- Sample women aged 45 - 60 who have been discharged from the hospital following ischemic stroke, randomly assign some to experimental treatment.
- Sample women aged 45 - 60 who have been admitted to the hospital for ischemic stroke, randomly assign some to experimental treatment.

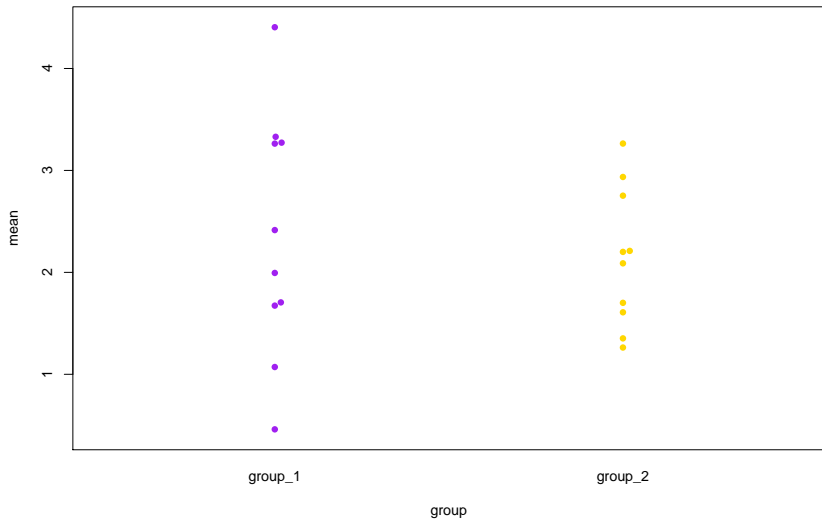
Which may suffer from selection bias?

Again, recall our example inferential question: **What is the effect of an experimental treatment on mortality following ischemic stroke in women age 45 - 60?**

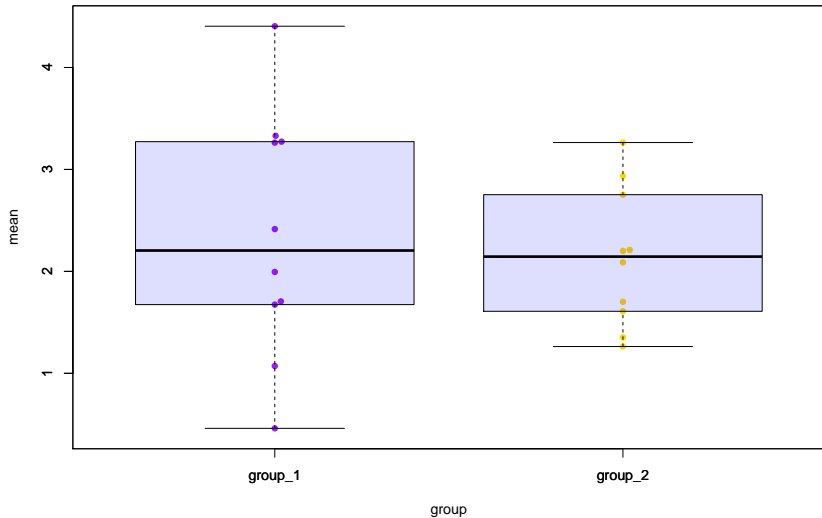
What if we do not assign the experimental treatment, but the physician decides whether or not to administer treatment to the patient?

For a gentle introduction to similar problems and techniques that control for such problems in observational study designs, check out [The Book of Why](#) by Judea Pearl.

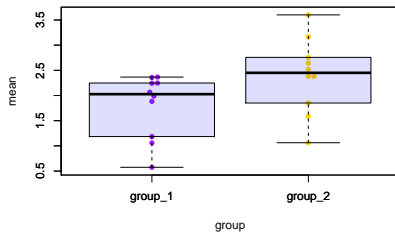
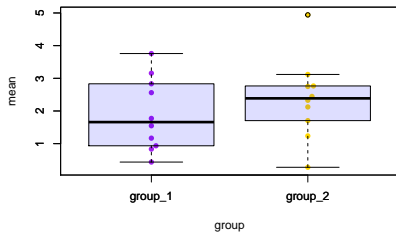
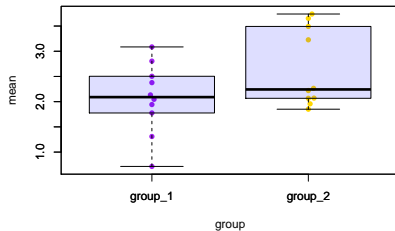
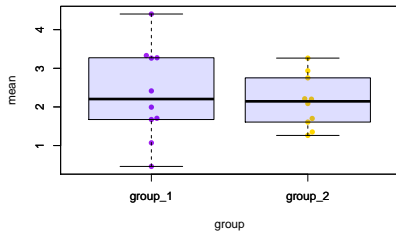
Inference: Simple Means



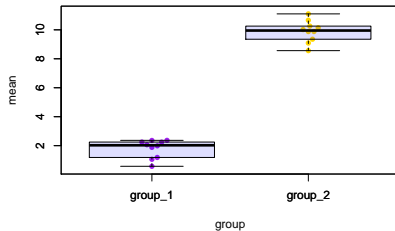
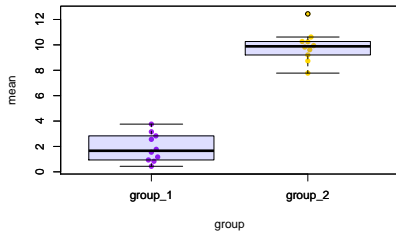
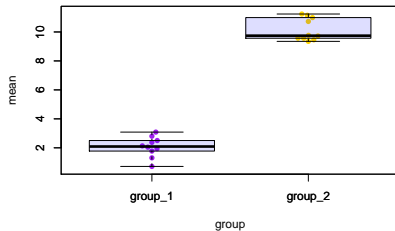
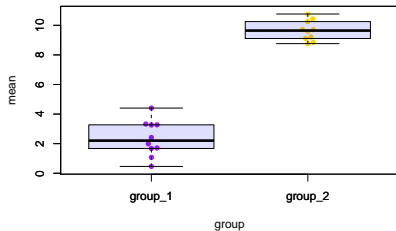
Inference: Simple Means



Inference: Small Effect Size, Small Sample Size



Inference: Large Effect Size, Small Sample Size



Inference: Small Effect Size, Large Sample Size

