

Validity

Kinds of statistics

- Descriptive (about the data)
- Inferential (about the world)

Neither is more important than the other!!

Kinds of statistics

- Exploratory (I don't have a hypothesis or theory. I don't know what's going to happen)
- Confirmatory (If theory X is true, then the data I collect should look like Y)

Neither is more important than the other!!

Exploratory ----- Confirmatory

Kinds of research

- Experimental (we introduce an intervention and look at the effects; researcher introduced assignment)
- Observational (we measure/survey our participants without trying to affect them; no researcher-introduced assignment)

Typically we pair some kinds of statistical tests with experimental work and other kinds of tests with observational work.

In reality, most statistical tests can be used with most kinds of research. It's not so much the kind of research that matters, but **which statistic helps to answer your question** and **what types of variables do you have?**

- We'll discuss the first point throughout the course
- Let's discuss variables now

Validity

Why statistics

- An essential aid to “signal detection”
- A universal language for communicating what we find.
- Required for competent evaluation of others’ work.

Our Goals

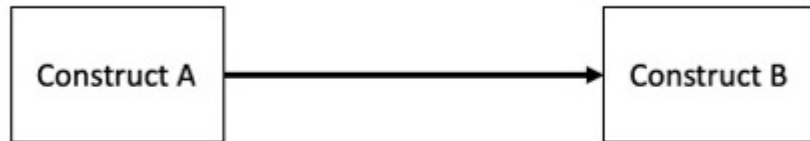
Advanced skill in quantitative methods carries with it the responsibility to use those skills carefully and ethically.

Today, we'll discuss methodological issues present in statistics.

- It can be tempting to use statistics to fix poor research design.
- These issues cannot be fixed quantitatively (even when it looks like they can).
- After today, we focus on what happens after you collect data. But it is still your job to study research design, data collection, and theoretical logic.

Constructs

- Our basic goal in science is to make inferences about the causal relations between constructs.



- We can't do that directly, so we rely on proxies for those constructs
 - | We are measuring the invisible

Measuring the Invisible

We can't do that directly, so we rely on proxies for those constructs.



In order to infer that $A \rightarrow B$, we have to make three assumptions:

- X is a good proxy for A
- Y is a good proxy for B
- X and Y are causally related

Measuring the Invisible

- When the first two assumptions are true, the relation between X and Y will provide a good estimate of the relation between A and B.
- What threatens our ability to carry out this seemingly simple task?
 - How do quantitative methods help us solve these problems?

Validity

Four kinds of validity in research threaten our ability to make valid causal inferences. Solving each problem either directly requires quantitative methods or makes use of principles that are central to quantitative methods.

- Internal validity
- Construct validity
- External validity

Internal validity

- **Definition:** the validity of the inference that X and Y are causally related.
 - Given that X and Y are correlated, can we validly infer that the relation is causal?

Threats to internal validity

- ambiguous temporal precedence
- selection
- attrition
- history
- maturation
- regression
- testing
- instrumentation

Threats to internal validity

- **ambiguous temporal precedence**
- selection
- attrition
- history
- maturation
- regression
- testing
- instrumentation

Temporal precedence can be established in an experiment because treatment precedes outcome.

But, when treatment is not possible, then logic and common sense can sometimes dictate temporal precedence.

- prenatal nutrition and cognitive development
- depression and cancer

Threats to internal validity

- ambiguous temporal precedence
- **selection**
- attrition
- history
- maturation
- regression
- testing
- instrumentation

Any systematic differences between groups that might account for an observed effect.

- Test scores of students who visit the Psychology tutoring center vs students who do not visit tutoring center.

How to combat this?

Threats to internal validity

- ambiguous temporal precedence
- selection
- **attrition**
- history
- maturation
- regression
- testing
- instrumentation

Even if random assignment is used, participants may drop out of the study, producing unequal groups, a situation that has the same inferential problems as selection.

Threats to internal validity

- ambiguous temporal precedence
- selection
- attrition
- **history**
- maturation
- regression
- testing
- instrumentation

History refers to any event that occurs between the beginning of treatment and the measurement of outcome that might have produced the observed effect.

- A marketing campaign intended to increase beer sales happens to coincide with other events that might have the same effect: a particularly hot period of weather, a long losing streak by the St. Louis Cardinals, etc.

Threats to internal validity

- ambiguous temporal precedence
- selection
- attrition
- history
- **maturation**
- regression
- testing
- instrumentation

Maturation refers to changes in the organism that occur regardless of treatment and that may masquerade as a treatment effect.

- A school-wide educational intervention is predicted to increase achievement test scores. The entire school must get the same curriculum, so a control group in the school is not possible.

Threats to internal validity

- ambiguous temporal precedence
- selection
- attrition
- history
- maturation
- **regression**
- testing
- instrumentation

Regression (to the mean) occurs when participants are selected because of their extreme scores and those scores are unreliable. The scores will regress toward the mean at the second assessment

- *Sports Illustrated* cover jinx
- Tall men father not-so-tall sons (Galton)

Threats to internal validity

- ambiguous temporal precedence
- selection
- attrition
- history
- maturation
- regression
- **testing**
- instrumentation

Testing refers to the possible change that may occur just because participants have been previously measured. These are often called practice or fatigue effects.

- Students do better on the first half of test compared to the second.
- Students do better in the second half of the term compared to the first.

Threats to internal validity

- ambiguous temporal precedence
- selection
- attrition
- history
- maturation
- regression
- testing
- **instrumentation**

Change may occur because the measurement changes over time, perhaps becoming more or less reliable.

Instrumentation reflects changes in the measurement; testing reflects changes in the object of measurement.

"When a measure becomes a target, it ceases to be a good measure."
(Goodhart)

The key point with internal validity is that something else besides the treatment is a plausible alternative explanation for any apparent treatment effect.

Solving threats to internal validity is a research design problem, not a statistics problem. Nonetheless, quantitative methods play a key role in making the case for internal validity.

Removing the influence of other variables

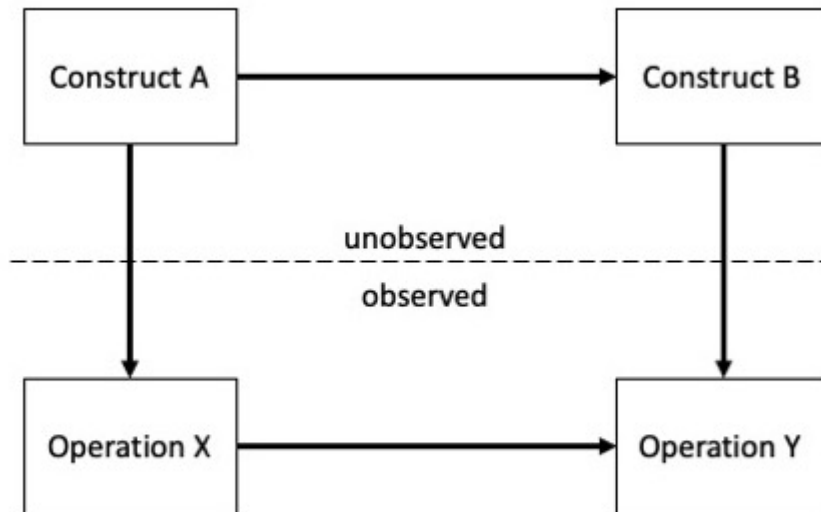
If the "other variables" can be measured, their influence can be statistically controlled so that the hypothesized relation can be detected more accurately.

However:

Statistical control should best be thought of as a method of last resort, to be used when design controls are not available or have failed.

Construct validity

- The validity of the inference that a given operationalization of a construct does a good job representing the construct.



- Construct validity refers to the correctness of the label that is applied to the operation. It depends on first demonstrating adequate reliability and then is bolstered by demonstrating relations of the target operation to other operations.

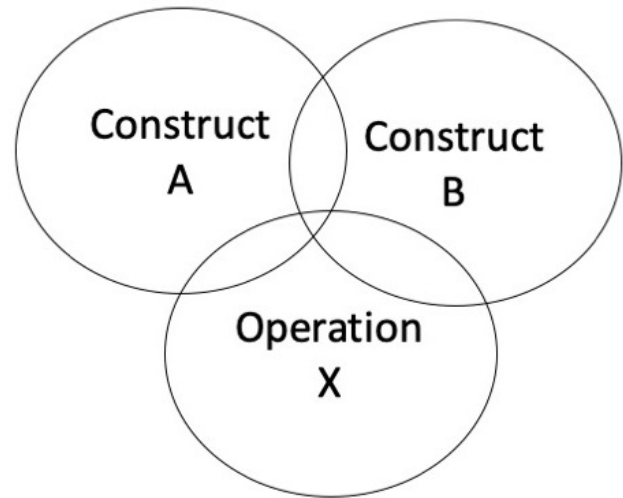
Threats to construct validity

- **inadequate explication of constructs**
- **construct confounding**
- confounding constructs with levels of constructs
- reactive self-report changes
- reactivity to the experimental situation
- experimenter expectancy
- novelty and disruption effects

Construct confounding

Operations usually tap more than one construct. Failure to recognize the full set of constructs embedded in the operation can lead to incorrect inferences about the constructs that are active.

A self-report of optimism might also reflect self-esteem or positive affect.



External validity

- **Definition:** The validity of the inference that a causal relation between operations *generalizes* to other units, treatments, observations, or settings.



Threats to external validity

- sampling bias
- experimenter effects
- Hawthorne effect
- testing effects
- situation effects

Threats to external validity

- **sampling bias**
- experimenter effects
- Hawthorne effect
- testing effects
- situation effects

We can't study an entire population, so instead we take samples. If your sample is not representative of the population however, then how on Earth can you generalize to that population?

Threats to external validity

- sampling bias
- **experimenter effects**
- Hawthorne effect
- testing effects
- situation effects

What if I told all my research participants that my job depended on the outcome of the study they are in? The participants might change their behaviors. And their behavior is what we are studying. So anything I find might not be generalizable (even if unintentional). Need to remain neutral!

Threats to external validity

- sampling bias
- experimenter effects
- **Hawthorne effect**
- testing effects
- situation effects

The Hawthorne Effect says that the fact that people know they are being observed might be enough to change their behavior.

Ex: if someone is in a study about stress, and they know they are in the study (informed consent), maybe they make themselves seem more stressed than they actually are. They might fill out surveys with scores reflecting higher levels of stress than what they might fill out if they didn't think they were being watched.

Threats to external validity

- sampling bias
- experimenter effects
- Hawthorne effect
- **testing effects**
- situation effects

Testing effects are especially critical in pre/post designs. At the pre-test, they are nervous. But at the post-test, they know what to expect and are less anxious etc. Really problematic when studying something like, say, anxiety.

Threats to external validity

- sampling bias
- experimenter effects
- Hawthorne effect
- testing effects
- **situation effects**

Situation effects can be things like time of day, the setting of the experiment/location etc.

What if you study recall memory. You have all your participants come in before 10am. You find an effect.

Now you repeat the study, but you have all your participants come in after 8pm. You don't find an effect.

Your Ethical Duty as a Scientist

Advanced skill in quantitative methods carries with it the responsibility to use those skills **carefully and ethically**.

- Know the shortcomings of your study
- REPORT the shortcomings of your study
- Let your readers understand the limitations
- Do NOT overstate your findings (or let the press overstate them) Confusing Qualitative and Quantitative
 - happens more than you think; esp in the machine learning world
 - can straight up get the wrong numbers (will come back to this with correlation!)