# PoliSci 4782    Political Analysis II

## Causality and Research Design

Seamus Wagner

The Ohio State University

# Two Orientations of Statistical Analysis

- **Causal inference**:
  - The goal is to identify a causal connection between two variables.
  - The unbiasedness of $\hat{\beta}$ is prioritized.
  - More relevant to scientific discovery.

- **Predictive inference**:
  - The goal is to forecast future observations of outcome variables based on the past observations.
  - The out-of-sample prediction accuracy is prioritized.
  - More relevant to industries and decision-making agencies.

# Causality

- Causality is not something can be directly computed by statistics.

- We only have correlation (e.g. correlation coefficient in 4781) and conditional correlation (e.g. regression coefficient in 4781 and 4782) for understanding relationships between two variables.

- At best, causality that we discuss in statistical analysis is *an assumed relationship, plus some evidence*.

# Understanding Causality by Counterfactual

Some math for understanding causality:

- Suppose that the potential outcomes of variable $Y$ [$E(Y_0)$ or $E(Y_1)$] depends on the existence of the cause $D$ (or the assignment of the "treatment"). $D$ takes 1 when the cause exists and 0 otherwise.

- We can use the following equation to present the causality between $D$ and $Y$ for individuals indexed by $i$:

$$E(Y_i) \equiv D_i \ E(Y_{i1}) + (1 - D_i) \ E(Y_{i0})$$

- $E(Y_{i0})$ and $E(Y_{i1})$ are the **counterfactual** to each other for each individual $i$

- In theory, the causal effect of $D$ on $Y$ is $\Delta = E(Y_{i1}) - E(Y_{i0})$

# Fundamental Problem of Causal Inference

- If $D_i = 1$, we only observe $Y_{i1}$ for individual $i$ ($Y_{i0}$?)
- If $D_i = 0$, we only observe $Y_{i0}$ for individual $i$ ($Y_{i1}$?)

That said, for every unit we can only observe one potential outcome depending on $D_i$—the counterfactual is purely hypothetical. Therefore, we are never able to compute $E(Y_{i1}) - E(Y_{i0})$ at the individual level.

## Average Treatment Effect

- Because of the fundamental problem at the individual level, we turn to the sample/group level instead.
- **The average treatment effect:**

$$ATE = \bar{E}(Y_{i1}|D = 1) - \bar{E}(Y_{i0}|D = 0)$$

- This average treatment effect can also be rewritten as the following, by a trick of adding and subtracting $\bar{E}(Y_{i0}|D = 1)$

$$\bar{E}(Y_{i1}|D = 1) - \bar{E}(Y_{i0}|D = 1) + [\bar{E}(Y_{i0}|D = 1) - \bar{E}(Y_{i0}|D = 0)]$$

- The formula shows that $ATE$ is the "true" causal effect $(\bar{E}(Y_{i1}|D = 1) - \bar{E}(Y_{i0}|D = 1))$ plus an extra component which people generally refer to as selection bias.

# Interpretation of ATE

Average treatment effect = true causal effect + selection bias

- Good causal inference requires *(1) an unbiased estimate on average treatment effect* and *(2) a successful effort to minimize selection bias.*

- Statistical models alone can only achieve (1) at best, but not (2).

# Why Is Regression Insufficient?

In the best-case scenario that our model is correctly specified with treatment $D$ and the known confounder $Z$:

$$Y = \alpha + \beta_1 D + \beta_2 Z + \epsilon$$

$$ATE = \beta_1 = E(Y|D = 1, Z) - E(Y|D = 0, Z)$$

- What about other potential confounders that we don't know or cannot measure?
- What if the sample is structurally different from population or other samples?
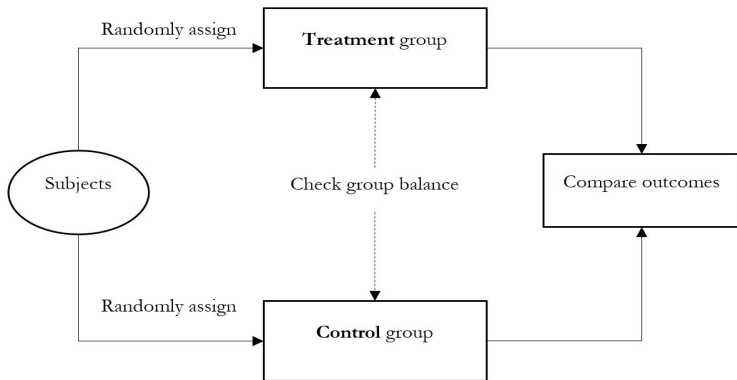
# How Can Randomized Experiment Do Better?

Random treatment assignment is the key:

- randomization
    - make sure each individual has an equal chance to get treated ($D = 1$)
    - such that individuals are impossible to self-select into either treatment group or control group.

- treatment manipulation
    - make treatment assignment completely independent of any other factors (both observed and unobserved in theory)
    - because assignment is dictated by random chance

# Randomized Experiment



If the procedure is completely random and covariates are perfectly balanced, ATE = TRUE causal effect

# Research Design

Even when experiment is impossible, researchers need to leverage research design to mimic the randomized experimental setting by

- taking advantage of "as-if random" or "exogenous shock" settings (natural experiment, instrumental variables, regression discontinuity)

- improving data balance (matching) or making data structurally more similar to the population (weighting)

- . . .

# Understanding Causal Inference Better

Imai, K., King, G. and Stuart, E.A., 2008. Misunderstandings between experimentalists and observationalists about causal inference. *Journal of the royal statistical society*, 171(2), pp.481-502.

Start with a few building blocks:

- *ATE* is the average treatment effect from sample
- *PATE* is the (true) treatment effect from population
- $\Delta$ is the estimation error of causal inference with sample
- subscript $S$ for sample selection and $T$ for treatment imbalance
- subscript $X$ for observed confounder and $U$ for unobserved confounder

# Decomposing Estimation Error

$$\Delta \equiv PATE - ATE$$
$$= \Delta_S + \Delta_T$$
$$= (\Delta_{S_X} + \Delta_{S_U}) + (\Delta_{T_X} + \Delta_{T_U})$$

- $S$ for sample selection and $T$ for treatment imbalance
- $X$ for observed confounder and $U$ for unobserved confounder

# $\Delta_{S_X}$

Estimation error caused by sample selection with respect to the known confounder:

- vanishes, if sample = population (e.g. in census)
- vanishes, if random sampling within strata (such that sample represents population)
- does not matter, if we don't care about population

# $\Delta_{S_U}$

Estimation error caused by sample selection with respect to the unknown confounder:

- vanishes, if random sampling from well-defined population
- vanishes, assuming that the distribution of $U$ in sample is identical to that in population)
- unverifiable in nature

# $\Delta_{T_X}$

Estimation error caused by treatment imbalance with respect to the known confounder:

- equals to 0, when $X$ balanced between the treated and control
  - This can be done by either *ex ante* blocking (split units on $X$ evenly) or *ex post* matching (pruning unmatched units)

# $\Delta_{T_U}$

Estimation error caused by treatment imbalance with respect to the unknown confounder:

- vanishes, by assumptions or random treatment assignment
- unverifiable in nature

# Evaluating Clinical Trials

nonrandom selection, small $n$, some blocking on covariates (age, race, etc.), random treatment assignment

- $\Delta_{S_X} \neq 0$
- $\Delta_{S_U} \neq 0$
- $\Delta_{T_X} \rightsquigarrow 0$ (depending on how well blocking works)
- $E(\Delta_{T_U}) = 0$

# Observational Study

no stratification, nonrandom selection, large $n$, no blacking or matching, nonrandom treatment assignment (no intervention of researchers)

- $\Delta_{S_X} \approx 0$, if representative, or corrected by weighting or matching
- $\Delta_{S_U} \neq 0$
- $\Delta_{T_X} \neq 0$
- $E(\Delta_{T_U}) \neq 0$, except by assumptions

# Observational Study, Well-matched

no stratification, nonrandom selection, large *n*, no blacking or matching, nonrandom treatment assignment (no intervention of researchers)

- $\Delta_{S_X} \approx 0$, if representative, or corrected by weighting or matching
- $\Delta_{S_U} \neq 0$
- $\Delta_{T_X} \approx 0$
- $E(\Delta_{T_U}) \neq 0$, except by assumptions

# Pros and Cons between Experiments and Observational Study

- Randomized experiment:
  - good at $\Delta_T$
  - relatively weak at $\Delta_S$ ("external validity")

- Observational study:
  - good at $\Delta_S$, if large $n$
  - weak at $\Delta_T$
  - need more effort in understanding and adjusting for $X$ (via matching or better modeling)

# Coming Up

- Model dependence (variable imbalance)

- Matching (techniques to improve treatment balance and causal inferences)