Economics 1960 – Literture review, modelling and identification strategies

Steven Lee Brown University

October 17 & 21, 2025 October 24 & 28, 2025

Why do I have to do literature review?

- All research builds on past research
 - "Standing on the shoulders of giants"
- Read papers: So you aren't duplicating research already done!
- Read papers: So you know how to convince readers your paper is worth reading!
 - Readers expect you to directly address how your findings relate to other papers in the topic/field

Lots of online resources available

- Google Scholar: https://scholar.google.com/
- For econ research specifically
 - IDEAS/RePEc: https://ideas.repec.org/
 - NBER Working paper series: https://www.nber.org/papers
 - CEPR Discussion papers: https://cepr.org/publications/discussion-papers

Google Scholar

How to search

- Distill your research question into keywords
- Don't necessarily read the whole paper
 - Read the abstract... then triage
 - Read the introduction... then triage
 - Read the whole paper
- If an author shows up many times in searches, helpful to go to their research page directly
- Know the "top" journals in your field

How to use papers to find more papers

- Use survey/review papers
 - Survey/review papers give an overview of the literature for a particular topic
 - Lots of papers will be cited here!
 - Top survey/review journals in econ: Journal of Economic Perspectives and Journal of Economic Literature
- Use the references of key related papers
 - Any high-quality paper will have its own literature review
 - If this paper is highly related to yours, many of its citations will be relevant for you!

Models and Identification Strategies

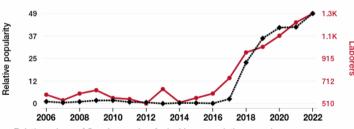
Most economics papers want to say something about causality

- We don't want to know if X just happens to be associated with Y
- We want to say X "causes" Y
- Or, in theory papers: Under a, b, c assumptions, we can say X causes Y

Popularity of the 'spiderman pointing' meme

correlates with

The number of veterinary technologists and technicians in Kansas



- ◆ Relative volume of Google searches for 'spiderman pointing meme' (without quotes, in the United States) · Source: Google Trends
- BLS estimate of veterinary technologists and technicians in Kansas · Source:
 Bureau of Larbor Statistics

2006-2022, r=0.977, r²=0.954, p<0.01 · tylervigen.com/spurious/correlation/31294

https://tylervigen.com/spurious-correlations

Convincing readers of causality can fail for a number of reasons

- Not underpinned by a theoretical model
- Not supported by an identification strategy

Models

What is a model?

- Simplified representation of reality used to explain how individuals/firms/groups make decisions
- You've encountered many models
 - Model of supply and demand: $Q^D = a bp$, $Q^s = c + dp$: $Q^D(p^*) = Q^S(p^*)$
 - Labor-leisure model: $T = h + \ell$, C = w, $U = U(C, \ell)$: max $U = U(C^*, \ell^*)$
 - Model of firm behavior: Q = f(K, L), $\pi = p \cdot Q wL rK$: $\max \pi = \pi(K^*, L^*)$

What are the building blocks of a model?

- Agents who is making the decision?
- Preferences/objectives what do the agents want?
- Constraints why can't they get everything they want?
- Decisions how can agents change their outcomes? what's in their control?
- (Sometimes) Equilibrium condition how do everyone's decisions tie together?

Every paper (outside of econometrics) has a model lurking in the background

- What is the effect of Medicaid expansion on mortality?
 - More money → better health? Intuitive?
- For individuals
 - Preferences: Maximize utility which is function of consumption and health
 - Constraints: Everyone can't spend more than they earn; stricter for low SES individuals
 - Decisions: Spend on health or other consumption goods?
 - Prediction: If Medicaid expansion goes to people that are very budget constrained, mortality goes down; otherwise, ambiguous

Bottom line

For theory thesis writers

- Helpful to work with existing models
- Every model has a set of simplifying restrictions, within each building block
- Loosen some restrictions
 - Does that change the implications of the model?
 - Does that better represent how agents in the real world operate?

For empirical writers

- Can I concisely describe what economic decision-making process underlies my research question?
- Even if it doesn't make it into your draft, can be helpful to clarify why your research is interesting or which additional analyses to conduct

Identification strategies

Causal inference

- Causal inference is the process of determining the causal relationship of one variable/phenomena with another
- An economics empirical strategy tries to uncover this causal relationship

Causal inference

- Causal inference is the process of determining the causal relationship of one variable/phenomena with another
- An economics empirical strategy tries to uncover this causal relationship
- Gold standard: Randomization from an experiment

Causal inference

- Causal inference is the process of determining the causal relationship of one variable/phenomena with another
- An economics empirical strategy tries to uncover this causal relationship
- Gold standard: Randomization from an experiment
- Alternative: Identitification strategy from statistical research design

Crash course on "potential outcomes"

- We care about some outcome Y and we think treatment $D \in \{0,1\}$ can influence Y
- What are the natural questions?
 - What happens to Y if agent i receives the treatment (D=1), or not (D=0)
 - What's the difference between the resulting Y's

Crash course on "potential outcomes"

- We care about some outcome Y and we think treatment $D \in \{0,1\}$ can influence Y
- What are the natural questions?
 - What happens to Y if agent i receives the treatment (D=1), or not (D=0)
 - What's the difference between the resulting Y's
- Let's define *potential* outcome $Y_i(D)$
 - $Y_i(D=1)$: outcome if unit i receives treatment D=1
 - $Y_i(D=0)$: outcome if unit i receives treatment D=0 (no treatment)
- The causal effect of the treatment:

$$\tau_i = Y_i(1) - Y_i(0)$$

and the population average

$$\tau = E[Y_i(1) - Y_i(0)]$$

Crash course on "potential outcomes"

- We care about some outcome Y and we think treatment $D \in \{0,1\}$ can influence Y
- What are the natural questions?
 - What happens to Y if agent i receives the treatment (D=1), or not (D=0)
 - What's the difference between the resulting Y's
- Let's define *potential* outcome $Y_i(D)$
 - $Y_i(D=1)$: outcome if unit i receives treatment D=1
 - $Y_i(D=0)$: outcome if unit i receives treatment D=0 (no treatment)
- The causal effect of the treatment:

$$\tau_i = Y_i(1) - Y_i(0)$$

and the population average

$$\tau = E[Y_i(1) - Y_i(0)]$$

- ... but we can't observe both $Y_i(1)$ and $Y_i(0)$

Empirical strategy tries to get at a version of τ after some assumptions

- Randomization (Experiment)
 - Randomly split sample into group A and group B
 - Only A gets treatment
 - Compare group A against group B since A and B are similar through randomization

Empirical strategy tries to get at a version of au after some assumptions

- Randomization (Experiment)
 - Randomly split sample into group A and group B
 - Only A gets treatment
 - Compare group A against group B since A and B are similar through randomization
- Difference-in-Differences
 - I observe groups A and B before treatment happens
 - Only A gets treatment at time t
 - Compare group A against group B after time t and before time t

Empirical strategy tries to get at a version of au after some assumptions

- Randomization (Experiment)
 - Randomly split sample into group A and group B
 - Only A gets treatment
 - Compare group A against group B since A and B are similar through randomization
- Difference-in-Differences
 - I observe groups A and B before treatment happens
 - Only A gets treatment at time t
 - Compare group A against group B after time t and before time t
- Regression Discontinuity
 - I observe groups A and B that are cleanly delineated at some cutoff of characteristic X
 - Only A gets treatment
 - Compare group A against group B very close to the cutoff of X

Empirical strategy tries to get at a version of au after some assumptions

- Randomization (Experiment)
 - Randomly split sample into group A and group B
 - Only A gets treatment
 - Compare group A against group B since A and B are similar through randomization

Difference-in-Differences

- I observe groups A and B before treatment happens
- Only A gets treatment at time t
- Compare group A against group B after time t and before time t

Regression Discontinuity

- I observe groups A and B that are cleanly delineated at some cutoff of characteristic X
- Only A gets treatment
- Compare group A against group B very close to the cutoff of X

Instrumental variables

- I observe some phenomena that pushes people to get treatment but is unimportant for the outcome
- Compare outcomes for treated vs. untreated only for those responsive to the phenomena

Randomization

$$y_i = \alpha + \beta \cdot D_i + X_i + \varepsilon_i \tag{1}$$

- Across your sample, you treated some individuals and some you did not
- $-y_i$ is the outcome of interest for agent i
- D_i is an indicator function $D_{it} = \mathbb{1}_{i \in \{\mathit{Treat}\}}$ (assigned treatment status for i)
- X_i are a vector of controls $X_i = \{x_i^1, x_i^2, \dots\}$

Randomization

$$y_i = \alpha + \beta \cdot D_i + X_i + \varepsilon_i \tag{1}$$

- Across your sample, you treated some individuals and some you did not
- $-y_i$ is the outcome of interest for agent i
- D_i is an indicator function $D_{it} = \mathbb{1}_{i \in \{Treat\}}$ (assigned treatment status for i)
- X_i are a vector of controls $X_i = \{x_i^1, x_i^2, \dots\}$

Key assumptions

- SUTVA: Treatment status for person i doesn't affect person j's outcome
- Random assignment: D_i randomly assigned
 - Balance test

$$x_i^n = \delta \cdot D_i + \eta_i$$

You should find $\delta \approx 0$

Difference-in-differences (DiD)

$$y_i = \alpha_i + \gamma_t + \beta \cdot D_{it} + X_{it} + \varepsilon_{it}$$
 (2)

- Observe multiple units across multiple time periods
- Some units treated starting $t=t^*$, called post-period, $Post=\mathbb{1}_{t\geq t^*}$
- Treatment is then $D_{it} = \mathbb{1}_{i \in \{Treat\}} \cdot Post$
- Unit fixed effects, α_i ; time fixed effects, γ_t
 - Each unit has different starting points
 - All units may face common shocks in a time period

Difference-in-differences (DiD)

$$y_i = \alpha_i + \gamma_t + \beta \cdot D_{it} + X_{it} + \varepsilon_{it}$$
 (2)

- Observe multiple units across multiple time periods
- Some units treated starting $t=t^*$, called post-period, $Post=\mathbb{1}_{t\geq t^*}$
- Treatment is then $D_{it} = \mathbb{1}_{i \in \{\mathit{Treat}\}} \cdot \mathit{Post}$
- Unit fixed effects, α_i ; time fixed effects, γ_t
 - Each unit has different starting points
 - All units may face common shocks in a time period

Key assumptions

- Parallel trends: w/o D_{it} , treated units have same change in outcomes as untreated units
 - Think about potential outcomes! If $i \notin \{Treat\}$ is proxy for treated group potential outcomes $Y_{i \in \{Treat\}}(0)$, need to be sure the two groups trend together when $t \leq m$
- No anticipation: units do not react in anticipation of treatment arrival in t^st

Dynamic DiD (Event Study)

$$y_i = \alpha_i + \gamma_t + \sum_{m=-K}^{L} \beta_m \cdot z_{i,t+m} + X_{it} + \varepsilon_{it}$$
(3)

- Suppose you observe outcomes for $t \in [t^* K, t^* + L]$, where treatment starts at $t = t^*$
- Define $m = t t^*$: time from treatment
- $z_{i,t+m} = \mathbb{1}_{i \in \{\mathit{Treat}\}} \cdot \mathbb{1}_{t=t^*-m}$, dummy variable for "leads" and "lags" of treatment status

Purpose

- Show how the effect of treatment evolves over course of the treatment
- Placebo test for the plausibility of parallel trends

Regression discontinuity (RD)

$$y_i = \alpha + \beta \cdot \mathbb{1}_{r_i \ge \tau} + f(r_i) + X_i + \varepsilon_i \tag{4}$$

- Suppose there's an attribute observed for all units and a value for that attribute which determines treatment
- $-r_i$ is the attribute, called the "running variable"
- $-\tau$ is the cutoff value
- Then $D_i = \mathbb{1}_{r_i \geq \tau}$

Regression discontinuity (RD)

$$y_i = \alpha + \beta \cdot \mathbb{1}_{r_i \ge \tau} + f(r_i) + X_i + \varepsilon_i \tag{4}$$

- Suppose there's an attribute observed for all units and a value for that attribute which determines treatment
- r_i is the attribute, called the "running variable"
- $-\tau$ is the cutoff value
- Then $D_i = \mathbb{1}_{r_i > \tau}$

Key assumptions

- Continuity at cutoff: At the limit, agents just below and just above au are comparable
- No manipulation/sorting: Agents aren't placing themselves on either side of au on purpose

Instrumental Variables/Two-stage least squares (2SLS)

First-stage:

$$D_i = \pi_0 + \pi_1 Z_i + \pi_2 X_i + u_i \tag{5}$$

Second-stage:

$$Y_i = \beta_0 + \beta_1 \hat{D}_i + \beta_2 X_i + \varepsilon_i \tag{6}$$

- Suppose you have some attribute from "nature" that changes treatment but not the outcome
- D_i is the "endogenous treatment" which is influenced by Z_i , the instrumental variable
- $-\hat{D}_i$ is the predicted treatment stemming from Z_i

Instrumental Variables/Two-stage least squares (2SLS)

First-stage:

$$D_i = \pi_0 + \pi_1 Z_i + \pi_2 X_i + u_i \tag{5}$$

Second-stage:

$$Y_i = \beta_0 + \beta_1 \hat{D}_i + \beta_2 X_i + \varepsilon_i \tag{6}$$

- Suppose you have some attribute from "nature" that changes treatment but not the outcome
- D_i is the "endogenous treatment" which is influenced by Z_i , the instrumental variable
- \hat{D}_i is the predicted treatment stemming from Z_i

Key assumptions

- Independence: Z_i is assigned as-good-as randomly
- Instrument relevance: Z_i actually predicts D_i
- Exclusion restriction: Z_i only affects Y_i through D_i

Common threats to identification

- Omitted variable bias: Treatment D is correlated with some factor W. W also affects Y
 - Ex: Does college $(D) \uparrow \text{ wages } (Y)$? High school grades (W) are correlated with D and Y
- Reverse causality: Y causes treatment D
 - Ex: Does democracy (D) cause economic growth (Y)? Well, growth could lead to democracy
- Selection bias: Agents self-select into treatment D based on expected outcomes
 - Ex: Does tutoring improve grades? More motivated students sign up for tutoring
- Spillovers: Agent i's treatment affects agent j's outcomes
 - Ex: Does door-to-door can vassing increase votes? Knocking on i's door $\to i$ speaks to neighbor j

No approach is truly assumption free

- You are almost always assuming something in your model, empirical strategy, etc.
- Best papers:
 - State the key assumptions
 - Show the results
 - Provide evidence or robustness checks to make assumptions plausible
- Last step is the hardest, downfall of most projects
- Still, being upfront about assumptions shows you want to be serious economist

Some tips

- Start by imagining the ideal randomized experiment
 - Then think of ways your setting differs from this ideal
- Study the institutional context carefully
- Many ID strategies exploit "quirks" in rules, policies, and human norms
 - Counties have different policies
 - Oversubscribed program ightarrow random lottery determines access
 - Eligibility for benefit/program must be "fair" and "objective" \rightarrow eligibility based on clear (e.g. age, income, test scores)