

Making Causal Critiques

Day 4 - How much are we Learning?

Jonathan Phillips

January 24, 2019

How much are we Learning?

- ▶ Everything we have discussed so far has been about the **accuracy** of a causal claim
- ▶ But not every study is as valuable to political science
- ▶ We *learn* more from some studies than from others
 1. Reliability of the claim
 2. Reproducibility of the claim
 3. Scope (*generalizability*) of the claim

Robustness

- ▶ For simplicity, we publish a paper with a 'final' result
 - ▶ 1% extra GDP growth increases the President's chance of re-election by 5%
 - ▶ One more year of drought increases the risk of conflict by 10%
- ▶ But how **confident** are we in these figures?
- ▶ Good studies include estimates of uncertainty
 - ▶ 1% extra GDP growth increases the President's chance of re-election by 5% with a standard deviation of 0.2%
 - ▶ One more year of drought increases the risk of conflict by 10%, with a standard deviation of 30%
- ▶ But these confidence intervals are usually for a *single* methodology and set of assumptions

Robustness

- ▶ What if our assumptions were wrong?
- ▶ How much would our results change if we used a different methodology?
 - ▶ Including different controls
 - ▶ Including alternative measures of the variables
 - ▶ Including or excluding outliers
- ▶ If we can change all these things and still get the same answers, our result is **reliable** and **robust**

Robustness

- ▶ For example, Michalopoulos and Papaioannou (2013) show that more centralized pre-colonial societies in Africa have more economic activity today
- ▶ Robustness tests include:
 - ▶ Extra controls for disease, land, natural resources
 - ▶ Alternative model for spatial autocorrelation
 - ▶ Country fixed effects to focus only on within-country variation
 - ▶ Comparing only neighbouring societies
 - ▶ Alternative codings of centralized pre-colonial societies
 - ▶ Alternative measures of economic activity (nightlights etc.)
 - ▶ Different units of analysis - grid squares instead of ethnic territories

Robustness

- ▶ Robustness tests help avoid **researcher bias**
 - ▶ Running 200 models with different covariates
 - ▶ Only reporting one that is significant
 - ▶ But even if there was **no causal effect** in the data, *by chance* we would expect 20 models to produce significant effects

Reproducibility

1. If we take the same data and apply the same method, do we get the same result?
 - ▶ Often, no! Only 35% in Brazilian political science journals (Avelino and Desposato 2018)
2. If we take another sample of data and apply the same method, do we get the same result?
 - ▶ Very rarely done

Reproducibility

- ▶ A big problem for reproducibility is **publication bias**
 - ▶ Lots of researchers perform lots of studies
 - ▶ Some find positive results, some negative, many 'null' findings
 - ▶ But journals want readers, and readers like positive results
 - ▶ So only the positive results get published
- ▶ If you're reading a paper, think of the ten other papers you're *not* reading that tried the same thing and found no effect

Reproducibility

- ▶ One solution is **Pre-registration**
 - ▶ Submit your study design to a website - what you will analyse and how
 - ▶ Everyone knows who is researching what, and if they published or not
 - ▶ Researchers are also less tempted to 'pick' their preferred analysis after seeing the data
 - ▶ Eg. EGAP Pre-Registration

Generalizability

- ▶ But even if studies are robust and reproducible, **how much** are we learning?
- ▶ We can learn very little even from a precise, bias-free study:
 - ▶ IgNobel Prize
 - ▶ "Suicide rates are linked to the amount of country music played on the radio"
 - ▶ "Is using voodoo dolls effective?"
 - ▶ "Why do old men have big ears?"
 - ▶ "How exposure to a crocodile encourages people to gamble"

Generalizability

► **Internal Validity**

- Are the conclusions of the study accurate *within* the sample?
- Are the assumptions valid, is our causal effect biased?
- Is the conclusion reliable if we use slightly different assumptions?

► **External Validity**

- How far can the results 'travel' outside of the study design?
 1. Does the study reflect a wider population?
 2. How big, representative and interesting is that wider population?

Generalizability

- ▶ For example, Chattopadhyay and Duflo (2004) argue that women leaders invest more in education using data from 265 villages in two states in India (West Bengal and Rajasthan)
- ▶ But does the conclusion hold if we get more data from:
 1. 265 different villages?
 2. Different states?
 3. Different countries?
 4. Different years?

Generalizability

- ▶ Most studies are designed with generalizability in mind:
 - ▶ Representative Samples are drawn from a target population
 - ▶ So we can use **statistical inference** to extend our conclusions from the sample to the population
 - ▶ Note this only works if we know all the units (hidden tribes etc.)
 - ▶ But Chattopadhyay and Duflo (2004) did not take a representative sample of villages - they surveyed all villages in one district
 - ▶ Their widely-cited paper *only* applies to Birbhum and Udaipur districts
 - ▶ We have no evidence of how women leaders govern elsewhere in India or the world

Generalizability

- ▶ Specific causal research designs also restrict the scope of our findings
 - ▶ Precisely because we had to restrict our sample to find appropriate counterfactuals
 - ▶ The new comparisons are often less representative or interesting
- ▶
- ▶ Instead of an **Average Treatment Effect (ATE)** they represent a **Local Average Treatment Effect (ATE)**
 - ▶ A treatment effect applicable only to those people who were affected by the methodology's treatment: **compliers**

Field Experiments

- ▶ Field experiments require lots of compromises and assumptions
 - ▶ Costs
 - ▶ Ethical restrictions
 - ▶ Consent
- ▶ Implementation is limited to a small sample, often non-representative
- ▶ And the findings *only* apply to that sample
- ▶ Or maybe only to a sub-group of that sample

Field Experiments

- ▶ External Validity Limitations of Field Experiments:

Field Experiments

- ▶ External Validity Limitations of Field Experiments:
 - ▶ What theory are we testing? Can't accumulate knowledge without theory. The causal mechanisms are still a black box.
 - ▶ Limited portability of findings - context matters for the treatment effects:
 - ▶ Eg. CCTs improve child health only where clinics are available
 - ▶ Average effects may not apply to any individual
- s
- ▶ Naive application of policy implications
- ▶ How much do the results depend on researcher oversight?

Lab Experiments

- ▶ Problems generalizing from the lab:
 - ▶ **Hawthorne effect:** Lab context influences behaviour, social desirability bias
 - ▶ **Context effects:** The real-world always provides more information, more history
 - ▶ **Process effects:** People care *how* decisions are made
 - ▶ **Selection effects:** Actors in specific roles are rarely representative samples, 'WEIRD' or pro-social lab subjects

Lab Experiments

- ▶ The lab differs from the field:

Lab Experiments

- ▶ The lab differs from the field:
 - ▶ The stakes
 - ▶ The norms (specific norms of being an experimental subject)
 - ▶ The degree of scrutiny
 - ▶ The sample of individuals
 - ▶ The degree of anonymity

Lab Experiments

- ▶ Levitt and List 2006 argue lab experiments are *inherently* flawed because the decisions we want to measure are likely to change depending on the degree of **scrutiny**

Lab Experiments

- ▶ Levitt and List 2006 argue lab experiments are *inherently* flawed because the decisions we want to measure are likely to change depending on the degree of **scrutiny**
- ▶ “You tip more when you’re on a date”
- ▶ Social norms are activated, eg. treating one-shot games like repeated games
- ▶ Scrutiny alters who wants to make a decision as well as the decision they make
- ▶ Subjects use cues (heuristics) to draw on ‘similar’ situations from the real world

Lab Experiments

- ▶ Many studies find more cooperation in the lab than in the real world

Lab Experiments

- ▶ Many studies find more cooperation in the lab than in the real world
 - ▶ Scrutiny increases cooperation
 - ▶ Anonymity reduces cooperation
 - ▶ That's interesting in itself! We can manipulate the degree of scrutiny/anonymity etc.

Lab Experiments

- ▶ Lab experiments may be generalizable where norms/morality are less important

Lab Experiments

- ▶ Lab experiments may be generalizable where norms/morality are less important
 - ▶ ???

Survey Experiments

- ▶ Treatment occurs *within* the survey questionnaire
 - ▶ Different versions of the questionnaire randomly applied
 - ▶ Not a field experiment: Still an artificial context
 - ▶ Not a lab experiment: People not brought to a single location or interacting

Conjoint Survey Experiments

- ▶ How do people make choices between many options?
- ▶ Treatments are often 'bundles', but which aspect matters most?

Conjoint Survey Experiments

- ▶ Hainmueller et al 2013 - How do attitudes to immigrants depend on immigrant characteristics?

Conjoint Survey Experiments

- ▶ Hainmueller et al 2013 - How do attitudes to immigrants depend on immigrant characteristics?
- ▶ Vary education, profession, language, gender, national origin, etc.
- ▶ Profiles
 - ▶ Attributes
 - ▶ Values
- ▶ Randomize attribute order to prevent bias

| | Immigrant 1 | Immigrant 2 |
|--------------------------------|---|---|
| Prior Trips to the U.S. | Entered the U.S. once before on a tourist visa | Entered the U.S. once before on a tourist visa |
| Reason for Application | Reunite with family members already in U.S. | Reunite with family members already in U.S. |
| Country of Origin | Mexico | Iraq |
| Language Skills | During admission interview, this applicant spoke fluent English | During admission interview, this applicant spoke fluent English |
| Profession | Child care provider | Teacher |
| Job Experience | One to two years of job training and experience | Three to five years of job training and experience |
| Employment Plans | Does not have a contract with a U.S. employer but has done job interviews | Will look for work after arriving in the U.S. |
| Education Level | Equivalent to completing two years of college in the U.S. | Equivalent to completing a college degree in the U.S. |
| Gender | Female | Male |

On a scale from 1 to 7, where 1 indicates that the United States should absolutely not admit the immigrant and 7 indicates that the United States should definitely admit the immigrant, how would you rate immigrant 1?



Using the same scale, how would you rate immigrant 2?



choice outcomes hereafter. Second, in "rating-based conjoint analysis," respondents give a numerical rating to each profile which represents their degree of preference for the profile. This format is preferred by some analysts who contend that such ratings provide more direct, finely grained information about respondents' preferences. We call this latter type of outcome a *rating outcome*.

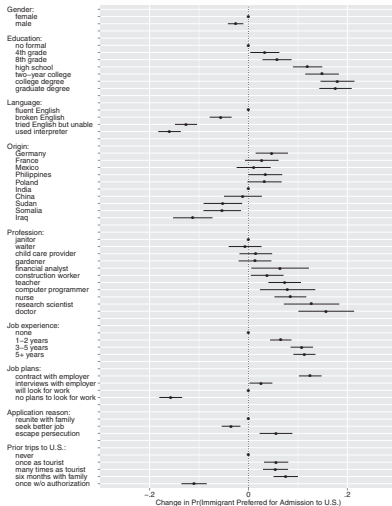


Fig. 3 Effects of immigrant attributes on preference for admission. This plot shows estimates of the effects of the randomly assigned immigrant attributes on the probability of being preferred for admission to the United States. Estimates are based on the regression estimators with clustered standard errors; bars represent 95% confidence intervals. The points without horizontal bars denote the attribute value that is the reference category for each attribute.

Conjoint Survey Experiments

- ▶ How realistic are the responses?
 - ▶ Not a behavioural measure; nothing at stake
 - ▶ Social desirability bias
 - ▶ Not like real-world preference-formation process
- ▶ Hainmueller et al 2014 - compare conjoint responses to a Swiss referendum
- ▶ Citizens voted on specific naturalization applicants (Really!)

Figure S11: Effects of Applicant Attributes on Opposition to Naturalization Request (Un-weighted Survey Sample)

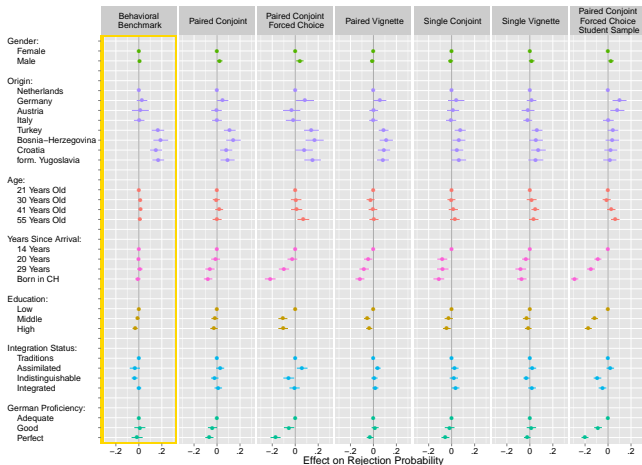


Figure shows point estimates (dots) and corresponding, cluster-robust 95 % confidence intervals (horizontal lines) from ordinary least squares regressions. The dots on the zero line without confidence intervals denote the reference category for each applicant attribute.

Conjoint Survey Experiments

- ▶ But note the conjoint method still hugely under-estimated the overall rejection rate
- ▶ 21% versus 37% in reality

Regression Discontinuity

- ▶ The LATE is for those people who were so close to the discontinuity that whether they were treated or not is basically random
 - ▶ Even though those cases are rare (eg. tied elections)
 - ▶ Even though we use data from a lot more people to estimate the LATE
- ▶ Do we care about those people at the discontinuity?

Regression Discontinuity

- ▶ For example, Titiunik et al (2011) use a regression discontinuity on close elections in Brazilian municipalities to show that incumbent Mayors are more likely to lose (negative incumbency effect)
 - ▶ But this does **not** mean that there is a negative incumbency effect in most Brazilian municipalities
 - ▶ Only about 500 out of 5,570 municipalities had 'close' elections (within $\pm 3\%$)
 - ▶ Those municipalities were more urban, southern and wealthy than the rest
 - ▶ We do not learn anything about places where the result was a landslide (70-80%)
 - ▶ But these are the places where incumbents probably benefitted a lot!

Regression Discontinuity

- ▶ Similarly, geographic regression discontinuities only tells us the effect of living on one side of the border *for people who live by the border*
 - ▶ But who chooses to live by a border? People who like rural areas, migrants etc.
 - ▶ Self-selection bias has come back!

Instrumental Variables

- ▶ Critique (Deaton 2009):

Instrumental Variables

- ▶ Critique (Deaton 2009):
 - ▶ Our causal models need to represent a theory, not just be an arbitrary equation

Instrumental Variables

- ▶ Critique (Deaton 2009):
 - ▶ Our causal models need to represent a theory, not just be an arbitrary equation
 - ▶ If we use 'convenient' instruments, our causal effect and complier population are out of our control and might not be interesting

Instrumental Variables

- ▶ Critique (Deaton 2009):
 - ▶ Our causal models need to represent a theory, not just be an arbitrary equation
 - ▶ If we use 'convenient' instruments, our causal effect and complier population are out of our control and might not be interesting
 - ▶ LATE causal estimates are not a good guide to policy effects

Instrumental Variables

- ▶ Critique (Deaton 2009):
 - ▶ Our causal models need to represent a theory, not just be an arbitrary equation
 - ▶ If we use 'convenient' instruments, our causal effect and complier population are out of our control and might not be interesting
 - ▶ LATE causal estimates are not a good guide to policy effects
 - ▶ 'External' to our model is not the same as 'Exogenous', and we can't test exogeneity

Instrumental Variables

- ▶ Critique (Deaton 2009):
 - ▶ Our causal models need to represent a theory, not just be an arbitrary equation
 - ▶ If we use 'convenient' instruments, our causal effect and complier population are out of our control and might not be interesting
 - ▶ LATE causal estimates are not a good guide to policy effects
 - ▶ 'External' to our model is not the same as 'Exogenous', and we can't test exogeneity
 - ▶ Where the instrument is an arbitrary rule, there is often sorting as people re-adjust

Learning in Political Science

- ▶ So how much can we learn?
 - ▶ We have to make careful judgments based on internal and external validity
 - ▶ Ideally combining multiple methodologies to compare low-bias low-generalizability evidence with high-bias high-generalizability evidence