

Session III

External Validity

Thomas J. Leeper

Government Department
London School of Economics and Political Science

Share your TESS Examples

In groups of 4–5, share your TESS examples

- What was the researcher's question?
- How did they test it experimentally?
- What was interesting or surprising about the designs?

Take about 7–8 minutes.

1 External Validity of a Sample

- Design-based
- Model-based

2 Other Notions of External Validity

- Settings
- Unit
- Treatments
- Outcomes

3 Participant Recruitment

1 External Validity of a Sample

- Design-based
- Model-based

2 Other Notions of External Validity

- Settings
- Unit
- Treatments
- Outcomes

3 Participant Recruitment

Think–Pair–Share

Consider the following question:

What makes an experiment (or any research study) generalizable?
What does it mean for a study's results to “generalize”?

- 1 Write or think to yourself for 90 seconds
- 2 Then, discuss with the person next to you

“The Gold Standard”

a population-based experiment uses survey sampling methods to produce a collection of experimental subjects that is representative of the target population of interest for a particular theory . . . the population represented by the sample should be representative of the population to which the researcher intends to extend his or her findings. In population-based experiments, experimental subjects are randomly assigned to conditions by the researcher

p2. from Mutz, Diana. 2011. *Population-Based Survey Experiments*. Princeton University Press.

Surveys Start with an Inference Population

- We want to speak to a population
- But what population is it?
 - A national population?
 - Adults in Western, industrialized democracies?
 - All human beings?
- This is rarely specified, but is important when we think about whether a sample is appropriate

A Hypothetical Census

- Advantages
 - Perfectly representative
 - Sample statistics are population parameters
- Disadvantages
 - Costs
 - Feasibility
 - Need

1 External Validity of a Sample

- Design-based
- Model-based

2 Other Notions of External Validity

- Settings
- Unit
- Treatments
- Outcomes

3 Participant Recruitment

Sampling Considerations . . .

- Design-based survey samples all work from the premise of each unit having a *known, non-zero* probability of being sampled
 - SRS is representative per se
 - Non-self-weighting samples representative when weighted
- Random sampling ensures that samples are, *in expectation*, representative of the population *in all respects*
 - Demographics
 - Covariances
 - Potential outcomes

Representativeness

What does it mean for a sample to be representative?

- Census?
- Probability-based sampling?
- Quota fulfillment?
- Others?

Which of these matter?

Combining Probability Sampling and Experimental Design

- Sample is representative of population in every respect (in expectation)
- Sample Average Treatment Effect (SATE) is the average of the sample's individual-level treatment effects
 - Unbiased estimate of PATE
 - Not necessarily any unit's individual treatment effect
 - Blocking might reduce variance
 - Optimized for estimating SATE

Credibility of all of this is based on *design* only

Sampling aspect only works in a world of perfect coverage and no response bias

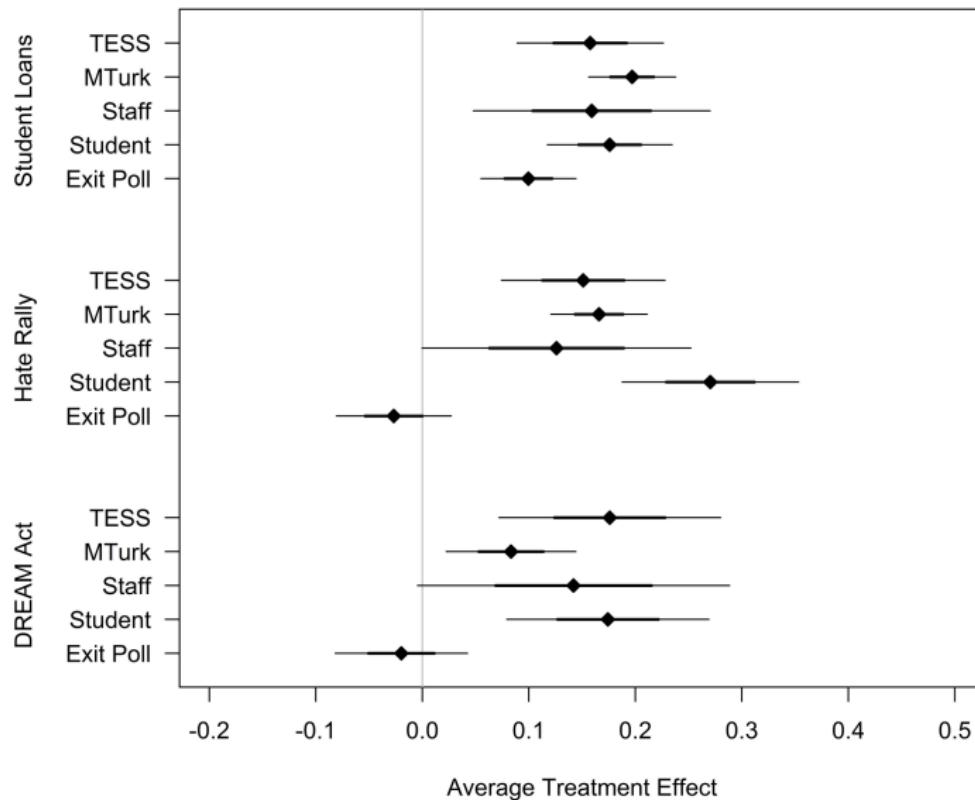
My View

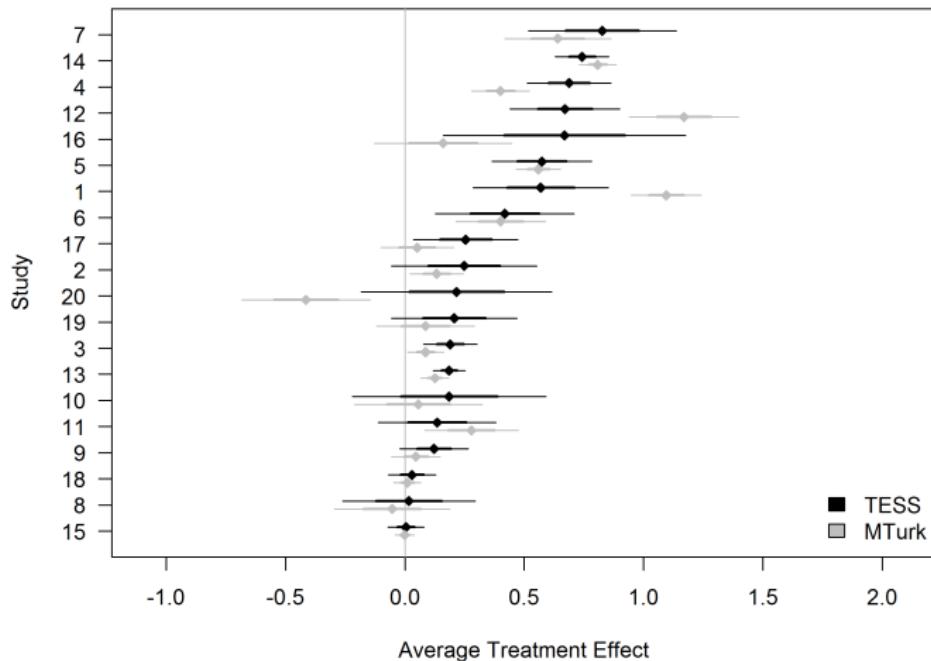
100% design-based inference does not exist!

- All survey designs involve reweighting adjustments
- Representativeness is a more complex issue than demographic comparisons
- Randomization gives us clear causal inference about a *local* effect
 - Sacrifice representativeness for causal inference
 - Try to figure nature of the *localness*

| | GfK | Poll | Student | Staff | MTurk | Ads | ANES |
|------------------|------|-------|---------|-------|-------|-------|-------|
| Dem. (%) | 51.3 | 86.1 | 75.7 | 66.4 | 62.1 | 72.1 | 46.2 |
| Rep. (%) | 46.0 | 7.7 | 17.8 | 16.4 | 20.3 | 14.7 | 39.3 |
| Lib. (%) | 27.8 | 75.4 | 68.5 | 62.7 | 60.4 | 66.2 | 23.8 |
| Con. (%) | 35.3 | 9.4 | 14.7 | 19.8 | 19.1 | 17.7 | 36.1 |
| Fem. (%) | 51.1 | 60.8 | 56.4 | 50.8 | 41.7 | 65.3 | 51.9 |
| White (%) | 77.9 | 67.6 | 62.9 | 60.2 | 76.0 | 53.8 | 80.4 |
| Age | 49.4 | 40-49 | 18-24 | 25-34 | 25-34 | 25-34 | 50-54 |
| Interest | 2.8 | 3.5 | 3.2 | 2.8 | 2.7 | 3.0 | 3.0 |
| N | 593 | 741 | 299 | 128 | 1024 | 80 | – |

Mullinix et al. 2015. "The Generalizability of Survey Experiments." *Journal of Experimental Political Science*.





My View

100% design-based inference does not exist!

- All survey designs involve reweighting adjustments
- Representativeness is a more complex issue than demographic comparisons
- Randomization gives us clear causal inference about a *local* effect
 - Sacrifice representativeness for causal inference
 - Try to figure nature of the *localness*

Common Differences

- Most common thing to focus on is demographic representativeness
 - Sears (1986): “students aren’t real people”
 - Western, educated, industrialized, rich, democratic (WEIRD) psychology participants
- But do those characteristics actually matter?

Common differences

Shadish, Cook, and Campbell tell us to think about:

- Surface similarities
- **Ruling out irrelevancies**
- **Making discriminations**
- Interpolation/extrapolation

Focus on effect heterogeneity

- **Constant effects:** TE_i and Y_{i0} are same for all observations
- **Homogeneous effects:** TE_i is same for all observations
- **Heterogeneous effects:** TE_i is different for all observations

Focus on effect heterogeneity

- Think about and make an evidence-based argument for why you think there are (or are not) heterogeneous effects
- If you think there is heterogeneity, then we probably do not care about the SATE anyway
- Conditional Average Treatment Effect:
$$E[Y_{1i}|X = 1, Z = z] - E[Y_{0i}|X = 0, Z = z]$$

1 External Validity of a Sample

- Design-based
- Model-based

2 Other Notions of External Validity

- Settings
- Unit
- Treatments
- Outcomes

3 Participant Recruitment

Stratification/Blocking

As soon as we care about heterogeneous effects, it makes sense to stratify and block on factors that might moderate the treatment effect.

As soon as we identify all sources of heterogeneity, it doesn't matter what sample we use because effects are *by definition* homogeneous within such strata.

But, we never know when we've reached that point!

If we acknowledge and start thinking about effect heterogeneity, does this mean we can use any convenient group of participants as if they were probability samples?

No. Of course not.

Not All “Samples” Are Alike

- Different types:
 - Passive/opt-in/“river sampling”
 - Sample of convenience
 - Snowball sample
 - Students
 - Crowdsourcing
- Differ in numerous ways
 - Cost
 - “Experience”
 - Attentiveness
 - Demographics

Costs per participant

From one of my studies:

| Sample | Cost | n | Cost/participant |
|-----------|---------|------|------------------|
| National | \$13200 | 593 | \$22.26 |
| Exit Poll | \$3000 | 741 | \$4.05 |
| Students | \$0 | 299 | \$0 |
| Staff | \$1280 | 128 | \$10.00 |
| MTurk | \$550 | 1024 | \$0.54 |
| Ads | \$636 | 80 | \$7.95 |

Participant Experience

- A lot of growing concern about experience
- Larger literature on “panel conditioning”
 - Inconclusive evidence
- Some numbers:
 - MTurk workers are doing 100+ studies per month
 - Numbers are the same for YouGov panelists

Reweighting

- If effects are heterogeneous, it may be possible to *reweight* unrepresentative data to match a population
- Any method for this is “model-based” (rather than “design-based”)
- Not widely used or evaluated (yet)
- All techniques build on the idea of stratification

Review of Stratification

- 1 Define population
- 2 Construct a sampling frame
- 3 Identify variables we already know about units in the sampling frame
- 4 Stratify sampling frame based on these characteristics
- 5 Collect an SRS within each stratum
- 6 Aggregate our results

Post-Stratification

- Used to correct for nonresponse, coverage errors, and sampling errors
- Reweight sample data to match population distributions
 - Divide sample and population into strata
 - Weight units in each stratum so that the weighted sample stratum contains the same proportion of units as the population stratum does
- There are numerous related techniques

Post-Stratification: Example

- Imagine our sample ends up skewed on immigration status and gender relative to the population

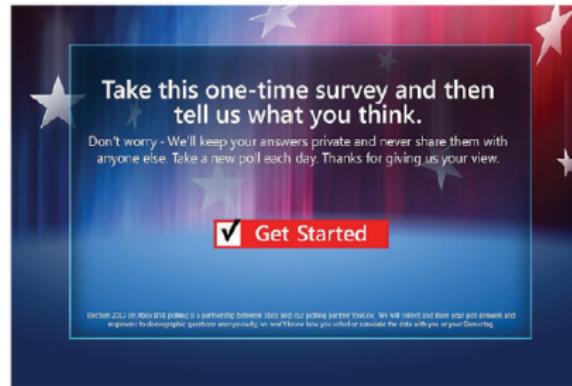
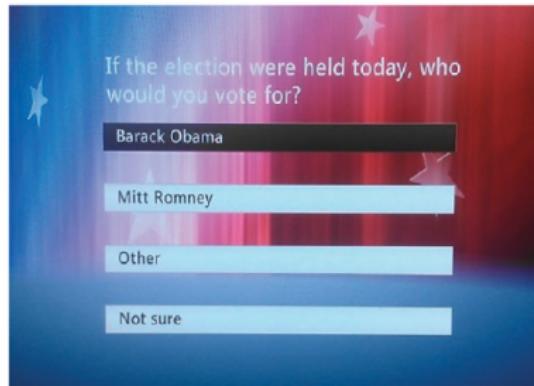
| Group | Pop. | Sample | Rep. | Weight |
|---------------------|------|--------|-------|--------|
| Native-born, Female | .45 | .5 | Over | 0.900 |
| Native-born, Male | .45 | .4 | Under | 1.125 |
| Immigrant, Female | .05 | .07 | Over | 0.714 |
| Immigrant, Male | .05 | .03 | Under | 1.667 |

- PS weight is just $w_{ps} = N_I/n_I$

Post-Stratification

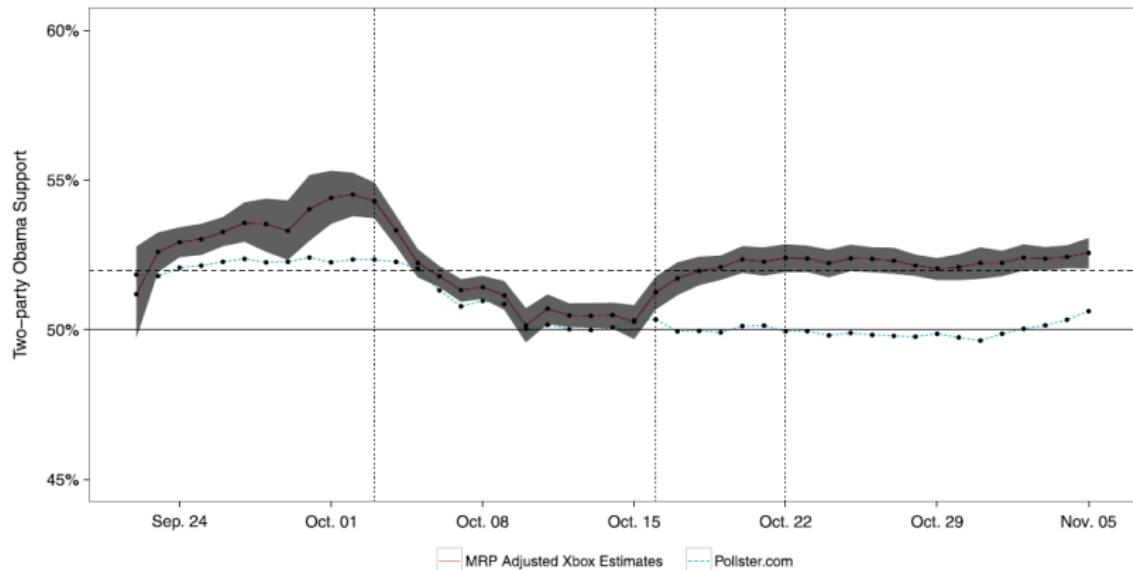
- This is the basis for inference in non-probability samples
 - *Demographic* representativeness
- Online panels will reweight sample based on age, sex, education, etc.
- Purely design-based surveys are increasingly rare

The Xbox Study

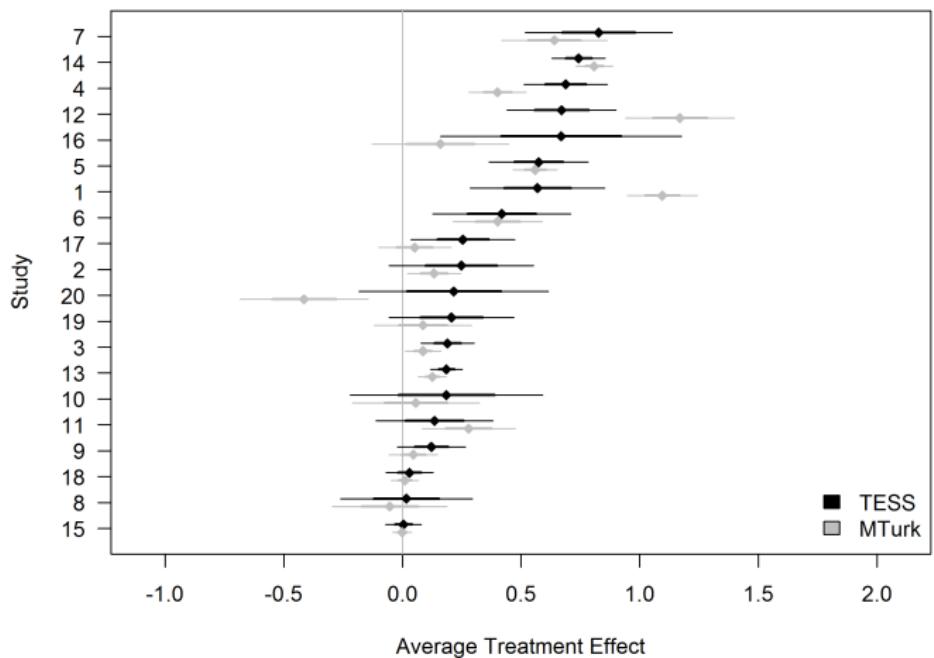


Wang et al. 2015. "Forecasting elections with non-representative polls." *International Journal of Forecasting*.

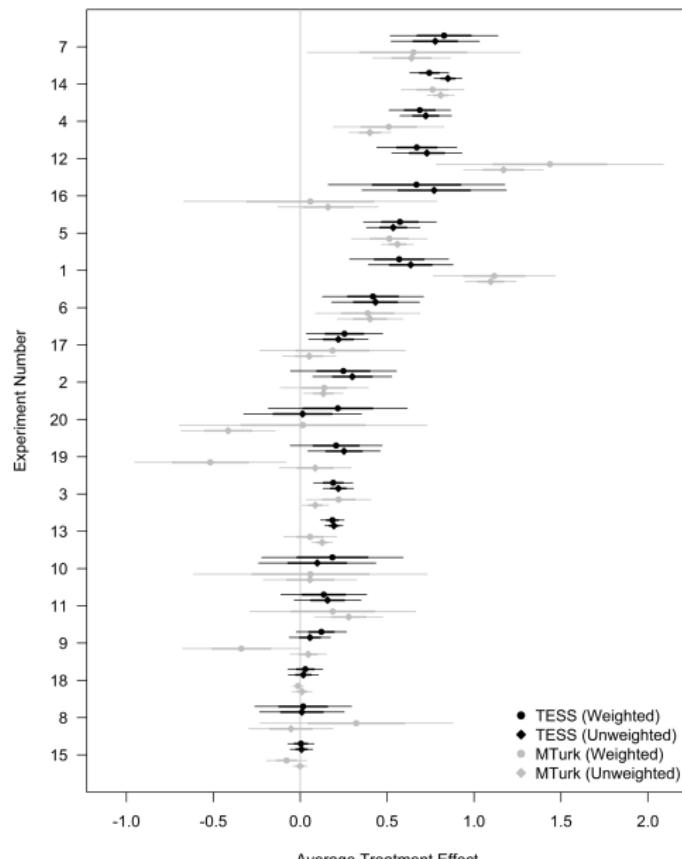
The Xbox Study



Wang et al. 2015. "Forecasting elections with non-representative polls."
International Journal of Forecasting.



Mullinix et al. 2015. "The Generalizability of Survey Experiments." *Journal of Experimental Political Science*.



Propensity Score Approach

- 1 Define a target population
- 2 Estimate a propensity score model
 - Pool experimental samples and target population units
 - Predict membership of all target and sample units in the experimental sample
- 3 Using fitted logits, divide population & sample into strata
- 4 Estimate stratum-specific ATE
- 5 Calculate weighted average of stratum-level estimates

Propensity Score Approach

Target population average treatment effect:

$$\sum_{v=1}^5 p(v) T(v) \quad (1)$$

where $p(v)$ is the proportion of the target population in a given stratum, v , and $T(v)$ is the estimated effect from stratum v of the experimental sample

Propensity Score Approach

Effect variance:

$$\sum_{v=1}^5 p(v)^2 V(v), \quad (2)$$

where $V(v)$ is the variance of the estimated experimental sample effect for stratum v

Propensity Score Subclassification Estimator

| Stratum | Weights | | | Estimates | | |
|---------|---------|--------|-------------|--------------|--------------|-------------|
| | Nat'l | Sample | Loan | DREAM 1 | DREAM 2 | Rally |
| 1 | 0.20 | 0.83 | 0.94 (0.08) | 0.06 (0.11) | -0.22 (0.12) | 0.74 (0.10) |
| 2 | 0.20 | 0.11 | 0.99 (0.26) | 0.22 (0.37) | -0.28 (0.36) | 0.77 (0.29) |
| 3 | 0.20 | 0.04 | 1.28 (0.43) | -0.61 (0.58) | -1.76 (0.54) | 1.00 (0.45) |
| 4 | 0.20 | 0.01 | 1.99 (0.73) | 0.29 (1.12) | 0.56 (0.89) | 1.44 (0.79) |
| 5 | 0.20 | 0.00 | | | | |
| Sample | - | - | 1.04 (0.30) | -0.01 (0.44) | -0.34 (0.38) | 0.79 (0.33) |
| Nat'l | - | - | 1.14 (0.18) | 0.02 (0.22) | -0.94 (0.23) | 0.94 (0.19) |

So does reweighting solve everything forever?

- Need well-defined target population
 - and detailed covariate data
 - and large stratum sizes
- Purely model-based, so only as good as the model
 - What unobservables might there be?
 - What reweighting might worse bias?
- Non-coverage is a potential problem
- Not well-tested on experimental data

Questions?

1 External Validity of a Sample

- Design-based
- Model-based

2 Other Notions of External Validity

- Settings
- Unit
- Treatments
- Outcomes

3 Participant Recruitment

SUTO Framework

- Cronbach (1986) talks about generalizability in terms of UTO
- Shadish, Cook, and Campbell (2001) speak similarly of:
 - **S**ettings
 - **U**nits
 - **T**reatments
 - **O**utcomes
- External validity depends on all of these

Population

- Setting
- Units
- Treatments
- Outcomes

Your Study

- Setting
- Units
- Treatments
- Outcomes

In your study, how do these correspond?
how do these differ?
do these differences matter?

Heterogeneity due to Settings

- We should expect heterogeneity related to settings!
- How do we use/explore this?
 - Comparative research designs where experiments provide measures for each case
 - Over-time replications of the same design
 - Replication of a design across contexts with unknown sources of variability?
- Can we control for context?

Pretreatment Dynamics

"If the experiment explores a communication that regularly occurs in 'reality,' then reactions in the experiment might be contaminated by those 'regular' occurrences prior to the experiment."¹

¹p.875 from Druckman & Leeper. 2012. "Learning More from Political Communication Experiments: Pretreatment and Its Effects." *American Journal of Political Science* 56(4): 875–896.

Pretreatment Dynamics

- Pretreatment is a feature of an experimental setting, treatment, and sample, wherein the effect of the treatment has already occurred²
- Consequences:
 - Biased effect estimates
- Mitigation:
 - Measure pretreatment
 - Avoid “pretreated” treatments or contexts
 - Study units not already treated
 - Theorize repeated effects

²Or, units having already been treated are otherwise affected differently.

Questions?

Heterogeneity due to Units

Most commonly studied source of heterogeneity is covariate-related (i.e., characteristics of units).

If we think there might be covariate-related effect heterogeneity, what can we do?

- Best solution: manipulate the moderator
- Next best: block on the moderator
- Least best: post-hoc exploratory approaches

Block Randomization

- Basic idea: randomization occurs within strata defined before treatment assignment
- CATE is estimate for each stratum; aggregated to SATE
- But...
 - Blocked randomization only works in exactly the same situations where stratified sampling works
 - Need to observe covariates pre-treatment in order to block on them, so works in panels but not cross-sectional designs
 - More precise SATE estimate

Questions?

Three Post-hoc Approaches

- Suggestive evidence
- Regression using treatment-by-covariate interactions
- Automated approaches
- (Replication and meta-analysis)

Suggestive Evidence

We can never know $\text{Var}(TE_i)$! But...

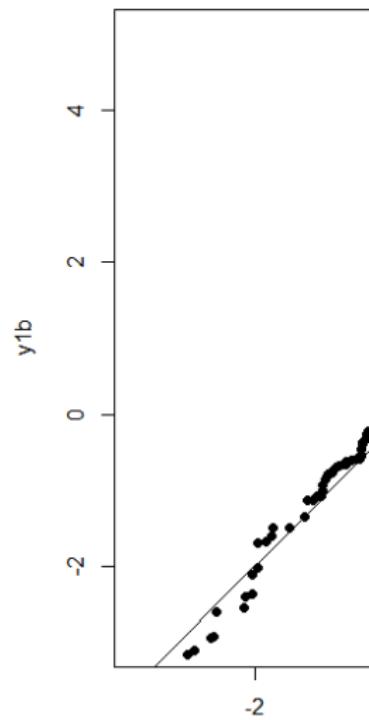
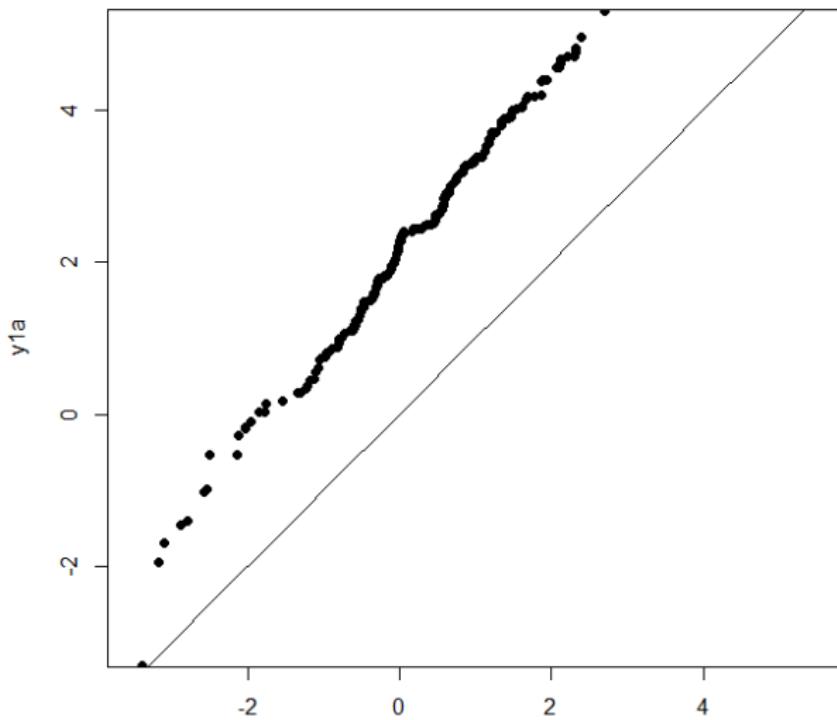
- Quantile-quantile plots
 - Compare the distribution of Y_0 's to distribution of Y_1 's
 - If homogeneity, a vertical shift in Y_1 's
 - If heterogeneity, a slope $\neq 1$
- Equality of variance tests
 - If homogeneity, variance should be equal
 - If heterogeneity, variances should differ

QQ Plots

```
# y_0 data
set.seed(1)
n <- 200
y0 <- rnorm(n) + rnorm(n, 0.2)

# y_1 data (homogeneous effects)
y1a <- y0 + 2 + rnorm(n, 0.2)
# y_1 data (heterogeneous effects)
y1b <- y0 + rep(0:1, each = n/2) + rnorm(n, 0.2)

qqplot(y0, y1a, pch=19, xlim=c(-3,5), ylim=c(-3,5), asp=1)
curve((x), add = TRUE)
qqplot(y0, y1b, pch=19, xlim=c(-3,5), ylim=c(-3,5), asp=1)
curve((x), add = TRUE)
```



Equality of Variance tests

```
> var.test(y0, y1a)

  F test to compare two variances

data: y0 and y1a
F = 0.60121, num df = 199, denom df = 199,
  p-value = 0.0003635
alternative hypothesis:
  true ratio of variances is not equal to 1
95 percent confidence interval:
  0.4549900 0.7944289
sample estimates:
ratio of variances
  0.6012131
```

Equality of Variance tests

```
> var.test(y0, y1b)

  F test to compare two variances

data: y0 and y1b
F = 0.53483, num df = 199, denom df = 199,
  p-value = 1.224e-05
alternative hypothesis:
  true ratio of variances is not equal to 1
95 percent confidence interval:
  0.4047531 0.7067133
sample estimates:
ratio of variances
  0.5348312
```

Questions?

Regression Estimation

Aside: Regression Adjustment in Experiments, Generally

- Recall the general advice that we do not need covariates in the regression to “control” for omitted variables (because there are none)
- Including covariates can reduce variance of our SATE by explaining more of the variation in Y

Scenario

Imagine two regression models. Which is correct?

- 1 Mean-difference estimate of SATE is “not significant”
- 2 Regression estimate of SATE, controlling for sex, age, and education, is “significant”

This is a small-sample dynamic, so make these decisions pre-analysis!

Treatment-Covariate Interactions

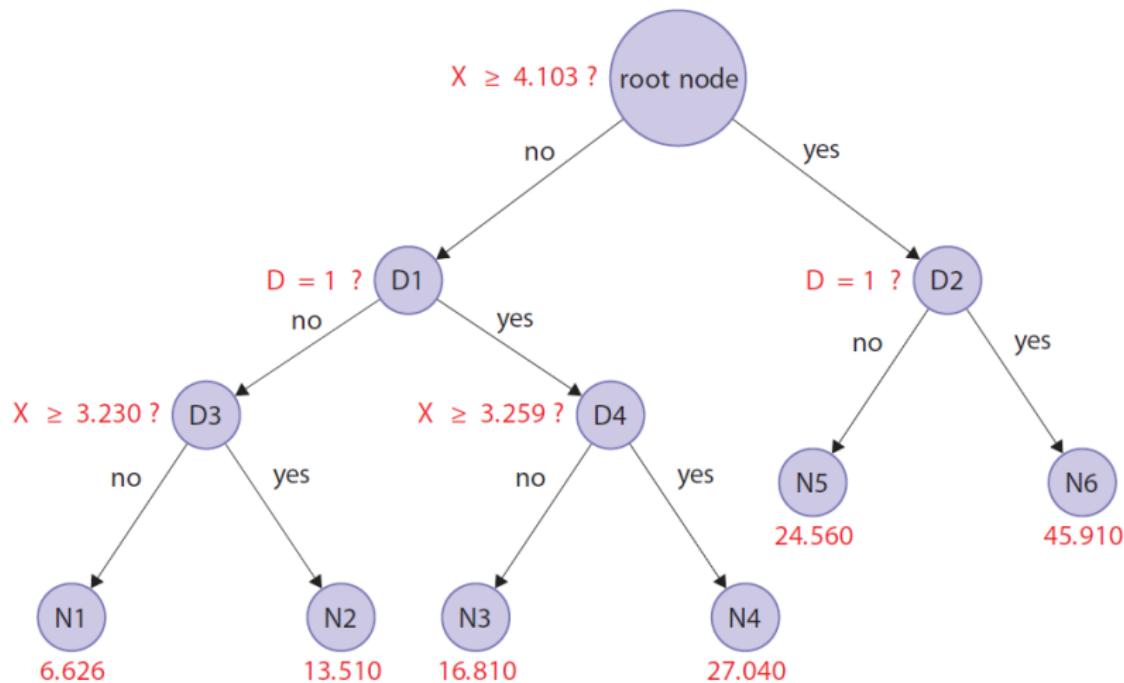
- The regression paradigm allows us to estimate CATEs using interaction terms
 - X is an indicator for treatment
 - M is an indicator for possible moderator
- SATE: $Y = \beta_0 + \beta_1 X + e$
- CATEs:

$$Y = \beta_0 + \beta_1 X + \beta_2 M + \beta_3 X * M + e$$

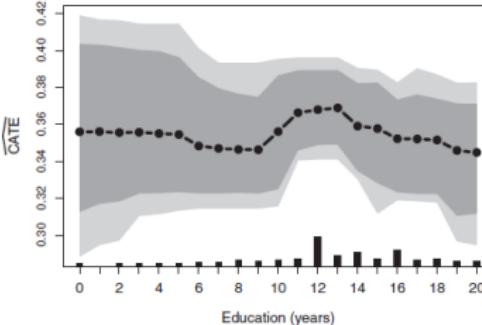
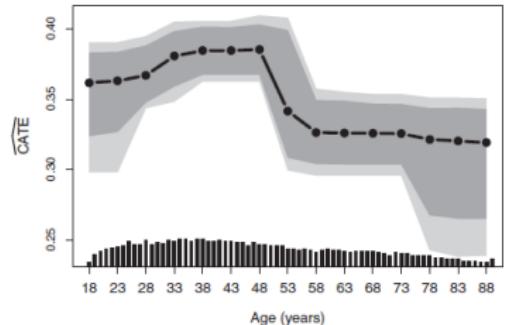
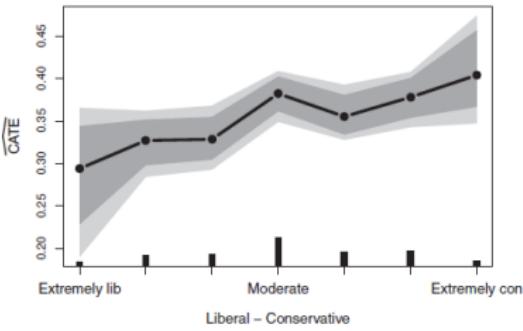
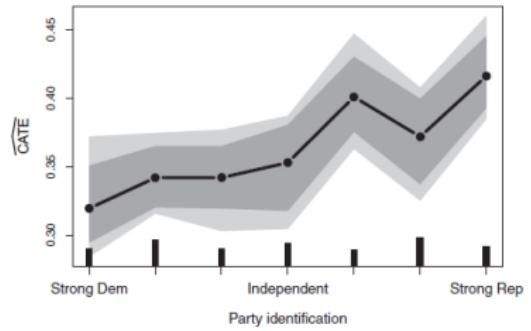
- Homogeneity: $\beta_3 = 0$
- Heterogeneity: $\beta_3 \neq 0$

BART

- Estimate CATEs in a fully automated fashion
- “Bayesian Additive Regression Trees”
 - Essentially an ensemble machine learning method
- Iteratively split a sample into more and more homogeneous groups until some threshold is reached using binary (cutpoint) decisions
- Repeat this a bunch of times, aggregating across results



Green & Kern. 2012. "Modeling Heterogeneous Treatment Effects in Survey Experiments with Bayesian Additive Regression Trees." *Public Opinion Quarterly* 76(3): 491–511.



Considerations

- BART is totally automated, conditional on the set of covariates used
- Only really works with dichotomous covariates
- Not widely used or tested
- Totally post-hoc and atheoretical

Considerations

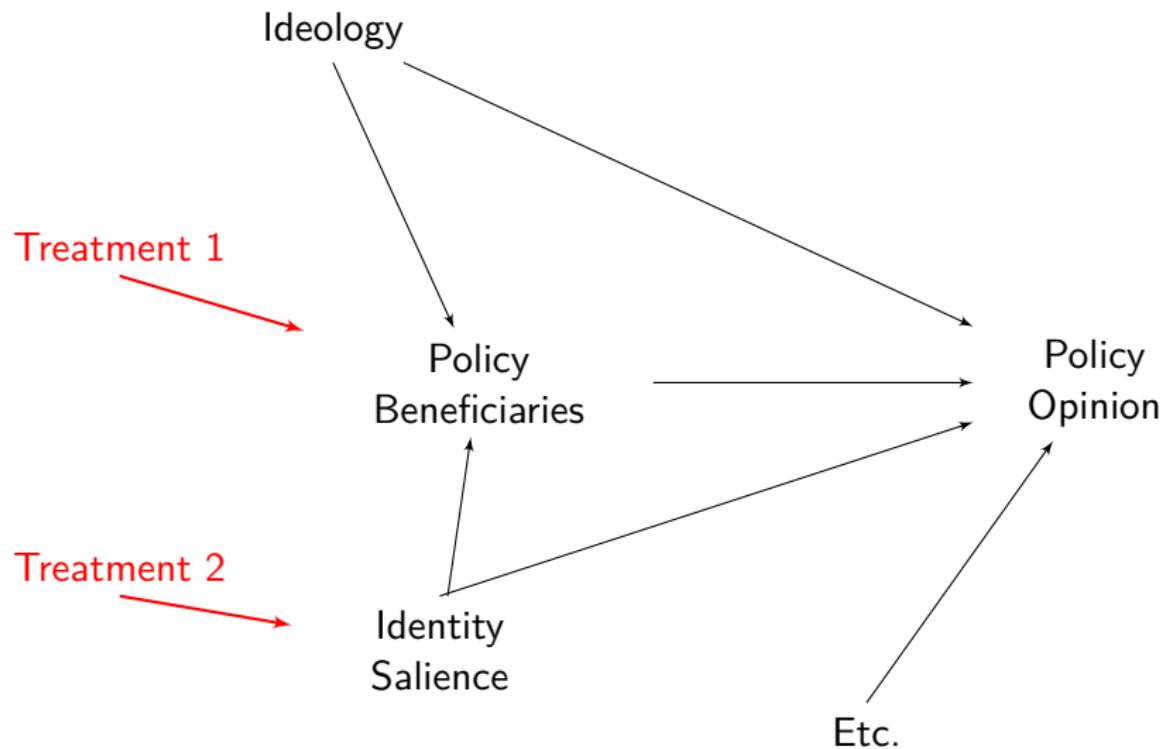
- Coefficients on moderators have no causal interpretation without further conditioning on observables
- Nearly unlimited potential moderators
 - First-order interactions with every covariate in dataset
 - Second-, third-order, etc. interactions
- Thus, multiple comparisons problem!
- Power (esp. if M is continuous)

Simply: Manipulating the moderator variable is the best way to estimate a heterogeneous effect!

Why is this true?

Complex Designs

- An experiment can have any number of conditions
 - Up to the limits of sample size
 - More than 8–10 conditions is typically unwieldy
- Typically analyze complex designs using ANOVA or regression, but we are still ultimately interested in pairwise comparisons to estimates SATEs
 - Treatment–treatment, or treatment-control
 - Without control group, we don't know which treatment(s) affected the outcome



Ex. Question-as-treatment³

- How close do you feel to your ethnic or racial group? How close do you feel to other Americans?
- Some people have said that taxes need to be raised to take care of pressing national needs. How willing would you be to have your taxes raised to improve education in public schools? Some people have said that taxes need to be raised to take care of pressing national needs. How willing would you be to have your taxes raised to improve educational opportunities for minorities?

³Transue. 2007. "Identity Salience, Identity Acceptance, and Racial Policy Attitudes: American National Identity as a Uniting Force." *American Journal of Political Science* 51(1): 78–91.

2x2 Factorial Design

| Condition | |
|----------------------|-------|
| Educ. for Minorities | Y_1 |
| Schools | Y_0 |

| Condition | Americans | Own Race |
|----------------------|-----------|-----------|
| Educ. for Minorities | $Y_{1,0}$ | $Y_{1,1}$ |
| Schools | $Y_{0,0}$ | $Y_{0,1}$ |

Two ways to *parameterize* this

Dummy variable regression (i.e., treatment-control CATEs):

$$Y = \beta_0 + \beta_1 X_{0,1} + \beta_2 X_{1,0} + \beta_3 X_{1,1} + \epsilon$$

Interaction effects (i.e., treatment-treatment CATEs):

$$Y = \beta_0 + \beta_1 X_{1,1} + \beta_2 X_{2,1} + \beta_3 X_{1,1} * X_{2,1} + \epsilon$$

Use `margins` to extract marginal effects

Considerations

- Need to have hypotheses about heterogeneity a priori
- Factorial designs can quickly become unwieldy and expensive
- Need to consider what CATEs are of theoretical interest
 - Treatment-control
 - Treatment-treatment
 - Marginal effects, averaging across other factors

Probably obvious, but . . .

| Factors | Conditions per factor | Total Conditions | <i>n</i> |
|---------|-----------------------|------------------|----------|
| 1 | 2 | 2 | 400 |
| 1 | 3 | 3 | 600 |
| 1 | 4 | 4 | 800 |
| 2 | 2 | 4 | 800 |
| 2 | 3 | 6 | 1200 |
| 2 | 4 | 8 | 1600 |
| 3 | 3 | 9 | 1800 |
| 3 | 4 | 12 | 2400 |
| 4 | 4 | 16 | 3200 |

Assumes power to detect a relatively small effect, but no consideration of multiple comparisons.

Considerations

- Need to have hypotheses about heterogeneity a priori
- Factorial designs can quickly become unwieldy and expensive
- Need to consider what CATEs are of theoretical interest
 - Treatment-control
 - Treatment-treatment
 - Marginal effects, averaging across other factors

Questions?

Treatment Preferences/Self-Selection

Bennett and Iyengar:⁴

manipulational control actually weakens the ability to generalize to the real world where exposure to stimuli is typically voluntary. Accordingly, it is important that experimental researchers use designs that combine manipulation with self-selection of exposure.

⁴p.724 from Bennett & Iyengar. 2008. "A new era of minimal effects? The changing foundations of political communication." *Journal of Communication* 58(4): 707-31.

Hovland: ⁵

It should be possible to assess what demographic and personality factors predispose one to expose oneself to particular communications and then to utilize experimental and control groups having these characteristics. Under some circumstances the evaluation could be made on only those who select themselves, with both experimental and control groups coming from the self-selected audience.

⁵p.16 from Hovland. 1959. "Reconciling conflicting results derived from experimental and survey studies of attitude change." *American Psychologist* 14(1): 8-17.

Treatment Preferences I

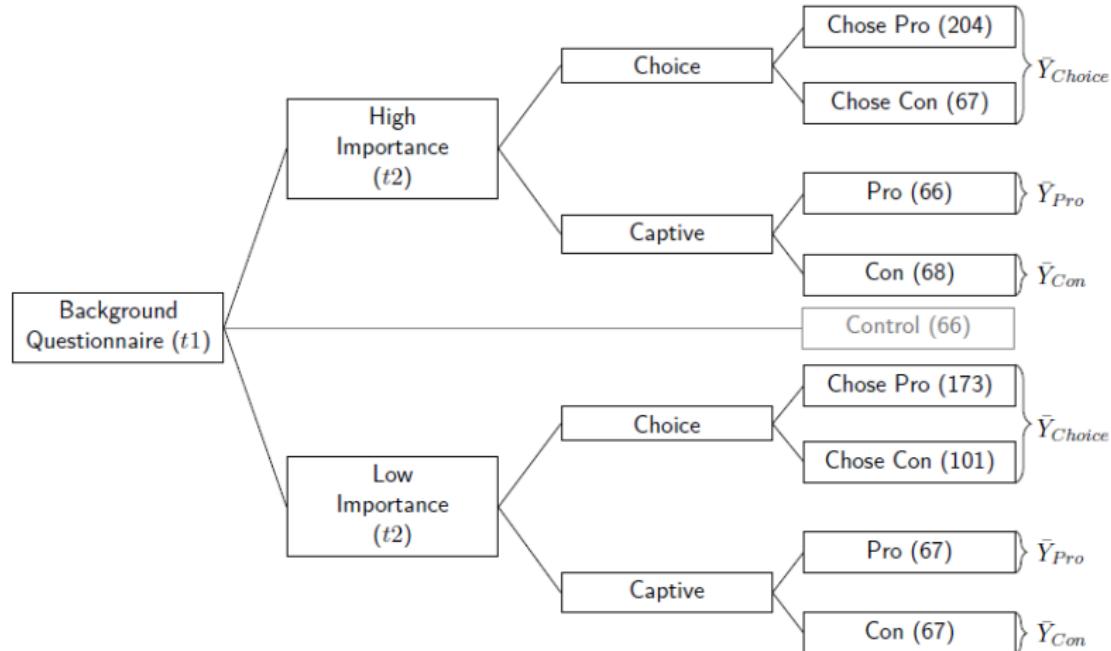
- Experiments are about inferring effect of X on Y
- Respondents may have preferences over whether they are treated or untreated⁶
- Origins of this discussion are in the medical literature⁷
- Closely related to the notion of placebo effects

⁶Rucker. 1989. "A Two-Stage Trial Design for Testing Treatment, Self-Selection, and Treatment Preference Effects." *Statistics in Medicine* 8: 477–485.

⁷Swift & Callahan. 2009. "The Impact of Client Treatment Preferences on Outcome: A Meta-Analysis." *Journal of Clinical Psychology* 65(4): 368–381.

Treatment Preferences I

- Treatment preferences may be an important factor in:
 - Compliance
 - Effect heterogeneity
- Depending on your treatments, you may want to measure preferences
 - 1 Stated preference measures
 - 2 Designs that reveal preferences



Analyzing 3-Group Preference Trials⁸

1 SATE: $\bar{Y}_T - \bar{Y}_C$

2 CATE (Prefer T): $\frac{\bar{Y}_{Choice} - \bar{Y}_C}{\hat{\alpha}}$

3 CATE (Prefer C): $\frac{\bar{Y}_T - \bar{Y}_{Choice}}{1 - \hat{\alpha}}$

Note: $\alpha = Pr(T|Choice)$

⁸GK2011 Package for R. <https://cran.r-project.org/package=GK2011>

Questions?

Attention and Satisficing

One final issue with unit-related sources of heterogeneity is how we handle or analyze survey-experimental data where we think participants “misbehaved”.

This falls into a couple of broad categories:

- 1 Noncompliance (discussed earlier)
- 2 Survey Satisficing
- 3 Apparent Inattention

Substantive Manipulation Check

- Two common approaches:
 - Information recall or understanding
 - Measure level of manipulated treatment variable
- Risky to remove cases based on this because it is a form of conditioning on post-treatment variables
- May be useful to consider either a mediator or effects

Attention Checking

- Online mode invites satisficing
- Attention checking can help, but is imperfect

Apparent Satisficing

- Filter out respondents based on response behavior
- Some common measures:
 - “Straightlining”
 - Non-differentiation
 - Acquiescence
 - Nonresponse
 - DK responding
 - Speeding
- Difficult to detect
- Difficult to distinguish from “real” responses

Metadata/Paradata

- Timing
 - Some survey tools will allow you to time page
 - Make a prior rules about dropping participants for speeding
- Mousetracking or eyetracking
 - Mousetracking is unobtrusive
 - Eyetracking requires participants opt-in
- Record focus/blur browser events

Direct Measures

- How closely have you been paying attention to what the questions on this survey actually mean?

- While taking this survey, did you engage in any of the following behaviors? Please check all that apply.
 - Use your mobile phone
 - Browse the internet
 - ...

Instructional Manipulation Check

Do you agree or disagree with the decision to send British forces to fight ISIL in Syria? We would like to know if you are reading the questions on this survey. If you are reading carefully, please ignore this question, do not select any answer below, and click “next” to proceed with the survey.

Strongly disagree

Somewhat disagree

Neither agree nor disagree

Somewhat agree

Strongly agree

Attention Checking

In summary...

- Attention checking can be useful
- Lots of options
- No obvious best metric
- Can be analytically consequential

How should we deal with respondents that appear to not be paying attention, not “taking” the treatment, or not responding to outcome measures?

- 1 Keep them
- 2 Throw them away

Best Practice: Protocol

- Excluding respondents based on survey behavior is one of the easiest ways to “p-hack” an experimental dataset
 - Inattention, satisficing, etc. will tend to reduce the size of the SATE
- So regardless of how you handle these respondents, these should be decisions that are made *pre-analysis*

When are you excluding participants?

Pre-Treatment

- Satisficing behaviors
- Inattention
- Covariate-based selection
- Pretreated

Post-Treatment

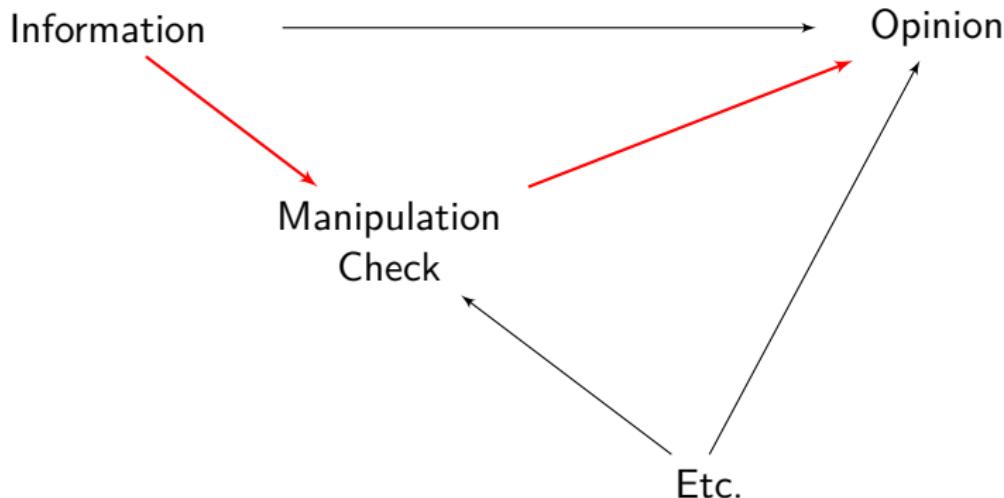
- Speeding on treatment
- “Failing” a manipulation check
- Drop-off

Pre-Treatment Exclusion

- This is totally fine from a causal inference perspective
- Advantages:
 - Focused on engaged respondents
 - Likely increase impact of treatment
- Disadvantages:
 - Changing definition of sample (and thus population)

Post-Treatment Exclusion

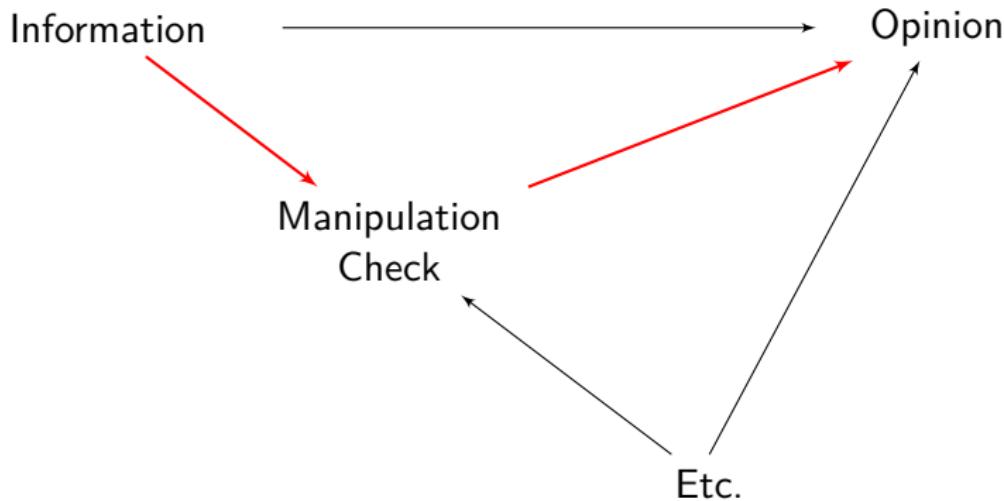
This is much more problematic because it involves controlling for a *post-treatment* variable



Risk that estimate of β_1 is diminished because effect is being carried through the manipulation check. Introduction of “collider bias” wherein values of the manipulation check are affected by other factors.

Post-Treatment Exclusion

- Any post-treatment exclusion is problematic and should be avoided
- Can estimate a LATE
 - Interpretation: Effect of manipulation check among those whose value of the check can be changed by the treatment manipulation
- Non-response or attrition is the same as researcher-imposed exclusion
 - Not problematic if MCAR
 - Nothing really to be done if caused by treatment



Risk that estimate of β_1 is diminished because effect is being carried through the manipulation check. Introduction of “collider bias” wherein values of the manipulation check are affected by other factors.

Post-Treatment Exclusion

- Any post-treatment exclusion is problematic and should be avoided
- Can estimate a LATE
 - Interpretation: Effect of manipulation check among those whose value of the check can be changed by the treatment manipulation
- Non-response or attrition is the same as researcher-imposed exclusion
 - Not problematic if MCAR
 - Nothing really to be done if caused by treatment

Protocol

Protocol is the complete planning document for how to design, implement, and analyze an experiment.⁹

⁹Thomas J. Leeper. 2011. "The Use of Protocol in the Design and Reporting of Experiments." *The Experimental Political Scientist*.

Protocol

1 Theory/hypotheses

2 Instrumentation

- Manipulation(s)
- Outcome(s)
- Covariate(s)
- Manipulation check(s)

3 Sampling

4 Implementation

5 Analysis

Why bother?

- Be clear to yourself what you're trying to do before you do it
- Assess the literature for best practices
- Highlight areas in need of pilot testing
- Economize questionnaire development
- Study preregistration

Questions?

Questions?

Heterogeneity due to Treatments

- We should expect this! Why?
- What can we do?
 - Pilot testing
 - Replication
 - More complex design
 - Conjoint experiments

Conjoint Designs I

- “Classic vignettes” taken to an extreme
 - Address heterogeneity w/r/t SUTO
- Example: Judge whether to admit an immigrant to your country
- Respondents see a series of vignettes that are fully randomized along any number of dimensions
 - Sex, Education, Language proficiency, etc.
- Outcome is judgment (binary or rating scale)

Conjoint Designs II

Why is this useful?

- Understand complex decision-making
- Within-subjects comparisons
- Heterogeneous effects across versions of treatment
- Pilot testing: Sensitivity of design to specification of *compound* vignette

Please read the descriptions of the potential immigrants carefully. Then, please indicate which of the two immigrants you would personally prefer to see admitted to the United States.

| | 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 |

Immigrant 1 Immigrant 2

If you had to choose between them, which of these two immigrants should be given priority to come to the United States to live?

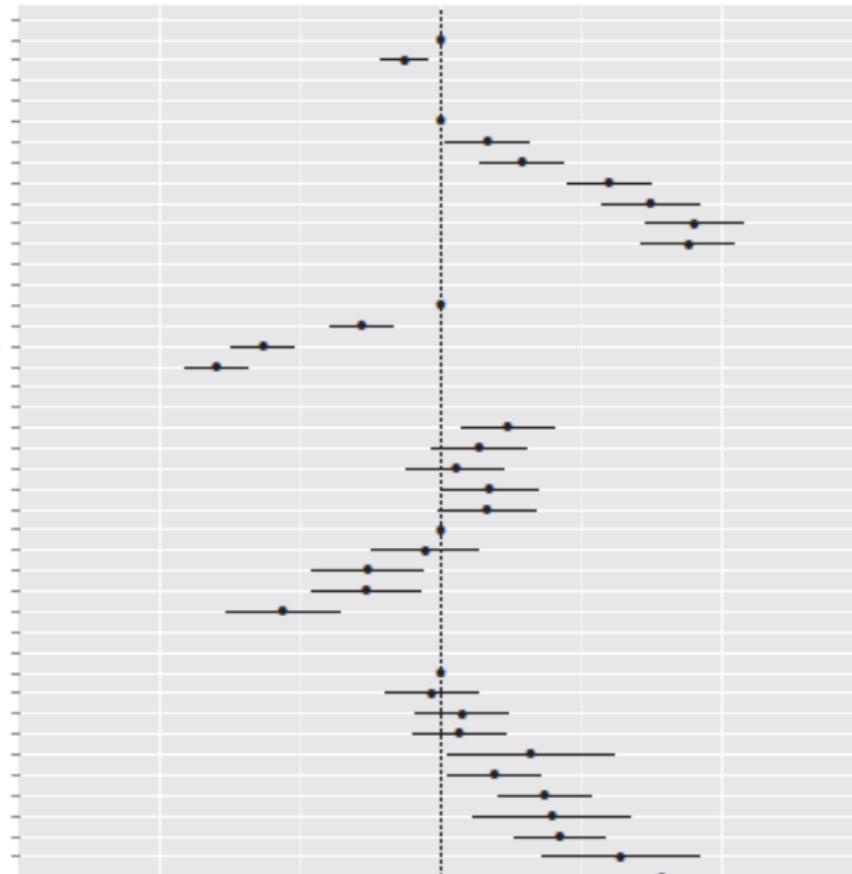
Gender:
female
male

Education:
no formal
4th grade
8th grade
high school
two-year college
college degree
graduate degree

Language:
fluent English
broken English
tried English but unable
used interpreter

Origin:
Germany
France
Mexico
Philippines
Poland
India
China
Sudan
Somalia
Iraq

Profession:
janitor
waiter
child care provider
gardener
financial analyst
construction worker
teacher
computer programmer
nurse
research scientist
doctor



Conjoint Designs III

- As long as profiles are randomized, this is just a complex factorial design where we can estimate *marginal effect* of each attribute
 - Treatment-control SATE, conditional on all other randomized factors
- Assumptions:
 - Fully randomized profiles
 - No “carry-over” effects
 - No profile order effects

Replication

- Conjoint solve one problem: they identify the relative size of sources of heterogeneity within a given treatment
- But how should we consider experiments testing the same theory using different treatments?
 - “Triangulation”
 - Consistent directionality
 - Consistent (standardized) effect sizes
- Big conclusion: replication is important and there's not enough of it.

Questions?

Heterogeneity due to *Outcomes*

- This is expected!
 - E.g., non-equivalent outcomes
- Reasonable to explore multiple outcomes
 - Multiple comparisons
 - Power considerations
 - Construct validity
- What outcomes you measure depend on your theory
- Lots of potential for behavioral measures!

Behavioural measures

Some behaviours that can be directly measured through survey questionnaires.

Three broad categories:

- 1 Behavioural measures that provide survey paradata
- 2 Behavioural measures that operationalize attitudes
- 3 Behavioural measures that operationalize behaviours

Behavioural Measures for Paradata

Why?

- Respondents use of the survey tells us something meaningful about their behaviour

What?

- Nonresponse
- Response latencies
- Reading times
- Answer switching
- Eye tracking
- Mouse tracking
- Smartphone metadata

Behavioural Measures for Attitudes

Why?

- Attitudinal self-reports might be “cheap talk”

What?

- Implicit Association Test
- Incentivized Survey questions

Behavioural Measures for Behaviour

Why?

- We want to observe or affect behaviour (e.g., in an experiment)

What?

- Directly measure or initiate a direct measure of a behaviour
- May be measured by something that occurs within the confines of the survey or something outside of the survey

Example 1: Active Information Choice

- “Followed link” identification¹⁰
- Information boards¹¹
- Video choice¹²
- Dynamic Process Tracing Environment¹³

¹⁰ Guess, AM. 2015. "Measure for Measure." *Political Analysis* 23: 59–75. doi:10.1093/pan/mpu010

¹¹ Leeper, TJ. 2014. "The Informational Basis for Mass Polarization." *Public Opinion Quarterly* 78(1): 27–46. doi:10.1093/poq/nft045

¹² Arceneaux, K & Johnson, M. 2012. *Changing Minds or Changign Channels*. Chicago: The University of Chicago Press.

¹³<https://dpte.polisci.uiowa.edu/dpte/>

Remember, please check **ALL** rows containing any links shown in **PURPLE**. Leave all other rows unchecked.

- [LINK](#) [LINK](#)
- [LINK](#) [LINK](#)
- [LINK](#) [LINK](#)
- [LINK](#) [LINK](#)
- [LINK](#) [LINK](#) [LINK](#)
- [LINK](#)
- [LINK](#) [LINK](#) [LINK](#)
- [LINK](#)
- [LINK](#)
- [LINK](#) [LINK](#)

Example 1: Active Information Choice

- “Followed link” identification¹⁰
- Information boards¹¹
- Video choice¹²
- Dynamic Process Tracing Environment¹³

¹⁰ Guess, AM. 2015. "Measure for Measure." *Political Analysis* 23: 59–75. doi:10.1093/pan/mpu010

¹¹ Leeper, TJ. 2014. "The Informational Basis for Mass Polarization." *Public Opinion Quarterly* 78(1): 27–46. doi:10.1093/poq/nft045

¹² Arceneaux, K & Johnson, M. 2012. *Changing Minds or Changign Channels*. Chicago: The University of Chicago Press.

¹³<https://dpte.polisci.uiowa.edu/dpte/>

Reports From the Hive,
Where the Swarm
Concurs

Pay for Performance
Improves Quality of
Health Care Through
Collaborative Medicine

Why are 3-D Movies so
Bad?

Physicians Group Says
Quality Will Improve
Under Outcome-based
Payments

Council Is Set to
Consider Increases in
Hotel and Property Taxes

Doctors Can Work
Together to Improve
Patient Health, But Need
Appropriate Incentives

Patients Better Served
When Providers Paid for
Health Outcomes

Improving America's
Health Requires Provider
Incentives, Not 'Fee-for-
Service'

When Paid for Outcomes,
Doctors Have Little
Reason to Treat Highest
Risk Patients

A Bowl of Chili with
Bragging Rights

SEC Vote Requires
Business Filings to Add
Environmental Risks to
Bottom Line

Anatomy of a Tear-
Jerker

Spammers Use the
Human Touch to Avoid
CAPTCHA

USDA Raises Corn
Export Outlook

Will a Standardized
System for Verifying
Web Identity Ever
Catch On?

Wellness, Rather
Than Illness, Is Focus
Under Outcome-
Accountable Care

Gender Differences in
Education Need
Innovative Solution

Heart Attack While
Dining at Heart Attack
Grill in Las Vegas

Out of the O.R., T.R.
Knight Back Onto the
Stage

Paying Doctors Based
on Outcomes Will
Lead to Rationing

Example 1: Active Information Choice

- “Followed link” identification¹⁰
- Information boards¹¹
- Video choice¹²
- Dynamic Process Tracing Environment¹³

¹⁰ Guess, AM. 2015. "Measure for Measure." *Political Analysis* 23: 59–75. doi:10.1093/pan/mpu010

¹¹ Leeper, TJ. 2014. "The Informational Basis for Mass Polarization." *Public Opinion Quarterly* 78(1): 27–46. doi:10.1093/poq/nft045

¹² Arceneaux, K & Johnson, M. 2012. *Changing Minds or Changign Channels*. Chicago: The University of Chicago Press.

¹³<https://dpte.polisci.uiowa.edu/dpte/>

Stage: Primary Election

Sub-stage: Early Primary

Time Remaining: 21:26

6:46

Andy Fischer's Political Experience

DELEGATE COUNT, END OF FEBRUARY

Republican Primary

Sam Green's Mother provides a Childhood Anecdote

Dana Turner's Picture

Terry Davis's Current Job Performance

Taylor Harris's Age

Iowa General Election

Hillary Clinton wins in South Dakota!



Stage: Pre-Election

Sub-stage: PE-2

Time Remaining: 0:00

0:00

Question 1 of 1

Primary elections require voters to choose the party they want to vote in. Before we begin the Iowa primary, please choose either the the Republican or Democrat Primary. You will see candidates for both parties but will be only able to vote in the party you choose.

- Republican
- Democrat

Select an answer, then click the End button to end the questionnaire.

End

Example 2: Sign-up/Enrolment

An extension of information choice behaviour would be explicit engagement in other kinds of (small) behaviours, such as:

- Entering an email address to receive information or join a mailing list^{14 15}
- Signing up for an appointment or further interaction

¹⁴ Leeper, TJ. 2017. "How Does Treatment Self-Selection Affect Inferences About Political Communication?" *Journal of Experimental Political Science*: In press.

¹⁵ Bolsen, Druckman, & Cook. 2014. "Communication and Collective Actions." *Journal of Experimental Political Science* 1(1): 24–38. doi:10.1017/xps.2014.2

Example 3: Incentivised Survey Questions

Definitions:

- A survey question is just a self-report
- An *incentivized* survey question attached financial gains or losses to the answer options

Paradigm could be applied to any measure of behavioural intentions to avoid cheap talk.

Mark your gamble selection with an **X** in the last column across from your preferred gamble.

| Gamble | Event | Payoff | Probabilities | Your Selection |
|--------|-------|--------|---------------|----------------|
| 1 | A | \$10 | 50% | |
| | B | \$10 | 50% | |
| 2 | A | \$18 | 50% | |
| | B | \$6 | 50% | |
| 3 | A | \$26 | 50% | |
| | B | \$2 | 50% | |
| 4 | A | \$34 | 50% | |
| | B | -\$2 | 50% | |
| 5 | A | \$42 | 50% | |
| | B | -\$6 | 50% | |

Example 3: Incentivised Survey Questions

Definitions:

- A survey question is just a self-report
- An *incentivized* survey question attached financial gains or losses to the answer options

Paradigm could be applied to any measure of behavioural intentions to avoid cheap talk.

Example 4: Purchasing Decisions

Common ways to study purchasing behaviour include:

- Direct attitudinal questions
- Retrospective and prospective self-reports
- Conjoint experiments

Another way is embedding a purchase in a survey.¹⁶

¹⁶Bolsen, T. 2011. "A Lightbulb Goes On." *Political Behavior* 35(1): 1–20. 10.1007/s11109-011-9186-5



Example 5: Donations

- Miller and Krosnick¹⁷ asked for charitable donations via cheque directly as part of a paper-and-pencil survey
- Klar and Piston¹⁸ offered respondents a survey incentive up-front for participation and then later offered them a chance to donate (a portion of payment) to a charity

¹⁷ Miller, Krosnick, & Lowe. N.d. "The Impact of Policy Change Threat on Financial Contributions to Interest Groups." Working paper.

¹⁸ Klar & Piston. 2015. "The influence of competing organisational appeals on individual donations." *Journal of Public Policy* 35(2): 171–91. doi:10.1017/S0143814X15000203

Example 6: Web Tracking Data

- 1 Active installation of a tracking app, such as YouGov Pulse¹⁹
²⁰
- 2 Post-hoc collection of web history files using something like
Web Historian²¹

¹⁹<https://yougov.co.uk/find-solutions/profiles/pulse/>

²⁰Guess, AM. N.d. "Media Choice and Moderation." Working paper,
<https://dl.dropboxusercontent.com/u/663930/GuessJMP.pdf>.

²¹<http://www.webhistorian.org/>

Other Possibilities

- Coordination tasks
 - Synchronous group tasks²²
 - Game play
 - Simulations
- Offering incentives to perform future behaviour (tracked elsewhere)
- OAuth/API integrations w/ other platforms
 - Merging website usage data w/ survey data
 - Treating website sign-up or usage as behavioural outcomes
 - Linking with smartphone metadata

²²Mao, Mason, Suri, Watts. 2016. "An Experimental Study of Team Size and Performance on a Complex Task." *PLoS ONE* 11(4): e0153048. doi:10.1371/journal.pone.0153048

Other Possibilities

- Coordination tasks
 - Synchronous group tasks²²
 - Game play
 - Simulations
- Offering incentives to perform future behaviour (tracked elsewhere)
- OAuth/API integrations w/ other platforms
 - Merging website usage data w/ survey data
 - Treating website sign-up or usage as behavioural outcomes
 - Linking with smartphone metadata

²²Mao, Mason, Suri, Watts. 2016. "An Experimental Study of Team Size and Performance on a Complex Task." *PLoS ONE* 11(4): e0153048. doi:10.1371/journal.pone.0153048

Some principles for survey measures of behaviour

- 1 Know why you are collecting a behavioural measure!
- 2 Know whether you are studying a past, present, or future behaviour.
- 3 Be creative! Recognise possibilities and limitations of any given survey mode.
- 4 Validate, validate, validate!

Activity!

With a partner, brainstorm how one or more these behavioural measures might be applied to a survey experiment (either as outcome, treatment, covariate, or behavioural check) relevant to your own work or your organisation.

SUTO Punchline 1: Replication!

- If we think effects are homogeneous (across SUTO), then replications in other SUTO conditions should provide us the same SATE (within sampling error)
- If we think effects are heterogeneous, then replications should give *systematically* different SATE (or CATE) estimates
 - Identify those patterns of heterogeneity using meta-analysis
 - Regress effect estimates from multiple studies on SUTO features of each study

SUTO Punchline 2: **What do you want to know?**

- Do we want to know SATE, CATE(s), or both?
- Decide in advance
 - Include in protocol
 - Design study to estimate CATE(s)
- Estimation of unit-related CATEs
 - Block randomization
 - Post-hoc procedures

Questions?

1 External Validity of a Sample

- Design-based
- Model-based

2 Other Notions of External Validity

- Settings
- Unit
- Treatments
- Outcomes

3 Participant Recruitment

Recruitment Considerations

- Recruitment
 - Sampling
 - Opt-in
 - A mix of each
- Incentives
- Frequency of participation
 - MTurk panelists do 100+ studies per month
 - YouGov panelists do nearly as many
- “Profile” variables
- Quotas, post-stratification, weighting
- Respondent “quality”

Professional Panels

- Big players: SSI, YouGov, GfK, TNS/Gallup
- Online panels of respondents
- Respondents participate for incentives
- Study costs are negotiated
 - Sample size
 - Study length (number of survey items)
 - Targeting
 - Timing

Opt-in (Crowdsourcing) Sites

- Not exactly a panel (fully opt-in)
- Incentivized participation
- Prominent examples
 - MTurk
 - Crowdflower
 - Microworkers
 - Prolific Academic
 - Google Surveys

“River Sampling”

- Not using an existing subject pool
 - Link sharing or posting on websites
 - Using email list
 - Online advertising (Google, Facebook)
- My advice: don't do this unless you have no other choice!

Custom Panels

- Creating your own panel is great
 - Carefully sample on specific characteristics
 - Organize repeated interviewing or interaction
- Lots of additional issues
 - Attrition
 - Compensation
 - Panel Conditioning
- See Callegaro et al. 2014. *Online Panel Research: A Data Quality Perspective*. Wiley.

My Advice, Elaborated

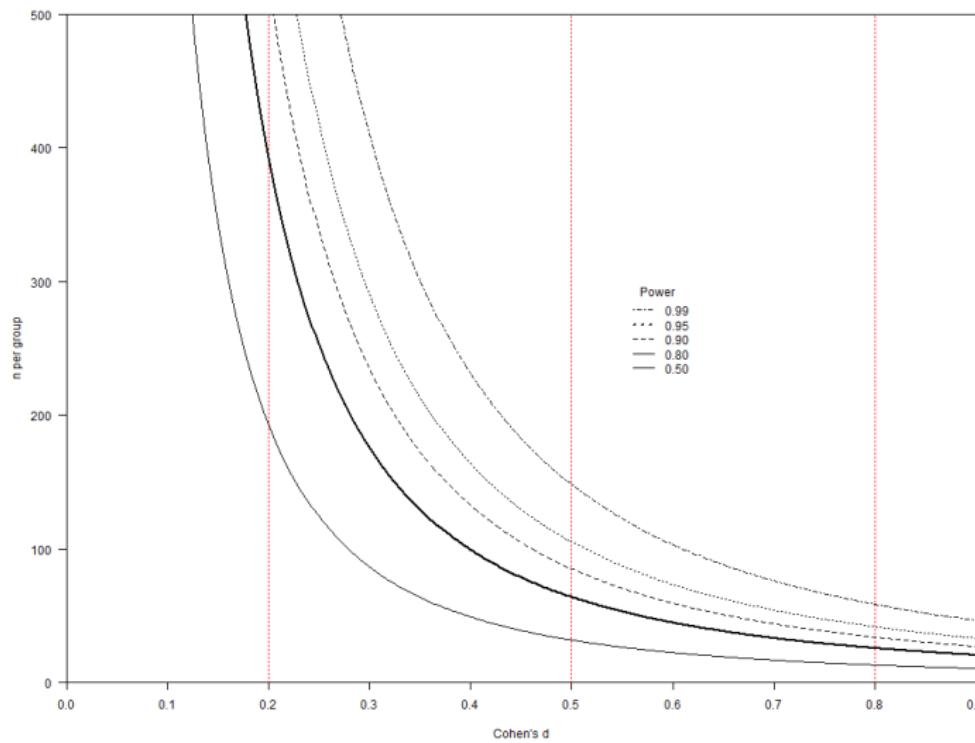
- Only work with populations where each unit is uniquely identifiable
- Without this, you risk many things:
 - Ambiguous eligibility
 - Retakes, treatment crossover
 - No way to evaluate response rates/bias
- Know something about your sample
 - How does it differ from your target of inference?
 - What theories or evidence would suggest those differences should matter?
 - What can you do to adjust or control for those *consequential* differences?

Measure, Measure, Measure

The only way to evaluate a sample is to know something about it.

The best way to convince reviewers is to rule out irrelevancies.

Don't forget statistical power...



And don't forget costs, either!

From one of my studies:

| Sample | Cost | n | Cost/participant |
|-----------|---------|------|------------------|
| National | \$13200 | 593 | \$22.26 |
| Exit Poll | \$3000 | 741 | \$4.05 |
| Students | \$0 | 299 | \$0 |
| Staff | \$1280 | 128 | \$10.00 |
| MTurk | \$550 | 1024 | \$0.54 |
| Ads | \$636 | 80 | \$7.95 |

External Validity

SUTO

Recruitment

