Unpacking the Black Box of Causality: Learning about Causal Mechanisms from Experimental and Observational Studies

### Kosuke Imai Princeton University

Talk at the School of Communication Research University of Amsterdam

June 10, 2015

Joint work with Luke Keele Dustin Tingley Teppei Yamamoto

# Project References (click the article titles)

#### • General:

- Unpacking the Black Box of Causality: Learning about Causal Mechanisms from Experimental and Observational Studies. *American Political Science Review*
- Theory:
  - Identification, Inference, and Sensitivity Analysis for Causal Mediation Effects. *Statistical Science*
- Extensions:
  - A General Approach to Causal Mediation Analysis. *Psychological Methods*
  - Experimental Designs for Identifying Causal Mechanisms. *Journal* of the Royal Statistical Society, Series A (with discussions)
  - Identification and Sensitivity Analysis for Multiple Causal Mechanisms: Revisiting Evidence from Framing Experiments. *Political Analysis*

#### Software:

• mediation: R Package for Causal Mediation Analysis. *Journal of Statistical Software* 

Kosuke Imai (Princeton)

# Identification of Causal Mechanisms

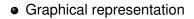
- Causal inference is a central goal of scientific research
- Scientists care about causal mechanisms, not just about causal effects
- Randomized experiments often only determine whether the treatment causes changes in the outcome
- Not how and why the treatment affects the outcome
- Common criticism of experiments and statistics:

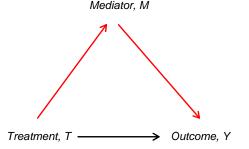
black box view of causality

• Question: How can we learn about causal mechanisms from experimental and observational studies?

Present a general framework for statistical analysis and research design strategies to understand causal mechanisms

- Show that the sequential ignorability assumption is required to identify mechanisms even in experiments
- Offer a flexible estimation strategy under this assumption
- Introduce a sensitivity analysis to probe this assumption
- Illustrate how to use statistical software mediation
- Solution Consider research designs that relax sequential ignorability



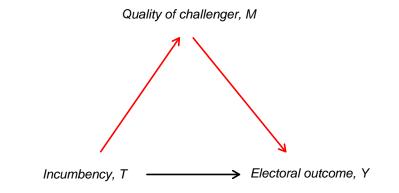


- · Goal is to decompose total effect into direct and indirect effects
- Alternative approach: decompose the treatment into different components
- Causal mediation analysis as quantitative process tracing

# Decomposition of Incumbency Advantage

- Incumbency effects: one of the most studied topics in American politics
- Consensus emerged in 1980s: incumbency advantage is positive and growing in magnitude
- New direction in 1990s: Where does incumbency advantage come from?
- Scare-off/quality effect (Cox and Katz): the ability of incumbents to deter high-quality challengers from entering the race
- Alternative causal mechanisms: name recognition, campaign spending, personal vote, television, etc.

# Causal Mediation Analysis in Cox and Katz

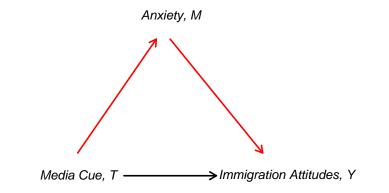


- How much of incumbency advantage can be explained by scare-off/quality effect?
- How large is the mediation effect relative to the total effect?

# Psychological Study of Media Effects

- Large literature on how media influences public opinion
- A media framing experiment of Brader et al.:
  - (White) Subjects read a mock news story about immigration:
    - Treatment: Hispanic immigrant in the story
    - Control: European immigrant in the story
  - Measure attitudinal and behavioral outcome variables:
    - Opinions about increasing or decrease immigration
    - Contact legislator about the issue
    - Send anti-immigration message to legislator
- Why is group-based media framing effective?: role of emotion
- Hypothesis: Hispanic immigrant increases anxiety, leading to greater opposition to immigration
- The primary goal is to examine how, not whether, media framing shapes public opinion

# Causal Mediation Analysis in Brader et al.



- Does the media framing shape public opinion by making people anxious?
- An alternative causal mechanism: change in beliefs
- Can we identify mediation effects from randomized experiments?

# The Standard Estimation Method

• Linear models for mediator and outcome:

$$Y_{i} = \alpha_{1} + \beta_{1} T_{i} + \xi_{1}^{\top} X_{i} + \epsilon_{1i}$$
  

$$M_{i} = \alpha_{2} + \beta_{2} T_{i} + \xi_{2}^{\top} X_{i} + \epsilon_{2i}$$
  

$$Y_{i} = \alpha_{3} + \beta_{3} T_{i} + \gamma M_{i} + \xi_{3}^{\top} X_{i} + \epsilon_{3i}$$

where  $X_i$  is a set of pre-treatment or control variables

- 1 Total effect (ATE) is  $\beta_1$
- 2 Direct effect is  $\beta_3$
- 3 Indirect or mediation effect is  $\beta_2 \gamma$
- Effect decomposition:  $\beta_1 = \beta_3 + \beta_2 \gamma$ .
- Some motivating questions:
  - What should we do when we have interaction or nonlinear terms?
  - What about other models such as logit?
    - In general, under what conditions can we interpret  $\beta_1$  and  $\beta_2 \gamma$  as causal effects?
  - What do we really mean by causal mediation effect anyway?

- Observed data:
  - Binary treatment:  $T_i \in \{0, 1\}$
  - Mediator:  $M_i \in \mathcal{M}$
  - Outcome:  $Y_i \in \mathcal{Y}$
  - Observed pre-treatment covariates:  $X_i \in \mathcal{X}$
- Potential outcomes model (Neyman, Rubin):
  - Potential mediators:  $M_i(t)$  where  $M_i = M_i(T_i)$
  - Potential outcomes:  $Y_i(t, m)$  where  $Y_i = Y_i(T_i, M_i(T_i))$
- Total causal effect:

$$\tau_i \equiv Y_i(1, M_i(1)) - Y_i(0, M_i(0))$$

• Fundamental problem of causal inference: only one potential outcome can be observed for each *i* 

- $M_i(1)$ :
  - Quality of her challenger if politician i is an incumbent
    - Level of anxiety individual *i* would report if he reads the story with Hispanic immigrant
- $Y_i(1, M_i(1))$ :
  - Election outcome that would result if politician *i* is an incumbent and faces a challenger whose quality is  $M_i(1)$
  - Immigration attitude individual *i* would report if he reads the story with Hispanic immigrant and reports the anxiety level  $M_i(1)$
- $M_i(0)$  and  $Y_i(0, M_i(0))$  are the converse

• Causal mediation (Indirect) effects:

$$\delta_i(t) \equiv Y_i(t, M_i(1)) - Y_i(t, M_i(0))$$

- Causal effect of the change in *M<sub>i</sub>* on *Y<sub>i</sub>* that would be induced by treatment
- Change the mediator from  $M_i(0)$  to  $M_i(1)$  while holding the treatment constant at *t*
- Represents the mechanism through M<sub>i</sub>
- Zero treatment effect on mediator  $\Longrightarrow$  Zero mediation effect
- Examples:
  - Part of incumbency advantage that is due to the difference in challenger quality induced by incumbency status
  - Difference in immigration attitudes that is due to the change in anxiety induced by the treatment news story

### • Direct effects:

$$\zeta_i(t) \equiv Y_i(1, M_i(t)) - Y_i(0, M_i(t))$$

- Causal effect of  $T_i$  on  $Y_i$ , holding mediator constant at its potential value that would realize when  $T_i = t$
- Change the treatment from 0 to 1 while holding the mediator constant at *M<sub>i</sub>(t)*
- Represents all mechanisms other than through *M<sub>i</sub>*
- Total effect = mediation (indirect) effect + direct effect:

$$\tau_i = \delta_i(t) + \zeta_i(1-t) = \frac{1}{2} \{ (\delta_i(0) + \zeta_i(0)) + (\delta_i(1) + \zeta_i(1)) \}$$

### Mechanisms

- Indirect effects:  $\delta_i(t) \equiv Y_i(t, M_i(1)) Y_i(t, M_i(0))$
- Counterfactuals about treatment-induced mediator values

### Manipulations

- Controlled direct effects:  $\xi_i(t, m, m') \equiv Y_i(t, m) Y_i(t, m')$
- Causal effect of directly manipulating the mediator under  $T_i = t$

### Interactions

- Interaction effects:  $\xi(1, m, m') \xi(0, m, m')$
- The extent to which controlled direct effects vary by the treatment

# What Does the Observed Data Tell Us?

- Recall the Brader et al. experimental design:
  - randomize T<sub>i</sub>
  - Provide the measure M<sub>i</sub> and then Y<sub>i</sub>
- Among observations with  $T_i = t$ , we observe  $Y_i(t, M_i(t))$  but not  $Y_i(t, M_i(1 t))$  unless  $M_i(t) = M_i(1 t)$
- But we want to estimate

$$\delta_i(t) \equiv Y_i(t, M_i(1)) - Y_i(t, M_i(0))$$

- For t = 1, we observe  $Y_i(1, M_i(1))$  but not  $Y_i(1, M_i(0))$
- Similarly, for t = 0, we observe  $Y_i(0, M_i(0))$  but not  $Y_i(0, M_i(1))$
- We have the identification problem 

  Need assumptions or better research designs

Kosuke Imai (Princeton)

**Causal Mechanisms** 

#### Incumbency advantage:

- An incumbent  $(T_i = 1)$  faces a challenger with quality  $M_i(1)$
- We observe the electoral outcome  $Y_i = Y_i(1, M_i(1))$
- We also want  $Y_i(1, M_i(0))$  where  $M_i(0)$  is the quality of challenger this incumbent politician would face if she is not an incumbent

## Media framing effects:

- A subject viewed the news story with Hispanic immigrant ( $T_i = 1$ )
- For this person,  $Y_i(1, M_i(1))$  is the observed immigration opinion
- *Y<sub>i</sub>*(1, *M<sub>i</sub>*(0)) is his immigration opinion in the counterfactual world where he still views the story with Hispanic immigrant but his anxiety is at the same level as if he viewed the control news story

In both cases, we can't observe  $Y_i(1, M_i(0))$  because  $M_i(0)$  is not realized when  $T_i = 1$ 

# Sequential Ignorability Assumption

• Proposed identification assumption: Sequential Ignorability (SI)

$$\{Y_i(t',m),M_i(t)\} \perp T_i \mid X_i = x, \qquad (1)$$

$$Y_i(t',m) \perp M_i(t) \mid T_i = t, X_i = x$$
 (2)

- In words,
  - $T_i$  is (as-if) randomized conditional on  $X_i = x$
  - 2  $M_i(t)$  is (as-if) randomized conditional on  $X_i = x$  and  $T_i = t$
- Important limitations:
  - In a standard experiment, (1) holds but (2) may not
  - X<sub>i</sub> needs to include all confounders
  - If a state of the state of
  - Randomizing *M<sub>i</sub>* via manipulation is not the same as assuming *M<sub>i</sub>(t)* is as-if randomized

Back to Brader *et al.*:

- Treatment is randomized  $\implies$  (1) is satisfied
- But (2) may not hold:
  - Pre-treatment confounder or X<sub>i</sub>: state of residence those who live in AZ tend to have higher levels of perceived harm and be opposed to immigration
  - Post-treatment confounder: alternative mechanism beliefs about the likely negative impact of immigration makes people anxious
- Pre-treatment confounders  $\implies$  measure and adjust for them
- Post-treatment confounders  $\implies$  adjusting is not sufficient

Under SI, both ACME and average direct effects are nonparametrically identified (can be consistently estimated without modeling assumption)

- ACME  $\overline{\delta}(t)$  $\int \int \mathbb{E}(Y_i \mid M_i, T_i = t, X_i) \left\{ dP(M_i \mid T_i = 1, X_i) - dP(M_i \mid T_i = 0, X_i) \right\} dP(X_i)$
- Average direct effects  $\bar{\zeta}(t)$

 $\int \int \left\{ \mathbb{E}(Y_i \mid M_i, T_i = 1, X_i) - \mathbb{E}(Y_i \mid M_i, T_i = 0, X_i) \right\} dP(M_i \mid T_i = t, X_i) dP(X_i)$ 

Implies the general mediation formula under any statistical model

Kosuke Imai (Princeton)

**Causal Mechanisms** 

• Linear structural equation model (LSEM):

$$\begin{aligned} \mathbf{M}_i &= \alpha_2 + \beta_2 \mathbf{T}_i + \xi_2^\top \mathbf{X}_i + \epsilon_{i2}, \\ \mathbf{Y}_i &= \alpha_3 + \beta_3 \mathbf{T}_i + \gamma \mathbf{M}_i + \xi_3^\top \mathbf{X}_i + \epsilon_{i3}. \end{aligned}$$

- Fit two least squares regressions separately
- Use product of coefficients  $(\hat{\beta}_2 \hat{\gamma})$  to estimate ACME
- Use asymptotic variance to test significance (Sobel test)
- Under SI and the no-interaction assumption  $(\bar{\delta}(1) \neq \bar{\delta}(0)), \hat{\beta}_2 \hat{\gamma}$  consistently estimates ACME
- Can be extended to LSEM with interaction terms
- Problem: Only valid for the simplest LSEM

#### • The procedure:

- Regress Y on T and show a significant relationship
- Regress M on T and show a significant relationship
- Regress Y on M and T, and show a significant relationship between Y and M

### • The problems:

- First step can lead to false negatives especially if indirect and direct effects in opposite directions
- The procedure only anticipates simplest linear models
- On't do star-gazing. Report quantities of interest

#### Model outcome and mediator

- Outcome model:  $p(Y_i | T_i, M_i, X_i)$
- Mediator model:  $p(M_i | T_i, X_i)$
- These models can be of any form (linear or nonlinear, semi- or nonparametric, with or without interactions)
- **2** Predict mediator for both treatment values  $(M_i(1), M_i(0))$
- Solution Predict outcome by first setting  $T_i = 1$  and  $M_i = M_i(0)$ , and then  $T_i = 1$  and  $M_i = M_i(1)$
- Compute the average difference between two outcomes to obtain a consistent estimate of ACME
- Monte-Carlo or bootstrap to estimate uncertainty

# **Example: Binary Mediator and Outcome**

• Two logistic regression models:

$$\Pr(M_i = 1 \mid T_i, X_i) = \log i t^{-1} (\alpha_2 + \beta_2 T_i + \xi_2^\top X_i)$$
  
$$\Pr(Y_i = 1 \mid T_i, M_i, X_i) = \log i t^{-1} (\alpha_3 + \beta_3 T_i + \gamma M_i + \xi_3^\top X_i)$$

- Can't multiply  $\beta_2$  by  $\gamma$
- Difference of coefficients  $\beta_1 \beta_3$  doesn't work either

$$\Pr(Y_i = 1 \mid T_i, X_i) = \operatorname{logit}^{-1}(\alpha_1 + \beta_1 T_i + \xi_1^{\top} X_i)$$

- Can use our algorithm (example:  $\mathbb{E}\{Y_i(1, M_i(0))\}$ )
  - Predict  $M_i(0)$  given  $T_i = 0$  using the first model
  - Compute  $Pr(Y_i(1, M_i(0)) = 1 | T_i = 1, M_i = \widehat{M}_i(0), X_i)$  using the second model

# Sensitivity Analysis

- Standard experiments require sequential ignorability to identify mechanisms
- The sequential ignorability assumption is often too strong
- Need to assess the robustness of findings via sensitivity analysis
- Question: How large a departure from the key assumption must occur for the conclusions to no longer hold?
- Parametric sensitivity analysis by assuming

$$\{Y_i(t',m),M_i(t)\}\perp T_i \mid X_i = x$$

but not

$$Y_i(t',m) \perp M_i(t) \mid T_i = t, X_i = x$$

Possible existence of unobserved pre-treatment confounder

# Parametric Sensitivity Analysis

- Sensitivity parameter:  $\rho \equiv \text{Corr}(\epsilon_{i2}, \epsilon_{i3})$
- Sequential ignorability implies  $\rho = 0$
- Result:

$$\overline{\delta}(\mathbf{0}) = \overline{\delta}(\mathbf{1}) = \frac{\beta_2 \sigma_1}{\sigma_2} \left\{ \widetilde{\rho} - \rho \sqrt{(1 - \widetilde{\rho}^2)/(1 - \rho^2)} \right\},$$

where  $\sigma_j^2 \equiv \operatorname{var}(\epsilon_{ij})$  for j = 1, 2 and  $\tilde{\rho} \equiv \operatorname{Corr}(\epsilon_{i1}, \epsilon_{i2})$ .

- When do my results go away completely?
- $\bar{\delta}(t) = 0$  if and only if  $\rho = \tilde{\rho}$
- Easy to estimate from the regression of *Y<sub>i</sub>* on *T<sub>i</sub>*:

$$Y_i = \alpha_1 + \beta_1 T_i + \epsilon_{i1}$$

Kosuke Imai (Princeton)

# Interpreting Sensitivity Analysis with R squares

- Interpreting ρ: how small is too small?
- An unobserved (pre-treatment) confounder formulation:

$$\epsilon_{i2} = \lambda_2 U_i + \epsilon'_{i2}$$
 and  $\epsilon_{i3} = \lambda_3 U_i + \epsilon'_{i3}$ 

- How much does U<sub>i</sub> have to explain for our results to go away?
- Sensitivity parameters: R squares
   Proportion of previously unexplained variance explained by U<sub>i</sub>

$$\mathcal{R}_M^{2*} \equiv 1 - rac{\operatorname{var}(\epsilon'_{i2})}{\operatorname{var}(\epsilon_{i2})}$$
 and  $\mathcal{R}_Y^{2*} \equiv 1 - rac{\operatorname{var}(\epsilon'_{i3})}{\operatorname{var}(\epsilon_{i3})}$ 

Proportion of original variance explained by U<sub>i</sub>

$$\widetilde{R}_M^2 \equiv \frac{\operatorname{var}(\epsilon_{i2}) - \operatorname{var}(\epsilon'_{i2})}{\operatorname{var}(M_i)} \quad \text{and} \quad \widetilde{R}_Y^2 \equiv \frac{\operatorname{var}(\epsilon_{i3}) - \operatorname{var}(\epsilon'_{i3})}{\operatorname{var}(Y_i)}$$

• Then reparameterize  $\rho$  using  $(R_M^{2*}, R_Y^{2*})$  (or  $(\tilde{R}_M^2, \tilde{R}_Y^2)$ ):

$$\rho = \operatorname{sgn}(\lambda_2 \lambda_3) R_M^* R_Y^* = \frac{\operatorname{sgn}(\lambda_2 \lambda_3) \widetilde{R}_M \widetilde{R}_Y}{\sqrt{(1 - R_M^2)(1 - R_Y^2)}},$$

where  $R_M^2$  and  $R_Y^2$  are from the original mediator and outcome models

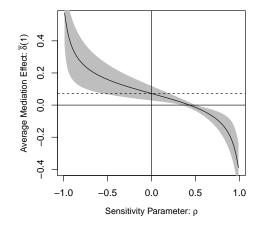
- $sgn(\lambda_2\lambda_3)$  indicates the direction of the effects of  $U_i$  on  $Y_i$  and  $M_i$
- Set  $(R_M^{2*}, R_Y^{2*})$  (or  $(\tilde{R}_M^2, \tilde{R}_Y^2)$ ) to different values and see how mediation effects change

# Reanalysis: Estimates under Sequential Ignorability

- Original method: Product of coefficients with the Sobel test
  - Valid only when both models are linear w/o T-M interaction (which they are not)
- Our method: Calculate ACME using our general algorithm

Outcome variables	Product of Coefficients	Average Causal Mediation Effect ( $\delta$ )
Decrease Immigration	.347	.105
ہں) Support English Only Laws	[0.146, 0.548] .204	[0.048, 0.170] .074
$\overline{\delta}(1)$ Request Anti-Immigration Information	[0.069, 0.339] .277	[0.027, 0.132] .029
$\bar{\delta}(1)$	[0.084, 0.469]	[0.007, 0.063]
Send Anti-Immigration Message $\overline{\delta}(1)$	.276 [0.102, 0.450]	.086 [0.035, 0.144]

# Reanalysis: Sensitivity Analysis w.r.t. p

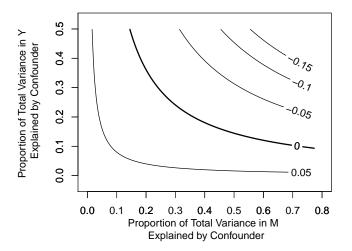


 ACME > 0 as long as the error correlation is less than 0.39 (0.30 with 95% CI)

Kosuke Imai (Princeton)

**Causal Mechanisms** 

Amsterdam (June 10, 2015) 30 / 57



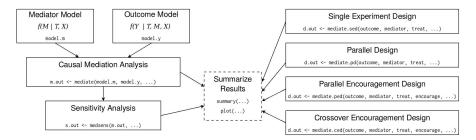
 An unobserved confounder can account for up to 26.5% of the variation in both Y<sub>i</sub> and M<sub>i</sub> before ACME becomes zero

Kosuke Imai (Princeton)

**Causal Mechanisms** 

Model-Based Inference

#### **Design-Based Inference**



Fit models for the mediator and outcome variable and store these models

```
> m <- lm(Mediator ~ Treat + X)
```

```
> y <- lm(Y ~ Treat + Mediator + X)
```

Mediation analysis: Feed model objects into the mediate() function. Call a summary of results

```
> summary(m.out)
```

Sensitivity analysis: Feed the output into the medsens () function. Summarize and plot

- > s.out <- medsens(m.out)</pre>
- > summary(s.out)
- > plot(s.out, "rho")
- > plot(s.out, "R2")

	Outcome Model Types							
Mediator Model Types	Linear	GLM	Ordered	Censored	Quantile	GAM	Survival	
Linear (lm/lmer)	$\checkmark$	$\checkmark$	√*	$\checkmark$	$\checkmark$	√*	$\checkmark$	
GLM (glm/bayesglm/glmer)	$\checkmark$	$\checkmark$	√*	$\checkmark$	$\checkmark$	√*	$\checkmark$	
Ordered (polr/bayespolr)	$\checkmark$	$\checkmark$	√*	$\checkmark$	$\checkmark$	√*	$\checkmark$	
Censored (tobit via vglm)	-	-	-	-	-	-	-	
Quantile (rq)	√*	√*	√*	√*	√*	√*	$\checkmark$	
GAM (gam)	√*	√*	√*	√*	√*	√*	√*	
Survival (survreg)	$\checkmark$	$\checkmark$	√*	$\checkmark$	$\checkmark$	√*	$\checkmark$	

Types of Models That Can be Handled by mediate. Stars (\*) indicate the model combinations that can only be estimated using the nonparametric bootstrap (i.e. with boot = TRUE).

- Treatment/mediator interactions, with formal statistical tests
- Treatment/mediator/pre-treatment interactions and reporting of quantities by pre-treatment values
- Factoral, continuous treatment variables
- Cluster standard errors/adjustable CI reporting/p-values
- Support for multiple imputation
- Multiple mediators
- Multilevel mediation (NEW!)

Please read our vignette file here.

Based on the same algorithm

Hicks, R, Tingley D. 2011. Causal Mediation Analysis. Stata Journal. 11(4):609-615.

ssc install mediation

More limited coverage of models (just bc. of time though!)

medeff (equation 1) (equation 2) [if] [in] [[weight]] ,
[sims(integer) seed(integer) vce(vcetype) Level(#)
interact(varname)] mediate(varname) treat(varname)

Where "equation 1" or "equation 2" are of the form (For equation 1, the mediator equation):

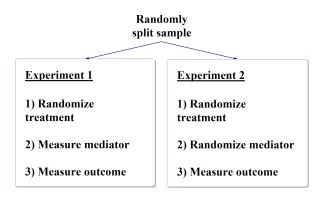
probit M T x

or

regress M T x

# **Beyond Sequential Ignorability**

- Without sequential ignorability, standard experimental design lacks identification power
- Even the sign of ACME is not identified
- Need to develop alternative experimental designs for more credible inference
- Possible when the mediator can be directly or indirectly manipulated
- All proposed designs preserve the ability to estimate the ACME under the SI assumption
- Trade-off: statistical power
- These experimental designs can then be extended to natural experiments in observational studies



- Must assume no direct effect of manipulation on outcome
- More informative than standard single experiment
- If we assume no T-M interaction, ACME is point identified

# Why Do We Need No-Interaction Assumption?

#### • Numerical Example:

Prop.	$M_{i}(1)$	$M_{i}(0)$	$Y_{i}(t, 1)$	$Y_{i}(t, 0)$	$\delta_i(t)$
0.3	1	0	0	1	-1
0.3	0	0	1	0	0
0.1	0	1	0	1	1
0.3	1	1	1	0	0

• 
$$\mathbb{E}(M_i(1) - M_i(0)) = \mathbb{E}(Y_i(t, 1) - Y_i(t, 0)) = 0.2$$
, but  $\bar{\delta}(t) = -0.2$ 

- The Problem: Causal effect heterogeneity
  - T increases M only on average
  - M increases Y only on average
  - T M interaction: Many of those who have a positive effect of T on M have a negative effect of M on Y (first row)
- A solution: sensitivity analysis (see Imai and Yamamoto, 2013)
- Pitfall of "mechanism experiments" or "causal chain approach"

Why study brain?: Social scientists' search for causal mechanisms underlying human behavior

• Psychologists, economists, and even political scientists

**Question**: What mechanism links low offers in an ultimatum game with "irrational" rejections?

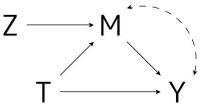
• A brain region known to be related to fairness becomes more active when unfair offer received (single experiment design)

Design solution: manipulate mechanisms with TMS

• Knoch et al. use TMS to manipulate — turn off — one of these regions, and then observes choices (parallel design)

## **Encouragement Design**

- Direct manipulation of mediator is difficult in most situations
- Use an instrumental variable approach:



- Advantage: allows for unobserved confounder between M and Y
- Key Assumptions:
  - Z is randomized or as-if random
  - No direct effect of Z on Y (a.k.a. exclusion restriction)

# Example: Social Norm Experiment on Property Taxes

- Lucia Del Carpio. "Are Neighbors Cheating?"
- Treatment: informing average rate of compliance
- Outcome: compliance rate obtained from administrative records
- Large positive effect on compliance rate  $\approx$  20 percentage points
- Mediators:
  - social norm (not measured; direct effect)
  - M<sub>1</sub>: beliefs about compliance (measured)
  - M<sub>2</sub>: beliefs about enforcement (measured)
- Instruments:
  - Z<sub>1</sub>: informing average rate of enforcement
  - Z<sub>2</sub>: payment-reminder
- Assumptions:
  - $Z_1$  affects Y only through  $M_1$  and  $M_2$
  - 2  $Z_2$  affects Y only through  $M_1$
- Results:
  - Average direct effect is estimated to be large
  - The author interprets this effect as the effect of social norm

- Recall ACME can be identified if we observe  $Y_i(t', M_i(t))$
- Get  $M_i(t)$ , then switch  $T_i$  to t' while holding  $M_i = M_i(t)$
- Crossover design:
  - Round 1: Conduct a standard experiment
  - 2 Round 2: Change the treatment to the opposite status but fix the mediator to the value observed in the first round
- Very powerful identifies mediation effects for each subject
- Must assume no carryover effect: Round 1 must not affect Round 2
- Can be made plausible by design

# Example: Labor Market Discrimination

- EXAMPLE Bertrand & Mullainathan (2004, AER)
  - Treatment: Black vs. White names on CVs
  - Mediator: Perceived qualifications of applicants
  - Outcome: Callback from employers
  - Quantity of interest: Direct effects of (perceived) race
  - Would Jamal get a callback if his name were Greg but his qualifications stayed the same?
  - Round 1: Send Jamal's actual CV and record the outcome
  - Round 2: Send his CV as Greg and record the outcome
  - Assumption: their different names do not change the perceived qualifications of applicants
  - Under this assumption, the direct effect can be interpreted as blunt racial discrimination

- Key difference between experimental and observational studies: treatment assignment
- Sequential ignorability:
  - Ignorability of treatment given covariates
  - Ignorability of mediator given treatment and covariates
- Both (1) and (2) are suspect in observational studies
- Statistical control: matching, propensity scores, etc.
- Search for quasi-randomized treatments: "natural" experiments
- How can we design observational studies?
- Experiments can serve as templates for observational studies

#### EXAMPLE Back to incumbency advantage

- Use of cross-over design (Levitt and Wolfram)
  - 1st Round: two non-incumbents in an open seat
  - 2nd Round: same candidates with one being an incumbent
- Assume challenger quality (mediator) stays the same
- Estimation of direct effect is possible
- Redistricting as natural experiments (Ansolabehere et al.)
  1st Round: incumbent in the old part of the district
  2nd Round: incumbent in the new part of the district
- Challenger quality is the same but treatment is different
- Estimation of direct effect is possible

# **Multiple Mediators**



- Quantity of interest = The average indirect effect with respect to M
- W represents the alternative observed mediators
- Left: Assumes independence between the two mechanisms
- Right: Allows *M* to be affected by the other mediators *W*
- Applied work often assumes the independence of mechanisms
- Under this independence assumption, one can apply the same analysis as in the single mediator case
- For causally dependent mediators, we must deal with the heterogeneity in the *T* × *M* interaction as done under the parallel design ⇒ sensitivity analysis

## Unpacking the Standard Path-Analytic Approach

• Applied social scientists often use the following model:

$$\begin{aligned} M_i &= \alpha_M + \beta_M T_i + \xi_M^\top X_i + \epsilon_{iM} \\ W_i &= \alpha_W + \beta_W T_i + \xi_W^\top X_i + \epsilon_{iW} \\ Y_i &= \alpha_3 + \beta_3 T_i + \gamma M_i + \theta^\top W_i + \xi_3^\top X_i + \epsilon_{i3} \end{aligned}$$

- The mediation effects are then estimated as  $\hat{\beta}_M \hat{\gamma}$  for *M* and  $\hat{\beta}_W \hat{\theta}$  for *W*
- We can show that these are consistent for  $\bar{\delta}^M_i$  and  $\bar{\delta}^W_i$  under the above assumption and linearity
- However, because of the assumed independence between mechanisms, analyzing one mechanism at a time will also be valid, e.g.,

$$\begin{aligned} \mathbf{M}_i &= \alpha_2 + \beta_2 \mathbf{T}_i + \xi_2^\top \mathbf{X}_i + \epsilon_{i2} \\ \mathbf{Y}_i &= \alpha_3 + \beta_3 \mathbf{T}_i + \gamma \mathbf{M}_i + \xi_3^\top \mathbf{X}_i + \epsilon_{i3} \end{aligned}$$

### Identification of Causally Related Mechanisms

• Consider the (weak) sequential ignorability assumption:

for any *t*, *m*, *w*, *x*.

- Unconfundedness of M<sub>i</sub> conditional on both pre-treatment (X<sub>i</sub>) and observed post-treatment (W<sub>i</sub>) confounders
- Corresponds to sequential randomization unlike Assumption 1
- The no  $T \times M$  interaction assumption required for the identification of  $\overline{\delta}(t)$  under Assumption 2:

 $Y_i(1, m, W_i(1)) - Y_i(0, m, W_i(0)) = Y_i(1, m', W_i(1)) - Y_i(0, m', W_i(0))$ 

# The Proposed Framework

- Problem: The no interaction assumption is often too strong (e.g. Does the effect of perceived issue importance invariant across frames?)
- We use a varying-coefficient linear structural equations model to:
  - Allow for homogeneous interaction for point identification
     Develop a sensitivity analysis in terms of the degree of heterogeneity in the interaction effect
- Consider the following model:

$$\begin{aligned} M_i(t, \boldsymbol{w}) &= \alpha_2 + \beta_{2i}t + \xi_{2i}^\top \boldsymbol{w} + \mu_{2i}^\top t \boldsymbol{w} + \lambda_{2i}^\top \boldsymbol{x} + \epsilon_{2i}, \\ Y_i(t, \boldsymbol{m}, \boldsymbol{w}) &= \alpha_3 + \beta_{3i}t + \gamma_i \boldsymbol{m} + \kappa_i t \boldsymbol{m} + \xi_{3i}^\top \boldsymbol{w} + \mu_{3i}^\top t \boldsymbol{w} + \lambda_{3i}^\top \boldsymbol{x} + \epsilon_{3i}, \end{aligned}$$

where  $\mathbb{E}(\epsilon_{2i}) = \mathbb{E}(\epsilon_{3i}) = 0$ 

- Allows for dependence of *M* on *W*
- Coefficients are allowed to vary arbitrarily across units

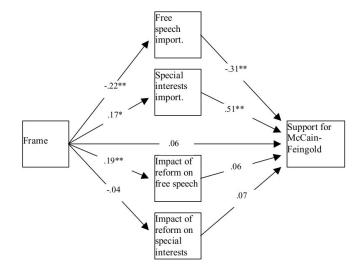
1

**Example:** Druckman and Nelson (2003) (N = 261)

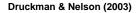
- Treatment: News paper article on a proposed election campaign finance reform, emphasizing either its positive or negative aspect
- Outcome: Support for the proposed reform
- Primary mediator: Perceived importance of free speech
- Alternative (confounding) mediator: Belief about the impact of the proposed reform
- Original analysis finds the importance mechanism to be significant, implicitly assuming its independence from beliefs

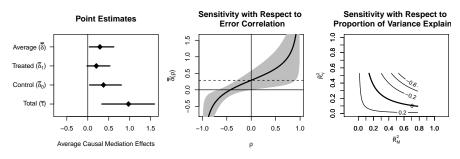
# **Original Analysis Assumes Independent Mechanisms**

#### Druckman and Nelson, p.738



# Analysis with the Independence Assumption



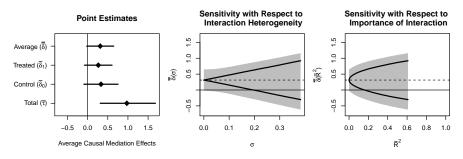


- Weakly significant average indirect effects ([0.025, 0.625]), accounting for 28.6% of the total effect
- Moderate degree of sensitivity to the mediator exogeneity ( $\bar{\delta} = 0$  when  $\rho = -0.43$  or  $\tilde{R}_M^2 \tilde{R}_Y^2 = 0.078$ )
- Concern: the importance mechanism may be affected by the belief content mechanism

Kosuke Imai (Princeton)

# Analysis without the Independence Assumption

#### Druckman & Nelson (2003)



- Similar results with slightly wider CI ([-0.021, 0.648])
- Lower bound on  $\bar{\delta}$  is zero when  $\sigma=$  0.195, or 51% of its upper bound
- This translates to the interaction heterogeneity explaining 15.9% of the variance of the outcome variable

Kosuke Imai (Princeton)

Causal Mechanisms

# **Concluding Remarks**

- Even in a randomized experiment, a strong assumption is needed to identify causal mechanisms
- However, progress can be made toward this fundamental goal of scientific research with modern statistical tools
- A general, flexible estimation method is available once we assume sequential ignorability
- Sequential ignorability can be probed via sensitivity analysis
- More credible inferences are possible using clever experimental designs
- Insights from new experimental designs can be directly applied when designing observational studies
- Multiple mediators require additional care when they are causally dependent

The project website for papers and software:

http://imai.princeton.edu/projects/mechanisms.html

Email for questions and suggestions:

kimai@princeton.edu