# Causal Mediation Analysis with Non-linear Models, Multiple Mediators, and DAGs

James Long
Department of Biostatistics
MD Anderson Cancer Center

SMU Statistics Seminar
November 13, 2020

# Outline

# Outline

# Causation and Mediation

Causation in science and medicine:

- ▶ Statins reduce incidence of heart attack

- ▶ SNP increases protein expression

- ▶ BRCA mutations cause breast cancer

# Causation and Mediation

Causation in science and medicine:

- Statins reduce incidence of heart attack

  **by lowering cholesterol.**

- SNP increases protein expression

  **by increasing mRNA expression.**

- BRCA mutations cause breast cancer

  **by inhibiting DNA repair.**

## Causation and Mediation

Causation in science and medicine:

- ▶ Statins reduce incidence of heart attack

     **by lowering cholesterol.**

- ▶ SNP increases protein expression

     **by increasing mRNA expression.**

- ▶ BRCA mutations cause breast cancer

     **by inhibiting DNA repair.**

**Causal mediators:** variables that facilitate a causal relation.
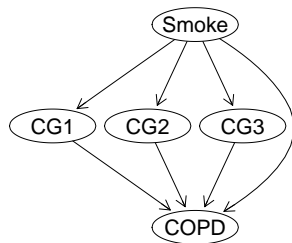


$x$ changes $y$ by altering $m$.

# Smoking → Methylation → COPD

- ▶ Chronic obstructive pulmonary disease (COPD): chronic inflammation of lungs, obstructs airflow
- ▶ Smoking is leading risk factor.
- ▶ Hypothesis: Smoking causes COPD by **methlyating** genes. Gaynor et al. [2018], Wan et al. [2012]

**Data:**
- ▶ Smoke = 1 for smoker, 0 no
- ▶ CG1,CG2,CG3 = methylation at 3 sites
- ▶ COPD = 1 or 0 disease status



**Scientific Questions:**
- ▶ Does methylation mediate smoking-COPD causal relation?
- ▶ Which sites?
- ▶ How much of the total effect is mediated by methylation?

# Overview of Talk

▶ Causal Modeling with DAGs

   ▶ Is $m$ a mediator?

   ▶ What percentage of total effect of $x$ on $y$ is mediated by $m$? Are there other causal pathways?

   ▶ Control for confounders

▶ Our Framework: `mediateR` package

   ▶ Estimation with non–linear models

   ▶ Multiple mediators (e.g. many gene expressions)

   ▶ Effect Scales

▶ Example: TCGA Kidney Renal Cell Carcinoma Data
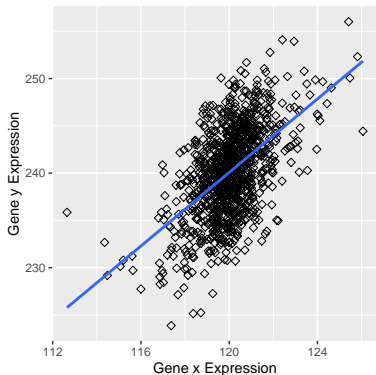
▶ Discussion

# Outline

# Statistical and Causal Models

- ► x = Gene x Expression
- ► y = Gene y Expression
- ► Fit linear model:

$$y \approx 1.94x + 6.88$$

- ► How do we interpret model predictions at x=124?

$$y \approx 1.94 \times 124 + 6.88 \approx 248$$



<u>Association:</u> If I **see** x=124, y $\approx$ 248. Passive observers of nature.
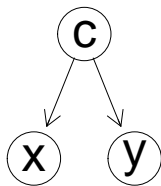<u>Causation:</u> If I set (or **do**) x=124, y $\approx$ 248. Intervene in nature.

**The statistical model (i.e. set of assumptions) does not permit the causal conclusion. Need a causal model for data.**
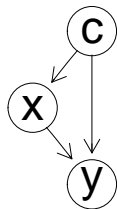
# DAG Definition

A **directed acyclic graph** (DAG) is a set of edges (directed arrows) between vertices (variables) such that there are no loops (acyclic).
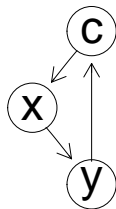
Example: Vertices = variables = $(c, x, y)$



DAG          DAG          Not a DAG (loop)

# Causal DAG

**Causal DAGs** encode causal assumptions by specifying the process by which variables are generated.

## Example

(c,x,y) are three gene expression levels



$$c \leftarrow f_c(\epsilon_c)$$
$$x \leftarrow f_x(c, \epsilon_x)$$
$$y \leftarrow f_y(c, \epsilon_y)$$

for some functions $f_c, f_x, f_y$ and independent[a] random variables $(\epsilon_c, \epsilon_x, \epsilon_y)$.

Not true:  $y \leftarrow f_y(x, c, \epsilon_y)$ because no $x$ to $y$ arrow.

a. This is the Markovian assumption. More discussion later in talk.

## External Interventions with Causal DAGs
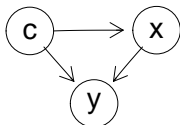
**Definition:** $y^{x=0}$ (counterfactual) is the value of $y$ if $x$ is set to 0

**Example of Computing $y^{x=0}$ using Causal DAG:**

Original Causal DAG:



$$c \leftarrow f_c(\epsilon_c)$$
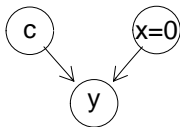$$x \leftarrow f_x(c, \epsilon_x)$$
$$y \leftarrow f_y(x, c, \epsilon_y)$$

Mutilated DAG (intervention on $x$):



$$c \leftarrow f_c(\epsilon_c)$$
$$x \leftarrow 0$$
$$y \leftarrow f_y(x = 0, c, \epsilon_y)$$

**Summarizing the Causal Effect:** (Pearl [2009])

- $p(y^{x=0} = y) \equiv p(y|do(x = 0))$
- $\mathbb{E}[y^{x=0}] \equiv \mathbb{E}[y|do(x = 0)] = \int yp(y|do(x = 0))dy$

# $y^{x=2}$ and $x^{y=2}$ in Simple Model

Assumptions:

(x, y) is jointly normal



$x \leftarrow \mu + \epsilon_x$

$y \leftarrow \beta_0 + \beta_x x + \epsilon_y$

where $\epsilon_x \sim N(0, \sigma_x^2)$ and
$\epsilon_y \sim N(0, \sigma_y^2)$, independent

$y^{x=2}$ (value of y when x set to 2):



$x \leftarrow 2$

$y \leftarrow \beta_0 + 2\beta_x + \epsilon_y$

$\mathbb{E}[y|do(x = 2)] = \mathbb{E}[y|x = 2] = \beta_0 + 2\beta_x$
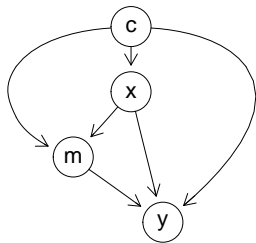
$x^{y=2}$ (value of x when y set to 2):



$x \leftarrow \mu + \epsilon_x$

$y \leftarrow 2$

$\mathbb{E}[x|do(y = 2)] = \mathbb{E}[x] = \mu$

# More Complex Causal Model

- c = SNP
- x = gene expression (binarized to high / low)
- m = protein expression (continuous)
- y = phenotype (continuous)



$$c \leftarrow f_c(\epsilon_c)$$
$$x \leftarrow f_x(c, \epsilon_x)$$
$$m \leftarrow f_m(x, c, \epsilon_m) = \alpha_x x + \alpha_c c + \epsilon_m$$
$$y \leftarrow f_y(x, c, m, \epsilon_y) = \beta_x x + \beta_c c + \beta_m m + \epsilon_y$$
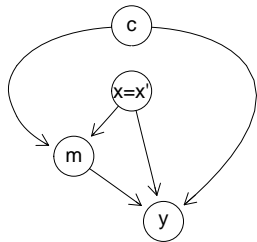
orange: causal assumptions from DAG
blue: statistical assumptions

**Goal:** Compute expected change in phenotype when changing gene expression from low ($x = 0$) to high ($x = 1$)

$$TE(0, 1) \equiv \mathbb{E}[y^{x=1}] - \mathbb{E}[y^{x=0}] = \mathbb{E}[y|do(x = 1)] - \mathbb{E}[y|do(x = 0)]$$

# Computation of Total Effect

For $x' = 0$ or $1$, mutilated graph is:



$$c \leftarrow f_c(\epsilon_c)$$
$$x \leftarrow x'$$
$$m \leftarrow f_m(x = x', c, \epsilon_m) = \alpha_x x' + \alpha_c c + \epsilon_m$$
$$y \leftarrow f_y(x = x', c, m, \epsilon_y) = \beta_x x' + \beta_c c + \beta_m m + \epsilon_y$$

After some algebra

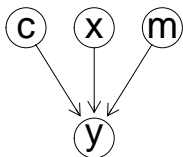$$\mathbb{E}[y \mid do(x = x')] = (\beta_x + \alpha_x \beta_m) x' + z$$

where $z$ does not depend on $x'$. So

$$TE(0, 1) = \mathbb{E}[y \mid do(x = 1)] - \mathbb{E}[y \mid do(x = 0)]$$
$$= \underbrace{\beta_x}_{\text{direct effect}} + \underbrace{\alpha_x \beta_m}_{\substack{\text{indirect effect} \\ \text{passing through } m}}$$

# Variables to Control For

- "Controlling for everything" gives **wrong** result:
  - regress $y$ on $(x, m, c)$
  - coefficient on $x$ is $\beta_x$
  - $\beta_x$ ignores $x \to m \to y$ effect
  - gives correct TE for:



- "Controlling for nothing" gives **wrong** result:
  - $TE \neq \mathbb{E}[y|x = 1] - \mathbb{E}[y|x = 0]$
- "Controlling for c" gives **correct** result:
  - regress $y$ on $(x, c)$
  - coefficient on $x$ is $\beta_x + \beta_m \alpha_x$

**Message:** Determining variables to control for requires assumptions on the causal structure of data. DAG can express these assumptions formally.

# Summary of Causal DAGs

▶ DAGs encode causal assumptions

▶ Enable computation of effects of interventions

▶ Judea Pearl pioneered this approach, see Pearl [2009]

# Outline

# Mediation in Linear Models

Developed by Baron and Kenny [1986], Sobel [1982], others.



- ▶ Direct Effect $\equiv \beta_x$
- ▶ Indirect Effect $\equiv \alpha_x \beta_m$
- ▶ Total Effect $\equiv \beta_x + \alpha_x \beta_m$

$\equiv$ is "by definition"

Statistical Inference: Fit two least squares models:

1. m | x,c
2. y | x,m,c

- ▶ Point estimators, confidence intervals, tests, based on parameter estimates e.g.

$$H_0 : \alpha_x \beta_m = 0 \text{ i.e. m is not a mediator}$$

$$H_a : \alpha_x \beta_m \neq 0 \text{ i.e. m is a mediator}$$

**Limitation: Only applies to linear models, single mediator.**

# Direct Effect: General Model Definition

The direct effect* of changing $x$ from 0 to 1 is

$$y^{x=1, m^{x=0}} - y^{x=0}$$

where

$$y^{x=1, m^{x=0}} = \text{the value of } y \text{ when setting } x = 1 \text{ for the direct path}$$
$$\text{from } x \to y \text{ while generating } m \text{ as if } x = 0$$
$$y^{x=0} = \text{the value of } y \text{ when setting } x = 0$$



* Based on Pearl [2001] and VanderWeele and Vansteelandt [2010] definitions.

# Natural Direct Effect Interpretation

- $x = 1$ if smoker, 0 no
- $m =$ gene methylation
- $y = 1$ if COPD, 0 if no

The direct effect

$$y^{x=1, m^{x=0}} - y^{x=0}$$

is the change in COPD if someone smokes but takes a methylation prevention drug which blocks any potential change in COPD induced via methylation.

# Natural Direct Effect

The **natural direct effect** is:

$$DE(0,1) = \mathbb{E}[y^{x=1,m^{x=0}}] - \mathbb{E}[y^{x=0}]$$

where <u>assuming causal DAG structure</u>[*]:

$$\mathbb{E}[y^{x=1,m^{x=0}}] = \int_{m,c} \mathbb{E}[y|x=1,m,c]p(m|x=0,c)p(c)dmdc$$

$$\mathbb{E}[y^{x=0}] = \int_{m,c} \mathbb{E}[y|x=0,m,c]p(m|x=0,c)p(c)dmdc$$

[*] See Pearl [2001].

# Indirect Effect: General Model Definition

The indirect effect* of changing $x$ from 0 to 1 is
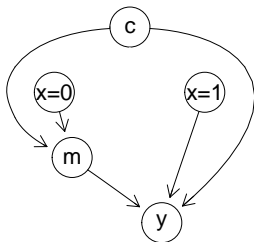
$$y^{x=1} - y^{x=1, m^{x=0}}$$

where

$$y^{x=1, m^{x=0}} = \text{the value of } y \text{ when setting } x = 1 \text{ for the direct path}$$
$$\text{from } x \to y \text{ while generating } m \text{ as if } x = 0$$

$$y^{x=1} = \text{the value of } y \text{ when setting } x = 1$$

$y^{x=1}$



$y^{x=1, m^{x=0}}$

*Based on VanderWeele and Vansteelandt [2010] and others. See Pearl [2001] for alternative definition.

# Natural Indirect Effect

The **natural indirect effect** is:

$$IE(0,1) = \mathbb{E}[y^{x=1}] - \mathbb{E}[y^{x=1,m^{x=0}}]$$

where <u>assuming causal DAG structure</u>:

$$\mathbb{E}[y^{x=1,m^{x=0}}] = \int_{m,c} \mathbb{E}[y|x=1,m,c]p(m|x=0,c)p(c)dmdc$$

$$\mathbb{E}[y^{x=1}] = \int_{m,c} \mathbb{E}[y|x=1,m,c]p(m|x=1,c)p(c)dmdc$$

# Mediation Formula

**Mediation formula**:

$$TE(x', x'') = \mathbb{E}[y|do(x = x'')] - \mathbb{E}[y|do(x = x')]$$
$$= DE(x', x'') + IE(x', x'')$$

- Previously x'=0 and x"=1.
- For linear models, effects only depend on $(x'' - x')$:
  - $DE(x', x'') = \beta_x(x'' - x')$
  - $IE(x', x'') = \alpha_x \beta_m(x'' - x')$

# Outline

# Causal Mediation for Genomic Data Sets

- ▶ Many exposures and potential mediators
    - ▶ mRNA expression of many genes
    - ▶ Protein expression of many genes
- ▶ Non–linear response models
    - ▶ Logistic model for disease status
    - ▶ Cox PH model for survival time

# Summary of Our Framework



- ▶ Numerical approximation of DE and IE integrals
- ▶ Multiple exposures $x$, Multivariate mediators $m$, Non–linear response models $y$
- ▶ Various Effect Scales: Mean, Odds, Restricted Mean

    `mediateR` https://github.com/longjp/mediateR

# Estimation

Require estimate of:

$$\mathbb{E}[y^{x=1, m^{x=0}}] = \int_{m,c} \mathbb{E}[y|x=1, m, c]p(m|x=0, c)p(c)dmdc$$

Approaches to estimating effects:

- ▶ Rare disease assumption with logistic response model
  - ▶ Assume rare disease ($P(y=1)$ small)
  - ▶ Effects are approximately linear combinations of path coefficients [VanderWeele and Vansteelandt, 2010, Huang et al., 2014]
- ▶ Common disease assumption with logistic response model
  - ▶ Use probit approximation to logistic [Gaynor et al., 2018]
- ▶ **Numerically approximate integral: Discuss now**
  - ▶ Use empirical distribution of $c$ to approximate $p(c)$
  - ▶ Draw $\widehat{m}$ from plug–in estimate of $p(m|x=0, c)$
  - ▶ Use plug–in estimate of $\mathbb{E}[y|x=1, m, c]$

# Numerical Integration in Logistic Model

<u>Statistical Model:</u> $(c_i, x_i, m_i, y_i)$ i.i.d. from

- $m = \alpha^{(x)}x + \alpha^{(c)}c + \alpha^{(0)} + \epsilon$
- $y \sim Bern((1 + e^{-(x^T\beta^{(x)} + m^T\beta^{(m)} + c^T\beta^{(c)} + \beta^{(0)})})^{-1})$

$\alpha^{(x)} \in \mathbb{R}^r$, $\alpha^{(c)} \in \mathbb{R}^{r \times q}$, $\alpha^{(0)} \in \mathbb{R}^r$, $\epsilon \sim N(0, \Sigma_\epsilon)$, $\Sigma_\epsilon \in \mathbb{R}^{r \times r}$, $\beta^{(x)} \in \mathbb{R}^1$, $\beta^{(c)} \in \mathbb{R}^q$, $\beta^{(m)} \in \mathbb{R}^r$, and $\beta^{(0)} \in \mathbb{R}^1$.

<u>Then:</u>

$$\mathbb{E}[y^{x=1, m^{x=0}}] = \int_{m,c} \mathbb{E}[y|x=1, m, c]p(m|x=0, c)p(c)dmdc$$

$$\approx \frac{1}{n}\sum_{i=1}^{n}\widehat{\mathbb{E}}[y|x=1, \widehat{m}_i, c_i]$$

where

$$\widehat{m}_i \sim N(m; \widehat{\alpha}^{(x)}0 + \widehat{\alpha}^{(c)}c_i + \widehat{\alpha}^{(0)}, \widehat{\Sigma_\epsilon})$$

$$\widehat{\mathbb{E}}[y|x=1, \widehat{m}_i, c_i] = \frac{1}{1 + e^{-\widehat{\beta}^{(x)}(1) - \widehat{m}_i^T\widehat{\beta}^{(m)} - c_i^T\widehat{\beta}^{(c)} - \widehat{\beta}^{(0)}}}$$

# Notes on Numerical Integration

- Calculations for $\mathbb{E}[y^{x=1}]$ and $\mathbb{E}[y^{x=0}]$ similar

- More generally, $x = 0$ and $x = 1$ can be replaced by any $x = x'$ and $x = x''$

- Similar strategy employed for $y$ with Cox proportional hazards model

- Bootstrap sample $(c_i, x_i, m_i, y_i)$ to obtain sampling distributions of $DE$, $IE$, $TE$
  - CI and hypothesis tests based on bootstrap quantiles

# Effect Scale for Logistic Response Model

▶ Odds ratios useful for summarizing effects with binary responses

▶ For logistic model, effects may be defined on odds scale:

$$DE^o(0,1) \equiv \frac{\frac{\mathbb{E}[y^{x=1,m^{x=0}}]}{1-\mathbb{E}[y^{x=1,m^{x=0}}]}}{\frac{\mathbb{E}[y^{x=0}]}{1-\mathbb{E}[y^{x=0}]}} \quad IE^o(0,1) \equiv \frac{\frac{\mathbb{E}[y^{x=1}]}{1-\mathbb{E}[y^{x=1}]}}{\frac{\mathbb{E}[y^{x=1,m^{x=0}}]}{1-\mathbb{E}[y^{x=1,m^{x=0}}]}}$$

▶ Mediation Formula on the Odds Scale:

$$TE^o(0,1) \equiv \frac{\frac{\mathbb{E}[y^{x=1}]}{1-\mathbb{E}[y^{x=1}]}}{\frac{\mathbb{E}[y^{x=0}]}{1-\mathbb{E}[y^{x=0}]}} = DE^o(0,1) \times IE^o(0,1)$$

Proposed by VanderWeele and Vansteelandt [2010]

# Effect Scale for Survival Response

The mean of $y$ restricted to $L$ is

$$\mathbb{E}[\min(y, L)]$$

where $L$ is some constant

- ▶ Restricted mean scale popular in survival applications because mean estimate has high variance

- ▶ Direct and Indirect Effects on restricted mean scale

$$DE^R(0,1) \equiv \mathbb{E}[\min(y^{x=1, m^{x=0}}, L)] - \mathbb{E}[\min(y^{x=0}, L)]$$
$$IE^R(0,1) \equiv \mathbb{E}[\min(y^{x=1}, L)] - \mathbb{E}[\min(y^{x=1, m^{x=0}}, L)]$$

- ▶ Mediation Formula on Restricted Mean Scale:

$$TE^R(0,1) \equiv \mathbb{E}[\min(y^{x=1}, L)] - \mathbb{E}[\min(y^{x=0}, L)]$$
$$= DE^R(0,1) + IE^R(0,1)$$

# Outline

# TCGA Kidney Renal Cell Carcinoma

$n = 470$ patients with:

- ▶ Clinical features, e.g. survival time
- ▶ RPPA Protein Expression
- ▶ mRNA Expression



Questions:

1. What gene expressions are associated with / causing changes in survival?

2. Are these changes being mediated at the protein level?

Data available: https://portal.gdc.cancer.gov/

# Gene Expression–Survival Correlations

- ▶ Network et al. [2013] identified 5 core metabolic pathways: PTEN, TCA cycle, Fatty acid synthesis, AMPK, Pentose phosphate
- ▶ Summarize each pathway at mRNA level with 1st principal component
- ▶ Fit Cox model on survival given pathway scores
- ▶ Survival curves by pathway score risk:

# Mediation Model

- mRNA expression can alter protein expression (via increased translation by ribosomes)
- Treat 5 proteins as mediators

  AMPKA alpha, AMPK pT172, ACC pS79, ACC, and PTEN
- For each mRNA pathway
  - Compute $DE^R$, $IE^R$, $TE^R$ with $L = 2000$ at
    - $x' = $ 5th percentile of pathway score
    - $x'' = $ 95th percentile of pathway score
- Interpretation of $DE^R$ for Pentose Phosphate
  - 1) $y^{x=x'', m^{x=x'}} = $ generate the mediators from Pentose phosphate at 5th percentile and then survival with Pentose phosphate at 95th percentile.
  - 2) $y^{x=x'} = $ Generate both mediators and response with Pentose phosphate at 5th percentile.
  - Take expected difference (lifetime restricted to 2000) of 1 - 2

# Results

| Pathway | Indirect | | Direct | | Total | |
|---|---|---|---|---|---|---|
| PTEN | -29 | [-146,77] | 203 | [-49,409] | 174 | [-74,384] |
| TCA cycle | 40 | [-23,120] | 289 | [42,494] | 329 | [117,525] |
| Fatty acid synthesis | -156 | [-339,20] | -290 | [-537,-91] | -446 | [-654,-268] |
| AMPK | 23 | [-74,126] | 8 | [-292,316] | 30 | [-245,328] |
| Pentose phosphate | -94 | [-247,79] | -181 | [-511,91] | -274 | [-574,-29] |

Table 1: Indirect, Direct, and Total effects and 95% confidence intervals (in days) of metabolomic mRNA expression as mediated by protein expression.

Under model assumptions:

▶ TCA cycle, Fatty Acid Synthesis, Pentose phosphate have significant total effects (at $\alpha = 0.05$)
▶ Fatty acid synthesis largest indirect effect estimate, but not significant (at $\alpha = 0.05$)

# Outline

# Summary

- ▶ Causal DAGs graphically represent a set of assumptions about the causal structure of the data.

- ▶ Causal DAG + Statistical model $\rightarrow$ predict interventions

- ▶ Assumptions implied by causal DAGs may be implausible, especially for observational genetic data.

- ▶ Manuscript describing work: https://arxiv.org/abs/2011.06061

- ▶ Software Resources

  - ▶ mediateR (https://github.com/longjp/mediateR): package I codeveloped for implementing models. Joint work with Kim-Anh Do, Min Jin Ha, Ehsan Irajizad (MDACC Biostatistics), and James Doecke (CSIRO).

  - ▶ mediation (https://cran.r-project.org/web/packages/mediation/index.html): package with partially overlapping functionality (Imai et al. [2010])

# Inferring Causal Network from Data

**Causal Discovery:** Inferring causal structure (e.g. DAG) from data.

- ▶ PC algorithm: Estimates skeleton of DAG

  (Kalisch and Bühlmann [2007], Spirtes et al. [2000], Ha et al. [2016])

- ▶ Invariant Causal Prediction: Use mix of observational / experimental data to determine entire causal structure
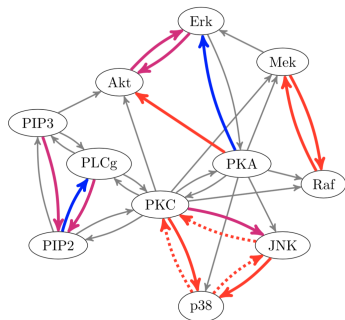
  Peters et al. [2016]



Figure source Meinshausen et al. [2016]

# Potential Outcomes Framework

- **DAG Framework**:
    - Graph specifies causal assumptions.
    - $y^{x=0}$ and $y^{x=1}$ (counterfactuals) distributions inferred from graph.
    - Developed / advocated by: Pearl [2009], Spirtes et al. [2000], Wright [1934]

- **Potential Outcomes Framework:**
    - $y^{x=0}$ and $y^{x=1}$ (potential outcomes) are primitive notions.
    - Causal knowledge conveyed via conditional independence assumptions, e.g.

    $$y^{x=x'} \perp\!\!\!\perp x | c \ \forall \ x'$$

    - Developed / advocated by: Rubin [2005], Splawa-Neyman et al. [1990]
    - "Direct and indirect causal effects via potential outcomes" Rubin [2004]

Thank you. Questions?

# Bibliography I

R. M. Baron and D. A. Kenny. The moderator–mediator variable distinction in social psychological research: Conceptual, strategic, and statistical considerations. *Journal of personality and social psychology*, 51(6):1173, 1986.

S. M. Gaynor, J. Schwartz, and X. Lin. Mediation analysis for common binary outcomes. *Statistics in medicine*, 2018.

M. J. Ha, W. Sun, and J. Xie. Penpc: A two-step approach to estimate the skeletons of high-dimensional directed acyclic graphs. *Biometrics*, 72(1):146–155, 2016.

Y.-T. Huang, T. J. VanderWeele, and X. Lin. Joint analysis of snp and gene expression data in genetic association studies of complex diseases. *The annals of applied statistics*, 8(1):352, 2014.

K. Imai, L. Keele, and T. Yamamoto. Identification, inference and sensitivity analysis for causal mediation effects. *Statistical science*, pages 51–71, 2010.

M. Kalisch and P. Bühlmann. Estimating high-dimensional directed acyclic graphs with the pc-algorithm. *Journal of Machine Learning Research*, 8(Mar):613–636, 2007.

N. Meinshausen, A. Hauser, J. M. Mooij, J. Peters, P. Versteeg, and P. Bühlmann. Methods for causal inference from gene perturbation experiments and validation. *Proceedings of the National Academy of Sciences*, 113 (27):7361–7368, 2016.

C. G. A. R. Network et al. Comprehensive molecular characterization of clear cell renal cell carcinoma. *Nature*, 499 (7456):43, 2013.

J. Pearl. Direct and indirect effects. In *Proceedings of the seventeenth conference on uncertainty in artificial intelligence*, pages 411–420. Morgan Kaufmann Publishers Inc., 2001.

J. Pearl. *Causality*. Cambridge university press, 2009.

J. Peters, P. Bühlmann, and N. Meinshausen. Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(5): 947–1012, 2016.

D. B. Rubin. Direct and indirect causal effects via potential outcomes. *Scandinavian Journal of Statistics*, 31(2): 161–170, 2004.

# Bibliography II

D. B. Rubin. Causal inference using potential outcomes: Design, modeling, decisions. *Journal of the American Statistical Association*, 100(469):322–331, 2005.

M. E. Sobel. Asymptotic confidence intervals for indirect effects in structural equation models. *Sociological methodology*, 13:290–312, 1982.

P. Spirtes, C. N. Glymour, R. Scheines, D. Heckerman, C. Meek, G. Cooper, and T. Richardson. *Causation, prediction, and search*. MIT press, 2000.

J. Splawa-Neyman, D. M. Dabrowska, and T. Speed. On the application of probability theory to agricultural experiments. essay on principles. section 9. *Statistical Science*, pages 465–472, 1990.

T. J. VanderWeele and S. Vansteelandt. Odds ratios for mediation analysis for a dichotomous outcome. *American journal of epidemiology*, 172(12):1339–1348, 2010.

E. S. Wan, W. Qiu, A. Baccarelli, V. J. Carey, H. Bacherman, S. I. Rennard, A. Agusti, W. Anderson, D. A. Lomas, and D. L. DeMeo. Cigarette smoking behaviors and time since quitting are associated with differential dna methylation across the human genome. *Human molecular genetics*, 21(13):3073–3082, 2012.

S. Wright. The method of path coefficients. *The annals of mathematical statistics*, 5(3):161–215, 1934.