

# The mechanism-based approach for age-period-cohort analysis

Arvid Sjölander

Department of Medical Epidemiology and Biostatistics  
Karolinska Institutet

# Age-period-cohort models

- ▶ Age-period-cohort models have a long history in epidemiology, social science and econometrics
  - ▶ Fannon and Nielsen (2019), Fosse and Winship (2019), and Murphy and Yang (2018)
- ▶ The purpose is to assess how an outcome of an individual is 'related' to three different time variables:
  - ▶ the time since the individual was born (age)
  - ▶ the calendar date at which the outcome is observed (period)
  - ▶ the calendar date at which the individual was born (cohort)

# A determinisic relation

- ▶ The three time variables are deterministically related:

$$\text{age} = \text{period} - \text{cohort}$$

- ▶ e.g., an individual who is born in year 1950 is 20 years old in 1970:

$$20 = 1970 - 1950$$

# Statistical associations

- ▶ The deterministic age-period-cohort relation is not problematic if we only care about statistical associations
- ▶ We can then:
  - ▶ model how the outcome is associated with, say, cohort and period, and
  - ▶ determine the outcome's association with age from fitted model by converting either cohort or period to age
- ▶ Usual goodness-of-fit tests relevant here as well...

# Causal effects

- ▶ The typical aim is more ambitious: we want to estimate the 'independent' causal effects of age, period and cohort
  - ▶ rarely stated explicitly though
- ▶ An identifiability problem: cannot contrast different values of one time variable while holding the two others fixed
  - ▶ somewhat similar to attempting to adjust for a confounder that is perfectly correlated with the exposure

# The constrained effects approach

- ▶ By far most common, traditionally
- ▶ Imposes parametric constraints on the three effects to make them identifiable
  - ▶ Firebaugh and Davis (1988), Glenn (1994), and Myers and Lee (1998): assume one of the effects is 0
  - ▶ more elaborate proposals as well Knoke and Hout (1974), Mason et al. (1973), and Nakamura (1986)
- ▶ Problems:
  - ▶ the constraints are often artificial; no *a priori* reason to believe that they hold
  - ▶ results are typically sensitive to the choice of constraints

# The mechanism-based approach

- ▶ Utilizes mediators on causal pathways between age-period-cohort and the outcome
- ▶ A sufficiently informative set of mediators may 'explain' the age, period and cohort effects and make them identifiable
  - ▶ Heckman and Robb (1985), O'Brien (2000), and Winship and Harding (2008): mainly informal arguments, and no explicit causal estimand
  - ▶ Bijlsma et al. (2017) and Sjölander and Gabriel (2025): formal development with modern causal inference methods

# Outline

Motivating example (Winship & Harding, 2008)

Nonparametric identification

Parametric estimation

Motivating example, revisited



# Outline

Motivating example (Winship & Harding, 2008)

Nonparametric identification

Parametric estimation

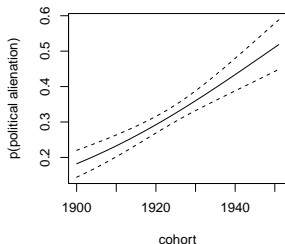
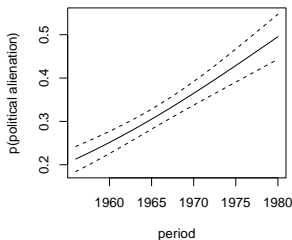
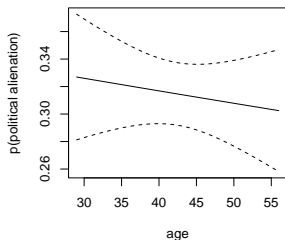
Motivating example, revisited

# Data and aim

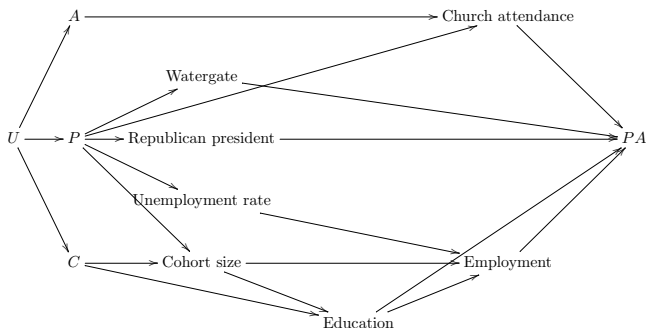
- ▶ Data from the American National Election Studies (ANES)
  - ▶ academically run national surveys of voters in the US, conducted around every presidential election
  - ▶ publicly available at <https://electionstudies.org/>
- ▶ Inclusion criterion: married white males age 29 to 56 surveyed in 1956, 1960, 1964, 1968, 1976, and 1980
  - ▶  $n = 1605$  (questionnaires, not individuals)
- ▶ Outcome: political alienation
  - ▶ binary indicator of the respondent agreeing with the statement '*I don't think officials care much what people like me think*'.

# Statistical associations: one-by-one

$$\text{logit}\{p(\text{political alienation}|X)\} = \beta_0 + \beta_1 X, \quad X \in \{\text{age, period, cohort}\}$$



# Assumed causal diagram (Winship & Harding, 2008)



- ▶ The role of  $U$  is to enforce the relation  $A = P - C$
- ▶ Key assumptions
  - ▶ **no unmeasured confounding**: no common causes of  $(A, P, C)$  and other variables
  - ▶ **partial exclusion**: no variable (mediator or outcome) is directly caused by both  $A, P$  and  $C$  – will modify this definition later to cover more general scenarios

## Analysis (Winship & Harding, 2008)

- ▶ Regression of each variable (mediators and outcome) on its direct causes
  - ▶ no identifiability problem when fitting models due to partial exclusion assumption
- ▶ Path-specific effects through product-of-coefficient method
  - ▶ standard in linear structural equation modeling (Bollen, 1989)
  - ▶ presumed causal interpretation of estimated effects due to no unmeasured confounding assumption

# Results (Winship & Harding, 2008)

	Estimate	95 Percent CI
Period effect (1976 vs. 1960)		
P → Watergate → PA	.4939	.3337, .6629
P → Republican president → PA	.3576	.2010, .5205
P → unemployment rate → employment → PA	.0021	-.0017, .0067
P → cohort size → employment → PA	.0018	-.0015, .0059
P → cohort size → education → PA	-.0099	-.0238, .0043
P → cohort size → education → employment → PA	-.0001	-.0003, .0001
P → church attendance → PA	.0487	.0163, .0973
Total	.8940	.7252, 1.0633
Cohort effect (1936-1939 vs. 1908-1911)		
C → cohort size → employment → PA	.0018	-.0015, .0059
C → cohort size → education → PA	-.0099	-.0238, .0043
C → cohort size → education → employment → PA	-.0001	-.0003, .0001
C → education → PA	-.1470	-.2321, -.0676
C → education → employment → PA	-.0012	-.0044, .0009
Total	-.1565	-.2446, -.0791
Age effect (37-40 vs. 49-52)		
A → church attendance → PA	.0077	-.0157, .0310
Total	.0077	-.0157, .0310
Grand total	.7453	.5590, .9214

# Hmm...

- ▶ No causal estimand or proofs
  - ▶ not clear what is being estimated, i.e., how to interpret the obtained estimates
- ▶ Product-of-coefficient method is only valid for linear models
  - ▶ Winship and Harding (2008) used logistic and probit models
- ▶ All measured mediators used in the analysis
  - ▶ requires extensive regression modeling, which makes the analysis vulnerable to model misspecification bias

# Outline

Motivating example (Winship & Harding, 2008)

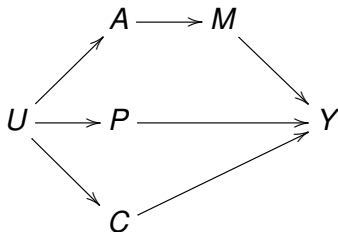
Nonparametric identification

Parametric estimation

Motivating example, revisited



## A simple example



- ▶ The role of  $U$  is to enforce the relation  $A = P - C$
- ▶ Key assumptions
  - ▶ **no unmeasured confounding**: no common causes of  $(A, P, C)$  and  $(M, Y)$
  - ▶ **partial exclusion**:

$$M \perp_d \{P, C\} | A$$

$$Y \perp_d A | \{P, C, M\}$$

where  $V_1 \perp_d V_2 | V_3$  denotes  $V_1$  and  $V_2$  d-separated by conditioning on  $V_3$

## Causal estimand (Bijlsma et al., 2017)

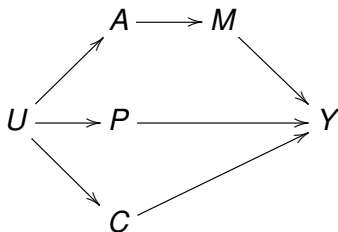
- ▶ Let  $Y(a, p, c)$  be the potential outcome for a given individual, if age, period and cohort were set to  $a$ ,  $p$  and  $c$ 
  - ▶ not necessarily obeying the deterministic relation  $a = p - c$
- ▶ Let  $E[Y(a, p, c)]$  be the mean potential outcome, if age, period and cohort were set to  $a$ ,  $p$  and  $c$  for all individuals
- ▶ Causal effect of taking  $A$  from  $a$  to  $a'$ , while holding  $(P, C)$  fixed at  $(p, c)$ :

$$E[Y(a', p, c)] - E[Y(a, p, c)]$$

# Quite controversial!

- ▶ An intervention that sets  $A$ ,  $P$ , and  $C$  to fixed values is highly hypothetical
  - ▶ e.g., we have no time machines
- ▶ Is this a problem? A longstanding debate in causal inference!
  - ▶ Hernán (2005) and Holland (1986): one should only consider a counterfactual as meaningful if one can specify a practically feasible intervention that would make the counterfactual observable
  - ▶ Pearl (2018): *'counterfactuals and causal effects are defined independently of those [practically feasible] interventions and therefore, are not to be denied existence, or rendered "inconsistent" by the latter's imperfections'*; similar views in Glymour and Glymour, 2014
- ▶ What is the alternative?

# Nonparametric identification (Sjölander & Gabriel, 2025)



$$\begin{aligned} E[Y(a, p, c)] &\stackrel{(1)}{=} E^*(Y|A = a, P = p, C = c) \\ &= \sum_m \left\{ E^*(Y|A = a, P = p, C = c, M = m) \right. \\ &\quad \left. \times p^*(M = m|A = a, P = p, C = c) \right\} \\ &\stackrel{(2)}{=} \sum_m \underbrace{E(Y|P = p, C = c, M = m)p(M = m|A = a)}_{\text{identifiable from } p(A, P, C, M, Y)} \end{aligned}$$

- ▶  $E^*(\cdot)$ ,  $p^*(\cdot)$ : hypothetical world where  $(A, P, C)$  vary freely
- ▶ (1): no unmeasured confounding
- ▶ (2): partial exclusion

# General result: data and assumptions

- ▶ Data:  $A, P, C, \mathbf{M}_K = (M_1, \dots, M_K), Y$
- ▶ Assumptions:
  - ▶ **no unmeasured confounding**: no common causes of  $(A, P, C)$  and  $(\mathbf{M}_K, Y)$
  - ▶ **partial exclusion**: there exists proper subsets  $\{R_1(A, P, C), \dots, R_K(A, P, C)\}$  and  $R_Y(A, P, C)$ , with complements  $\{R'_1(A, P, C), \dots, R'_K(A, P, C)\}$  and  $R'_Y(A, P, C)$ , such that

$$M_k \perp_d R'_k(A, P, C) | \{R_k(A, P, C), \mathbf{M}_{k-1}\} \text{ for } k = 1, \dots, K$$

and

$$Y \perp_d R'_Y(A, P, C) | \{R_Y(A, P, C), \mathbf{M}_K\}$$

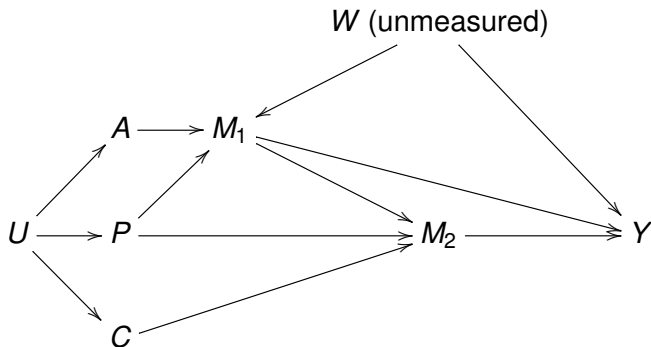
## General result: the APC-formula

- ▶ Under the no unmeasured confounding and partial exclusion assumptions,  $E[Y(a, p, c)]$  is identified as

$$E[Y(a, p, c)] = \sum_{\mathbf{m}_K} \left\{ E[Y | R_Y(a, p, c), \mathbf{M}_K = \mathbf{m}_K] \right. \\ \left. \times \prod_{k=1}^K p[M_k = m_k | R_k(a, p, c), \mathbf{M}_{k-1} = \mathbf{m}_{k-1}] \right\}$$

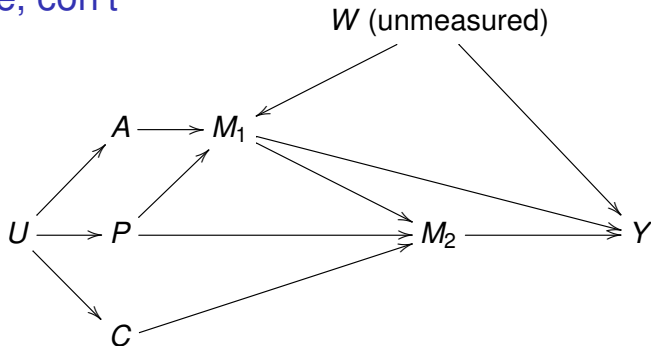
- ▶ Each term contains at most two of  $(a, p, c)$ , so there is no identification problem
- ▶ Similar to the nonparametric G-formula in longitudinal studies with time-varying exposures and confounders
  - ▶ Robins (1986)

## Example: identification with two mediators and unmeasured mediator-outcome confounding



- ▶ Cannot use  $M_1$  alone for identification, since  $Y$  is not d-separated from any of  $(A, P, C)$  by conditioning on  $M_1$
- ▶ Cannot use  $M_2$  alone for identification, since  $M_2$  is not d-separated from any of  $(A, P, C)$

## Example, con't



- Can use  $M_1$  and  $M_2$  together for identification since

$$M_1 \perp_d C$$

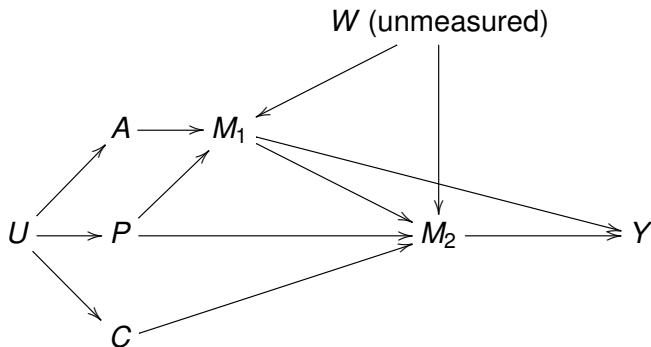
$$M_2 \perp_d A | \{P, C, M_1\}$$

$$Y \perp_d C | \{A, P, M_1, M_2\}$$

$$\begin{aligned}
 E[Y(a, p, c)] &= \sum_{m_1, m_2} \left\{ E(Y | \underbrace{A = a, P = p, M_1 = m_1, M_2 = m_2}_{=R_Y(a, p, c)}) \right. \\
 &\quad \times \left. p(M_2 = m_2 | \underbrace{P = p, C = c, M_1 = m_1}_{=R_2(a, p, c)}) p(M_1 = m_1 | \underbrace{A = a, P = p}_{R_1(a, p, c)}) \right\}.
 \end{aligned}$$

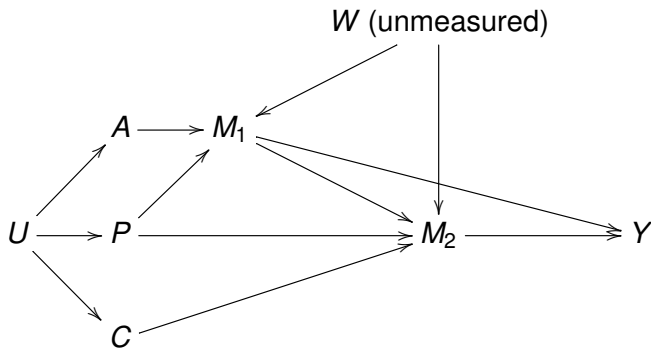


## Example: non-identifiability due to mediator-mediator confounding



- ▶ Cannot use  $M_1$  alone for identification, since  $Y$  is not d-separated from any of  $(A, P, C)$  by conditioning on  $M_1$
- ▶ Cannot use  $M_2$  alone for identification, since  $M_2$  is not d-separated from any of  $(A, P, C)$

## Example: non-identifiability due to mediator-mediator confounding, con't



- ▶ Cannot use  $M_1$  and  $M_2$  together for identification, since
  - ▶  $M_2$  is not d-separated from any of  $(A, P, C)$  by conditioning on  $M_1$
  - ▶  $M_1$  is not d-separated from any of  $(A, P, C)$  by conditioning on  $M_2$

# Outline

Motivating example (Winship & Harding, 2008)

Nonparametric identification

Parametric estimation

Motivating example, revisited

# Parametric models

$$p[M_k = m_k | R_k(a, p, c), \mathbf{M}_{k-1} = \mathbf{m}_{k-1}; \alpha_k] \quad \text{for } k = 1, \dots, K$$

$$E[Y | R_Y(a, p, c), \mathbf{M}_K = \mathbf{m}_K; \beta]$$

$$\begin{aligned} E[Y(a, p, c)] &= \sum_{\mathbf{m}_K} \left\{ E[Y | R_Y(a, p, c), \mathbf{M}_K = \mathbf{m}_K; \hat{\beta}] \right. \\ &\quad \times \left. \prod_{k=1}^K p[M_k = m_k | R_k(a, p, c), \mathbf{M}_{k-1} = \mathbf{m}_{k-1}; \hat{\alpha}_k] \right\} \end{aligned}$$

- ▶ Similar to parametric G-formula estimator of causal effects in longitudinal studies with time-varying confounding
  - ▶ Taubman et al. (2009) and Westreich et al. (2012)
- ▶ Low-dimensional  $\mathbf{M}$ : direct summation
- ▶ High-dimensional  $\mathbf{M}$ : Monte Carlo simulation

# Monte Carlo simulation

1. Fit models to obtain estimates  $\hat{\alpha}_1, \dots, \hat{\alpha}_K, \hat{\beta}$ , and let  $N$  be large integer, e.g.,  $N = 10,000$ .
2. For  $i = 1, \dots, N$ , repeat steps 3-5.
3. Define

$$\hat{m}_0^i(a, p, c) = \emptyset$$

4. For  $k = 1, \dots, K$ , define

$$\hat{\mathbf{m}}_{k-1}^i(a, p, c) = \{\hat{m}_0^i(a, p, c), \dots, \hat{m}_{k-1}^i(a, p, c)\}$$

and generate a prediction  $\hat{m}_k^i(a, p, c)$  as a random draw from the fitted model

$$p[M_k = m_k | R_k(a, p, c), \mathbf{M}_{k-1} = \hat{\mathbf{m}}_{k-1}^i(a, p, c); \hat{\alpha}_k]$$

5. Define the prediction

$$\hat{y}^i(a, p, c) = E[Y | R_Y(a, p, c), \mathbf{M}_K = \hat{\mathbf{m}}_K^i(a, p, c); \hat{\beta}]$$

6. Estimate  $E[Y(a, p, c)]$  as

$$\hat{E}[Y(a, p, c)] = \sum_{i=1}^N \hat{y}^i(a, p, c) / N$$

# Outline

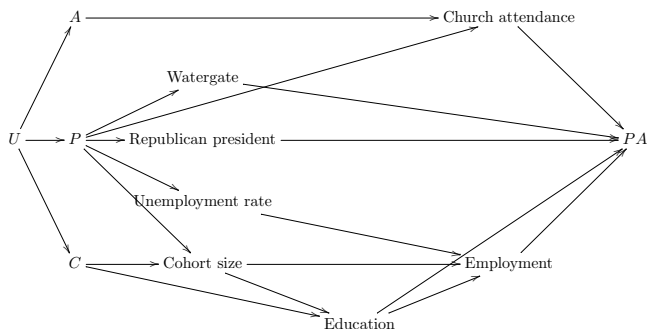
Motivating example (Winship & Harding, 2008)

Nonparametric identification

Parametric estimation

Motivating example, revisited

# Nonparametric identification



- It suffices to use church attendance for identification, since

$$\text{church attendance} \perp_d C | \{A, P\}$$

$$PA \perp_d A | \{P, C, \text{church attendance}\}$$

$$\begin{aligned} E[PA(a, p, c)] &= \sum_m \left\{ E(PA | P = p, C = c, \text{church attendance} = m) \right. \\ &\quad \times \left. p(\text{church attendance} = m | A = a, P = p) \right\}. \end{aligned}$$

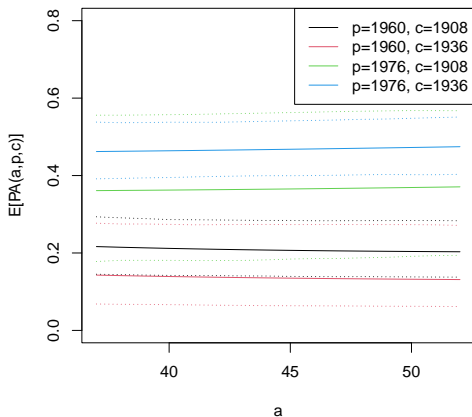
## Parametric estimation

$$E[PA(a, p, c)] = \sum_m \left\{ E(PA|P = p, C = c, \text{church attendance} = m) \right. \\ \left. \times p(\text{church attendance} = m|A = a, P = p) \right\}.$$

- ▶ Age: continuous (29-56)
- ▶ Period: categorical (1956, 1960, 1964, 1968, 1976, 1980)
- ▶ Cohort: continuous (1900-1951)
- ▶ Church attendance: categorical (never, seldom, regularly, often); multinomial logistic regression
- ▶ Political alienation: binary; ordinary (binomial) logistic regression
- ▶ Main effects and all two-way interactions in both regression models
- ▶ Direct summation over  $m$ , 95% confidence intervals with a nonparametric bootstrap, 1000 resamples

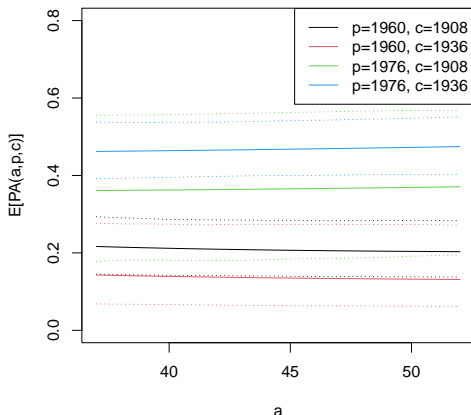


# Results



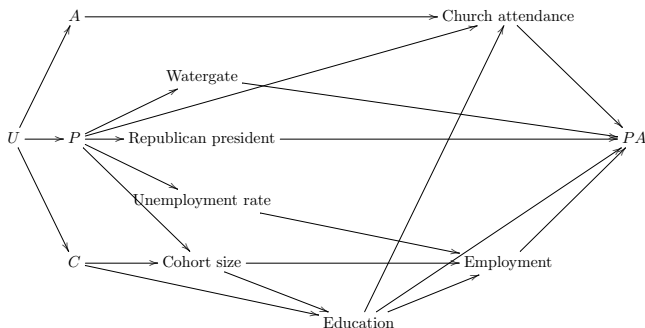
- ▶ Virtually no age effect
- ▶ Strong positive period effect:  
 $E[PA(a, 1976, c)] > E[PA(a, 1960, c)]$  for all  $\{a, c\}$
- ▶ Qualitatively similar to Winship and Harding (2008)

# Results



- ▶ Strong cohort effect
  - ▶ negative when  $p = 1960$ :  
 $E[PA(a, 1960, 1936)] < E[PA(a, 1960, 1908)]$  for all  $a$
  - ▶ positive when  $p = 1976$ :  
 $E[PA(a, 1976, 1936)] < E[PA(a, 1976, 1908)]$  for all  $a$
  - ▶ unnoticed by Winship and Harding (2008) since they had no period-cohort interactions

# Nonparametric identification, cont'd



$$\text{education} \perp_d A | \{P, C\}$$

$$\text{church attendance} \perp_d C | \{A, P, \text{education}\}$$

$$PA \perp_d A | \{P, C, \text{education}, \text{church attendance}\}$$

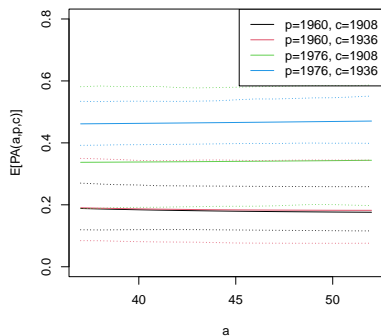
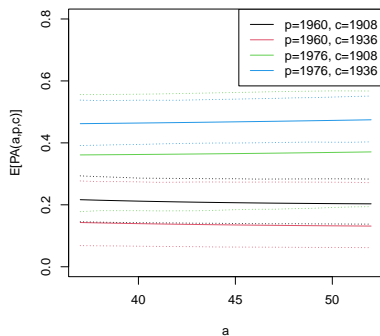
$$E[PA(a, p, c)]$$

$$= \sum_{m_1, m_2} \left\{ E(PA | P = p, C = c, \text{education} = m_1, \text{church attendance} = m_2) \right.$$

$$\times p(\text{church attendance} = m_2 | A = a, P = p, \text{education})$$

$$\left. \times p(\text{education} = m_1 | P = p, C = c) \right\}$$

# Results



► No cohort effect when  $p = 1960$

# Summary

- ▶ If we only care about statistical associations, then the deterministic age-period-cohort relation is not a problem
- ▶ If we care about causal effects, then it poses both conceptual and identifiability problems
- ▶ Under no unmeasured confounding and partial exclusion, nonparametric identification is given by the APC-formula
- ▶ Estimation can be carried out with parametric models, similar to parametric G-formula estimation

## Future work

- ▶ Less controversial causal estimand, without interventions on age-period-cohort?
- ▶ How to find the minimal sufficient set of mediators for identification?
  - ▶ work with Chihao Yan (student at biostat master program)
- ▶ Parametric model for  $E[Y(a, p, c)]$ , estimation with IPW?
  - ▶ work with Patrick Schnell, visiting associated professor from Ohio State University

# References I

- 
- Bijlsma, M., Daniel, R., Janssen, F., & De Stavola, B. (2017). An assessment and extension of the mechanism-based approach to the identification of age-period-cohort models. *Demography*, 54(2), 721–743.
- 
- Bollen, K. A. (1989). *Structural equations with latent variables*. John Wiley & Sons.
- 
- Fannon, Z., & Nielsen, B. (2019). Age-period-cohort models. In *Oxford research encyclopedia of economics and finance*.
- 
- Firebaugh, G., & Davis, K. E. (1988). Trends in antiblack prejudice, 1972-1984: Region and cohort effects. *American Journal of Sociology*, 94(2), 251–272.
- 
- Fosse, E., & Winship, C. (2019). Analyzing age-period-cohort data: A review and critique. *Annual Review of Sociology*, 45, 467–492.
- 
- Glenn, N. D. (1994). Television watching, newspaper reading, and cohort differences in verbal ability. *Sociology of Education*, 216–230.
- 
- Glymour, C., & Glymour, M. R. (2014). Commentary: Race and sex are causes. *Epidemiology*, 25(4), 488–490.
- 
- Heckman, J., & Robb, R. (1985). Using longitudinal data to estimate age, period and cohort effects in earnings equations. In *Cohort analysis in social research: Beyond the identification problem* (pp. 137–150). Springer.
- 
- Hernán, M. A. (2005). Invited commentary: Hypothetical interventions to define causal effects—afterthought or prerequisite? *American Journal of Epidemiology*, 162(7), 618–620.

# References II



Holland, P. W. (1986). Statistics and causal inference. *Journal of the American statistical Association*, 81(396), 945–960.



Knoke, D., & Hout, M. (1974). Social and demographic factors in american political party affiliations, 1952-72. *American Sociological Review*, 700–713.



Mason, K. O., Mason, W. M., Winsborough, H., & Poole, K. W. (1973). Some methodological issues in cohort analysis of archival data. *American Sociological Review*, 242–258.



Murphy, C. C., & Yang, C. Y. (2018). Use of age-period-cohort analysis in cancer epidemiology research. *Current Epidemiology Reports*, 5, 418–431.



Myers, D., & Lee, S. W. (1998). Immigrant trajectories into homeownership: A temporal analysis of residential assimilation. *International Migration Review*, 32(3), 593–625.



Nakamura, T. (1986). Bayesian cohort models for general cohort table analyses. *Annals of the Institute of Statistical Mathematics*, 38, 353–370.



O'Brien, R. M. (2000). Age period cohort characteristic models. *Social Science Research*, 29(1), 123–139.



Pearl, J. (2018). Does obesity shorten life? Or is it the soda? On non-manipulable causes.. *Journal of Causal Inference*, 6(2).



# References III



Robins, J. (1986). A new approach to causal inference in mortality studies with a sustained exposure period—application to control of the healthy worker survivor effect. *Mathematical Modelling*, 7(9-12), 1393–1512.



Sjölander, A., & Gabriel, E. E. (2025). A generalization of the mechanism-based approach for age–period–cohort models. *Epidemiology*, 36(2), 227–236.



Taubman, S., Robins, J., Mittleman, M., & Hernán, M. (2009). Intervening on risk factors for coronary heart disease: An application of the parametric g-formula. *International Journal of Epidemiology*, 38(6), 1599–1611.



Westreich, D., Cole, S., Young, J., Palella, F., Tien, P., Kingsley, L., Gange, S., & Hernán, M. (2012). The parametric g-formula to estimate the effect of highly active antiretroviral therapy on incident aids or death. *Statistics in Medicine*, 31(18), 2000–2009.



Winship, C., & Harding, D. J. (2008). A mechanism-based approach to the identification of age–period–cohort models. *Sociological Methods & Research*, 36(3), 362–401.