

1 **Abstract**

2 Many methods have been proposed to solve the age-period-cohort (APC) linear identification
3 problem, but most are not theoretically informed and may lead to biased estimators of APC
4 effects. One exception is the mechanism-based approach recently proposed and based on *Pearl's*
5 *front door criterion*; it ensures consistent APC effect estimators in the presence of a complete set
6 of intermediate variables between one of age, period, cohort and the outcome of interest, as long
7 as the assumed parametric models for all the relevant causal pathways are correct. Through a
8 simulation study mimicking APC data on cardiovascular mortality, we demonstrate possible
9 pitfalls that users of the mechanism-based approach may encounter under realistic conditions,
10 namely when 1) the set of available intermediate variables is incomplete; 2) intermediate
11 variables are affected by two or more of the APC variables, but this feature is not acknowledged
12 in the analysis 3) unaccounted confounding is present between intermediate variables and the
13 outcome. Furthermore, we show how the mechanism-based approach can be extended beyond the
14 originally proposed linear and probit regression models to incorporate all generalized linear
15 models, as well as non-linearities in the predictors, using Monte Carlo simulation. Based on the
16 observed biases that resulted from departures from underlying assumptions we formulated
17 guidelines for the application of the mechanism-based approach (extended or not).

18

19

20

1 **Introduction**

2 Demographers, epidemiologists, sociologists and others have attempted to break down outcomes
3 of interest into constituent effects caused by, or associated with, age, calendar time, and time of
4 birth. This is also known as age-period-cohort (APC) analysis. Age effects refer to changes in the
5 outcome as the age of individuals in the study population progresses. For example, as individuals
6 age, cardiovascular function declines and hence older individuals tend to have worse
7 cardiovascular health than younger individuals. Period effects refer to changes that occur in an
8 outcome as calendar time progresses. They can represent sudden changes or temporary changes
9 in an outcome, such as spikes in death rates due to war or famine, but may also represent gradual
10 changes such as those produced by the accumulation of minor improvements in public health
11 infrastructure over time influencing mortality rates in all age-groups. Finally, birth cohort effects
12 represent differences between generations that aren't attributable to differences in age or calendar
13 time. Conceptually, they commonly represent the effects of shared formative experiences of
14 individuals in a birth cohort, either *in utero* or during other critical phases in the life course (Ben-
15 Shlomo and Kuh 2002). The effect of these formative years would then remain largely constant
16 in that cohort throughout the remaining life course, and are therefore independent of age and
17 calendar time. For example, the cohort that was *in utero* during the Dutch Hunger Winter in
18 1944-1945 had worse health even later in life (Ekamper et al. 2015), compared to other cohorts.
19 Furthermore, birth cohort has been found to be strongly tied to smoking behaviour in various
20 Western countries (Preston et al. 2006; Verlato et al. 2006).

21 Unfortunately, decomposing an outcome (Y) into the separate effects of age (A), period
22 (P) and cohort (C) using for example a linear regression model (e.g. $E(Y|A,P,C) = \eta + \alpha \cdot A + \beta \cdot$
23 $P + \theta \cdot C$) imposes an identification problem. Since $A = P - C$ there is linear dependency between
24 the three variables, and consequently any linear model involving these three variables cannot

1 have a unique solution. To circumvent this, various techniques have been introduced that
2 constrain the model specification (e.g. Clayton & Schifflers 1987; Held and Riebler 2012;
3 Holford 2006; Yang et al. 2008), so that a solution can be found. However, the technical
4 constraints that have been proposed are arbitrary and do not lead to meaningful measures of
5 effect (Bell and Jones 2014; Luo 2013). Estimation of the parameters may be unbiased, but only
6 under the constraints that have been imposed, and hence do not reflect the ‘true effects’ (see
7 below) of age, period and cohort that we seek (Luo 2013). Luo and Fienberg both argue in favour
8 of a ‘paradigm shift’ in a recent discussion in *Demography* (Fienberg 2013; Kuhn 1996; Luo
9 2013); APC analysis, they argue, needs to become more theoretically informed. Simply fitting a
10 regression model to an outcome given age, period and cohort, without any forethought or
11 theoretical reasoning, cannot result in meaningful effect estimates for these variables.

12 Few authors explicate what they mean by ‘true’ or ‘meaningful’ effect estimates. Viewed
13 from one perspective, since the relationship $A = P - C$ always holds, then a regression model with
14 age, period and cohort as covariates truly has infinitely many solutions, and thus there is no
15 problem to be solved. However, since those who write about this problem talk about there being
16 one special solution out of those infinitely many solutions which is
17 ‘correct’/‘true’/‘valid’/‘meaningful’, it must be that they are (albeit implicitly) thinking of a
18 hypothetical world, different from the actual world, in which age, period and cohort can be
19 manipulated such that the identity $A = P - C$ is broken.

20 More formally, one could take an explicitly causal perspective using potential outcomes
21 with age, period and cohort as ‘independent’ exposures. Let $Y(a,p,c)$ be the potential outcome
22 which would occur if A were set to a , P to p and C to c , without necessarily abiding by the
23 relationship $a = p - c$ (Rubin 1974). Then the causal model

1
$$E(Y(a, p, c)) = \eta^* + \alpha^* \cdot a + \beta^* \cdot p + \theta^* \cdot c \quad (1)$$

2 has one solution, and this is presumably the ‘true’ solution to which the various authors on this
3 topic refer.

4 In fact, imagining a hypothetical world in which time can be manipulated is difficult
5 enough; contemplating one in which three different ‘aspects’ of time, namely age, period and
6 cohort, can be independently manipulated requires an even wilder imagination, and is therefore
7 unlikely to be truly of interest. More realistically, we can view Eq. 1. as being a shorthand for

8
$$E\left(Y(c_a, c_p, c_c)\right) = \eta' + \alpha' \cdot c_a + \beta' \cdot c_p + \theta' \cdot c_c \quad (2)$$

9 where c_a , c_p and c_c are the set of all immediate *consequences* of age, period and cohort,
10 respectively. So if being born in a particular cohort meant that you would be born during a
11 famine, it is this famine that we imagine we could manipulate (and, say, prevent), rather than the
12 cohort of birth itself. But since we may not have all these consequences at our disposal, model (2)
13 is replaced (as a shorthand) by model (1).

14 Given this reframing of the model of interest as a causal model, it makes sense to consider
15 methods from causal inference (Pearl 2000) to analyse data from APC studies. This was done by
16 Winship & Harding (Winship and Harding 2008) in 2008, and was called the mechanism-based
17 approach. In particular, their approach uses *Pearl’s front door criterion* to identify the APC
18 causal parameters α^* , β^* , and θ^* (Pearl 2000). In short, the mechanism-based approach uses
19 intermediate variables on the path between one of the three APC variables and the outcome in
20 order to estimate the effect of one of these three variables on the outcome indirectly, and to
21 estimate the effect of the remaining two APC variables directly (with the method generalizable to
22 modelling intermediate variables for two of the three APC variables). The approach naturally
23 leads to drawing a directed acyclic graph (Glymour 2006; VanderWeele et al. 2008) depicting the

1 assumed relationships between A, P, C, the intermediate variables being considered, and the
2 outcome. It thus motivates researchers to be explicit about their substantive assumptions.

3 The method requires that a complete set of intermediate variables can be found for at least
4 one of the three APC variables. By a ‘*complete* set of intermediate variables’ for A, for example,
5 we mean a set of variables M_1, M_2, \dots, M_K that are affected by age and which themselves affect
6 the outcome Y in such a way that *all* of the effect of A on Y is via this set of intermediate
7 variables.

8 However, in a realistic setting, finding a complete set of intermediate variables even for
9 just one of the three APC variables is unlikely. Also, the partial set of intermediate variables that
10 may be available could be dependent on more than one APC variable. Furthermore there may be
11 variables that affect both the intermediate variable(s) and the outcome. All these settings (if they
12 cannot somehow be accounted for) threaten the mechanism-based approach with bias, and one of
13 the aims of this paper is to demonstrate these potential sources of bias and their magnitude in
14 these realistic scenarios.

15 Another challenge also arises; the mechanism-based approach has been developed for use
16 in linear and probit regression models for Y, and linear and probit regression models for the
17 intermediate variables M_1, M_2, \dots, M_K . While some analytical solutions (e.g. Winship and Mare
18 1983) could be adopted to extend this approach to deal with logistic regression models for
19 outcome and/or mediators, they are complex to implement. Moreover, only approximate methods
20 are available to deal with settings where the variables included in the outcome model interact or
21 have some other non-linear effects, even when Y, M_1, M_2, \dots, M_K are all continuous and
22 modelled using linear models (Preacher and Hayes 2008; VanderWeele 2015; Jiang and
23 VanderWeele 2015).

1 In this paper, in order to illustrate possible pitfalls one may encounter using mechanism-
2 based APC models, we assess their performance under realistic settings, namely when 1) only a
3 partial set of mediators is available; 2) some of the intermediate variables are affected by two or
4 more of the APC variables, but this feature is not acknowledged in the analysis; 3) unmeasured
5 confounding affects the intermediate variables and the outcome. Furthermore, we extend the
6 mechanism-based approach to settings with any fully-parametric model for the outcome and
7 intermediate variables by approximating the estimation of the APC parameters by Monte Carlo
8 simulation. R code demonstrating the mechanism-based approach, and its extension, is available
9 as supplementary material.

10

11

12

1 **Methods**

2 *The mechanism-based approach*

3 The mechanism-based approach is based on exploiting the fact that age, period and cohort affect
4 the outcome through intermediate variables (hereafter called ‘mediators’) (Winship and Harding
5 2008). The key idea is that while age, period and cohort are deterministically related, the
6 intermediate variables along the paths from these to the outcome will be affected by other (APC-
7 independent) causes, and hence can be used (if measured) to circumvent the identification
8 problem. We now discuss this in more detail.

9 Consider for simplicity the setting depicted in Figure 1, which shows a causal directed
10 acyclic graph (DAG). Causal DAGs are formal graphical representations of the assumed causal
11 relationships between the variables under study (Glymour 2006). Here the number of mediators
12 K is equal to 2, and the mediators M_1 and M_2 being considered lie on causal pathways from P to
13 Y . Note that there is no arrow from P to Y in the DAG, representing the assumption that all of
14 the effect of P on Y is via M_1 and M_2 . Also note that there is no arrow from either of A or C to
15 the mediators, nor shared common causes of the mediators and any other variables in the DAG.
16 Finally, the two paths from P to Y are ‘separate’, in the sense that M_1 does not affect M_2 nor vice
17 versa, and neither does any variable along either path affect a variable on the other path, nor share
18 any common causes: M_1 and M_2 are assumed to be conditionally independent given P . These
19 strong structural assumptions concerning the roles of M_1 and M_2 (some of which can be relaxed –
20 more on this later) allow the identification of the APC effects in a two-stage procedure as we now
21 describe. In the first step, separate models for each mediator on P are fitted. In the second step, a
22 model for the outcome Y on A , C , M_1 and M_2 is fitted. If Figure 1 is correct and the outcome and
23 the mediators are continuous variables and modelled using linear regression, or the outcome and
24 mediators are binary and modelled using probit regression models, and none of these models in

1 truth includes product terms or other non-linearities, then the effect of A and C on Y (α^* and θ^*)
2 are equal to their regression coefficients in the outcome (second-step) model, while the effect of
3 P on Y (β^*) is equal to the sum of the effects along the two pathways involving M_1 and M_2 . The
4 effect along a pathway is equal to the product of two regression coefficients; for the P – M_1 – Y
5 pathway, it is the product of the coefficient for P in the (first-step) regression of M_1 on P and the
6 coefficient for M_1 in the (second-step) regression of Y (that also includes M_2 , A and C as
7 covariates). Similar calculations apply to the P – M_2 – Y pathway, and the effects along the two
8 pathways are then summed to obtain the effect of P on Y. These calculations are an application of
9 the path tracing rules that are widely used in structural equation modelling (Mulaik 2009; Wright
10 1934), with standard errors for the estimated effect of P estimated using the delta method (in
11 simple settings) or more generally the bootstrap (MacKinnon 2008). See the supplemental
12 material for an applied example of the path tracing rule.

13

14 *Complications*

15 In a real life setting, a number of situations may occur that make mechanism-based estimation of
16 APC effects less straightforward. Firstly, if a complete set of mediators is not available for the
17 selected APC variable(s), then the effect estimators of the three APC variables described above
18 will be biased for α^* , β^* and θ^* because the required assumption that (at least) one of the three
19 APC variables is fully mediated by a set of measured mediators would not be met. Secondly, a
20 variable that we believe to be a mediator for one of the three APC variables may actually be a
21 mediator for more than one (for example, daily caloric intake is probably affected by age as well
22 as period). In this case, the regression coefficient for the APC variable that we did not believe to
23 be mediated by any of the mediators for P, M_1 , M_2 , ... M_K , will only capture the component of its
24 effect that is not mediated. Thirdly, the relationship between mediators and outcome may be

1 confounded, i.e. there may be a variable (either measured or unmeasured) that has a causal effect
2 on both the mediator(s) and the outcome. If this confounding is not controlled for in the outcome
3 model, the effect of the mediator on the outcome will be estimated with bias, and consequently
4 also the effect of the APC variable that is assumed to be mediated by it. Finally, the outcome
5 and/or the mediators may not be of a type that can be modelled by linear or probit regression, or
6 even if they can, when the models require product terms or other non-linearities. In this situation,
7 the path tracing rule needed to derive the causal effect of the mediated APC variable cannot be
8 used (Mulaik 2009).

9

10 *Simulation: approach*

11 We assess the mechanism-based approach through simulations. In our simulations, we attempt to
12 recreate a realistic setting in which APC analyses are performed, namely the study of
13 cardiovascular mortality. However, in order to demonstrate particular pitfalls we isolate sources
14 of bias in the APC effect estimates, and therefore simplify this real world setting into three
15 scenarios.

16 In each scenario, the study population and outcome are the same. The individuals in the
17 study population are generated to be aged between 40 and 95 years during the calendar years
18 1990 to 2015, and hence the whole data consists of birth cohorts ranging from 1896 to 1975. Age
19 and period for each record are generated according to uniform distributions but then categorised
20 into 5-year groups (for A and P), and cohort dependent on these categories ($C = P - A$). The
21 outcome is mortality due to cardiovascular disease (CVD), coded as 1 (death due to CVD) or 0
22 (alive or dead from other causes). It was generated in all scenarios according to either a probit or
23 logistic regression model, with the probability of CVD death generated as a function of age,
24 cohort and the period mediators. Age is set to account for 70% of the effect of the APC variables

1 on cardiovascular mortality, period (via its mediators) for 20%, and birth cohort for 10%. We
2 believe these percentages approximately correspond to realistic effects of age, period and cohort
3 in the period 1990 to 2015 in Western countries; the difference in incidence of CVD death
4 between individuals aged 40 and those aged 95 is very large, whereas the difference in incidence
5 of CVD death in these age groups between the year 1990 and 2015 is much smaller (Peeters et al.
6 2011). Due to the linear dependency phenomenon, it is unknown what part of these differences is
7 in truth attributable to each dimension.

8

9 *Simulation: mediators and confounders*

10 We simulate settings where P is the variable which has measured mediating variables. Results
11 however easily generalise to the alternative scenarios where A or C play this role (with due
12 numerical differences given their unequal assumed strength of effects). Four mediators on the
13 path from P to Y are included in the simulation study. They are body mass index (BMI),
14 smoking, statin therapy, and ‘unmeasured’. We choose the former three variables because they
15 are commonly described variables that are believed to affect CVD mortality, but many other
16 variables are also believed to affect CVD mortality (e.g. Blackmore et al. 2015; Capewell et al.
17 2000), which are represented by the ‘unmeasured’ variable. Together, these four variables
18 account for the entire period effect on the outcome. The ‘unmeasured’ and BMI variables are
19 continuous, whereas smoking and statin therapy are binary. We set each of the ‘measured’
20 mediators to account for ~20% of the period effect on cardiovascular mortality, while we set the
21 ‘unmeasured’ mediator to account for ~40%. Table 1 shows the direction of the effects of period
22 on the mediators, and of each mediator on the outcome; the effect of period on the ‘unmeasured’
23 mediators is linear, while its effect on the measured mediators is non-linear but monotonically
24 increasing or decreasing. Initially, these variables are made to act as mediators only on the path

1 between P and CVD mortality, but in some scenarios they are also affected by A or C, in which
2 case the total effects of these latter variables changes. In the scenario that includes confounding
3 (see below), the confounder is presence or absence of a particular gene, randomly assigned to be
4 present in 50% of individuals and set to have a positive effect on both the mediators and the
5 outcome (Sabol et al. 1998; Smith et al. 2015).

6

7 *Simulation: scenarios and variants*

8 In all simulations, we generate data for 100,000 individuals, each measured once. We simulate
9 three scenarios, with each scenario simulated 1000 times. In the data generating process for the
10 first scenario ('simple'), A and C have a direct effect on Y while only the effect of P is mediated.
11 In the second scenario ('more causes'), there is an additional (negative) effect of A on the
12 mediator BMI, which amounts to roughly 30% of the total age effect, and an additional (positive)
13 effect of C on smoking, which amounts to roughly 30% of the total cohort effect. Finally, in the
14 third scenario ('confounding'), genotype confounds the relationship between BMI and Y and
15 between smoking and Y (Figure 2); genotype has a positive effect on both BMI and smoking, and
16 a positive effect on CVD mortality. Genotype accounts for roughly 33% of the association
17 between age and CVD death and 35% of association between cohort and CVD death.

18 Each scenario has two variants. In the first variant, we generate Y and the binary
19 mediators using probit regression models, while in the second variant logistic regression models
20 are used instead. In both variants, linear regression models are used to generate continuous
21 variables. Since probit and logistic regression models transform parameters into probabilities in a
22 different way, the probabilities of cardiovascular mortality somewhat differ between the two
23 variants. We varied the value for the intercept in each variant so that the age-specific probabilities
24 of cardiovascular mortality were similar to those found in high income countries. However, due

1 to these differences in transformation, the extent of the bias found with the variants is not directly
2 comparable.

3 Finally, to demonstrate how the unequally distributed strength of age, period and cohort
4 affect the bias, we perform two sets of simulations in which we vary these strengths. In the first
5 set, we vary the size of the period effect from 0 to 100% in 20% increments, while
6 correspondingly reducing the size of the age effect, and keeping the size of the cohort dimension
7 constant at 10% (excluding the last increment, where cohort is necessarily set to 0%). The second
8 set is identical, but then the cohort effect size is varied and period kept constant at 20%. In both
9 sets, bias is generated by removing the first mediator ('unmeasured') from the estimation model
10 in the 'simple' scenario. These simulations were done with probit, logit and linear regression
11 variants. The logit and linear regression variants, plus a third set of simulations in which the age
12 effect size is varied, are described in the supplemental material.

13

14 *Simulation: estimation*

15 Estimation of the APC effects according to the mechanism-based approach consists of fitting
16 separate regression models for CVD death as function of A, C and the P-mediators and of each of
17 the mediators as functions of P. In all estimation models A, P and C are treated as categorical
18 variables through dummy coding (10 dummies for age, 4 dummies for period and 14 dummies
19 for cohort because cohort categories were forced to overlap in order to maintain the linear
20 identity, $C = P - A$). By generating cohort in this way, we maintain the linear dependency
21 between age, period and cohort and therefore follow the equal interval width definition (Luo and
22 Hodges 2016).

23 Following how we generated the variables in each variant, in the probit variant we use
24 probit regression models for CVD death and logistic regression models for the logistic variant.

1 The effects of P (and A and C when appropriate) on continuous mediators are estimated using
2 linear regression and for binary mediators using logistic or probit regression models. In all three
3 scenarios, we first perform our estimation under entirely correct assumptions; i.e. in the second
4 scenario we also model the additional paths from A and C to their mediators (as described in the
5 supplementary material), and in the third scenario we also control for the confounder. Moreover
6 in all scenarios the parametric forms used to fit these models is the same as that used to generate
7 the data. Then, to investigate the effect of incorrect assumptions, in the second scenario we omit
8 to model the path from A and C to the mediators, and in the third scenario we do not control for
9 confounding. Additionally, to explore the effect of including an incomplete set of mediators, in
10 all three scenarios we first remove the ‘unmeasured’ mediators from the estimation model, then
11 ‘BMI’, followed by ‘smoking’, and finally we remove all period mediators (i.e. fit an age-cohort
12 model). For completeness we also report the results of adopting a more traditional APC-approach
13 in the supplemental material.

14 Results obtained from each scenario, variant and model specification are summarised as
15 means of each parameter’s estimates over the 1000 simulations, which we compare to the
16 estimated values obtained from the correctly specified models to estimate the bias. Due to the
17 very large sample size, these estimates obtained from the correct models can be interpreted as the
18 true values.

19

20 *Extending the mechanism-based approach through Monte Carlo integration*

21 When a model with a general non-linear link function is used for a mediator or for the outcome
22 (or both), e.g. Poisson or logistic regression, or if the models include product terms or other non-
23 linearities, then the path tracing method cannot be used (Mulaik 2009). A different approach is
24 then required.

1 The basic intuition of our approach is as follows. First, similar to the traditional approach,
2 we estimate individual relations between age, period, cohort and mediators, and then between
3 mediators and outcome (step 1 and step 2, respectively). The difference now being that the
4 statistical models used for these steps are allowed to have non-linear functional forms. However,
5 this enhanced flexibility comes at a price; the traditional multiplication of coefficients along
6 pathways is no longer possible. Therefore, Monte Carlo integration is used instead (Robert and
7 Casella 2004); the coefficients from step 1 and 2 are used in step 3 and 4 to generate a new
8 dataset that does not suffer from APC linear dependency (*Pearl's front door* criterion makes this
9 possible) but that reflect the original data structure (more details below). In step 5, an APC model
10 is then fitted to this newly created dataset to provide estimates of the effects of age, period and
11 cohort. Our approach is, in many ways, analogous to the Monte Carlo estimation of the
12 parametric G-formula (Keil et al. 2014, Hernán and Robins 2013). A similar approach has also
13 been suggested in mediation analysis (VanderWeele 2015).

14 We treat age, period and cohort as continuous variables to simplify the presentation
15 below, but in our simulations we model age, period and cohort as categorical variables.
16 Furthermore, we describe here only the case where there is one mediator for period, and where
17 the mediator has only one cause. See the supplementary material for a description of a more
18 general setting. We proceed as follows:

19 Step 1) Mediator estimation: fit a model for the mediator. If it is continuous, we can use
20 linear regression, e.g.

$$21 \qquad M = \gamma_0 + \gamma_1 \cdot P + \nu$$

22 Where we assume $\nu \sim N(0, \sigma_M^2)$ Note that the assumption on the distribution of the error terms in
23 the linear regression model is non-trivial if there are non-linearities involving M in the model for
24 Y. If instead a mediator is binary, we can use logistic regression, e.g.

1
$$\text{logit}\{E(M|P)\} = \gamma_0 + \gamma_1 \cdot P$$

2 Let $(\hat{\gamma}_0, \hat{\gamma}_1)$ be the estimates of (γ_0, γ_1) from the appropriate model. If the mediator is continuous,
 3 also save the estimate of the error variance $\hat{\sigma}_M^2$. These estimates will be used in step 3.

4 Step 2) Outcome estimation: fit a model for the outcome. If the outcome is continuous, fit
 5 a model using linear regression, e.g.

6
$$Y = \delta_0 + \delta_1 \cdot A + \delta_2 \cdot C + \delta_3 \cdot M + \xi$$

7 where $\xi \sim N(0, \sigma_Y^2)$. Whereas if the outcome is binary, use logistic regression, e.g.

8
$$\text{logit}\{E(Y|A, M, C)\} = \delta_0 + \delta_1 \cdot A + \delta_2 \cdot C + \delta_3 \cdot M$$

9 Let $(\hat{\delta}_0, \hat{\delta}_1, \hat{\delta}_2, \hat{\delta}_3)$ be the estimates of $(\delta_0, \delta_1, \delta_2, \delta_3)$ from the appropriate model. If the outcome
 10 is continuous, also save the estimate of the error variance $\hat{\sigma}_Y^2$. These estimates will be used in step
 11 4.

12 Step 3) Mediator simulation: for each of a range of period values \tilde{p} , simulate the mediator
 13 $\tilde{M}(\tilde{p})$. The values \tilde{p} could be randomly generated for example using a discrete uniform
 14 distribution, but its range should be equal to the range empirically observed in the data that were
 15 used for estimation. For example, if we have data ranging from 1990 to 2015, then that would be
 16 the range of values for \tilde{p} that we use. If instead we categorised this into 5-year periods (1990-
 17 1994, 1995-1999, etc.), then we generate values of \tilde{p} corresponding to these categories. If a
 18 mediator is continuous, use the estimates of the linear regression model in step 1 to simulate

19
$$\tilde{M}(\tilde{p}) = \hat{\gamma}_0 + \hat{\gamma}_1 \cdot \tilde{p} + \hat{v}$$

20 Where \hat{v} is randomly drawn from $N(0, \hat{\sigma}_M^2)$. If instead the mediator is binary, use the estimates
 21 from the logistic regression model in step 1 to simulate $\tilde{M}(\tilde{p})$ from a Bernoulli distribution with
 22 mean

23
$$\hat{\mu}_{m,\tilde{p}} = \frac{\exp(\hat{\gamma}_0 + \hat{\gamma}_1 \cdot \tilde{p})}{1 + \exp(\hat{\gamma}_0 + \hat{\gamma}_1 \cdot \tilde{p})}$$

1 The number of $\tilde{M}(\tilde{p})$ values to be simulated need not be equal to the number of observations in
 2 the data, as long as the entire empirical range is covered, but the more values we simulate, the
 3 less our final estimates will be affected by Monte Carlo error. The values of $\tilde{M}(\tilde{p})$ will be used in
 4 step 4.

5 Step 4) Outcome simulation: for each of a range of age, period and cohort values $(\tilde{a}, \tilde{p}, \tilde{c})$,
 6 simulate the potential outcome $\tilde{Y}(\tilde{a}, \tilde{p}, \tilde{c})$. Since \tilde{p} is already generated in step 3, they can be re-
 7 used instead of re-generated. As previously, the range of these values should be equal to the
 8 range empirically observed in age, period and cohort in the data respectively. However, we
 9 choose \tilde{a} , \tilde{p} and \tilde{c} independently, i.e. the identity $\tilde{a} = \tilde{p} - \tilde{c}$ should not hold. If the outcome is
 10 continuous, then using the linear regression estimates of step 2, and the simulated mediator
 11 values of step 3, simulate

$$12 \quad \tilde{Y}(\tilde{a}, \tilde{p}, \tilde{c}) = \hat{\delta}_0 + \hat{\delta}_1 \cdot \tilde{a} + \hat{\delta}_2 \cdot \tilde{c} + \hat{\delta}_3 \cdot \tilde{M}(\tilde{p}) + \hat{\xi}$$

13 Where $\hat{\xi}$ is randomly drawn from $N(0, \hat{\sigma}_Y^2)$ and $\tilde{M}(\tilde{p})$ is as generated in step 3. If instead the
 14 outcome is binary, use the estimates from the logistic regression model in step 2, and the
 15 simulated mediator values of step 3, to simulate $\tilde{Y}(\tilde{a}, \tilde{p}, \tilde{c})$ from a Bernoulli distribution with
 16 mean

$$17 \quad \hat{\mu}_{y, \tilde{a}, \tilde{p}, \tilde{c}} = \frac{\exp(\hat{\delta}_0 + \hat{\delta}_1 \cdot \tilde{a} + \hat{\delta}_2 \cdot \tilde{M}(\tilde{p}) + \hat{\delta}_3 \cdot \tilde{c})}{1 + \exp(\hat{\delta}_0 + \hat{\delta}_1 \cdot \tilde{a} + \hat{\delta}_2 \cdot \tilde{M}(\tilde{p}) + \hat{\delta}_3 \cdot \tilde{c})}$$

18 Step 5) Estimate age, period and cohort effects using the simulated values $\tilde{Y}(\tilde{a}, \tilde{p}, \tilde{c})$. If
 19 $\tilde{Y}(\tilde{a}, \tilde{p}, \tilde{c})$ was generated as a continuous variable, use linear regression with \tilde{a} , \tilde{p} and \tilde{c} as
 20 covariates. If instead $\tilde{Y}(\tilde{a}, \tilde{p}, \tilde{c})$ was generated as a binary variable, use logistic regression with \tilde{a} ,
 21 \tilde{p} and \tilde{c} as covariates. These models will be identifiable because \tilde{a} , \tilde{p} and \tilde{c} should have been

1 chosen independently. The estimated parameters can be interpreted as the age, period and cohort
2 effects from the causal model (Eq. 1).

3 Step 6) Use the non-parametric bootstrap to estimate the standard errors for the parameter
4 estimates (Efron and Tibshirani 1993). This step consists of re-sampling with replacement from
5 the original data, a dataset of equal size, and then repeating steps 1 through 5, and saving the
6 parameter estimates of the effects of \tilde{a} , \tilde{p} and \tilde{c} at the end of step 5. The sample standard
7 deviations of the distributions of the \tilde{a} , \tilde{p} and \tilde{c} effect estimates can be used as estimates of the
8 standard errors, or for example the empirical 2.5% and 97.5% quantiles of these distributions can
9 be used to derive 95% confidence intervals directly (with improvements such as ‘bias-corrected
10 and accelerated’ intervals to be recommended).

11 The supplemental material includes R-code demonstrating the application of this
12 technique. In linear settings, this technique results in findings identical to the traditional path
13 tracing method, as long as a sufficiently large number of Monte Carlo simulations are carried out.

14

1 **Linear dependency and expected bias**

2 While age, period and cohort could be pairwise independent, the three variables together have a
3 deterministic relationship: $A = P - C$. This linear dependency will determine the (direction of the)
4 bias when the mechanism-based approach is used and not all the mediators are included in the
5 estimation model. We demonstrate this here.

6 In a linear setting, the causal model described in Eq. 1. could be expressed as

7
$$Y(a, p, c) = \eta^* + \alpha^* \cdot a + \beta^* \cdot p + \theta^* \cdot c + \varepsilon^*$$

8 Where α^* , β^* and θ^* represent the effects of age, period and cohort, η^* the intercept and ε^* the
9 mean-zero error terms. This model corresponds to an (unidentified) associational model:

10
$$Y = \eta + \alpha \cdot A + \beta \cdot P + \theta \cdot C + \varepsilon \quad (3)$$

11 Since $P = A + C$, we can rewrite Eq. 3. as

12
$$Y = \eta + \alpha \cdot A + \beta \cdot (A + C) + \theta \cdot C + \varepsilon \quad \Leftrightarrow \quad Y = \eta + (\alpha + \beta) \cdot A + (\theta + \beta) \cdot C + \varepsilon \quad (4)$$

13 Eq. 4. shows that if we fit an age-cohort model (omitting period) and interpret the coefficients of
14 A and C as representing the age and cohort effect respectively, the period effect is attributed to
15 age and cohort in equal parts, thereby biasing the presumed age and cohort effects in the direction
16 of the period effect. In the same way, since $A = P - C$, and $C = P - A$, we have

17
$$Y = \eta + (\beta + \alpha) \cdot P + (\theta - \alpha) \cdot C + \varepsilon \quad (5)$$

18
$$Y = \eta + (\alpha - \theta) \cdot A + (\beta + \theta) \cdot P + \varepsilon \quad (6)$$

19 when fitting period-cohort and age-period models, respectively. As Eq. 5. and Eq. 6. show, in the
20 period-cohort model, the period parameter is biased in the direction of the age effect, while the
21 cohort parameter is biased in the opposite direction. Similarly, in the age-period model, the
22 period parameter is biased in the direction of the cohort effect, while the age parameter is biased

1 in the opposite direction. Of course, in all three of these models, the effect of the omitted variable
2 would also be biased (unless its effect is null), as its effect estimate is effectively set to 0.

3 The same logic applies in the situation where we use mediators to estimate the effects of
4 age, period or cohort, while estimating the other two effects directly. Consider the earlier
5 described example where we used two mediators on the period path, while the effects of age and
6 cohort are estimated directly (Figure 1). However, this time, we omit the variable M_2 from our
7 estimation model, perhaps because it was not measured (Figure 3, left). When fitting a model for
8 the outcome (Y) using A , C and the measured period mediator (M_1) as explanatory variables,
9 there will be an association between A and the period mediators and between C and the period
10 mediators due to the aforementioned linear dependency (Figure 3, right). Since the model for Y
11 includes M_1 (i.e. we condition on M_1) the paths from A or C to Y via M_1 are blocked. However,
12 since we do not condition on the unmeasured mediator (M_2), the pathways from A and C to Y via
13 M_2 are not blocked, and hence their regression parameters will be biased by the contribution of
14 the additional paths from respectively A and C via M_2 (i.e. $\delta_4 \cdot \gamma_3$ for both).

15

1 **Results**

2 *Scenario 1: 'Simple'*

3 Removing mediators from the estimation of the period path introduces bias (Figure 4 for probit
4 variant, see supplementary material eFigure 3 for logistic variant). Initially, the magnitude of the
5 bias is relatively weak, but becomes stronger as more mediators are removed. The relative
6 magnitude appears strongest for birth cohort, which had weak negative parameter estimates when
7 the model was correctly specified, but these estimates become strongly negative when mediators
8 are removed from the estimation model. The directions of the bias are as expected (see section
9 'linear dependency and expected bias'). On average, removing the 'unmeasured' period
10 mediators has a negative effect on the estimated age and cohort effects, while it introduces a
11 positive bias on the period estimates; the estimated parameters for the age and cohort effects
12 become less positive, while those for the period effect become less negative. The same occurs
13 when smoking and statins are additionally removed. This is as expected, because the paths via the
14 'unmeasured' period mediators, smoking, and statins all have a negative effect on the outcome
15 (Table 1). Removing BMI leads to the opposite, i.e. it introduces a positive bias on the estimated
16 effects of age and cohort, and a negative one on that of period. Again, this is as expected because
17 the path from P to Y through BMI is positive (Table 1). Finally, removing all mediators (i.e.
18 fitting an AC model) results in the largest amount of bias (Figure 4).

19

20 *Scenario 2: 'More causes'*

21 In the scenario where age, in addition to period, is a cause of BMI, and cohort, in addition to
22 period, is a cause of smoking, we find that not including these relationships in the estimation
23 model also results in bias. As expected, this bias is largest for age, where the age effect is
24 overestimated when the negative effect of age on BMI is not included in the estimation model

1 (Figure 5). A similar, but much weaker bias is found for cohort when the effect of C on smoking
2 is not included (Figure 5 for the probit variant; supplementary material eFigure 4 for logistic
3 variant). There is no bias in the estimation of the period effect (Figure 5).

4 These results follow our expectation. The effect of age on the outcome via BMI is
5 negative. By not including age as a cause of BMI in the estimation model, this negative effect is
6 subtracted from the total age effect, thereby resulting in an overestimation of the age effect. The
7 effect of not modelling birth cohort as a cause of smoking follows the same logic, but the bias is
8 weaker because the effect size from birth cohort to the outcome via smoking is also smaller. In
9 this scenario, the period effect estimates are not biased because period is correctly modelled as a
10 cause of the four mediators. Removing mediators from the estimation model results in biases in
11 the same direction as found in the ‘simple’ scenario (see supplementary material eFigures 1 and
12 4), and the same explanations for these directions apply.

13

14 *Scenario 3: ‘Confounding’*

15 In the scenario where there is a variable (genotype) confounding the relationships between BMI
16 and CVD death, and between smoking and CVD death (both mediators on the P-Y path), we find
17 that failing to control for this confounder results in bias in all three age, period and cohort effect
18 estimates. The age parameter estimates become somewhat negatively biased when we do not
19 control for genotype in our outcome model (Figure 6 for the probit variant; supplementary
20 material Figure 5 for the logistic variant). The same occurs in the cohort effect, while the period
21 effect suffers from a small positive bias (Figure 6).

22 The bias observed here is caused by collider stratification (Cole et al. 2010; Elwert and
23 Winship 2014). Collider stratification bias occurs when two variables both have causal effects on
24 a third variable (the ‘collider’), and the collider is conditioned upon. In the scenario considered

1 here, genotype has a positive causal effect on BMI, on smoking, and on CVD mortality.
2 Therefore, BMI and smoking are colliders in the paths between P and genotype. By including
3 BMI and smoking (together with the other two mediators) in our outcome model, we create an
4 additional spurious association between A and Y via A – P – BMI – genotype – Y, and one
5 between C and Y via the path C – P – smoking – genotype – Y. Both of these spurious
6 associations are negative because those induced by collider stratification, P-BMI-genotype and P-
7 smoking-genotype, are both negative while A – P, C – P and genotype – CVD are all positive.

8

9 *Varying effect sizes*

10 Increasing the size of the period effect, and keeping the size of the cohort effect constant, resulted
11 in increased bias when the ‘unmeasured’ mediator was removed from the estimation model
12 (Figure 7 for the probit variant, eFigures 6 and 7 in the supplemental material for the logistic and
13 continuous variants). Keeping the period effect constant and increasing the size of the cohort
14 effect, resulted in roughly equal amounts of bias in the cohort effect estimate (Figure 7 and
15 eFigures 6 and 7 in the supplemental material); the bias was equal in the continuous variant, as
16 that estimation is done without link function transformations.

17

1 **Discussion**

2 We assessed the performance of mechanism-based APC models in realistic settings, and extended
3 the method to incorporate all generalized linear models, as well as non-linearities in the
4 predictors, using Monte Carlo simulation. We found that, in simple scenarios where there is a
5 single set of mediators for P which is not affected by unmeasured confounding, the mechanism-
6 based approach (extended or not) performed reasonably well for the estimation of all three age,
7 period and cohort effects, especially if the mediators that were not included in the estimation had
8 opposite signs (e.g. ‘unmeasured’ and ‘BMI’) so that their unmeasured contributions (partially)
9 cancelled out. The bias we observed was in line with our expectations derived from the APC
10 linear dependency. In the scenarios with additional complications, i.e. either when mediators
11 were caused by more than one of three APC components, or there was unmeasured confounding
12 of the mediator-outcome relationships, we found additional bias. This, again, was in line with our
13 expectations. Findings were similar for the probit and logistic variants, though we did not directly
14 compare these methods due to their differences in transforming effects into probabilities.

15

16 *(Un)testable assumptions*

17 APC models solve the linear dependency problem by imposing modelling constraints by fiat
18 (Fienberg 2013). The mechanism-based approach does not differ from this, as it is (commonly,
19 see next paragraph) not possible to identify from the data whether a variable is a mediator for
20 age, period, cohort, a combination of these two, or all three variables. This is, of course, the same
21 problem that occurs in any conventional APC analysis when attempting to decompose some
22 outcome into age, period and cohort effects. Therefore, the untestable assumptions that are made
23 in conventional APC analysis move to the mediator stage of the modelling procedure.

1 It is possible to statistically test the addition of mediators to the APC model (Winship and
2 Harding 2008). We have omitted this test from our assessment because the test is conditional on
3 having (at least) a full set of mediators for one of the three APC variables. We consider it likely
4 that in the majority of real-life applications it will not be possible to find a full set of mediators
5 for even one of the APC variables, and hence have chosen to focus our assessment on possible
6 biases that may be encountered, such as due to missing mediators.

7

8 *Simulating bias*

9 In our simulations, we considered 4 mediators on the period path, and kept the ‘more causes’
10 scenario separate from the ‘confounding’ scenario. This was done in order to illustrate,
11 separately, possible pitfalls which may be encountered when the mechanism-based approach is
12 used. In a real application of the method, many more mediators may exist, which may also have
13 more than one cause, and their relation with the outcome may be confounded. However, a more
14 complicated scenario need not necessarily result in more biased estimates because biases of
15 opposed sign may cancel each other out, such as when BMI was removed from the estimation
16 model in our simulations (since BMI was the only positive mediator). Most importantly, when
17 assessing potential bias, it is the size of the omitted pathways that matters, relative to the size of
18 the included pathways. This was demonstrated in the analysis where we varied the size of the
19 period effect and omitted a period mediator from the estimation model; a larger period effect
20 resulted in a larger bias. Whereas keeping the size of the period effect constant and increasing the
21 size of the cohort effect resulted in roughly the same magnitude of bias and smaller relative bias.
22 Analogous simulations where instead age or cohort mediators had been removed would have
23 yielded the same conclusions; only the sign of the bias would differ as shown in the section on
24 linear dependency and expected bias.

1

2 *Confounding and colliding*

3 Because the mechanism-based approach uses a mediation approach to APC analysis, it is also
4 subject to biases that are not present in traditional APC analysis. These are confounding of the
5 mediator-outcome relationship, and collider stratification bias (Elwert and Winship, 2014,
6 Greenland 2003).

7 Arguably, traditional APC analysis is not affected by confounding because age, period
8 and cohort are time dimensions and therefore not causally affected by other variables. In our
9 confounding scenario, the effect of BMI (and smoking) on CVD mortality was confounded by
10 Genotype, which thereby also affected the estimation of the period effect. Ideally, such a
11 confounder would be controlled. However, if the confounder is unmeasured, and therefore cannot
12 be controlled, a difficult choice has to be made. Removing the mediator induces bias via an
13 omitted pathway, whereas including the mediator induces confounding bias. Confounding of the
14 mediator-outcome relationship (and hence collider stratification bias) resulted in more bias than
15 omitting the relevant pathway in our simulations, but this is dependent on the relative strengths of
16 mediation and confounding.

17 Collider stratification bias occurs when two variables both have causal effects on a third
18 variable (the ‘collider’), and the collider is conditioned upon (Elwert and Winship, 2014,
19 Greenland 2003). This is also known as endogenous selection. Doing so creates an artificial
20 association between the two causes of the collider which is of sign opposite to the product of the
21 signs of the effects into the collider. In traditional APC analysis, collider stratification does not
22 occur because APC effects on an outcome are estimated without adjusting for confounders or
23 mediators. In the “Confounding” simulations, BMI (and smoking) was affected by Period and by
24 Genotype. Due to the linear dependency of period with age and cohort, collider stratification bias

1 also affected the age and cohort estimates. Also in this scenario, conditioning on Genotype would
2 have prevented this bias. If conditioning on both causes of the collider is not possible, a choice
3 has to be made between removing and keeping the mediator in the analysis. Greenland's (2003)
4 second 'rule of thumb' regarding endogenous selection states that including mediators which are
5 also colliders induces larger bias than the size of the mediator's dependence on its causes (Elwert
6 and Winship, 2014). This suggests removing such mediators from the analysis, despite inducing
7 bias via omitted pathways.

8

9 *Directed Acyclic Graphs*

10 We encourage investigators to draw causal DAGs of the relations between APC variables,
11 mediators and outcome. To represent the deterministic relationship between age, period and
12 cohort in the causal DAG, we used bold arrows, which is in line with conventions in causal
13 inference (Spiegelhalter et al. 2002; Spirtes et al. 2001) but differs from the representation used
14 by Winship and Harding (2008). We used these arrows because the relationship between age,
15 period and cohort is fundamentally different from that of other relationships in the DAG (which
16 are stochastic and causal, rather than deterministic).

17 By drawing the relations between APC and mediators, we are clear about the assumptions
18 underlying our analyses. This sets the mechanism-based approach apart from other APC
19 approaches, where transparency about constraints can be lacking (Luo 2013). Drawing causal
20 DAGs, and re-drawing once one of the three age, period or cohort variables is removed (as was
21 done in Figure 3) helps explicate otherwise hidden assumptions about the relationships between
22 age, period, cohort and their mediators, and can help identify possible biases such as collider
23 stratification bias.

24

1 *Guidelines for application of the mechanism-based approach*

2 Based on our assessment, when estimating age, period and cohort effects by applying the
3 mechanism-based approach (extended or not), we suggest considering the following questions:

4 a) Which APC variable is believed to have the weakest effect on the outcome in question?

5 b) Which APC variable has putative mediators available in the data?

6 c) How well-measured are these mediators?

7 d) How exhaustive are these mediators of all the pathways linking the respective APC
8 variable and the outcome?

9 e) Are there common causes of these mediators (e.g. is more than one APC variable
10 affecting the mediators)?

11 f) Are mediator-outcome relationships potentially confounded, and how severely?

12 To minimize bias, select the APC variable for which the answers to questions (a) through (f) are
13 most favourable. For this variable, model the mediated pathways while the other two effects can
14 be modelled directly. The degree of bias that remains in the analysis is dependent on the answers
15 to the above questions for the chosen APC variable, and dependent on ‘ordinary’ modelling
16 concerns such as correct modelling of the functional form of the relationship between APC
17 variables, mediators and outcome.

18 Answering questions (a) through (f) will likely require making assertions which cannot be
19 tested in the data, and is thereby similar to setting constraints by fiat like in various other APC
20 approaches. However, the difference is that the answers to these questions, particularly questions
21 (a), (d), (e), and (f), can be motivated based on theoretical reasoning. As described in our
22 introduction, various authors have argued that age-period-cohort analysis needs to become more
23 theoretically informed. A large suite of APC methods exist that solve the linear identification
24 problem non-transparently and without substantive theoretical justification. A strength of the

1 mechanism-based approach is that it motivates researchers to be explicit about their (substantive
2 theoretical) assumptions.

3

4 *Conclusion*

5 We demonstrated the performance of the mechanism-based approach to APC modelling in non-
6 ideal circumstances. Biases occurred when the assumed causal relations did not coincide with the
7 truth, such as when paths from causes to mediators were omitted from the estimation model, or if
8 there was unmeasured confounding of the mediator-outcome relationship. The direction of the
9 bias followed our expectations based on APC linear dependency. Size of bias is dependent on the
10 size of the effects involving confounders, omitted intermediate variables or pathways. Our
11 extension of the mechanism-based approach increases its utility, by allowing it to be easily
12 useable for models with non-linear link functions and parameterizations of any complexity. Our
13 brief guidelines, aided by causal DAGs, offer a useful tool for researchers who wish to implement
14 this approach.

15

1 **References**

2 Bell, A., Jones, K. (2014). Another 'futile quest'? A simulation study of Yang and Land's
3 Hierarchical Age-Period-Cohort model. *Demographic Research*, 30, 333-360.

4
5 Ben-Shlomo, Y., Kuh, D. (2002). A life course approach to chronic disease epidemiology:
6 conceptual models, empirical challenges and interdisciplinary perspectives. *International Journal*
7 *of Epidemiology*, 31 (2), 285-293.

8
9 Blackmore, H.L., Ozanne, S.E. (2015). Programming of cardiovascular disease across the life-
10 course. *Journal of Molecular and Cellular Cardiology*, 83, 122-130.

11
12 Capewell, S., Beaglehole, R., Seddon, M., McMurray, J. (2000). Explanation for the decline in
13 coronary heart disease mortality rates in Auckland, New Zealand, between 1982 and 1993.
14 *Circulation*, 102 (13), 1511-1516.

15
16 Clayton, D, Schifflers, E. (1987). Models for temporal variation in cancer rates. II: age-period-
17 cohort models. *Statistics in Medicine*, 6, pp. 469–481.

18
19 Cole, S.R., Platt, R.W., Schisterman, E.F., Chu, H., Westreich, D., Richardson, D., Poole, C.
20 (2010). Illustrating bias due to conditioning on a collider. *International Journal of*
21 *Epidemiology*, 39 (2), 417-420.

22
23 Efron B., Tibshirani R.J. (1994). An introduction to the Bootstrap. CRC Press: Boca Raton.

24

1 Ekamper, P., van Poppel F., Stein, A.D., Lumey, L.H. (2014). Independent and additive
2 association of prenatal famine exposure and intermediary life conditions with adult mortality
3 between age 18-63 years. *Social Science and Medicine*, 119, 232-239.
4
5 Elwert, F., Winship C. (2014). Endogenous Selection Bias: The Problem of Conditioning on a
6 Collider Variable. *Annual Review of Sociology*, 40, 31-53.
7
8 Fienberg, SE. (2013). Cohort Analysis' Unholy Quest: A Discussion. *Demography*, 2013, 50 (6),
9 1981-1984.
10
11 Glymour, MM. (2006). Chapter 16: Using causal diagrams to understand common problems in
12 social epidemiology. In J.M. Oaks and J.S. Kaufman (eds.), *Methods in Social Epidemiology* (pp.
13 387-422). San Francisco, CAL: Jossey-Bass.
14
15 Greenland S. (2003). Quantifying biases in causal models: classical confounding versus collider-
16 stratification bias. *Epidemiology* 14,300–306.
17
18 Hernán M.A., Robins J.M. Causal Inference. 2013 Boca Raton, Fla Chapman & Hall/CRC
19
20 Holford, T.R. (2006). Approaches to fitting age-period-cohort models with unequal intervals.
21 *Statistics in Medicine*, 25 (6), 977-993.
22
23 Held, L., Riebler, A. (2012). A conditional approach for inference in multivariate age-period-
24 cohort models. *Statistical Methods in Medical Research*, 21,311-329.

1

2 Jiang, Z., VanderWeele T.J. (2015). Jiang and VanderWeele Respond to “Bounding Natural
3 Direct and Indirect Effects”. *American Journal of Epidemiology*, 182 (2), 115-117.

4

5 Keil, A.P, Edwards, J.K., Richardson, D.B., Naimi, A.I., Cole, S.R. (2014). The parametric g-
6 formula for time-to-event data: intuition and a worked example. *Epidemiology*, 25(6),889-97.

7

8 Kuhn, T. S. (1996). *The structure of scientific revolutions* (3rd ed.). Chicago, IL: The University
9 of Chicago Press.

10

11 Luo, L. (2013). Assessing Validity and Application Scope of the Intrinsic Estimator Approach to
12 the Age-Period-Cohort Problem. *Demography*, 50 (6), 1945-1967.

13

14 Luo, L., Hodges, J.S. (2016), Block Constraints in Age-Period-Cohort models with Unequal-
15 Width Intervals. *Sociological Methods & Research*, 45(4), 700-726.

16

17 MacKinnon, D.P., Lockwood, C.M., Williams, J. (2004), Confidence Limits for the Indirect
18 Effect: Distribution of the Product and Resampling Methods. *Multivariate Behavioral Research*,
19 39(1), 99-128

20

21 Mulaik, S. (2009). Structural equation models. In: *Linear Causal Modeling with Structural*
22 *Equations* (pp. 119–138). Boca Raton, FL: CRC Press.

23

1 Pearl (2000). *Causality: Models, Reasoning, and Inference*. Cambridge: Cambridge University
2 Press.

3

4 Peeters, A., Nusselder, W.J., Stevenson, C., Boyko, E.J., Moon, L., Tonkin, A. (2011). Age-
5 specific trends in cardiovascular mortality rates in the Netherlands between 1980 and 2009.
6 *European Journal of Epidemiology* 26 (5): 369–373.

7

8 Preacher, K. J., & Hayes, A. F. (2008). Asymptotic and resampling strategies for assessing and
9 comparing indirect effects in multiple mediator models. *Behavior Research Methods*, 40, 879-
10 891.

11

12 Preston, S., Wang, H. (2006). Sex mortality differences in the United States: The role of cohort
13 smoking patterns. *Demography*, 43, 631–646.

14

15 Robert, C, Casella, G (2004). Monte Carlo Integration. In Robert, C. (ed.), *Monte Carlo*
16 *Statistical Methods* (pp. 79-122). Springer: New York.

17

18 Rubin, D.B. (1974). Estimating Causal Effects of Treatments in Randomized and Nonrandomized
19 Studies. *Journal of Educational Psychology*, 66 (5), 688–701.

20

21 Sabol, S.Z., Nelson, M.L., Fisher, C., Gunzerath, L., Brody, C.L., Hu, S., Sirota, L.A., Marcus,
22 S.E., Greenberg, B.D., Lucas, F.R., Benjamin, J., Murphy, D.L., Hamer, D.H. (1999). A genetic
23 association for cigarette smoking behavior. *Health Psychology*, 18 (1), 7-13.

24

1 Spiegelhalter, D.J., Thomas A, Best N.G., Lunn D. (2002). Winbugs User Manual. Version 1.4.
2 Cambridge. England: MRC Biostatistics Unit.
3
4 Spirtes P., Glymour C., Scheines R. (2001). Causation, Prediction and Search. The MIT Press:
5 Massachusetts: Cambridge.
6
7 Smith, J.G., Newton-Cheh, C. (2015). Genome-wide association studies of late-onset
8 cardiovascular disease. *Journal of Molecular and Cellular Cardiology*, 83, 131-14.
9
10 VanderWeele, T.J. (2015). *Explanation in Causal Inference: Methods for Mediation and*
11 *Interaction*, Oxford University Press.
12
13 VanderWeele, T.J., Hernán, M.A., Robins J.M. (2008). Causal directed acyclic graphs and the
14 direction of unmeasured confounding bias. *Epidemiology*. 19 (5), 720-8.
15
16 Verlato, G., Melotti, R., Corsico, A.G., Bugiani, M., Carrozzi, L., Marinoni, A., Dallari, R.,
17 Pirina, P., Struzzo P, Olivieri, M., de Marco, R.; ISAYA Study Group. (2006). Time trends in
18 smoking habits among Italian young adults. *Respiratory Medicine*, 100 (12), 2197-2206.
19
20 Winship, C, Harding, D.J. (2008). A Mechanism-Based Approach to the Identification of Age-
21 Period-Cohort Models. *Sociological Methods & Research*, 36 (3), 362-401.
22
23 Winship, C., Mare, R.D. (1983). Structural Equations and Path Analysis for Discrete Data. *The*
24 *American Journal of Sociology*, 89(1), 54-110.

1 Wright, S. (1934). The Method of Path Coefficients. *The Annals of Mathematical Statistics*, 5(3),
2 161-215.

3

4 Yang, Y., Schulhofer-Wohl, S., Fu, W. J., Land, K. C. (2008). The intrinsic estimator for age-
5 period-cohort analysis: What it is and how to use it. *American Journal of Sociology*, 113, 1697–
6 1736.

7

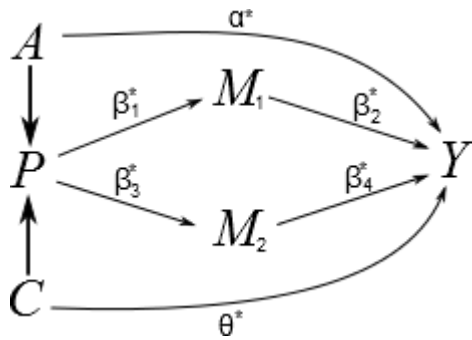
1

Mediator	Direction of effect		
	Period on mediator	Mediator on outcome	Period on outcome
BMI	Positive	Positive	Positive
Smoking	Negative	Positive	Negative
Statin	Positive	Negative	Negative
Unmeasured	Negative	Positive	Negative

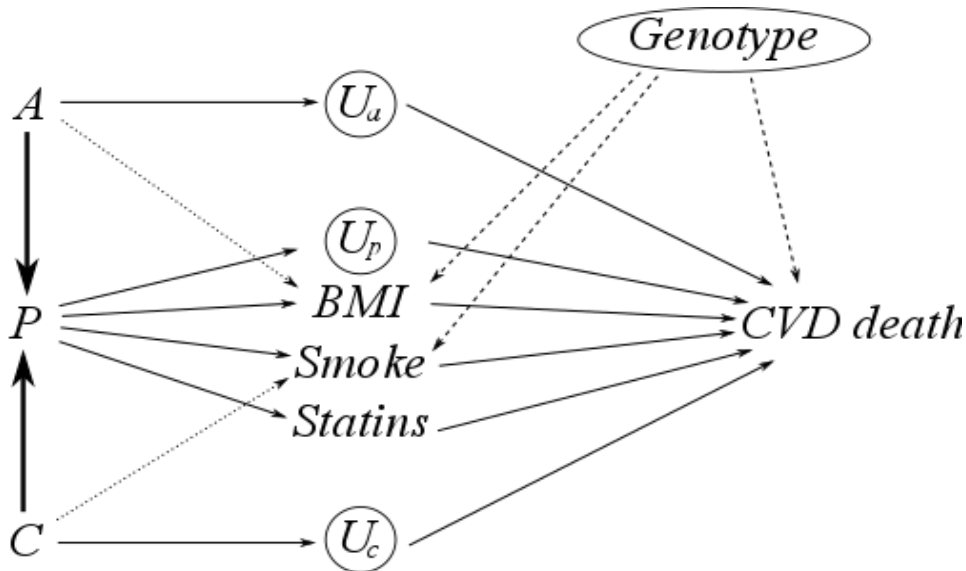
2 **Table 1.** Direction of effect of period on the mediators, mediators on cardiovascular mortality,
3 and the direction along the entire path of period on cardiovascular mortality.

4

5

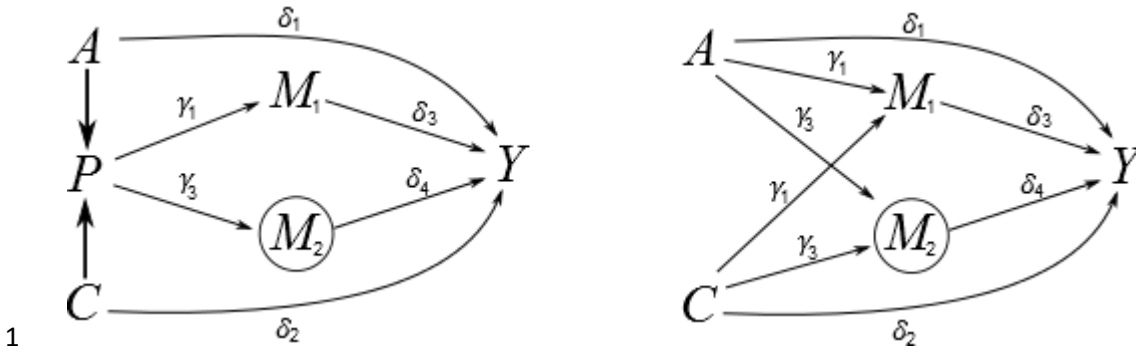


1
 2 **Figure 1.** Causal Directed Acyclic Graph, showing the age effect (α^*), cohort effect (θ^*), and the
 3 period effect (the β^* 's). The bold arrows represent deterministic relationships, while the non-bold
 4 arrows represent stochastic relationships.



Relations in scenario 'simple': ———
 Relations in scenario 'more causes': ——— &
 Relations in scenario 'confounding': ——— &

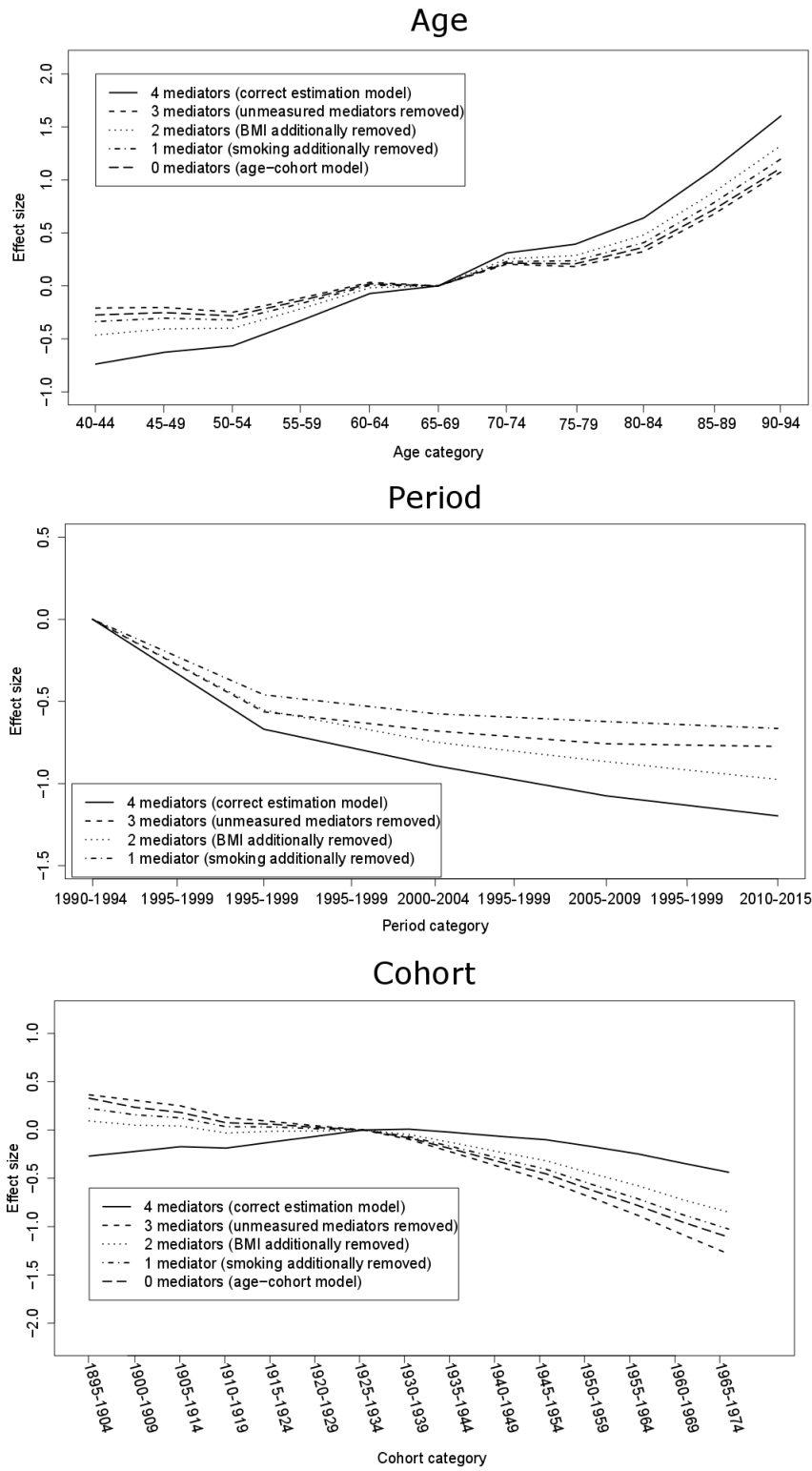
5
 6 **Figure 2.** Causal Directed Acyclic Graph of the 3 scenarios investigated by simulation. Bold
 7 arrows between age (A), period (P) and cohort (C) represent a deterministic relationship, whereas
 8 the remaining arrows represent stochastic causal relationships. Circled variables represent
 9 variables that are omitted from the estimation model (in some simulation setups).



2 **Figure 3.** Causal Directed Acyclic Graph when the age effect (δ_1) and cohort effect (δ_2) are
 3 estimated directly, and the period effect is estimated using mediators as per *Pearl's front door*
 4 *criterion*, while one period mediator is unmeasured. Left: effect estimates if M_2 is measured and
 5 P is included in the estimation process. Right: relationships that form due to linear dependency
 6 when period is excluded from estimation.

7

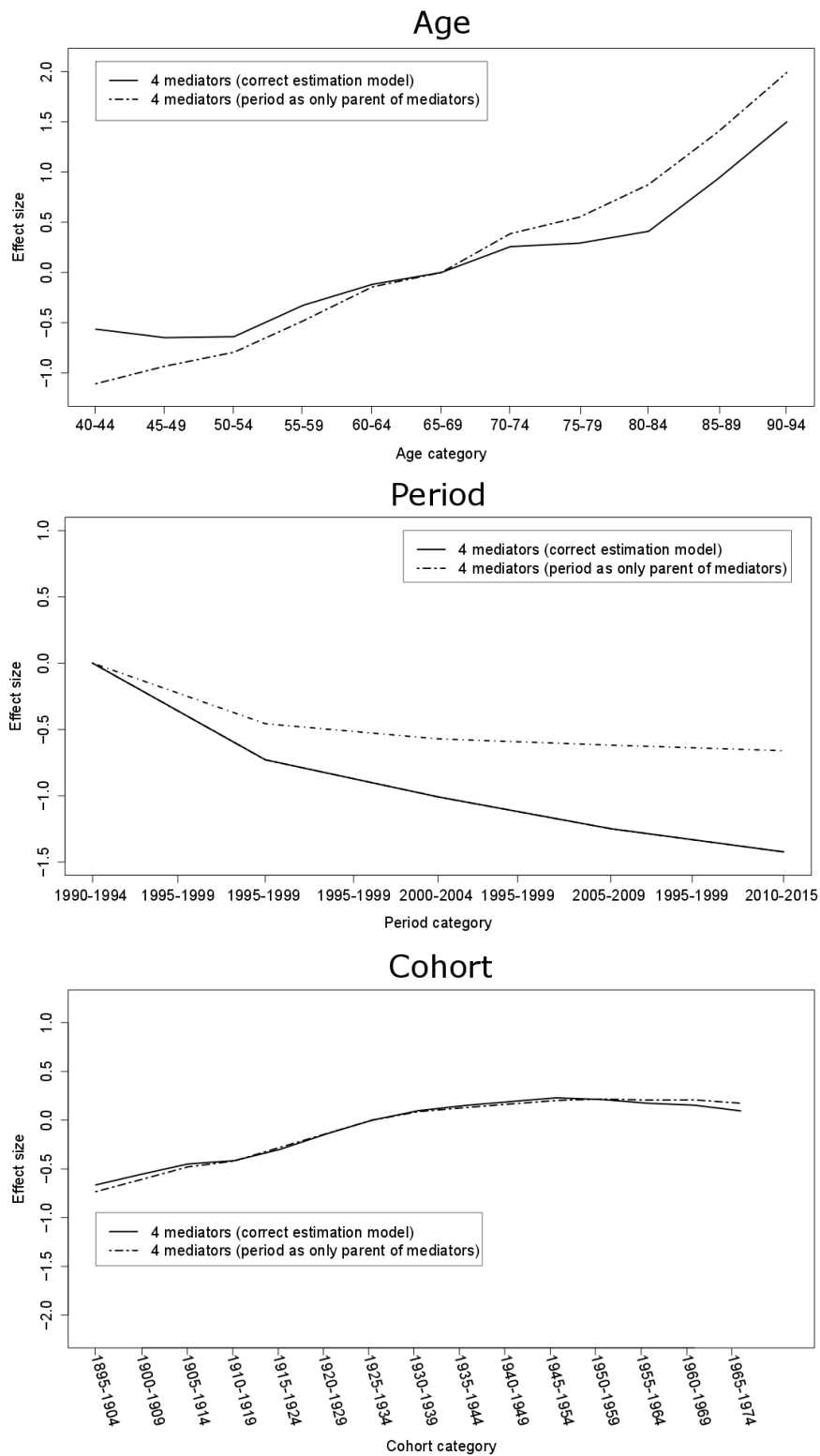
Probit Results: Scenario 1 'Simple'



1

2 **Figure 4.** Average estimated parameters for the APC effects in scenario 1 (“Simple”) using the
 3 mechanism-based approach: summary of 1000 simulations.

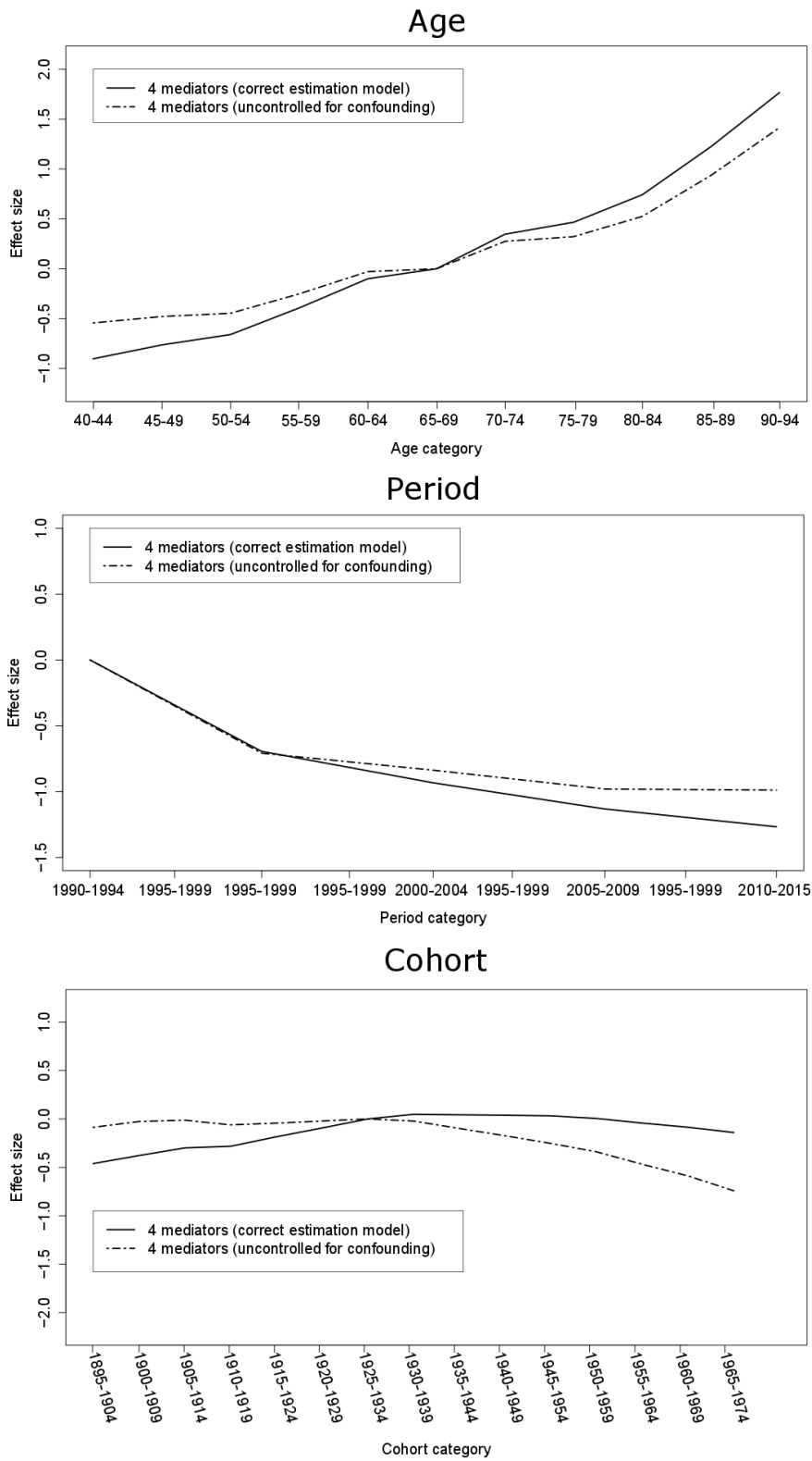
Probit Results: Scenario 2 'More Causes'



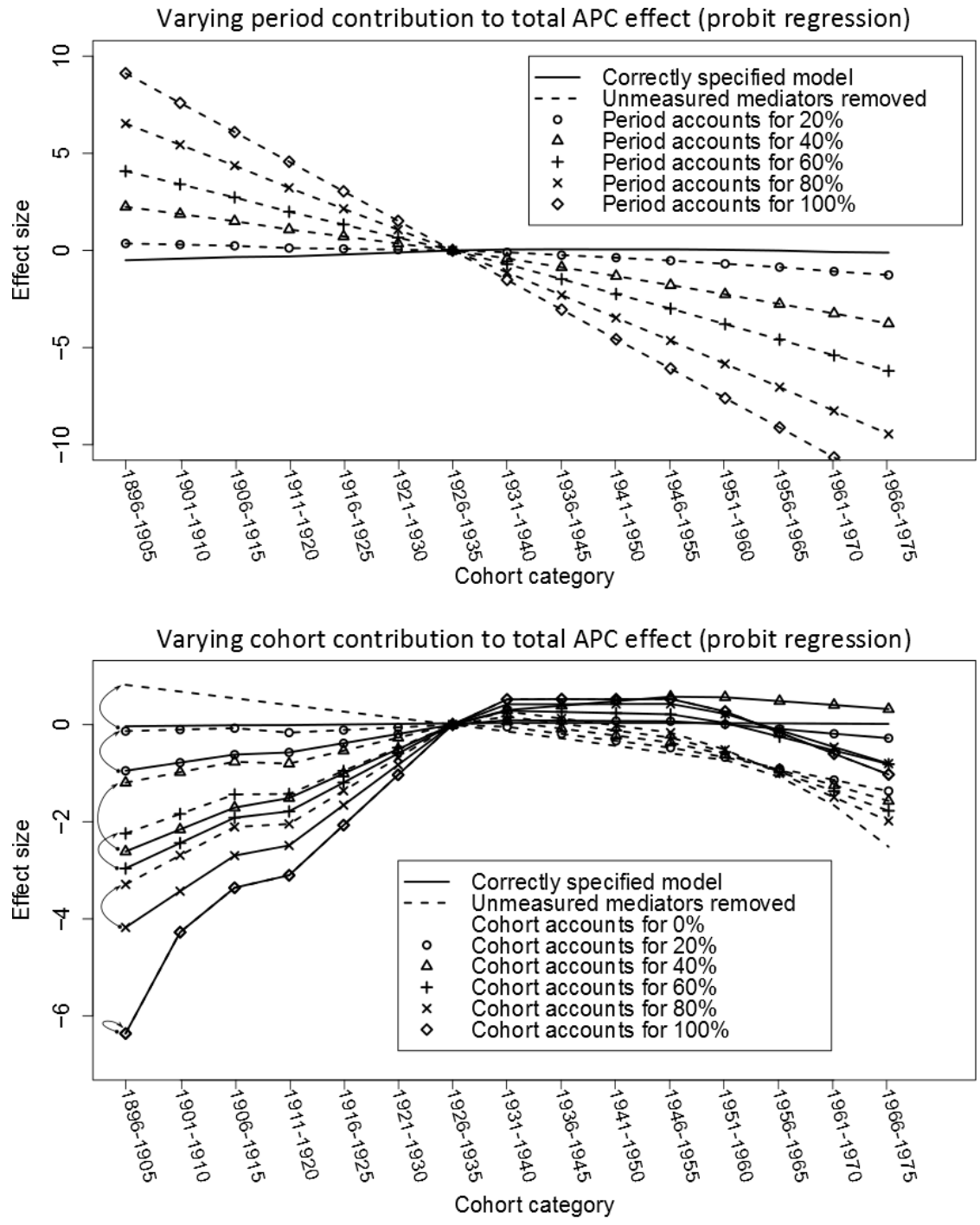
1

2 **Figure 5.** Average APC estimates in scenario 2 (“More causes”) using the mechanism-based
 3 approach: summary of 1000 simulations.

Probit Results: Scenario 3 'Confounding'



1
 2 **Figure 6.** Average APC estimates in scenario 3 (“Confounding”) using the mechanism-based
 3 approach: summary of 1000 simulations.



1
 2 **Figure 7.** Varying the effect sizes of period 0 to 100% of the total APC effect in 20% increments
 3 while keeping cohort effect constant (upper), varying the effect sizes of cohort 0 to 100% of the
 4 total APC effect in 20% increments while keeping period effect constant (lower). When period
 5 accounts for 100% the correctly specified cohort trend is a horizontal line at $y=0$. No bias when
 6 cohort accounts for 100% because then the period effect (source of bias) accounts for 0%. Only
 7 cohort figures (probit variant) shown. Arrows in lower figure indicates the size and direction of
 8 the bias.