

Confounding adjustment through front-door blocking in longitudinal studies

Arvid Sjölander⁽¹⁾, Rino Bellocco⁽²⁾

Abstract

A common aim of epidemiological research is to estimate the causal effect of a particular exposure on a particular outcome. Towards this end, observed associations are often 'adjusted' for potential confounding variables. When the potential confounders are unmeasured, explicit adjustment becomes unfeasible. It has been demonstrated that causal effects can be estimated even in the presence of unmeasured confounding, utilizing a method called 'front-door blocking'. In this paper we generalize this method to longitudinal studies. We demonstrate that the method of front-door blocking poses a number of challenging statistical problems, analogous to the famous problems associated with the method of 'back-door blocking'.

Keywords: Causal inference; Confounding; Directed Acyclic Graph

DOI: 10.2427/8757

1 INTRODUCTION

A common aim of epidemiological research is to estimate the causal effect of a particular exposure on a particular outcome. In observational (i.e. nonrandomized) studies, the exposure-outcome association is often confounded by extraneous factors, and cannot be given a causal interpretation unless the confounding factors have been properly adjusted for. For stationary exposures (i.e. exposures which do not vary over time), adjustments can be made using standard methods, e.g. stratification or parametric regression modeling.

When the exposure and confounders vary over time, adjustments require special techniques. Robins [1] showed that if there are time-varying confounders that are affected by previous exposure levels, and have an effect on later exposure levels, then standard methods will produce biased estimates, even if all confounders are observed and adjusted for in the analysis. He derived an expression for the causal exposure effect, in the presence of observed time-varying confounding, as a function of the observed data distribution. This expression was called 'the G-functional'. Later, Pearl [2] gave a graphical interpretation to the G-functional. He demonstrated that the

⁽¹⁾ **Corresponding Author**, Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Nobels Väg 12A, 171 77, Stockholm, Sweden. *e-mail:* arvid.sjoland@ki.se

⁽²⁾ Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Nobels Väg 12A, 171 77, Stockholm, Sweden; Department of Statistics, University of Milano-Bicocca, Piazza dell'Ateneo Nuovo, 1-20126 Milan, Italy. *e-mail:* rino.bellocco@ki.se

G-functional produces the causal exposure effect when the observed confounders block all ‘back-door paths’ between the exposure and the outcome in the underlying Directed Acyclic Graph (DAG).

A straightforward way to estimate the G-functional from finite data is to fit a parametric regression model for each of its separate components. However, as shown by Robins [1], this explicit modeling strategy is problematic for two reasons. 1) the model parameters will not in general have simple interpretations in terms of the causal exposure effect, and 2) the joint model will most likely be incompatible with the null hypothesis of no causal effect, thus having the feature of ruling out the causal null hypothesis *a priori*. This non-attractive feature was called ‘the G-null paradox’ by Robins [1]. The paper by Robins [1] triggered an intensive research effort, which has last for over two decades. To bypass the problems which are associated with the explicit modeling strategy, a variety of novel statistical methods have been developed (e.g. [3]–[9]). All these methods rely on the assumption of no unmeasured confounding.

Often, the assumption of no unmeasured confounding is not tenable. In rare occasions though, it may be possible to identify the causal exposure effect, even in the presence of unmeasured confounding. Pearl [2] considered a stationary scenario involving an exposure, an outcome, and an arbitrary set of unmeasured confounders. He showed that if a variable can be found which a) completely mediates the effect of the exposure on the outcome, and b) is not affected by the unmeasured confounders, then the causal exposure effect can be identified using data on this variable. When a variable satisfies criteria a) and b), the variable is said to ‘block the front-door’ between exposure and outcome. Estimating causal effects through the blocking of front-doors is principally different from estimating causal effects through the blocking of back-doors; whereas the former strategy allows for unmeasured confounding, the latter strategy does not.

In this paper we generalize the results by Pearl [2] to longitudinal scenarios. We derive an expression for the causal exposure effect, in the presence of unobserved time-varying confounding, utilizing the blocking of front-doors. In analogy with the G-functional, we call this expression ‘the F-functional’. We demonstrate that the F-functional shares the same statistical challenges as the G-functional. That is, by explicitly modeling each separate component of the F-functional, we run into problems of interpretation and incompatibilities with the causal null hypothesis.

The paper is organized as follows. In Section 2 we establish notation and definitions. In Section 3 we review the G-functional, and the statistical challenges associated with its estimation. In Section 4 we derive the F-functional and demonstrate that it is associated with the same estimation challenges as the G-functional. To illustrate causal structures, and to motivate our arguments, we will use DAGs. We refer the readers to Pearl [2] for a thorough review of the methods and concepts associated with DAGs.

2 Notation and definitions

We consider a longitudinal scenario, where an exposure of interest, A , and a (set of) potential confounder(s), L , have been measured on repeated occasions $t \in (0, 1, \dots, T)$. Let A_t and L_t denote the exposure and confounder variables at time t , respectively. We use the convention that $A_t = L_t = \emptyset$ if $t < 0$. We define the temporal order so that L_t is realized just before A_t . The outcome of interest, Y , is measured after the end of follow-up ($t = T$). We use $p(\cdot)$ generically for both probability distributions and density functions, and we use integral signs generically for both integrations and summations. We use $E(\cdot)$ for expected values (population averages). We assume that the observed data consists of n iid observations from $p(L_0, A_0, L_1, A_1, \dots, L_T, A_T, Y)$. We define $\bar{A}_t = (A_0, A_1, \dots, A_t)$ and $\bar{L}_t = (L_0, L_1, \dots, L_t)$ as the observed exposure history and confounder history up to t , respectively.

We use the potential outcome framework Rubin [10] to define causal effects. Let $Y^{\bar{a}_t}$ denote the potential outcome for a given individual under the exposure history $\bar{A}_t = \bar{a}_t$. If the individual factually attains levels $\bar{A}_t = \bar{a}_t$, then $Y^{\bar{a}_t}$ is trivially observed and equal to the factual outcome Y , i.e.

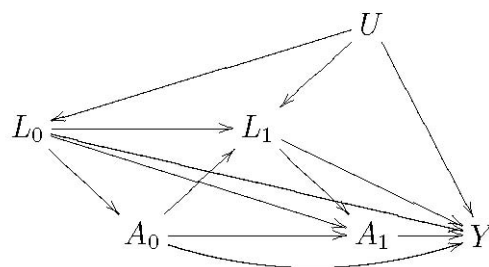
$$\bar{A}_t = \bar{a}_t \Rightarrow Y^{\bar{a}_t} = Y \quad \forall t. \quad (1)$$

The statement in equation (1) relates the potential outcomes to the observed variables, and it is often referred to as the ‘consistency assumption’ [2]. If the individual factually attains some other exposure levels than \bar{a}_t , then $Y^{\bar{a}_t}$ is unobserved or ‘counterfactual’ - it corresponds to a scenario that is contrary to fact. Let $E(Y^{\bar{a}_t})$ denote the average value of Y under the counterfactual scenario when the whole study population attains $\bar{A}_t = \bar{a}_t$. The joint causal effect of a shift in exposure from $\bar{A}_T = \bar{a}_T$ to $\bar{A}'_T = \bar{a}'_T$ on Y is defined as some comparison between $E(Y^{\bar{a}_T})$ and $E(Y^{\bar{a}'_T})$, e.g. $E(Y^{\bar{a}_T}) - E(Y^{\bar{a}'_T})$. When A is binary (0/1), we may for instance wish to compare the mean outcome when the population is unexposed at all times, $E(Y^{\bar{a}_T=0_T})$, with the mean outcome when the population is exposed at all times $E(Y^{\bar{a}_T=1_T})$.

3 The G-functional

The G-functional applies to scenarios in which each observed variable may be causally affected by each preceding variable. There may be an arbitrary set of unmeasured variables U which has a causal effect on \bar{L}_T and Y . However, it is assumed that U has no causal effect on A_t , apart from any effect mediated through \bar{L}_t . The DAG in Figure 1 displays the scenario for $T = 1$. In the setup depicted in Figure 1 there is no unmeasured confounding. Specifically, adjusting

Figure 1. A DAG in which $(\bar{L}_t, \bar{A}_{t-1})$ blocks all back-door paths between A_t and Y .



for $(\bar{L}_t, \bar{A}_{t-1})$ is sufficient to give the association between A_t and Y a causal interpretation. In graphical jargon we say that $(\bar{L}_t, \bar{A}_{t-1})$ blocks all back-door paths between A_t and Y . In terms of potential outcomes, Figure 1 implies

$$Y^{\bar{a}_T} \amalg A_t | \bar{L}_t, \bar{A}_{t-1} \quad \forall t, \quad (2)$$

which is often referred to as ‘sequential exchangeability’.

Robins [1] showed that under consistency (equation (1)) and sequential exchangeability (equation (2)), $E(Y^{\bar{a}_T})$ is identifiable and equal to

$$E(Y^{\bar{a}_T}) = \int_{\bar{l}_T} E(Y | \bar{a}_T, \bar{l}_T) \prod_{t=0}^T p(l_t | \bar{a}_{t-1}, \bar{l}_{t-1}) d\bar{l}_T. \quad (3)$$

He referred to the right handside of equation (3) as the ‘G-functional’. In practical scenarios, it is typically not feasible to estimate the G-functional non-parametrically, at a reasonable level of efficiency. A straightforward way to gain efficiency is to fit a parametric model for each component of the right handside of equation (3). The integration can be carried out using numerical techniques. However, this explicit modeling strategy is problematic for two reasons: 1) the model parameters will not in general have simple interpretations in terms of the causal exposure effect, and 2) the joint model will most likely not be compatible with the null hypothesis of no causal effect, thus having the non-attractive feature of ruling out the causal null hypothesis *a priori*. This feature was called ‘the G-null paradox’ by Robins [1].

To appreciate the first problem, consider the following example.

Example 1. Suppose that Y is binary (0/1), that L_t is continuous and unrestricted, and that we use the standard models

$$Y|\bar{A}_T, \bar{L}_T \sim \text{Ber} \left\{ \text{expit} \left(\beta_0 + \beta_1 \sum_{t=0}^T A_t + \beta_2 \sum_{t=0}^T L_t \right) \right\}, \quad (4)$$

and

$$L_t|\bar{A}_{t-1}, \bar{L}_{t-1} \sim N \left(\alpha_0 + \alpha_1 \sum_{k=0}^{t-1} A_k + \alpha_2 \sum_{k=0}^{t-1} L_k, \sigma^2 \right), \quad (5)$$

with $\text{expit}(x) = e^x/(1 + e^x)$. Replacing equation (4) and equation (5) into equation (3) gives

$$E(Y^{\bar{a}_T}) = \int_{\bar{L}_T} \text{expit} \left(\beta_0 + \beta_1 \sum_{t=0}^T a_t + \beta_2 \sum_{t=0}^T l_t \right) \prod_{t=0}^T \frac{\exp \left[-\frac{\{l_t - (\alpha_0 + \alpha_1 \sum_{k=0}^{t-1} a_k + \alpha_2 \sum_{k=0}^{t-1} l_k)\}^2}{2\sigma^2} \right]}{\sqrt{2\pi\sigma^2}} d\bar{l}_T. \quad (6)$$

In equation (6), $E(Y^{\bar{a}_T})$ depends on \bar{a}_T through the model parameters $(\alpha_0, \alpha_1, \alpha_2, \beta_0, \beta_1, \beta_2)$ in a rather complex fashion. Thus, to determine the causal effect of \bar{A}_T on Y under the models in equation (4) and equation (5), we would have to carry out the integration for each possible value of \bar{a}_T . This makes the explicit modeling strategy numerically untractable and highly non-transparent. \square

To appreciate the second problem, consider the following example.

Example 2. Suppose that the true causal structure is given by the DAG in Figure 2. In

Figure 2. A DAG in which $(\bar{L}_t, \bar{A}_{t-1})$ blocks all back-door paths between A_t and Y , and the causal null hypothesis holds.

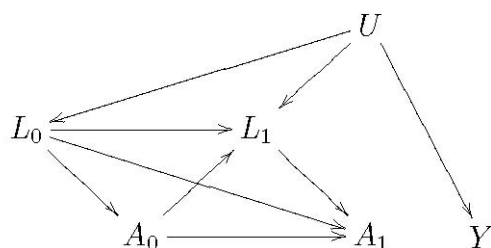


Figure 2, \bar{A}_T has no causal effect on Y , so that $E(Y^{\bar{a}_T})$ is not a function of \bar{a}_T ; we say that the causal null hypothesis holds. Nevertheless, L_t is conditionally associated with \bar{A}_{t-1} , given \bar{L}_{t-1} , due to the paths $A_k \rightarrow L_t$, $k \in (0, 1, \dots, t-1)$. Furthermore, Y is conditionally associated

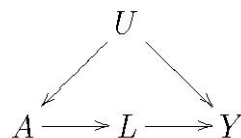
with \bar{A}_T , given \bar{L}_T , due to the paths $A_k \rightarrow L_t \leftarrow U \rightarrow Y$, $t \in (0, 1, \dots, T)$, $k \in (0, 1, \dots, t-1)$ being open by conditioning on L_t [2]. Thus, each component on the right handside of equation (3) depends on \bar{a}_T , even though the whole function does not. To specify parametric models for each component which satisfies this restriction is a very difficult task. Thus, ‘working models’ such as those in equation (4) and equation (5) are most likely misspecified so that the obtained parameter estimates are such that the estimated G-functional depends on \bar{a}_T , even though the causal null hypothesis in Figure 2 holds. \square

To bypass the problems which are associated with the explicit modeling strategy, a variety of novel statistical methods have been developed, including G-estimation and structural nested models [3], inverse probability weighting (IPW) and marginal structural models (MSMs) [4]–[7], and special techniques for the estimation of dynamic treatment regimes [8, 9]. All these methods rely on the assumption of no unmeasured confounding, or equivalently, on the assumption of sequential exchangeability.

4 The F-functional

Pearl [2] considered the stationary scenario depicted in Figure 3. In Figure 3, there is unmeasured

Figure 3. A DAG where L blocks the front-door between A and Y .



confounding for A and Y , through the common cause U . The covariate L satisfies the following important properties: a) it completely mediates the effect of A on Y , and b) it is not affected by U . When a) and b) hold, we say that L blocks the front-door between A and Y . Under the assumptions encoded in Figure 3, Pearl [2] showed that $E(Y^a)$ can be expressed as

$$E(Y^a) = \int_l p(l|a) \int_{a'} E(Y|a', l) p(a'|l) da' dl \quad (7)$$

Since the right handside of equation (7) only contains components which are directly computable from $\Pr(Y, L, A)$, it follows that the causal effect of A on Y is identifiable by observations on (Y, L, A) .

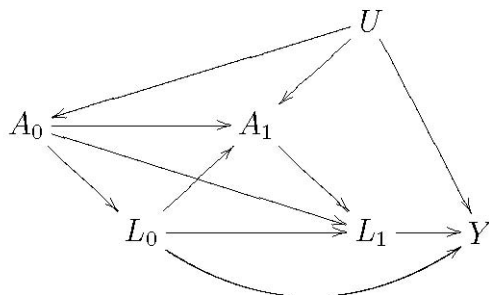
As a practical example, Pearl [2] considered a fictitious scenario where A , Y and L represent ‘smoking’, ‘lung cancer’, and the amount of tar deposited in a person’s lungs, respectively. Using hypothetical numbers, he demonstrated that smoking and lung cancer can be strongly associated, even in the absence of a causal effect, as computed through equation (7).

In reality, both exposure and covariates are often time-varying. For instance, both smoking habits and the amount of tar in the lungs are likely to vary over time. Thus, it is useful to extend Pearl’s method for time-varying data. Below, we generalize the identifying expression in equation (7) to longitudinal scenarios, analogously to the G-functional. We will call this generalized expression the ‘F-functional’. We show that the F-functional poses the same statistical challenges as the G-formula. For convenience, we will swap the temporal order so that A_t is realized before L_t .

We consider a longitudinal scenario where each observed variable may be causally influenced by any preceding variable. There may be an arbitrary set of unmeasured variables U which

has a causal effect on \bar{A}_T and Y . Thus, we allow for unmeasured confounding. We make two assumptions: a) all effect of A_t on Y is mediated through $(L_t, L_{t+1}, \dots, L_T)$, and b) U has no causal effect on L_t , apart from any effect mediated through \bar{A}_t . The DAG in Figure 4 displays the scenario for $T = 1$. When a) and b) hold, we say that $(L_t, L_{t+1}, \dots, L_T)$ completely blocks

Figure 4. A DAG where $(L_t, L_{t+1}, \dots, L_T)$ completely blocks the front-door between A_t on Y .



the front-door between A_t on Y . Assumptions a) and b) are clearly quite restrictive; in Section 5 we provide a discussion on practical scenarios where these assumptions can be expected to hold.

We now show that $E(Y^{\bar{a}_T})$ is identifiable by observations on $(Y, \bar{L}_T, \bar{A}_T)$, and given by

$$E(Y^{\bar{a}_T}) = \int_{\bar{l}_T} \prod_{t=0}^T p(l_t | \bar{a}_t, \bar{l}_{t-1}) \int_{\bar{a}'_T} E(Y | \bar{a}'_T, \bar{l}_T) \prod_{t=0}^T p(a'_t | \bar{a}'_{t-1}, \bar{l}_{t-1}) d\bar{a}'_T d\bar{l}_T. \quad (8)$$

Proof. Using standard graphical methods [2] it can be shown that Figure 4 implies the following conditional independencies:

$$L_t \perp\!\!\!\perp U | \bar{A}_t, \bar{L}_{t-1} \quad \forall t, \quad (9)$$

$$Y \perp\!\!\!\perp \bar{A}_t | U, \bar{L}_t \quad \forall t \quad (10)$$

$$Y^{\bar{a}_T} \perp\!\!\!\perp A_t | U, \bar{A}_{t-1}, \bar{L}_{t-1} \quad \forall t. \quad (11)$$

We have that

$$\begin{aligned} E(Y^{\bar{a}_T}) &= \int_u E(Y^{\bar{a}_T} | u) p(u) du \\ &\stackrel{\text{equation (11)}}{=} \int_u E(Y^{\bar{a}_T} | u, a_0) p(u) du \\ &= \int_{u, l_0} E(Y^{\bar{a}_T} | u, a_0, l_0) p(l_0 | u, a_0) p(u) du \\ &= \dots \\ &= \int_{u, \bar{l}_T} E(Y^{\bar{a}_T} | u, \bar{a}_T, \bar{l}_T) \prod_{t=0}^T p(l_t | u, \bar{a}_t, \bar{l}_{t-1}) p(u) du d\bar{l}_T \\ &\stackrel{\text{equation (1)}}{=} \int_{u, \bar{l}_T} E(Y | u, \bar{a}_T, \bar{l}_T) \prod_{t=0}^T p(l_t | u, \bar{a}_t, \bar{l}_{t-1}) p(u) du d\bar{l}_T \\ &\stackrel{\text{equation (9)}}{=} \int_{\bar{l}_T} \prod_{t=0}^T p(l_t | \bar{a}_t, \bar{l}_{t-1}) \int_u E(Y | u, \bar{a}_T, \bar{l}_T) p(u) du d\bar{l}_T, \end{aligned} \quad (12)$$

We can write $p(u)$ as

$$\begin{aligned}
 p(u) &= \int_{a'_0} p(u|a'_0)p(a'_0)da'_0 \\
 &\stackrel{\text{equation (9)}}{=} \int_{a'_0} p(u|a'_0, l_0)p(a'_0)da'_0 \\
 &= \dots \\
 &= \int_{\bar{a}'_T} p(u|\bar{a}'_T, \bar{l}_T) \prod_{t=0}^T p(a'_t|\bar{a}'_{t-1}, \bar{l}_{t-1}), \tag{13}
 \end{aligned}$$

We can now write $\int_u E(Y|u, \bar{a}_T, \bar{l}_T)p(u)du$ as

$$\begin{aligned}
 \int_u E(Y|u, \bar{a}_T, \bar{l}_T)p(u)du &\stackrel{\text{equation (13)}}{=} \int_{u, \bar{a}'_T} E(Y|u, \bar{a}_T, \bar{l}_T)p(u|\bar{a}'_T, \bar{l}_T) \prod_{t=0}^T p(a'_t|\bar{a}'_{t-1}, \bar{l}_{t-1})dud\bar{a}'_T \\
 &\stackrel{\text{equation (10)}}{=} \int_{u, \bar{a}'_T} E(Y|u, \bar{a}'_T, \bar{l}_T)p(u|\bar{a}'_T, \bar{l}_T) \prod_{t=0}^T p(a'_t|\bar{a}'_{t-1}, \bar{l}_{t-1})dud\bar{a}'_T \\
 &= \int_{\bar{a}'_T} E(Y|\bar{a}'_T, \bar{l}_T) \prod_{t=0}^T p(a'_t|\bar{a}'_{t-1}, \bar{l}_{t-1})d\bar{a}'_T, \tag{14}
 \end{aligned}$$

Combining equation (12) and equation (14) concludes the proof. \square

A straightforward way to estimate the F-functional from finite samples is to fit a regression model for each component on the right handside of equation (8), and carry out the integration numerically. However, this approach suffers from the same problems when applied to the F-functional as when applied to the G-functional. To see that the model parameters will not in general have simple interpretations in terms of the causal exposure effect, consider the following example.

Example 3. Suppose that Y and A_t are binary (0/1), that L_t is continuous and unrestricted, and that we use the standard models

$$Y|\bar{A}_T, \bar{L}_T \sim \text{Ber} \left\{ \text{expit} \left(\beta_0 + \beta_1 \sum_{t=0}^T A_t + \beta_2 \sum_{t=0}^T L_t \right) \right\}, \tag{15}$$

$$A_t|\bar{A}_{t-1}, \bar{L}_{t-1} \sim \text{Ber} \left\{ \text{expit} \left(\gamma_0 + \gamma_1 \sum_{k=0}^{t-1} A_k + \gamma_2 \sum_{k=0}^{t-1} L_k \right) \right\}, \tag{16}$$

and

$$L_t|\bar{A}_t, \bar{L}_{t-1} \sim N \left(\alpha_0 + \alpha_1 \sum_{k=0}^t A_k + \alpha_2 \sum_{k=0}^{t-1} L_k, \sigma^2 \right). \tag{17}$$

Replacing equation (15), equation (16), and equation (17) into equation (8) gives

$$E(Y^{\bar{a}_T}) = \int_{\bar{l}_T} \left[\prod_{t=0}^T \frac{\exp \left[-\frac{\{l_t - (\alpha_0 + \alpha_1 \sum_{k=0}^t a_k + \alpha_2 \sum_{k=0}^{t-1} l_k)\}^2}{2\sigma^2}}{\sqrt{2\pi\sigma^2}} \right] \right]$$

$$\int_{\bar{a}'_T} \expit \left(\beta_0 + \beta_1 \sum_{t=0}^T a'_t + \beta_2 \sum_{t=0}^T l_t \right) \prod_{t=0}^T I(a'_t = 1) \left\{ \expit \left(\gamma_0 + \gamma_1 \sum_{k=0}^{t-1} a'_k + \gamma_2 \sum_{k=0}^{t-1} l_k \right) \right\} I(a'_t = 0) \left\{ 1 - \expit \left(\gamma_0 + \gamma_1 \sum_{k=0}^{t-1} a'_k + \gamma_2 \sum_{k=0}^{t-1} l_k \right) \right\} d\bar{a}'_T d\bar{l}_T. \quad (18)$$

In equation (18), $E(Y^{\bar{a}_T})$ depends on \bar{a}_T through the model parameters $(\alpha_0, \alpha_1, \alpha_2, \beta_0, \beta_1, \beta_2, \gamma_0, \gamma_1, \gamma_2)$ in a rather complex fashion. Thus, to determine the causal effect of \bar{A}_T on Y under the models in equation (4) and equation (5), we would have to carry out the integration for each possible value of \bar{a}_T . This makes the explicit modeling strategy numerically untractable and highly non-transparent. \square

To see that the approach will suffer from a problem similar to the G-null paradox, consider the following example.

Example 4. Suppose that the true causal structure is given by the DAG in Figure 5. In

Figure 5. $(L_t, L_{t+1}, \dots, L_T)$ completely blocks the front-door between A_t on Y , and the causal null hypothesis holds.

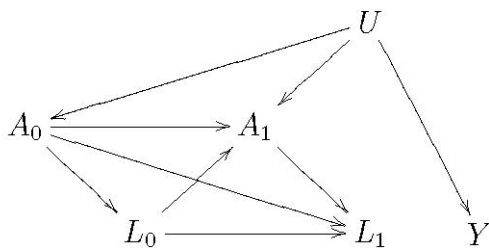


Figure 5, \bar{A}_T has no causal effect on Y , so that $E(Y^{\bar{a}_T})$ is not a function of \bar{a}_T . Nevertheless, L_t is conditionally associated with \bar{A}_t , given \bar{L}_{t-1} , due to the open paths $A_k \rightarrow L_t$, $k \in (0, 1, \dots, t)$. Y is conditionally associated with \bar{A}_T , given \bar{L}_T , due to the open paths $A_t \leftarrow U \rightarrow Y$, $t \in (0, 1, \dots, T)$. Finally, A_t is conditionally associated with \bar{A}_{t-1} , given \bar{L}_{t-1} , due to the open paths $A_k \rightarrow A_t$, $k \in (0, 1, \dots, t)$. Thus, each component on the right handside of equation (8) depends on \bar{a}_T , even though the whole function does not. \square

5 CONCLUSIONS

In this paper we have extended the method of ‘front-door blocking’ to longitudinal scenarios. We have derived an analytic expression for the causal exposure effect which we have called the ‘F-functional’. We have shown that although the F-functional can be used for non-parametric identification, several challenging problems remain to be solved to estimate the F-functional from finite data. The analogs to these problems for the G-functional have been the target of intense research, which has generated a variety of novel statistical methods. We believe that the F-functional has the potential of triggering a similar development.

A difference between back-door blocking and front-door blocking is that whereas it is easy to come up with realistic scenarios for which the former method can be used (e.g. [3]-[9]), it is

much more difficult to find scenarios for which the latter would be applicable. In particular, it is difficult to find variables which completely blocks the front-door between the exposure and the outcome. Pearl [2] considered a fictitious example involving smoking, lung cancer and tar deposited in the lungs. Another possible example could arise in nutrition epidemiology, where it has been established that diets rich in fruits and vegetables have a protective effect on several cardiovascular and cancer diseases. Among a number of mechanistic hypotheses, compounds with antioxidant properties have been proposed to explain these findings. For instance, it has been postulated that antioxidants, with additive and synergistic effects, completely mediates the protective effect of plant food intake on gastric cancer [11]. Furthermore, it may be reasonable to assume that any factors affecting the individual propensity towards a "healthy diet" and gastric cancer (e.g. education and socioeconomic status), have no effect on the antioxidant levels in the body. Thus, the F-functional could potentially be used to estimate the causal protective effect of plant food intake on gastric cancer, with antioxidant levels blocking all intermediate front-doors.

References

- [1] Robins JM. A new approach to causal inference in mortality studies with sustained exposure periods - Application to control of the healthy worker survivor effect. *Mathematical Modelling* 1986; 7: 1393-1512
- [2] Pearl J. *Causality: Models, Reasoning and Inference*. Cambridge: Cambridge University Press, 2000
- [3] Robins JM. The analysis of randomized and non-randomized AIDS treatment trials using a new approach to causal inference in longitudinal studies. In *Health Service Research Methodology: A focus on AIDS*. Edited by Sechrest L, Freeman H, Mulley A, pp 113-159, 2009
- [4] Robins JM. Marginal structural models. In *Proceedings of the American Statistical Association. Section on Bayesian Statistical Science*, 1-10, 1997
- [5] Robins J, Hernan M, Brumback B. Marginal structural models and causal inference in epidemiology. *Epidemiology* 2000; 11(5): 550-560
- [6] Hernan M, Brumback B, Robins J. Marginal structural models to estimate the causal effect of zidovudine on the survival of HIV-positive men. *Epidemiology* 2000; 11(5): 561-570
- [7] Hernan MA, Lanoy E, Costagliola D, Robins J. Comparison of dynamic treatment regimes via inverse probability weighting. *Basic and Clinical Pharmacology and Toxicology* 2006; 98: 237-242
- [8] Cain LE, Robins JM, Lanoy E, Logan R, Costagliola D, Hernan MA. When to start treatment? A systematic approach to the comparison of dynamic treatment regimes using observational data. *The International Journal of Biostatistics* 2010; 6: article 18.
- [9] Orellana L, Rotnitzky A, Robins JM. Dynamic regime structural mean models for estimation of optimal dynamic treatment regimes, Part I: Main content. *The International Journal of Biostatistics* 2010; 6: article 7
- [10] Rubin DB. Estimating causal effects of treatments in randomized and non-

randomized studies. *Journal of Educational Psychology* 1974; 66(5): 688-701

- [11] Serafini, M. Bellocco, R., Wolk, A., Ekstrom, A.M. Total antioxidant potential of fruit and vegetables and risk of gastric cancer. *Gastroenterology* 2002; 123(4): 985-991