

eAppendix of
 "On identification of natural direct effects when a
 confounder of the mediator is directly affected by
 exposure"
 by

Eric J. Tchetgen Tchetgen and Tyler J. VanderWeele

Identification for multiple independent exposure induced confounders
 of the mediator

Suppose that N consists of multiple binary variables $N = (N_1, \dots, N_k)$, and suppose that the nonparametric structural equations model (7) – (11) holds upon replacing equation (9) with the k equations:

$$N_j = g_{N_j}(C, E, \varepsilon_{N_j}), \quad j = 1, \dots, k$$

such that $\{\varepsilon_{N_j} : j = 1, \dots, k\}$ are mutually independent and are jointly independent of $\{\varepsilon_C, \varepsilon_E, \varepsilon_M, \varepsilon_Y\}$.

As mentioned in the text, for this assumption to hold, all common causes of each pair of variables in (N_1, \dots, N_k) would also need to be included in C , as illustrated in Figure 3. below.

Result 2: Assuming the nonparametric structural equations model (7) – (11), suppose that the $E - N$ Monotonicity Assumption holds for each of $(N_1, \dots, N_k) = N$, then $E\{Y(e, M(e^))\}$ is nonparametrically identified by the following formula:*

$$\sum_{m, n, n', c} \mathbb{E}(Y|e, m, n, c) \Pr(M = m | e^*, n', c) \prod_{j=1}^k f_j(n_j, n'_j, e, e^*, c) \Pr(C = c)$$

where

$$f_j(n_j, n'_j, e, e^*, c) = \begin{cases} \Pr\{N_j = 1|e^*, c\} & \text{if } n'_j = n_j = 1 \\ \Pr\{N_j = 1|e, c\} - \Pr\{N_j = 1|e^*, c\} & \text{if } n'_j = 0 \text{ and } n_j = 1 \\ 0 & \text{if } n'_j = 1 \text{ and } n_j = 0 \\ \Pr\{N_j = 0|e, c\} & \text{if } n'_j = n_j = 0 \end{cases}$$

Result 3: Assuming the nonparametric structural equations model (7) – (11), suppose that the $E - N$ Monotonicity Assumption holds for each of $(N_1, \dots, N_k) = N$, then the hazard function of $Y(e, M(e^))$ evaluated at y is nonparametrically identified and satisfies an additive hazards model of the form:*

$$\lambda_0(y) + \frac{\sum_{m,n,n',c} \gamma(y, e, m, n, c) \Gamma(y, e, m, n, c) \Pr(M = m | e^*, n', c) \prod_{j=1}^k f_j(n_j, n'_j, e, e^*, c) \Pr(C = c)}{\sum_{m,n,n',c} \Gamma(y, e, m, n, c) \Pr(M = m | e^*, n', c) \prod_{j=1}^k f_j(n_j, n'_j, e, e^*, c) \Pr(C = c)}$$

where $\Gamma(y, e, m, n, c) = \exp\{-\int_0^y \gamma(u, e, m, n, c) du\}$ and $f_j(n_j, n'_j, e, e^*, c)$ defined in Result 2.

Decomposition of $NDE(e, e^*)$

Consider again the setting of a binary recanting witness N . Recall that $NDE(e, e^*)$ captures the effects along the following two pathways: $E \rightarrow Y$ and $E \rightarrow N \rightarrow Y$. We note that

$$\begin{aligned} \mathbb{E}\{Y(e, M(e^*))\} - \mathbb{E}\{Y(e^*, M(e^*))\} &= \mathbb{E}\{Y(e, M(e^*), N(e)) - Y(e, M(e^*), N(e^*))\} \\ &+ \mathbb{E}\{Y(e, M(e^*), N(e^*)) - Y(e^*, M(e^*), N(e^*))\} \end{aligned}$$

$\mathbb{E}\{Y(e, M(e^*), N(e^*)) - Y(e^*, M(e^*), N(e^*))\}$ captures the pathway $E \rightarrow Y$, the portion of the direct effect not mediated by N , while $\mathbb{E}\{Y(e, M(e^*), N(e)) - Y(e, M(e^*), N(e^*))\}$ captures the pathway $E \rightarrow N \rightarrow Y$, the portion of the direct effect mediated by N . Under monotonicity of the effects of exposure on the recanting witness, we show next that $\mathbb{E}\{Y(e, M(e^*), N(e^*))\}$ is identified and therefore, both of these effects are nonparametrically identified.

Corollary 1: Assuming the nonparametric structural equations model (7) – (11), suppose that N is binary, and $E - N$ Monotonicity Assumption holds, then

$$\mathbb{E}\{Y(e, M(e^*), N(e^*))\} = \sum_{m, n, n', c} \mathbb{E}(Y|e, m, n', c) \Pr(M = m | e^*, n', c) f(n, n', e, e^*, c) \Pr(C = c)$$

with $f(n, n', e, e^*, c)$ given in Result 1.

Corollary 1 extends to the context of multivariate binary confounder N under the assumptions listed in Result 2. Details are omitted but are easily deduced from the presentation. Despite these important generalizations of Result 1, identification under monotonicity is still somewhat limited in that each N_j is restricted to be binary $j = 1, \dots, k$.

A data illustration using Proc NL MIXED in SAS

We briefly illustrate the methodology developed in this paper in the context of simulated data. We first generate data as would be observed in a randomized study of sample size 500 where E is randomized with probability 1/2, N is dichotomous with event probability

$$\Pr(N = 1|E) = \{1 + \exp(-0.5 - 0.75E)\}^{-1},$$

and M and Y are continuous:

$$Y = 60 + 2E + 3M + 1.5N + MN + \varepsilon_Y$$

$$M = 30 + 3E + 4N + \varepsilon_M$$

where ε_Y is normal with mean zero and variance 1.5 and ε_M is standard normal. There is no pre-exposure confounder C in these simulated data. Under monotonicity, it is easy to verify that $NDE(1, 0) = 6.877$ in these data. We describe how this direct effect can be estimated with Proc NLMIXED in SAS by providing sample code below. Using this sample code, we obtained an estimate of $NDE(1, 0)$ equal to 6.383 (95% confidence interval=(3.8328 – 8.9347)).

To further illustrate the methods developed in this paper, consider an alternative data generating mechanism mimicking an observational study with a single binary confounder

$$C \sim \text{Bernoulli}(1/3);$$

$$E|C \sim \text{Bernoulli}((1 + \exp(-0.5 - 0.6C))^{-1});$$

$$N = 50 - 0.75E + 0.5C + N(0, 1);$$

$$M = 30 + 3E + 3C + 2N + N(0, 1);$$

$$u \sim N(0, 2)$$

$$Y = 60 + 2E + 4C + 3M + 1.5N + EN + 2EM + N(0, 1.5)$$

We aim to estimate $NDE(1, 0, c)$, which under the assumption of no $M - N$ interaction, is equal

to

$$NDE(1, 0, 1) = 318.63,$$

$$NDE(1, 0, 0) = 310.13.$$

We further illustrate how this effect estimate can be obtained in Proc NLMIXED in SAS by providing sample code in the appendix. Using this sample code, we obtained the following estimates

$$NDE(1, 0, 1) = 317.52 : 95\% \text{ confidence interval} = (316.62 - 318.43)$$

$$NDE(1, 0, 0) = 309.60 : 95\% \text{ confidence interval} = (308.86 - 310.35)$$

SAS CODE

The data set 'data_example' contains variables E, N, M, Y from the simulated randomized study example.

The first sample code produced the maximum likelihood estimate of $NDE(1, 0)$ reported in the text.

```
proc nlmixed data=sample_example;

  parms alpha_0=1 alpha_e=2 alpha_m=3 alpha_n=4 alpha_mn=2

  theta_0=2 theta_e=3 theta_n=1

  eta_0=-1 eta_e=0.4 sigma_y=0.5 sigma_m=2;

  MuY= alpha_0+alpha_e*E+alpha_m*M+alpha_n*N

  +alpha_mn*M*N;

  ll_y=-((Y-MuY)**2)/(2*sigma_y)-0.5*log(sigma_y);
```

```

MUm=theta_0+theta_e*E+theta_n*N;

ll_m=-((M-MuM)**2)/(2*sigma_M)-0.5*log(sigma_M);

p_n=(1+exp(-(eta_0+eta_e*E)))**(-1);

ll_n= N*log (p_n)+(1-N)*log(1-p_n);

ll_o=ll_y+ll_m+ll_n;

omega= (1+exp(-(eta_0+eta_e)))**(-1)-(1+exp(-(eta_0)))**(-1);

model Y ~ general(ll_o);

estimate 'nde' alpha_e+(alpha_mn*theta_0+alpha_n)*omega;

run;

```

The data set 'data_example_2' contains variables C, E, N, M, Y from the simulated observation study example. The following sample code produces the maximum likelihood estimates of $NDE(1, 0, c)$, $c = 0, 1$, reported in the text.

```

proc nlmixed data=sample_example;

parms alpha_0=1 alpha_e=2 alpha_c=1 alpha_m=3 alpha_n=4

alpha_me=2 alpha_ne=1 theta_0=2 theta_e=3 theta_c=1 theta_n=0.5

eta_0=-1 eta_e=0.5 eta_c=1 sigma_y=0.5 sigma_m=2 sigma_n=1 ;

MuY= alpha_0+alpha_c*c+alpha_e*E+alpha_m*M+alpha_n*N+alpha_ne*E*N

+alpha_me*E*M;

ll_y=-((Y-MuY)**2)/(2*sigma_y)-0.5*log(sigma_y);

MUm=theta_0+theta_c*C+theta_e*E+theta_n*N;

ll_m=-((M-MuM)**2)/(2*sigma_M)-0.5*log(sigma_M);

MUUn=eta_0+eta_c*C+eta_e*E;

ll_n=      -((N-MuN)**2)/(2*sigma_N)-0.5*log(sigma_N);

```

```

ll_o=ll_y+ll_m+ll_n;

theta_00 = theta_0+theta_n*eta_0;

theta_cc = theta_c+theta_n*eta_c;

theta_ee = theta_e+theta_n*eta_e;

model N~general(ll_o);

estimate 'nde(1,0,1)' alpha_e+alpha_me*(theta_00+theta_cc)
+alpha_n*eta_e+alpha_ne*(eta_0+eta_e+eta_c);

estimate 'nde(1,0,0)' alpha_e+alpha_me*(theta_00)
+alpha_n*eta_e+alpha_ne*(eta_0+eta_e);

run;

```

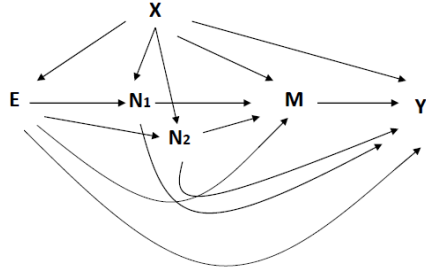


Figure 3. A setting with two independent confounders of M-Y relation affected by E.

PROOFS

Proof of Result 3: By Result 1, the log-survival curve of $Y(e, M(e^))$ at y is given by*

$$\log \sum_{m,n,n',c} \exp \left\{ - \int_0^y [\lambda_0(y) + \gamma(u, e, m, n, c)] du \right\} \\ \times f(M = m | E = e^*, N = n', C = c) \prod_{j=1}^k f_j(n_j, n'_j, e, e^*, c) f(C = c)$$

The result follows upon differentiation of this function with respect to y and multiplication by (-1) .

Proof of Result 4: Recall that $\mathbb{E}\{Y(e, M(e^))\}$*

$$\begin{aligned}
&= \sum_{m,n,n',c} \mathbb{E}(Y|e, m, n, c) f(M = m|E = e^*, N = n', c) \\
&\times f(N(e, c) = n, N(e^*, c) = n'|c) f(c) \\
&= \sum_{m,n,n',c} \left(\beta_m(e, m, c) + \beta_n(e, n, c) + \overbrace{\beta_{m,n}(e, m, n, c)}^{=0 \text{ by assumption}} + \bar{\beta}_{e,c}(e, c) \right) \text{ (by equation (??))} \\
&\times f(M = m|E = e^*, N = n', c) f(N(e, c) = n, N(e^*, c) = n'|c) f(c) \\
&= \sum_{m,n,n',c} (\beta_m(e, m, c) + \beta_n(e, n, c) + \bar{\beta}_{e,c}(e, c)) \\
&\times f(M = m|E = e^*, N = n', C = c) f(N(e, c) = n, N(e^*, c) = n'|c) f(c) \\
&= \sum_{m,n',c} \beta_m(e, m, c) f(M = m|E = e^*, N = n', c) f(N(e^*, c) = n'|c) f(c) \\
&+ \sum_{n,c} \beta_n(e, n, c) f(N(e, c) = n|c) f(c) + \sum_c \bar{\beta}_{e,c}(e, c) f(c) \\
&= \sum_{m,n',c} \beta_m(e, m, c) f(M = m|e^*, n', c) f(N = n'|e^*, c) f(c) \text{ (by NPSEM independence)} \\
&+ \sum_{n,c} \beta_n(e, n, c) f(N = n|e, c) f(c) + \sum_c \bar{\beta}_{e,c}(e, c) f(c) \text{ (by NPSEM independence)} \\
&= \sum_{m,c} \beta_m(e, m, c) f(M = m|e^*, c) f(c) \text{ (by marginalization over } n') \\
&+ \sum_{n,c} \beta_n(e, n, c) f(N = n|e, c) f(c) + \sum_c \bar{\beta}_{e,c}(e, c) f(c)
\end{aligned}$$

proving the result.

Proof of Result 1: Assuming the NPSEM (7) – (11),

$$\begin{aligned}
E \{Y(e, M(e^*))\} &= \sum_{y,m,n,n',c} yf(Y(e, m, n, c) = y, M(e^*, n', c) = m, N(e, c) = n, N(e^*, c) = n', C = c) \\
&= \sum_{y,m,n,n',c} yf(Y(e, m, n, c) = y | M(e^*, n', c) = m, N(e, c) = n, N(e^*, c) = n', C = c) \\
&\quad \times f(M(e^*, n', c) = m | N(e, c) = n, N(e^*, c) = n', C = c) \\
&\quad \times f(N(e, c) = n, N(e^*, c) = n' | C = c) \\
&\quad \times f(C = c) \\
&= \sum_{y,m,n,n',c} yf(Y(e, m, n, c) = y | E = e, M(e, n, c) = m, N(e, c) = n, C = c) \\
&\quad \times f(M(e^*, n', c) = m | E = e^*, N(e^*, c) = n', C = c) \\
&\quad \times f(N(e, c) = n, N(e^*, c) = n' | C = c) \\
&\quad \times f(C = c) \qquad \qquad \qquad \text{(by NPSEM independence)} \\
&= \sum_{y,m,n,n',c} yf(Y = y | E = e, M = m, N = n, C = c) \\
&\quad \times f(M = m | E = e^*, N = n', C = c) f(N(e, c) = n, N(e^*, c) = n' | C = c) f(C = c) \\
&\quad \text{(by consistency)}
\end{aligned}$$

then, under E – N Monotonicity Assumption,

$$\Pr(N(e) = 0, N(e^*) = 1 | c) = 0.$$

and

$$\Pr(N(e) = 1, N(e^*) = 1 | c) = \Pr(N(e^*) = 1 | c) = \Pr\{N = 1 | e^*, c\}$$

Similarly,

$$\Pr(N(e) = 0, N(e^*) = 0|c) = \Pr(N(e) = 0|c) = \Pr\{N = 0|e, c\}.$$

This implies

$$\begin{aligned} & \Pr(N(e) = 1, N(e^*) = 0|c) \\ &= 1 - \Pr(N(e) = 1, N(e^*) = 1|c) - \Pr(N(e) = 0, N(e^*) = 0|c) \\ &= \Pr\{N = 1|e, c\} - \Pr\{N = 1|e^*, c\} \end{aligned}$$

proving the result.

Proof of Result 2:

$$\begin{aligned} E\{Y(e, M(e^*))\} &= \sum_{y,m,n,n',c} yf(Y = y|E = e, M = m, N = n, C = c) \\ &\quad \times f(M = m|E = e^*, N = n^*, C = c) \\ &\quad \times f(N(e, c) = n, N(e^*, c) = n^*|C = c) \\ &\quad \times f(C = c) \\ &= \sum_{y,m,n,n',c} yf(Y = y|E = e, M = m, N = n, C = c) \\ &\quad \times f(M = m|E = e^*, N = n^*, C = c) \\ &\quad \times \prod_{j=1}^k f(N_j(e, c) = n_j, N_j(e^*, c) = n'_j|C = c) f(C = c) \end{aligned}$$

by independence of $(N_1(e, c), N_2(e^*, c)), \dots, (N_j(e, c), N_j(e^*, c))$. Next by monotonicity of the effect

of E on each N_j , one obtains as in the proof of Result 1

$$f\left(N_j(e, c) = n_j, N_j(e^*, c) = n'_j | C = c\right) = f_j\left(n_j, n'_j, e, e^*, c\right)$$

proving the result.