

## Bounds on Direct Effects in the Presence of Confounded Intermediate Variables

Zhihong Cai,<sup>1,\*</sup> Manabu Kuroki,<sup>2</sup> Judea Pearl,<sup>3</sup> and Jin Tian<sup>4</sup>

<sup>1</sup>Department of Biostatistics, Kyoto University School of Public Health, Yoshida-Konoe-cho, Sakyo-ku, Kyoto 606-8501, Japan

<sup>2</sup>Department of Systems Innovation, Osaka University, 1-3, Machikaneyama-cho, Toyonaka, Osaka 560-8531, Japan

<sup>3</sup>Department of Computer Science, UCLA, 4532 Boelter Hall, Los Angeles, California 90024-1596, U.S.A.

<sup>4</sup>Department of Computer Science, Iowa State University, 226 Atanasoff Hall, Ames, Iowa 50011, U.S.A.

\*email: cai@pbh.med.kyoto-u.ac.jp

**SUMMARY.** This article considers the problem of estimating the average controlled direct effect (ACDE) of a treatment on an outcome, in the presence of unmeasured confounders between an intermediate variable and the outcome. Such confounders render the direct effect unidentifiable even in cases where the total effect is unconfounded (hence identifiable). Kaufman et al. (2005, *Statistics in Medicine* **24**, 1683–1702) applied a linear programming software to find the minimum and maximum possible values of the ACDE for specific numerical data. In this article, we apply the symbolic Balke–Pearl (1997, *Journal of the American Statistical Association* **92**, 1171–1176) linear programming method to derive closed-form formulas for the upper and lower bounds on the ACDE under various assumptions of monotonicity. These universal bounds enable clinical experimenters to assess the direct effect of treatment from observed data with minimum computational effort, and they further shed light on the sign of the direct effect and the accuracy of the assessments.

**KEYWORDS:** Causal effect; Midpoint estimator; Potential response type; Stratified analysis.

### 1. Introduction

Estimation of the direct effect of a treatment on an outcome is a central concern in epidemiological and clinical research (Robins and Greenland, 1992; Buyse and Molenberghs, 1998; Wang and Taylor, 2002; Rubin, 2004; Kaufman et al., 2005; Taylor, Wang, and Thiebaut, 2005; Petersen, Sinisi, and van der Laan, 2006). Pearl (2001) gave a formal definition of the total effect decomposition into direct and indirect effects, and distinguished between the controlled direct effect and the natural direct effect, the former is obtained when intermediate variables are held constant at specific values. Kaufman et al. (2005) considered the problem of estimating the average controlled direct effect (ACDE) of a treatment on an outcome, in the presence of unmeasured confounders between an intermediate variable and the outcome. Such confounders render the direct effect unidentifiable even in cases where the total effect is unconfounded (hence identifiable). Kaufman et al. (2005) applied a linear programming software to find the minimum and maximum possible values of the ACDE for specific numerical data. They further proposed the midpoint between the minimum and maximum values as an estimator of the ACDE. However, they did not provide exact formulas of the bounds on the ACDE.

In this article, we apply the symbolic Balke–Pearl linear programming method (Balke, 1995; Balke and Pearl, 1997) to derive closed-form formulas of the upper and lower bounds

on the ACDE under various assumptions of monotonicity. In contrast to the numerical method of Kaufman et al. (2005), these symbolic formulas enable clinical experimenters to assess the direct effect of a treatment on an outcome from observed data with minimum computational effort, and they further shed light on the accuracy of the assessment. In addition, we derive bounds on the ACDE when covariate information is available. Moreover, we provide a formal formula for the midpoint estimator chosen by Kaufman et al. (2005), and propose a new stratified midpoint estimator that is more accurate when covariate measurements are available. In addition to the binary case, we further propose bounds on the ACDE in the case where observed variables are multicategorical. Finally, we illustrate our results through an empirical example in both binary and multicategorical cases.

### 2. Bounding Formulas

#### 2.1 Problem Description

To motivate our problem, we examine the data from the Lipid Research Clinics Coronary Primary Prevention Trial (LRC-CPPT) (shown in Table 1; LRT-CPPT group, 1984; Kaufman et al., 2005). The purpose of this study is to evaluate the efficacy of the cholesterol-lowering drug cholestyramine for the prevention of coronary heart disease (CHD) in 3806 hypercholesterolemia men. Our interest is to examine whether serum cholesterol level 1 year after initiation of cholestyramine was

**Table 1**  
*Definite CHD mortality or myocardial infarction events (Y) in the LRC-CPPT study according to randomized cholestyramine treatment group (X) and serum cholesterol (md/dl) at 1 year (Z) (Kaufman et al., 2005)*

	Placebo ( $x_1$ )			Cholestyramine treatment ( $x_0$ )		
	Cholesterol $\geq 280$ mg/dl ( $z_1$ )	Cholesterol $< 280$ mg/dl ( $z_0$ )	Total placebo	Cholesterol $\geq 280$ mg/dl ( $z_1$ )	Cholesterol $< 280$ mg/dl ( $z_0$ )	Total treatment
$y_1$	82	86	168	33	97	130
$y_0$	669	1081	1750	332	1426	1758
Total	751	1167	1918	365	1523	1888

an adequate surrogate endpoint (i.e., explanation) for the outcome of CHD.

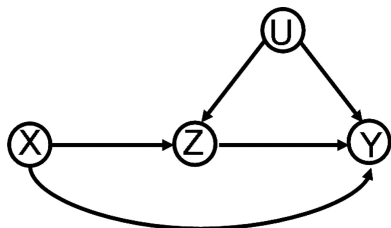
According to Freedman, Graubard, and Schatzkin (1992), a good surrogate endpoint is one that explains a large proportion of the total effect. A conventional approach to validate a surrogate endpoint is to estimate the relative contributions of the direct and indirect effects to the total effect. However, if there exist unmeasured confounding factors, for example, if there exist unmeasured genetic or life style factors that affect both cholesterol and CHD, estimating the direct effect requires careful causal analysis.

To model presence of unmeasured confounding, we consider the directed acyclic graph shown in Figure 1, where a treatment  $X$ , an intermediate  $Z$ , and an outcome  $Y$  are binary variables with values  $x, z$ , and  $y$ , respectively, ( $x \in \{x_0, x_1\}$ ,  $y \in \{y_0, y_1\}$ ,  $z \in \{z_0, z_1\}$ ), and  $U$  is a set of unmeasured variables, which is independent of  $X$ . In this figure, the treatment is assumed to be randomized, hence there is no confounder between  $X$  and  $Y$ , and the total effect of  $X$  on  $Y$  is identifiable. However, the set of unmeasured confounders  $U$  between  $Z$  and  $Y$  renders the direct effect of  $X$  on  $Y$  unidentifiable. In other words, it is impossible to estimate this direct effect without making further assumptions. The central aim of this article is to derive formulas of the bounds on the direct effect of  $X$  on  $Y$  in this causal structure.

The ACDEs are defined as

$$ACDE(z) = \text{pr}\{y_1 \mid \text{do}(x_1), \text{do}(z)\} - \text{pr}\{y_1 \mid \text{do}(x_0), \text{do}(z)\}, \tag{1}$$

for  $z \in \{z_0, z_1\}$ , where  $\text{do}(\cdot)$  denotes an imposed intervention (Pearl, 2000; Kaufman et al., 2005).  $\text{pr}\{y \mid \text{do}(x), \text{do}(z)\}$



**Figure 1.** A directed acyclic graph with a measured treatment  $X$ , an intermediate  $Z$ , and an outcome  $Y$ , and a set of unmeasured variables  $U$  (Kaufman et al., 2005).

indicates the probability of  $Y = y$  when we set  $X$  and  $Z$  to specific values  $x$  and  $z$ , respectively, by a joint intervention (Pearl, 2000).

Using the counterfactual notation of Neyman (1923) and Rubin (1974), equation (1) can be also written as

$$ACDE(z) = \text{pr}(Y_{x_1, z} = y_1) - \text{pr}(Y_{x_0, z} = y_1). \tag{2}$$

A formal translation from graphs to counterfactual models is given by Pearl (2000).

Equation (1) represents the average causal effect of  $X$  on  $Y$  when the causal path through  $Z$  is blocked by holding  $Z$  fixed at  $z_0$  or  $z_1$ . Note that equation (1) is different from the crude stratum-specific risk difference  $\text{pr}(y_1 \mid x_1, z) - \text{pr}(y_1 \mid x_0, z)$ . The latter stands for the observed conditional risk difference in stratum  $z$  (in our example, the subgroup with cholesterol  $< 280$  mg/dl or cholesterol  $\geq 280$  mg/dl), which represents the direct effect of  $X$  on  $Y$ , plus the spurious correlation between  $X$  and  $Y$  through the path  $X \rightarrow Z \leftarrow U \rightarrow Y$ . On the other hand, equation (1) represents the direct effect only, because the path  $X \rightarrow Z \leftarrow U \rightarrow Y$  had been blocked by an intervention on  $Z$ . If the ACDE equals 0 in our example, then we can judge that cholesterol ( $Z$ ) is a perfect surrogate endpoint for CHD ( $Y$ ), which suggests that lowering cholesterol level constitutes an adequate explanation for how the drug prevents the occurrence of CHD.

In order to derive bounds on the ACDE, we follow Kaufman et al. (2005) and define 64 potential response types. First, we consider  $X$  (cholestyramine) as a treatment and  $Z$  (cholesterol) as an outcome. Because  $X$  and  $Z$  are binary variables, there are four possible potential response types at the unit level: (1) a subject whose cholesterol increases regardless of taking cholestyramine or placebo (doomed), (2) a subject whose cholesterol decreases only by taking cholestyramine (causative), (3) a subject whose cholesterol decreases only by not taking cholestyramine (preventive), and (4) a subject whose cholesterol decreases regardless of taking cholestyramine or placebo (immune) (Greenland and Robins, 1986). We index these four types by a mapping variable  $r_z = 1, 2, 3, 4$ . Similarly, when we consider  $X$  (cholestyramine) as a treatment and  $Y$  (CHD) as an outcome with  $Z$  (cholesterol) fixed to  $z_0$  or  $z_1$ , there still exist doomed, causative, preventive, and immune potential response types. We denote these four types by a mapping variable  $r_{y|z_0} = 1, 2, 3, 4$  when  $Z$  is fixed to  $z_0$ , and another mapping variable  $r_{y|z_1} = 1, 2, 3, 4$  when  $Z$  is fixed to  $z_1$ . Therefore, any of the  $4 \times 4 \times 4$  index triples,

$(r_z, r_{y|z_0}, r_{y|z_1})$  represents a potential response type. The joint probability distribution of  $(r_z, r_{y|z_0}, r_{y|z_1})$  is defined by

$$q_{ijk} = \text{pr}(r_z = i, r_{y|z_0} = j, r_{y|z_1} = k),$$

for  $i, j, k = 1, 2, 3, 4$ , and the  $\{q_{ijk}\}$  represent the proportion of the 64 potential response types among the population. Thus, the population of interest is fully characterized by  $\{q_{ijk}\}$ , and we can rewrite equation (2) as

$$ACDE(z_1) = \sum_{i=1}^4 \sum_{j=1}^4 \left( \sum_{k \in \{1,2\}} q_{ijk} - \sum_{k \in \{1,3\}} q_{ijk} \right)$$

$$ACDE(z_0) = \sum_{i=1}^4 \sum_{k=1}^4 \left( \sum_{j \in \{1,2\}} q_{ijk} - \sum_{j \in \{1,3\}} q_{ijk} \right).$$

See Web Appendix A for a detail derivation.

Kaufman et al. (2005) applied a linear program software package to find the minimum and maximum possible values of the ACDE for specific numerical data. However, they did not provide exact formulas for the ACDE. Balke (1995) and Balke and Pearl (1997) describe a computer program that takes symbolic description of linear programming problems and returns symbolic expressions for the desired bounds. In this article, we apply this symbolic method to derive closed-form formulas for the ACDEs under three sets of assumptions. Details of this method are included in Web Appendix A.

### 2.2 No Assumption Case

When no assumption is made, there are 64 potential response types, whereas there are only eight observed conditional probabilities  $\text{pr}(y, z | x)$ . Using the symbolic Balke–Pearl method (Balke, 1995; Balke and Pearl, 1997), the formulas for the tightest lower and upper bounds on the ACDEs are given by

$$\text{pr}(y_0, z | x_0) + \text{pr}(y_1, z | x_1) - 1 \leq ACDE(z)$$

$$\leq 1 - \text{pr}(y_1, z | x_0) - \text{pr}(y_0, z | x_1), \quad (3)$$

for  $z \in \{z_1, z_0\}$ , which defines the range within which the ACDE must lie. It is remarkable that we get such a simple formula, consisting of only one additive expression in the lower bound and one additive expression in the upper bound.

To find when the lower bound coincides with the upper bound, we calculate their difference and obtain  $\text{pr}(z_{1-i} | x_0) + \text{pr}(z_{1-i} | x_1)$  for  $ACDE(z_i)$  ( $z_i \in \{z_1, z_0\}$ ). Hence, in order to make the lower bound equal the upper bound, both  $\text{pr}(z_{1-i} | x_0)$  and  $\text{pr}(z_{1-i} | x_1)$  must be zero. This indicates that the upper bound cannot coincide with the lower bound in both  $ACDE(z_0)$  and  $ACDE(z_1)$  at the same time, because the probabilities in all the cells must be zero in order to achieve it. That is, the bounding interval never vanishes, regardless of the observations.

In addition, it should be noted that equation (3) provides a simple testable criterion for the existence of a direct effect, that is, if  $\text{pr}(y_0, z | x_0) + \text{pr}(y_1, z | x_1) > 1$ , then we are assured that  $ACDE(z)$  is positive, and if  $\text{pr}(y_1, z | x_0) + \text{pr}(y_0, z | x_1) > 1$ ,  $ACDE(z)$  must be negative.

### 2.3 Monotonic Assumption Case

Kaufman et al. (2005) made two assumptions regarding the potential response types: (1) monotonic assumption, which means no unit-level causal effects of  $X$  on  $Z$  or of  $X$  on  $Y$  or

of  $Z$  on  $Y$  can be negative, and (2) no-interaction assumption, which means that, for all units, the response of  $Y$  to change in  $X$  does not depend on the level at which we hold  $Z$ . There are 18 potential response types that satisfy monotonic assumption, that is,  $\{q_{i11}, q_{i21}, q_{i22}, q_{i41}, q_{i42}, q_{i44} : i = 1, 2, 4\}$ , and 12 potential response types that satisfy both monotonic and no-interaction assumptions, that is,  $\{q_{i11}, q_{i22}, q_{i41}, q_{i44} : i = 1, 2, 4\}$  (Kaufman and Kaufman, 2006, personal communication). By applying the Balke–Pearl method, we derive closed-form formulas for the tightest bounds on the ACDEs in the two cases. The following equations give the lower and upper bounds under monotonic assumption:

$$\max \left\{ \begin{array}{c} 0 \\ \text{pr}(y_0, z_1 | x_0) - \text{pr}(y_0, z_1 | x_1) \end{array} \right\}$$

$$\leq ACDE(z_1) \leq \text{pr}(y_0 | x_0) - \text{pr}(y_0, z_1 | x_1). \quad (4)$$

$$\max \left\{ \begin{array}{c} 0 \\ \text{pr}(y_1, z_0 | x_1) - \text{pr}(y_1, z_0 | x_0) \end{array} \right\}$$

$$\leq ACDE(z_0) \leq \text{pr}(y_1 | x_1) - \text{pr}(y_1, z_0 | x_0). \quad (5)$$

It is seen that the interval collapses when  $\text{pr}(y_0, z_1 | x_0) = \text{pr}(y_0 | x_0)$  (or  $\text{pr}(y_0, z_0 | x_0) = 0$ ) in equation (4) and  $\text{pr}(y_1, z_0 | x_1) = \text{pr}(y_1 | x_1)$  (or  $\text{pr}(y_1, z_1 | x_1) = 0$ ) in equation (5). In these cases, the ACDEs can be evaluated by the lower or upper bound.

On the basis of the lower bounds of equations (4) and (5), we can judge whether there exist positive direct effects under the monotonic assumption. That is, if  $\text{pr}(y_0, z_1 | x_0) > \text{pr}(y_0, z_1 | x_1)$  and/or  $\text{pr}(y_1, z_0 | x_1) > \text{pr}(y_1, z_0 | x_0)$ , then we are assured that there exist positive direct effects.

Moreover, equations (4) and (5) provide a simple necessary test for the monotonic assumption. That is, if the monotonic assumption holds true, then the upper bounds should be no less than zero, because the lower bounds are nonnegative. Thus, if the observed quantities are  $\text{pr}(y_0 | x_0) < \text{pr}(y_0, z_1 | x_1)$  or  $\text{pr}(y_1 | x_1) < \text{pr}(y_1, z_0 | x_0)$ , then the upper bounds would be negative, which indicates that the monotonic assumption does not hold in this situation.

On the other hand, the following equation gives the lower and upper bounds under both monotonic and no-interaction assumptions:

$$\max \left\{ \begin{array}{c} 0 \\ \text{pr}(y_0, z_1 | x_0) - \text{pr}(y_0, z_1 | x_1) \\ \text{pr}(y_1, z_0 | x_1) - \text{pr}(y_1, z_0 | x_0) \\ \text{pr}(y_0, z_1 | x_0) - \text{pr}(y_0, z_1 | x_1) + \text{pr}(y_1, z_0 | x_1) - \text{pr}(y_1, z_0 | x_0) \end{array} \right\}$$

$$\leq ACDE(z) \leq \text{pr}(y_1 | x_1) - \text{pr}(y_1 | x_0), \quad (6)$$

for  $z \in \{z_1, z_0\}$ . It is seen that the upper bound is the total effect of  $X$  on  $Y$ . We can judge whether there exist positive direct effects from the lower bound of equation (6). That is, if either  $\text{pr}(y_0, z_1 | x_0) > \text{pr}(y_0, z_1 | x_1)$  or  $\text{pr}(y_1, z_0 | x_1) > \text{pr}(y_1, z_0 | x_0)$ , then we are assured that there exist positive direct effects.

Further, equation (6) provides a simple necessary test for both monotonicity and no-interaction assumptions. Because  $ACDE(z)$  must be nonnegative from the lower bounds of equation (6), then, if  $\text{pr}(y_1 | x_1) < \text{pr}(y_1 | x_0)$ , the upper bounds

would be negative, which indicates that at least one of the two assumptions is violated.

One more thing to be mentioned is that, because the monotonic assumption and no-interaction assumption add some constraints on the potential response types, the bounds under monotonic assumption should not be wider than those under no assumption, and similarly, the bounds under both monotonic and no-interaction assumptions should not be wider than those under monotonic assumption.

2.4 Estimation Accuracy

In Sections 2.2 and 2.3, we derive the estimators for the lower and upper bounds under three sets of assumptions. Another problem is the estimation accuracy of these estimators. For the no-assumption case, it is easy to obtain the exact variances for the lower and upper bounds. However, for the remaining two cases, it is very complicated to derive the variances for the lower bounds though it is easy to obtain the exact variances for the upper bounds. In Web Appendix B, we provide the variance estimators for the lower and upper bounds under the three cases. In addition, we evaluate the performance of the proposed variance estimators through simulation studies.

3. Stratified ACDE

The analysis of Section 2 applies to situations where all confounders between  $Z$  and  $Y$  are unmeasured. However, if some of these confounders are observed, this information is helpful in narrowing the bounds on direct effects. In this section, we consider the directed acyclic graph with the set of confounders  $U$  in Figure 1 being divided into two sets of variables: measured covariates  $S$  and unmeasured covariates  $W$ .

Then, we define the stratified ACDE as

$$ACDE(z | s) = \text{pr}\{y_1 | \text{do}(x_1), \text{do}(z), s\} - \text{pr}\{y_1 | \text{do}(x_0), \text{do}(z), s\} \tag{7}$$

for  $s \in \{s_1, \dots, s_k\}$  and  $z \in \{z_1, z_0\}$ , where

$$\text{pr}\{y | \text{do}(x), \text{do}(z), s\} = \frac{\text{pr}\{y, s | \text{do}(x), \text{do}(z)\}}{\text{pr}\{s | \text{do}(x), \text{do}(z)\}}.$$

Because  $S$  is a set of observed baseline covariates, it is not affected by  $X$  or  $Z$ , which implies  $\text{pr}\{s | \text{do}(x), \text{do}(z)\} = \text{pr}(s)$  (Pearl, 2000). Let  $k$  be the number of categories of  $S$ , then  $64^k$  potential response types are needed in order to obtain the tightest bounds on the ACDE. However, if we limit the potential responses to the case where  $S = s$  is observed, then there are in total 64 potential response types in each stratum. Therefore, we can apply the previous discussion to derive the bounds on the stratified ACDEs. When there exist 64 potential response types, the stratified ACDE is given by

$$\begin{aligned} &\text{pr}(y_0, z | x_0, s) + \text{pr}(y_1, z | x_1, s) - 1 \\ &\leq ACDE(z | s) \leq 1 - \text{pr}(y_1, z | x_0, s) - \text{pr}(y_0, z | x_1, s), \end{aligned}$$

for  $z \in \{z_1, z_0\}$ . Similarly, if the monotonic assumption holds true in stratum  $s$ , we can obtain

$$\begin{aligned} &\max \left\{ \begin{array}{c} 0 \\ \text{pr}(y_0, z_1 | x_0, s) - \text{pr}(y_0, z_1 | x_1, s) \end{array} \right\} \\ &\leq ACDE(z_1 | s) \leq \text{pr}(y_0 | x_0, s) - \text{pr}(y_0, z_1 | x_1, s), \\ &\max \left\{ \begin{array}{c} 0 \\ \text{pr}(y_1, z_0 | x_1, s) - \text{pr}(y_1, z_0 | x_0, s) \end{array} \right\} \\ &\leq ACDE(z_0 | s) \leq \text{pr}(y_1 | x_1, s) - \text{pr}(y_1, z_0 | x_0, s). \end{aligned}$$

In addition, if both the monotonic and no-interaction assumptions hold true in stratum  $s$ , we can obtain

$$\begin{aligned} &\max \left\{ \begin{array}{c} 0 \\ \text{pr}(y_0, z_1 | x_0, s) - \text{pr}(y_0, z_1 | x_1, s) \\ \text{pr}(y_1, z_0 | x_1, s) - \text{pr}(y_1, z_0 | x_0, s) \\ \text{pr}(y_0, z_1 | x_0, s) - \text{pr}(y_0, z_1 | x_1, s) + \text{pr}(y_1, z_0 | x_1, s) - \text{pr}(y_1, z_0 | x_0, s) \end{array} \right\} \\ &\leq ACDE(z | s) \leq \text{pr}(y_1 | x_1, s) - \text{pr}(y_1 | x_0, s), \end{aligned} \tag{8}$$

for  $z \in \{z_1, z_0\}$ . Thus, because the  $ACDE(z)$  can be obtained by

$$ACDE(z) = \sum_s ACDE(z | s)\text{pr}(s),$$

letting  $LB_s(z)$  and  $UB_s(z)$  be the lower bound and the upper bound in stratum  $s$ , respectively, the summarized bounds on the  $ACDE(z)$  by using covariate information can be evaluated by

$$\sum_s LB_s(z)\text{pr}(s) \leq ACDE(z) \leq \sum_s UB_s(z)\text{pr}(s). \tag{9}$$

These summarized bounds on direct effects are not wider than the bounds derived in Section 2, a simple proof of which is provided in Web Appendix C.

We would like to point out some practical requirements for the observed covariates  $S$ . First of all,  $S$  must be baseline covariates in order for the method to be valid. Moreover, we can divide such baseline covariates into the following three cases: (a)  $S$  is a confounder between  $Z$  and  $Y$ ; (b)  $S$  has an effect on  $Z$  but not on  $Y$ ; (c)  $S$  has an effect on  $Y$  but not on  $Z$ . If the measured covariate  $S$  satisfies any of the three cases, then the summarized bounds of equation (9) should not be wider than those provided in Section 2.

4. Midpoint Estimator

Kaufman et al. (2005) proposed the midpoint between the minimum and maximum values as an estimator of the ACDE, which is given as

$$\text{mRD}(z) = \frac{UB(z) + LB(z)}{2}, \quad z \in \{z_0, z_1\}, \tag{10}$$

where  $LB(z)$  and  $UB(z)$  are the linear programming minimum and maximum values for  $ACDE(z)$  derived from the observed probabilities using linear programming packages. With the derived formulas in Section 2, we can now present an exact formula of the midpoint estimator. For example, the midpoint estimator with no assumption is derived directly as

$$\begin{aligned} \text{mRD}(z) = &\frac{1}{2} \{ \text{pr}(y_0, z | x_0) + \text{pr}(y_1, z | x_1) \\ &- \text{pr}(y_1, z | x_0) - \text{pr}(y_0, z | x_1) \} \end{aligned} \tag{11}$$

**Table 2**

*A hypothetical example: proportion of potential response types in strata  $s_1$  and  $s_0$*

	$s_1$			$s_0$			
$q_{111}$	0.10	$q_{141}$	0.01	$q_{111}$	0.10	$q_{141}$	0.20
$q_{211}$	0.01	$q_{241}$	0.01	$q_{211}$	0.30	$q_{241}$	0.01
$q_{411}$	0.30	$q_{441}$	0.02	$q_{411}$	0.10	$q_{441}$	0.01
$q_{122}$	0.01	$q_{144}$	0.10	$q_{122}$	0.10	$q_{144}$	0.10
$q_{222}$	0.01	$q_{244}$	0.01	$q_{222}$	0.01	$q_{244}$	0.01
$q_{422}$	0.20	$q_{444}$	0.22	$q_{422}$	0.01	$q_{444}$	0.05

based on equation (3). The midpoint estimators for the remaining two cases can be derived in the same way. Thus, we can calculate the midpoint estimator from the observed data without using linear programming packages.

When covariate information is available, we propose a new stratified midpoint estimator, which is given by

$$\text{mRD}_s(z) = \sum_s \frac{LB_s(z) + UB_s(z)}{2} \text{pr}(s), \quad (12)$$

where  $S$  is a set of observed baseline covariates discussed in Section 3.

The new stratified midpoint estimator is superior to the midpoint estimator when some covariates are observed. To see this, we consider a hypothetical example when both monotonic and no-interaction assumptions hold true, and there is a binary observed covariate  $S$ . Table 2 shows the true proportion of 12 potential response types in each stratum, and Table 3 shows the observed conditional probabilities  $\text{pr}(y, z | x, s)$  induced from Table 2. Here,  $\text{pr}(s_1)$  is set to be 0.45.

Then, according to Kaufman et al.’s method, the bounds on the direct effect are (0.050, 0.175), and the midpoint estimate is 0.112. On the other hand, the bounds are (0.190, 0.230) in stratum  $s_1$  and (0.090, 0.130) in stratum  $s_0$  according to our formula (8). Then, we calculate the summarized lower and upper bounds according to our summarized formula (9), which are (0.135, 0.175), and the stratified midpoint estimator according to formula (12), which is 0.155. Here, we can calculate the true stratified ACDE from Table 2, which is 0.220 in stratum  $s_1$  and 0.120 in stratum  $s_0$ . In addition, the true ACDE is 0.165 from Table 2, which is included in both Kaufman et al.’s bounds and our bounds. However, it is seen that Kaufman et al.’s midpoint estimator is quite away from the true ACDE and outside our bounds, whereas the stratified midpoint estimator is close to the true ACDE.

**Table 3**

*Observed conditional probabilities  $\text{pr}(y, z | x, s)$  induced from Table 2*

		$s_1$		$s_0$	
		$y_1$	$y_0$	$y_1$	$y_0$
$x_1$	$z_1$	0.15	0.11	0.72	0.11
	$z_0$	0.50	0.24	0.11	0.06
$x_0$	$z_1$	0.11	0.11	0.30	0.20
	$z_0$	0.31	0.47	0.40	0.10

### 5. Extension to Multicategorical Case

In the discussion above, we consider the ACDE when observed variables are binary. In this section, we consider the case where  $X, Y$ , and  $Z$  are multicategorical variables. When the categorical treatment variable  $X$  is changed from  $x$  to  $x'$ , we define the ACDE as

$$\text{ACDE}(y, z, x, x') = \text{pr}\{y | \text{do}(x'), \text{do}(z)\} - \text{pr}\{y | \text{do}(x), \text{do}(z)\}.$$

where  $y$  and  $z$  are possible values of  $Y$  and  $Z$ , respectively. Then, we provide the lower and upper bounds on the ACDE under the multicategorical case:

$$\begin{aligned} -1 + \text{pr}(z | x) + \text{pr}(y, z | x') - \text{pr}(y, z | x) &\leq \text{ACDE}(y, z, x, x') \\ &\leq 1 - \text{pr}(z | x') + \text{pr}(y, z | x') - \text{pr}(y, z | x). \end{aligned}$$

The proof is given in Web Appendix D. When  $X, Y$ , and  $Z$  are binary variables, these bounds are consistent with equation (3). Kang and Tian (2006) provided a method to obtain the inequality constraint for causal effects from nonexperimental data in the presence of unobserved variables. The above bounds can also be obtained by using their method.

### 6. Empirical Example

#### 6.1 Binary Case

We illustrate our results through the example given in Section 2. Kaufman et al. (2005) collapsed the serum cholesterol values into two categories from five original categories, based on the data in Freedman et al. (1992). We will discuss the five categories in the next subsection.

Because treatment  $X$  is randomized, the total effect of  $X$  on  $Y$  can be estimated by the risk difference  $\text{pr}(y_1 | x_1) - \text{pr}(y_1 | x_0) = 0.0876 - 0.0689 = 0.0187$ . On the other hand, the observed stratum-specific risk difference is  $\text{pr}(y_1 | x_1, z_0) - \text{pr}(y_1 | x_0, z_0) = 0.0737 - 0.0637 = 0.0100$  in stratum  $z_0$ , and  $\text{pr}(y_1 | x_1, z_1) - \text{pr}(y_1 | x_0, z_1) = 0.1092 - 0.0904 = 0.0188$  in stratum  $z_1$ . Thus, as noted in Kaufman et al. (2005), there appears to be a direct causative effect of not receiving cholestyramine on the risk of CHD in each stratum of intermediate.

The bounds on the ACDE under no assumption are  $[-0.1999, 0.3850]$  in stratum  $z_0$ , and  $[-0.7814, 0.6337]$  in stratum  $z_1$ , which are relatively wide. Here, according to Kaufman et al. (2005), it is reasonable to assume that neither cholestyramine nor absence of hyperlipidaemia may elevate risk of the outcomes, nor may cholestyramine elevate serum cholesterol, leading to 18 potential response types for consideration. In addition, the necessary test for the monotonic assumption in Section 2 shows that  $\text{pr}(y_0 | x_0) - \text{pr}(y_0, z_1 | x_1) = 0.5823 > 0$ , and  $\text{pr}(y_1 | x_1) - \text{pr}(y_1, z_0 | x_0) = 0.0362 > 0$ , which suggests that the monotonic assumption holds for the data. Then, according to our formulas, the bounds are  $[0, 0.0362]$  in stratum  $z_0$ , and  $[0, 0.5823]$  in stratum  $z_1$ . The upper bound in stratum  $z_1$  can be as large as 0.5823, which is much larger than the total effect 0.0187. Even the midpoint estimator is 0.2912, larger than 0.0187. Therefore, it may not be helpful to calculate the relative contribution of the direct and indirect effects to the total effect, in order to validate the serum cholesterol level as a surrogate endpoint. One explanation is that there exists potential response type  $q_{442}$ , which

Table 4

Data from definite CHD mortality or myocardial infarction events ( $Y$ ) in the Lipid Research Clinics Coronary Primary Prevention Trial (Freedman et al., 1992) and the lower and upper bounds on the ACDE in five cholesterol categories

Cholesterol ( $Z$ )	Placebo ( $x_1$ )		Cholestyramine ( $x_0$ )		Bounds	
	$y_1$	$y_0$	$y_1$	$y_0$	Lower	Upper
<180	0	7	9	97	-0.949	0.992
180–230	8	83	34	641	-0.656	0.939
230–280	78	991	54	688	-0.595	0.455
280–330	64	572	23	281	-0.818	0.690
>330	18	97	10	51	-0.964	0.944

contributes the  $ACDE(z_1)$  value but does not contribute to the total effect.

When we restrict to 12 potential response types, again, the necessary test for no-interaction assumption holds, that is,  $\text{pr}(y_1 | x_1) - \text{pr}(y_1 | x_0) = 0.0187 > 0$ . The bounds on the  $ACDE(z)$  are  $[0, 0.0187]$  in both strata. The upper bound equals the total effect, because the interactive potential response type  $q_{442}$  does not exist. The midpoint estimator gives an estimate 0.0094, which indicates that there may exist a direct effect of cholestyramine treatment on CHD without mediating serum cholesterol.

Moreover, it is noted that the bounds under the monotonic assumption are narrower than those under no assumption, and the bounds under both monotonic and no-interaction assumptions are narrower than those under monotonic assumption. The reason is that these assumptions make some constraints on the potential response types.

### 6.2 Multicategorical Case

Freedman et al. (1992) provided the data from the LRC-CPPT study, where the serum cholesterol values ( $Z$ ) have five categories (shown in Table 4). Based on our formulas in Section 5, we calculate the lower and upper bounds when the serum cholesterol is fixed at each of the five categories, which are shown in Table 4. When we compare the bounds in binary case with those in Table 4, it is seen that with the number of categories of  $Z$  increases, the observed probabilities become smaller and the bounds become wider, which indicates that the width of the bounds is dependent on the sparsity of the observations. However, the bounds are helpful if one is interested in the ACDE under more detailed categories, which the bounds of binary case cannot provide.

### 7. Discussion

This article applied the symbolic Balke–Pearl method to derive closed-form formulas for the lower and upper bounds on the ACDEs under three sets of assumptions. We also considered extensions to situations where the treatment, the intermediate, and the outcome are multinomial rather than dichotomous variables, as well as situations in which the confounding factors are partially observed, so that covariate-adjusted bounds and midpoint estimators can be obtained.

Because our approach is nonparametric and mainly based on observed information, the proposed bounds define a range within which the direct effect must lie. On the basis of these deterministic bounds, one can narrow the bound width substantially by introducing subject matter constraints. There-

fore, these universal bounds are helpful for epidemiologists and clinical experimenters to assess the direct effect of treatment.

### 8. Supplementary Materials

Web Appendices and Tables referenced in Sections 2, 3, and 5 are available under the Paper Information link at the *Biometrics* website <http://www.biometrics.tibs.org>.

### ACKNOWLEDGEMENTS

We would like to thank Professor Sol Kaufman of University at Buffalo and Professor Jay S. Kaufman of University of North Carolina for their helpful communication. We are also grateful to Professor Tosiya Sato of Kyoto University for his valuable comments. This article was partially supported by the Japan Society for the Promotion of Science, the Ministry of Education, Culture, Sports, Science and Technology of Japan, the Kayamori Foundation of Informational Science Advancement, the Kurara Foundation, the Mazda Foundation, and the NSF grants IIS-0535223 and IIS-0347846.

### REFERENCES

- Balke, A. (1995). *Probabilistic counterfactuals: Semantics, computation, and applications*. Technical Report (R-242), UCLA Cognitive Systems Laboratory (Ph.D. Thesis).
- Balke, A. and Pearl, J. (1997). Bounds on treatment effects from studies with imperfect compliance. *Journal of the American Statistical Association* **92**, 1171–1176.
- Buyse, M. and Molenberghs, G. (1998). Criteria for the validation of surrogate endpoints in randomized experiments. *Biometrics* **54**, 1014–1029.
- Freedman, L. S., Graubard, B. I., and Schatzkin, A. (1992). Statistical validation of intermediate endpoints for chronic diseases. *Statistics in Medicine* **11**, 167–178.
- Greenland, S. and Robins, J. M. (1986). Identifiability, exchangeability, and epidemiological confounding. *International Journal of Epidemiology* **15**, 413–419.
- Kang, C. and Tian, J. (2006). Inequality constraints in causal models with hidden variables. In *Proceedings of the Twenty Second Conference on Uncertainty in Artificial Intelligence*, 411–420. Boston: UAI Press.
- Kaufman, S. and Kaufman, J. S. (2006). Personal Communication.
- Kaufman, S., Kaufman, J. S., MacLenose, R. F., Greenland, S., and Poole, C. (2005). Improved estimation of

- controlled direct effects in the presence of unmeasured confounding of intermediate variables. *Statistics in Medicine* **24**, 1683–1702. Correction, **25**, 3228.
- The Lipid Research Clinics Coronary Primary Prevention Trial Results. (1984). I. Reduction in incidence of coronary heart disease. *Journal of the American Medical Association* **251**, 351–364.
- Neyman, J. (1923). Justification of applications of the calculus of probabilities to the solutions of certain questions in agricultural experimentation. Excerpts English translation (1990). *Statistical Science* **5**, 463–472.
- Pearl, J. (2000). *Causality: Models, Reasoning, and Inference*. New York: Cambridge University Press.
- Pearl, J. (2001). Direct and indirect effects. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, 411–420. San Francisco: Morgan Kaufmann.
- Petersen, M. L., Sinisi, S. E., and van der Laan, M. J. (2006). Estimation of direct causal effects. *Epidemiology* **17**, 276–284.
- Robins, J. M. and Greenland, S. (1992). Identifiability and exchangeability for direct and indirect effects. *Epidemiology* **3**, 143–155.
- Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology* **66**, 688–701.
- Rubin, D. B. (2004). Direct and indirect causal effects via potential outcomes. *Scandinavian Journal of Statistics* **31**, 161–170.
- Taylor, J. M. G., Wang, Y., and Thiebaut, R. (2005). Counterfactual links to the proportion of treatment effect explained by a surrogate marker. *Biometrics* **61**, 1102–1111.
- Wang, Y. and Taylor, J. M. G. (2002). A measure of the proportion of treatment effect explained by a surrogate marker. *Biometrics* **58**, 803–812.

Received January 2007. Revised September 2007.  
Accepted September 2007.