

# Causal Inference from Indirect Experiments

Judea Pearl

Cognitive Systems Laboratory

Computer Science Department

University of California, Los Angeles, CA 90024

Phone: (310) 825-3243

Fax: (310) 825-2273

*judea@cs.ucla.edu*

## Abstract

An indirect experiment is a study in which randomized control is replaced by randomized *encouragement*, that is, subjects are encouraged, rather than forced, to receive a given treatment program. The purpose of this paper is to bring to the attention of experimental researchers simple mathematical results that enable us to assess, from indirect experiments, the strength with which causal influences operate among variables of interest. The results reveal that despite the laxity of the encouraging instrument, data from indirect experimentation can yield significant and sometimes accurate information on the impact

of a program on the population as a whole, as well as on the particular individuals who participated in the program.

Keywords: Causal reasoning, treatment evaluation, noncompliance, graphical models

## 1 Introduction

Standard experimental studies in the biological, medical, and behavioral sciences invariably invoke the instrument of randomized control, that is, subjects are assigned at random to various groups (or treatments or programs) and the mean differences between participants in different groups are regarded as measures of the efficacies of the associated programs. *Indirect experiments* are studies in which randomized control is either infeasible or undesirable, so randomized *encouragement* is instituted instead, that is, subjects are still assigned at random to various groups, but members of each group are encouraged, rather than forced, to participate in the program associated with the group; it is up to the individuals to select among the programs.

Recently, use of randomization in social and medical experimentation has been questioned. The objections seem to fall into three major categories:

1. Perfect control is hard to achieve or ascertain. Studies in which treatment is assumed to be randomized may be marred by uncontrolled imperfect compliance. For example, subjects experiencing adverse reactions to an experimental drug may decide to reduce the assigned dosage. Such imperfect compliance introduces appreciable bias into the conclusions that researchers draw from the data, and this bias cannot be corrected unless detailed models of compliance are constructed [9].
2. Denying subjects assigned to certain control groups the benefits of the best available treatment has moral and legal ramifications. For example, in AIDS

research it is difficult to justify placebo programs because those patients assigned to the placebo group would be denied access to potentially lifesaving treatment [16].

3. Randomization, by its very presence, may influence participation as well as behavior [11]. For example, eligible candidates may be wary of applying to a school once they discover that it deliberately randomizes its admission criteria. Likewise, as Kramer and Shapiro [13] note, subjects in drug trials may be less likely to participate in randomized trials than in nonexperimental studies, even when the treatments are equally nonthreatening.

Altogether, researchers are finally being forced to acknowledge that mandated randomization may undermine the reliability of experimental evidence and that experimentation with human subjects should include an element of *self-selection*.

This paper concerns the drawing of inferences from studies in which subjects are indeed given final choice of program, while randomization is confined to an indirect instrument that merely encourages or discourages participation in the various programs. For example, in evaluating the efficacy of a given training program, notices of eligibility may be sent to a randomly selected group of students or, alternatively, eligible candidates may be selected at random to receive scholarships for participating in the program. Similarly, in drug trials, subjects may be given randomly chosen advice on recommended dosage level, yet the final choice of dosage will be determined by the subjects to fit their individual needs.

The question we attempt to answer in this investigation is whether such indirect randomization can provide sufficient information to allow accurate assessment of the intrinsic merit of a program, as would be measured, for example, if the program were to be extended and mandated uniformly to the population. The analysis presented shows that, given a minimal set of assumptions, such inferences are indeed possible, albeit in the form of *bounds*, rather than precise point estimates, for the causal effect of the program or treatment. These bounds can be used by the analyst to guarantee that

the causal impact of a given program must be higher than one measurable quantity and lower than another.

Our most crucial assumption is that, for any given person, the encouraging instrument influences the treatment chosen by that person but has no effect on how that person would respond to the treatment chosen. The second assumption, one which is always made in experimental studies, is that subjects respond to treatment independently of one other. Other than these two assumptions, our model places no constraints on how tendencies to respond to treatments may interact with choices among treatments.

## 2 Problem Statement

The basic experimental setting associated with indirect experimentation is shown in Figure 1. To focus the discussion, we will consider a prototypical clinical trial with partial compliance although, in general, the model applies to any study in which a randomized instrument encourages subjects to choose one program over another.

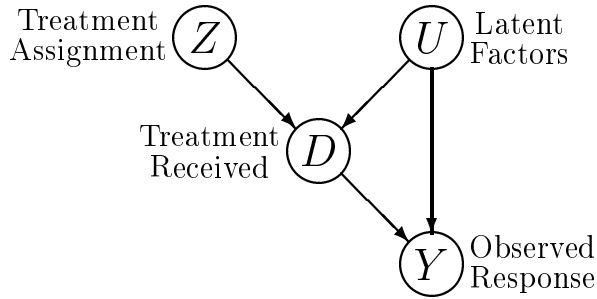


Figure 1: *Graphical representation of causal dependencies in a randomized clinical trial with partial compliance.*

We assume that  $Z$ ,  $D$ , and  $Y$  are observed binary variables where  $Z$  represents the (randomized) treatment assignment,  $D$  is the treatment actually received, and  $Y$  is the observed response.  $U$  represents all factors, both observed and unobserved, that influence the way a subject responds to treatments; hence, an arrow is drawn

from  $U$  to  $Y$ . The arrow from  $U$  to  $D$  denotes that the  $U$  factors may also influence the subject's choice of treatment  $D$ . For example, subjects who could benefit most from the treatment may be precisely those who decide, perhaps due to adverse initial reaction, to stop the treatment and switch to the control group. Thus, although  $D$  is shown to depend directly on  $Z$  and  $U$ , this dependence may represent a complex decision process standing between the assignment ( $Z$ ) and the actual treatment ( $D$ ).

To facilitate the notation, we let  $z$ ,  $d$ , and  $y$  represent, respectively, the values taken by the variables  $Z$ ,  $D$ , and  $Y$ , with the following interpretation:

$z \in \{z_0, z_1\}$ ,  $z_1$  asserts that treatment has been assigned ( $z_0$ , its negation);

$d \in \{d_0, d_1\}$ ,  $d_1$  asserts that treatment has been administered ( $d_0$ , its negation); and

$y \in \{y_0, y_1\}$ ,  $y_1$  asserts a positive observed response ( $y_0$ , its negation).

The domain of  $U$  remains unspecified and may, in general, combine the spaces of several random variables, both discrete and continuous.

The graphical model reflects two assumptions:

1. The assigned treatment  $Z$  does not influence  $Y$  directly but rather through the actual treatment  $D$ . In practice, any direct effect  $Z$  might have on  $Y$  would be adjusted for through the use of a placebo.
2.  $Z$  and  $U$  are marginally independent, as ensured through the randomization of  $Z$ , which rules out a common cause for both  $Z$  and  $U$ .

These assumptions impose on the joint distribution<sup>1</sup> the decomposition

$$P(y, d, z, u) = P(y|d, u) P(d|z, u) P(z) P(u) \tag{1}$$

---

<sup>1</sup>Only the expectation over  $U$  will enter our analysis, hence we take the liberty of denoting the distribution of  $U$  by  $P(u)$ , even though  $U$  may consist of continuous variables.

which, of course, cannot be observed directly because  $U$  is unobserved. However, the marginal distribution  $P(y, d, z)$  and, in particular, the conditional distribution  $P(y, d|z), z \in \{z_0, z_1\}$  are observed,<sup>2</sup> and the challenge is to assess from these distributions the average *change* in  $Y$  due to treatment.

If  $P(y|d, u)$  gives the probability that an individual with characteristic  $U = u$  will respond with  $Y = y$  to a treatment  $D = d$ , then taking the average over  $u$  gives the average response if the treatment  $d$  is applied uniformly to the population. Therefore, the *average causal effect* of  $D$  on  $Y$ ,  $\alpha$ , is defined as the difference

$$\alpha = E[P(y_1|d_1, u) - P(y_1|d_0, u)] \tag{2}$$

where  $E$  stands for the expectation taken over  $u$ .

When compliance is perfect,  $D$  and  $U$  are independent,  $P(y|d, u)$  can be written  $P(y, d|u)/P(d|u) = P(y, d|u)/P(d)$ , and  $E[P(y|d, u)]$  reduces to the observed conditional probability  $P(y|d)$ . Thus,  $\alpha$  would be measured by the observed mean difference between treated and untreated subjects,

$$\Delta(Y|D) = P(y_1|d_1) - P(y_1|d_0) \tag{3}$$

However, when compliance is not perfect, high values of  $\Delta(Y|D)$  may correspond to low or even negative values of  $\alpha$ . The discrepancy between  $\alpha$  and  $\Delta(Y|D)$  is known as *confounding bias*, which, clearly, cannot be eliminated by taking a larger sample size.

To circumvent this bias, researchers usually advocate use of *intent-to-treat analysis*, in which assignment groups are compared regardless of the treatment actually received. In other words, instead of estimating  $\alpha$ , we settle for

$$\Delta(Y|Z) = P(y_1|z_1) - P(y_1|z_0)$$

---

<sup>2</sup>In practice, of course, only a finite sample of  $P(y, d|z)$  will be observed, but since our task is one of identification, not estimation, we make the large-sample assumption and consider  $P(y, d|z)$  as given.

which represents the causal effect of the encouragement instrument  $Z$  on  $Y$ , since  $Z$  is randomized. Estimates based on intent-to-treat analysis are free of confounding bias as long as the experimental conditions perfectly mimic the conditions prevailing in the eventual usage of the treatment. In particular, the experiment should mimic subjects' incentives for receiving each treatment. In situations where field incentives are more compelling than experimental incentives, treatment effectiveness is determined by  $\alpha$ , and estimates based on  $\Delta(Y|Z)$  can be extremely misleading.<sup>3</sup> For example, imagine a study in which (a) the drug has an adverse effect on a large segment of the population and (b) only those members of the segment who drop from the treatment arm recover. The  $\Delta(Y|Z)$  measure will attribute these cases of recovery to the drug since they are part of the intent-to-treat arm, while in reality these cases have recovered by avoiding the treatment. The formulas reported in this paper should enable the analyst to determine the extent to which estimates based on intent-to-treat analysis deviate from the actual treatment effect. This information should be useful for assessing whether efforts to ensure population compliance have the potential to increase or decrease the overall benefit of the treatment.

Much of the statistical literature assumes that  $\alpha$  is the parameter of interest, because  $\alpha$  predicts the impact of applying the treatment uniformly (or randomly) over the population. However, if a policy maker is not interested in introducing new treatment policies but rather in deciding whether to maintain or terminate an existing program under its current incentive system, then the parameter of interest should measure the impact of the treatment *on the treated*, namely, the mean response of the treated subjects compared to the mean response of these same subjects had

---

<sup>3</sup>A similar weakness is inherent in the analysis of Angrist et al. [2]. They derive causal effect formulas for the unobservable subpopulation of “responsive” subjects, that is, subjects who would have changed treatment status if given a different assignment. This subpopulation cannot serve as a basis for policies involving the entire population because it is instrument dependent—individuals who are responsive in the study may not remain responsive in the field, where the incentives for obtaining treatment differ from those used in the study.

they not been treated [11]. The appropriate formula for this parameter is

$$\alpha^* = \sum_u [P(y_1|d_1, u) - P(y_1|d_0, u)]P(u|d_1) \quad (4)$$

which is similar to Eq. (2), save for replacing the unconditional expectation over  $u$  with the conditional expectation: given  $D = d_1$ . The formulas reported in this paper provide assessments for both  $\alpha$  and  $\alpha^*$ .

### 3 Summary of Results

The expression for  $\alpha$  (Eq. (2)) can be bounded by two simple formulas, each made up of observed parameters of  $P(y, d|z)$  ([22, 15, 17]; see Appendix I):

$$\begin{aligned} \alpha &\geq P(y_1|z_1) - P(y_1|z_0) - P(y_1, d_0|z_1) - P(y_0, d_1|z_0) \\ \alpha &\leq P(y_1|z_1) - P(y_1|z_0) + P(y_0, d_0|z_1) + P(y_1, d_1|z_0) \end{aligned} \quad (5)$$

Due to their simplicity and wide range of applicability, the bounds given by Eq. (5) were named the *natural bounds* [3]. The natural bounds guarantee that the causal effect of the actual treatment cannot be smaller than that of the encouragement by more than the sum of two measurable quantities,  $P(y_1, d_0|z_1) + P(y_0, d_1|z_0)$ ; they also guarantee that the causal effect of the treatment cannot exceed that of the encouragement by more than the sum of two other measurable quantities,  $P(y_0, d_0|z_1) + P(y_1, d_1|z_0)$ . The width of the natural bounds, not surprisingly, is given by the rate of noncompliance,  $P(d_1|z_0) + P(d_0|z_1)$ .

This width can be narrowed further using linear programming [3], which shows that, even under conditions of imperfect compliance, some experimental data (i.e.,  $P(x, y|z)$ ) can permit the precise evaluation of  $\alpha$ . These narrower bounds, which are



in fact the narrowest possible, are given by the following inequalities:

$$E[P(y_1|d, u)] \geq \max \left\{ \begin{array}{l} P(y_1, d|z_0) \\ P(y_1, d|z_1) \\ P(y_1|z_0) - P(y_1, \bar{d}|z_1) - P(y_0, d|z_1) \\ P(y_1|z_1) - P(y_1, \bar{d}|z_0) - P(y_0, d|z_0) \end{array} \right\} \triangleq L(d)$$

$$E[P(y_1|d, u)] \leq \min \left\{ \begin{array}{l} P(y_0|z_1) + P(y_1, \bar{d}|z_1) \\ P(y_0|z_0) + P(y_1, \bar{d}|z_0) \\ P(y_1|z_1) + P(y_0, \bar{d}|z_0) + P(y_1, d|z_0) \\ P(y_1|z_0) + P(y_0, \bar{d}|z_1) + P(y_1, d|z_1) \end{array} \right\} \triangleq U(d) \quad (6)$$

where  $d$  and its complement  $\bar{d}$  range over  $\{d_0, d_1\}$ . These inequalities induce the following bounds on  $\alpha$ :

$$U(d_1) - L(d_0) \geq \alpha \geq L(d_1) - U(d_0) \quad (7)$$

where  $U(d)$  and  $L(d)$  are defined in Eq. (6).

Although more complicated than the natural bounds of Eq. (5), these expressions are nevertheless easy to assess once we have the frequency data in the eight cells of  $P(y, d|z)$ . It can also be shown [3] that the natural bounds are optimal where we can safely assume that no subject is *contrarian*, that is, that no subject would consistently choose a treatment arm contrary to the one assigned.

Note that if the response  $Y$  is continuous, one can associate  $y_1$  and  $y_0$  with the binary events  $Y > t$  and  $Y \leq t$ , respectively, and let  $t$  vary continuously over the range of  $Y$ . Eqs. (6) would then provide bounds on the entire distribution of the treatment effect  $E[P(Y \leq t|d, u)]$ .

The analysis of the  $\alpha^*$  measure also produces informative bounds and, more remarkably, under conditions of *no intrusion* (namely,  $P(d_1|z_0) = 0$ , as in most clinical trials),  $\alpha^*$  can be identified precisely [5, 1]. The bounds governing  $\alpha^*$  are (see Ap-

pendix I)

$$\begin{aligned}\alpha^* &\geq \frac{P(y_1|z_1) - P(y_1|z_0)}{P(d_1)/P(z_1)} - \frac{P(y_0, d_1|z_0)}{P(d_1)} \\ \alpha^* &\leq \frac{P(y_1|z_1) - P(y_1|z_0)}{P(d_1)/P(z_1)} + \frac{P(y_1, d_1|z_0)}{P(d_1)}\end{aligned}\tag{8}$$

Clearly, in situations where treatment may only be obtained by those encouraged (by assignment),  $\alpha^*$  is identifiable and is given by

$$\alpha^* = \frac{P(y_1|z_1) - P(y_1|z_0)}{P(d_1|z_1)}\tag{9}$$

if  $P(d_1|z_0) = 0$ . Unlike the  $\alpha$  measure,  $\alpha^*$  is not an intrinsic property of the treatment, as it varies with the encouraging instrument. As noted in Section 2, the significance of the  $\alpha^*$  measure emerges primarily in studies where it is desired to evaluate the efficacy of an *existing* program on its current participants. In such studies, under the assumption that encouragement is randomized,  $\alpha^*$  is given simply by the mean response difference between the encouraged and nonencouraged populations, divided by the rate of participation,  $P(d_1|z_1)$ .

## 4 Example 1: The Effects of Cholestyramine

To demonstrate by example how the bounds for  $\alpha$  can be used to provide meaningful information about causal effects, consider the Lipid Research Clinics Coronary Primary Prevention Trial data (see [14]). A portion (337 subjects) of this data was analyzed in [9] and is the focus of this example. Subjects were randomized into two treatment groups of roughly equal size; in one group, all subjects were prescribed cholestyramine ( $z_1$ ), while the subjects in the other group were prescribed a placebo ( $z_0$ ). Over several years of treatment, each subject's cholesterol level was measured many times, and the average of these measurements was used as the post-treatment cholesterol level (continuous variable  $C_F$ ). The compliance of each subject was determined by tracking the quantity of prescribed dosage consumed (a continuous quantity).

In order to apply the bounds of Eq. (5) to data from this study, the continuous data is first transformed, using thresholds, to binary variables representing treatment assignment ( $Z$ ), received treatment ( $D$ ), and treatment response ( $Y$ ). The threshold for dosage consumption was selected as roughly the midpoint between minimum and maximum consumption, while the threshold for cholesterol level reduction was set at 28 units.

After thresholding, the data samples give rise to the following eight probabilities<sup>4</sup>:

$$\begin{aligned} P(y_0, d_0|z_0) &= 0.919 & P(y_0, d_0|z_1) &= 0.315 \\ P(y_0, d_1|z_0) &= 0.000 & P(y_0, d_1|z_1) &= 0.139 \\ P(y_1, d_0|z_0) &= 0.081 & P(y_1, d_0|z_1) &= 0.073 \\ P(y_1, d_1|z_0) &= 0.000 & P(y_1, d_1|z_1) &= 0.473 \end{aligned}$$

These data represent a compliance rate of

$$P(d_1|z_1) = 0.139 + 0.473 = 0.61$$

a mean difference (using  $P(z_1) = 0.50$ ) of

$$\Delta(Y|D) = P(y_1|d_1) - p(y_1|d_0) = \frac{0.473}{0.473 + 0.139} - \frac{0.073 + 0.081}{1 + 0.315 + 0.073} = 0.662$$

and an encouragement effect of

$$\Delta(Y|Z) = P(y_1|z_1) - P(y_1|z_0) = 0.073 + 0.473 - 0.081 = 0.465$$

According to Eq. (5),  $\alpha$  can be bounded by

$$\alpha \geq 0.465 - 0.073 - 0.000 = 0.392$$

$$\alpha \leq 0.465 + 0.315 + 0.000 = 0.780$$

---

<sup>4</sup>We make the large-sample assumption and take the sample frequencies as representing  $P(y, d|z)$ . To account for sample variability, all bounds should be supplemented with confidence intervals and significance levels, as in traditional analysis of controlled experiments. This account, however, is beyond the scope of the present paper.

These are remarkably informative bounds: although 38.8% of the subjects deviated from their treatment protocol, the experimenter can categorically state that when applied uniformly to the population, the treatment is guaranteed to improve by at least 39.2% the probability of reducing the level of cholesterol by 28 points or more. This guarantee is purely mathematical and does not rest on any assumed model of subject behavior.

The impact of treatment “on the treated” is equally revealing. Using Eq. (9),  $\alpha^*$  can be evaluated precisely (since  $P(d_1|z_0) = 0$ ):

$$\alpha^* = \frac{0.465}{0.610} = 0.762$$

In other words, those subjects who stayed in the program are much better off than they would have been if not treated: the treatment can be credited with reducing cholesterol levels by at least 28 units in 76.2% of these subjects.

## 5 Example 2: The Paradoxical Episode of Dr. Pearson

### 5.1 The Story

The school board decided to hire an expert statistician, Dr. Pearson, to examine the question of whether there is a real difference between the educational effectiveness of the town’s two high schools,  $A$  and  $B$ . After extensive discussions with all parties involved, Dr. Pearson decided to resolve the issue by conducting a clean randomized experiment. He obtained the list of all students about to enroll in next year’s freshman classes and selected at random two groups of 100 students each. He then went to the homes of all 200 students and persuaded parents to enroll their children in school  $A$  or school  $B$ , depending on the group the child happened to be in. Dr. Pearson then went back to his work at the university, after arranging to come back in the spring to analyze the test scores of the two groups.

When he returned, Dr. Pearson found that something had gone terribly wrong immediately after he had left town. Some kids found their assigned school too boring, others complained about school being too hard, and still others were plain contrarians. Altogether, only 50% of those assigned to school *A* actually remained in school *A* after the first week of classes, and exactly the same thing had happened with those assigned to school *B*. Dr. Pearson was totally depressed. While it is true that the test scores showed a marked performance difference between the two groups—50% of the students attending school *B* and none of those attending school *A* managed to pass the state exam—he knew very well that by all statistical standards the study was worthless.

As Dr. Pearson disgustedly prepared to quit town, his young assistant, Alex, asked to have one more look at the data. “What is there to look at?” snarled Dr. Pearson. “I can tell a classic case of noncompliance when I see one! Even when only 10% of the subjects switch groups, I have a hell of a time convincing my colleagues that the study is worth a nickel. Can you imagine how they would react if I were to show them data with a 50% crossover rate? They would say that exactly those students who would have performed better in school *A* decided to switch over to school *B* and that those who switched from school *B* to school *A* were precisely the students who would not have performed well no matter where they went. You know how convoluted statisticians can be”.

Surprisingly, after spending some time on the computer, Alex came back with a smile on his face: “Dr. Pearson, it wasn’t a total waste after all. No matter how you slice it, school *B* is a clear winner over school *A*. Barring sampling errors, school *B* gives students a 50% greater chance of passing the state exam”.

“What do you mean ‘winner’? You talk like one of those computer freaks. In our profession, we do not talk about winners or losers. We talk about the data and let our clients jump to their own conclusions. Let’s stick to the data, Alex”.

“I am talking about the data”, Alex replied. “What I mean is that if all the

subjects we selected had remained in their assigned schools, the percentage of success in group  $B$  (school  $B$ ) would have been exactly 50% higher than that of group  $A$  (school  $A$ ).

“You are a hopeless case, Alex. What’s all this talk about ‘had remained’ and ‘would have been’? Where the hell did you take your statistics 1 class? Next you are going to pretend that the data tells you how many would pass the exam if we were to convince all our subjects to enroll in school  $B$ , right?”

“50%”, said Alex.

“Here you are, a perfect oracle. Now how about if they enrolled in school  $A$ ?” asked Dr. Pearson, somewhat amused.

“Hmm, they would all fail”, said Alex, as he scanned the computer printout.

“This is too much. I give up”, said Dr. Pearson. “All this talk about ‘if we were’ and ‘if they would’ is getting on my nerves. I am getting out of this place”.

## 5.2 The Story Unfolds

As strange as it may sound, Dr. Pearson’s desperation was premature, and Alex’s confidence well justified. While Alex’s numerical predictions were based on observing the history of those 50 students who passed the exam (i.e., on how many of them switched schools), his confidence that such precise predictions should be feasible is grounded in the fact that (a) compliance was the same in both groups and (b) performance in one treatment arm (say school  $B$ ) was perfectly correlated with compliance. The specific observation that gave rise to Alex’s predictions was that none of the students who passed the exam had switched schools.

Under such an observation, the data read:

$$\begin{aligned} P(y_0, d_0|z_0) &= P(y_0, d_1|z_0) = P(y_0, d_0|z_1) = P(y_1, d_1|z_1) = 0.50 \\ P(y_0, d_1|z_0) &= P(y_1, d_1|z_0) = P(y_0, d_1|z_1) = P(y_1, d_0|z_1) = 0.00 \end{aligned} \tag{10}$$

where  $z_1$ ,  $d_1$ , and  $y_1$  stand for “ $B$ -assigned”, “ $B$ -attended”, and “exam passed”,

respectively, and  $z_0$ ,  $d_0$ , and  $y_0$  stand for their complements. Substituting these figures in Eqs. (6) and (7) yields equality between  $U(d)$  and  $L(d)$ , with

$$E[P(y_1|d_0, u)] = 0.00$$

$$E[P(y_1|d_1, u)] = 0.50$$

$$\alpha = 0.50$$

Thus, despite the huge rate of noncompliance (50%), the effectiveness of both programs can be determined precisely, fully supporting Alex's predictions.

The next subsection will provide an intuitive explanation for the reasoning behind these predictions.

### 5.3 The Story Explained

The data consist of four groups of students, denoted  $(A, A)$ ,  $(A, B)$ ,  $(B, A)$ , and  $(B, B)$ , where the first member of each pair denotes the school assigned and the second, the school actually attended. Assume that the data show all four groups to be of equal size and that all students in group  $(B, B)$  and no students in the other groups passed the exam.

That these data support Alex's claims can be seen from the following:

1. First we will show that those who switched from school  $A$  to school  $B$  would switch from  $B$  to  $A$  had they been assigned to  $B$ . In principle, group  $(A, B)$  may consist of two types of students: those who would stay in  $B$  if assigned to  $B$ , and those who would switch over to  $A$  if assigned to  $B$ . Call the former *contrarians* (always switching) and the latter *B-bound* (selecting  $B$  regardless of assignment). We argue that, barring sampling errors, there can be no *B-bound* students in  $(A, B)$ , only contrarians. Indeed, if there were a fraction  $p$  of *B-bound* students in the population, then, due to randomization, the same fraction  $p$  is expected to exist in the two assignment arms. However, all *B-bound*

students assigned to  $B$  must (by definition) belong to  $(B, B)$ , whose members all passed the exam. In contrast, all  $B$ -bound students in  $(A, B)$  failed the exam despite having been exposed to the same education (school  $B$ ) as the  $B$ -bound students from  $(B, B)$ . This disparity in performance can only be reconciled by assuming that  $p = 0$ , namely, that there are no  $B$ -bound students at all and that all members of  $(A, B)$  are contrarians.<sup>5</sup>

2. Next we will show that, barring sampling errors, members of group  $(A, B)$  would all fail the exam had they stayed in school  $A$ . Indeed, if all members of  $(A, B)$  are contrarians, then there are 50% contrarians among those assigned to  $A$ , and, by randomization, the same percentage of contrarians must exist among those assigned to  $B$ , which is perfectly consistent with the size of group  $(B, A)$ . Thus, all members of  $(B, A)$  must be contrarians, not  $A$ -bound, of precisely the same stock as members of  $(A, B)$ . Now, to find out how effective school  $A$  is for contrarians, we observe that all members of  $(B, A)$  failed in school  $A$ . We conclude, therefore, that members of  $(A, B)$  will fail as well, being no different from members of  $(B, A)$ . This substantiates Alex's first claim, that no change of performance would be expected if all students were persuaded to remain in their assigned schools.
  
3. It remains to be shown that students in group  $(A, A)$  are identical to those in  $(B, B)$  and, hence, would pass the exam had they been assigned to school  $B$ . Here we invoke the same arguments as given before to show that  $(A, A)$  could not consist of any  $A$ -bound students, that is,  $(A, A)$  students would stay in  $B$

---

<sup>5</sup>Note that the argument does not rest on the assumption that all  $B$ -bound students perform equally as a group. Rather, we argue that if there are factors that influence the performance of the  $B$ -bound group, the same sort of factors must act on  $B$ -bound students assigned to school  $A$  and on  $B$ -bound students assigned to school  $B$ . Hence, the average performance of the two subgroups must be the same. This follows from the assumption that the assignment in itself does not change the way a student would react to the educational program of the school attended by the student.



had they been assigned  $B$ . If there were any  $A$ -bound students in  $(A, A)$  they must also be present in  $(B, A)$ , but we have shown that the entire  $(B, A)$  group is made up of contrarians, which yields a contradiction. Thus, we conclude that the entire  $(A, A)$  group must consist of pure compliers, and they should perform just like students in  $(B, B)$  if assigned to  $B$ . This completes our explanation of Dr. Pearson's episode.<sup>6</sup>

The reader may appreciate that the reasoning behind our explanation, although logically valid, is not simple; we would not expect Dr. Pearson to do this kind of reasoning in his head on a daily basis. It is to save such labor that the human race has invented mathematical analysis. Fortunately, clinicians can now replace mental exercises with one simple formula, as given in Eq. (6).

## 6 Can the Model Be Validated?

It is well known that one cannot infer causation from nonrandomized studies unless one is willing to make some causal assumptions. This paper appears to be claiming that causal effects can be inferred (or at least bounded) from nonrandomized studies, which raises the following two questions. What are the causal assumptions that have entered into our analysis? Can these assumptions be validated experimentally?

As defined in Section 2, our model rests on two assumptions:  $Z$  is randomized, and  $Z$  has no side effect on  $Y$ . These two assumptions are equivalent to stating that  $Z$  is independent of  $U$ , a condition that economists call *exogeneity* and which qualifies  $Z$  as an *instrumental variable* [6] relative to the relation between  $D$  and  $Y$ .<sup>7</sup> For a long

---

<sup>6</sup>We leave it to the reader to show that precise determination of treatment effectiveness is feasible whenever (a) the percentage of subjects switching from  $A$  to  $B$  is the same as those switching from  $B$  to  $A$  and (b) in at least one treatment arm,  $A$  or  $B$ , performance is perfectly correlated with the assignment.

<sup>7</sup>Instrumental variables is a technique invented by the geneticist Sewal Wright [25] to help economists identify elasticities of supply and demand [10]. The key idea is that the coefficient  $b$

time, experimental verification of whether a variable  $Z$  is exogenous has been thought to be impossible, since the definition involves unobservable factors (or disturbances, as they are usually called) such as those represented by  $U$ . Imbens and Angrist [12], for example, state specifically that the model presented in Figure 1 is not testable even when  $Z$  is randomized. The notion of exogeneity, like that of causation itself, has been viewed as a product of subjective modeling judgment, not as an objective property that can be tested against the data.

The bounds derived in this paper tell a different story. Despite its elusive nature, exogeneity can be given an empirical test. The test is not guaranteed to detect all violations of exogeneity, but it can, in certain circumstances, screen out very bad would-be instruments.

The empirical test dictated by our analysis can be obtained from Eq. (6) (or Eq. (19) of Appendix I). By insisting that each upper bound be higher than the corresponding lower bound, we obtain the following testable constraints on the observed distribution:

$$\begin{aligned}
 P(y_0, d_0|z_0) + P(y_1, d_0|z_1) &\leq 1 \\
 P(y_0, d_1|z_0) + P(y_1, d_1|z_1) &\leq 1 \\
 P(y_1, d_0|z_0) + P(y_0, d_0|z_1) &\leq 1 \\
 P(y_1, d_1|z_0) + P(y_0, d_1|z_1) &\leq 1
 \end{aligned}
 \tag{11}$$

If any of these inequalities is violated, the investigator can safely deduce that at least one of the assumptions underlying our model is violated as well. If the assignment is carefully randomized, then any violation of these inequalities must be attributed to some direct influence that the assignment process has on subjects' responses (e.g., a traumatic experience). Alternatively, if direct effects of  $Z$  on  $Y$  can be eliminated, say in the (causal) equation  $Y = bX + U$  cannot be identified if  $X$  and  $U$  are correlated. However, if we can find a third variable  $Z$  that is correlated with  $X$  and (judged) uncorrelated with  $U$ , then  $b$  can be determined from the correlations between  $Z$ ,  $X$ , and  $Y$ , which yields  $b = R_{yz}/R_{xz}$ .

through an effective use of a placebo, then any observed violation of the inequalities can safely be attributed to spurious correlation between  $Z$  and  $U$ , namely, to selection bias.

The inequalities in Eq. (11), when generalized to multivalued variables, assume the form

$$\max_d \sum_y [\max_z P(y, d|z)] \leq 1 \quad (12)$$

which in [18] was called the *instrumental inequality*. A proof is given in Appendix II. We see that the instrumental inequality is violated when the controlling instrument  $Z$  manages to produce significant changes in the response variable  $Y$  while the treatment  $D$  remains constant. Although such changes could in principle be explained by strong correlations between  $U$ ,  $D$ , and  $Y$  (since  $D$  does not screen off  $Z$  from  $Y$ ), the instrumental inequality sets a limit on the magnitude of the changes.

The similarity of the instrumental inequality to Bell's inequality in quantum physics [8, 24] is not accidental; both inequalities delineate a class of observed correlations that cannot be explained by hypothesizing latent common causes. The instrumental inequality can, in a sense, be viewed as a generalization of Bell's inequality for cases where direct causal connection is permitted to operate between the correlated observables,  $D$  and  $Y$ .

The instrumental inequality can be tightened appreciably if we are willing to make additional assumptions about subjects' behavior—for example, that no individual can be discouraged by the encouragement instrument, or, mathematically, that for all  $u$  we have

$$P(d_1|z_1, u) \geq P(d_1|z_0, u)$$

Such an assumption amounts to having no contrarians in the population, namely, no subjects who will act in a spiteful manner, contrary to their assignment. Under this assumption, the inequalities in Eq. (11) can be tightened [3] to give

$$P(y, d_1|z_1) \geq P(y, d_1|z_0)$$

$$P(y, d_0|z_0) \geq P(y, d_0|z_1) \tag{13}$$

for all  $y \in \{y_0, y_1\}$ . Violation of these inequalities now means either selection bias or direct effect of  $Z$  on  $Y$  or the presence of defiant subjects.

## 7 Concluding Remarks

Intimidated by philosophical debates, undisciplined practice, and inadequate notation, many statisticians have been reluctant to tackle problems involving causal considerations. As a consequence, statistical problems connected with treatment evaluation, liability determination, and policy decisions have often been left to antiquated methods and ad hoc analysis. Intent-to-treat analysis of noncompliance in clinical trials is one example of such antiquated methods, but, unfortunately, it continues to dominate current practices in treatment evaluation.

The results reported in this paper<sup>8</sup>, a fallout from artificial intelligence research in causal reasoning, supplement existing methodologies in several ways. First, the bounds provided should allow traditional intent-to-treat analysts to evaluate how far  $\Delta(Y|Z)$  is from the actual treatment effect and whether efforts to enforce compliance are likely to decrease or increase overall treatment effectiveness. Second, thousands of extant clinical databases that were either, as in the case of Dr. Pearson, abandoned or improperly analyzed can now be filtered through the bounds of Eqs. (5) through (7), which should yield valuable clinical knowledge. Last, and not the least promising, the availability of these bounds should encourage deliberate design of indirect experiments, using self-selection, where randomized controlled experiments are infeasible

---

<sup>8</sup>Equations (5) and (9) were reported previously [1, 5, 15, 22], but have remained largely unnoticed, presumably because the derivations used counterfactual variables and esoteric notation. In comparison, the graphical formulation shown in Figure 1 appeals to a natural, process-based conceptualization of the problem and uses standard probabilistic notation in the derivation. These features have given rise to many new results (e.g., Eqs. (6), (8), (12), and (14)) and to a general bounding method [4] that promises to render causal analysis more accessible to rank-and-file researchers.

or undesirable for any of the three reasons described in the introduction.

Another set of results of possible interest to experimental researchers are those concerning the deduction of causal effects from purely observational studies, where no randomized instrument can be identified. Given an arbitrary causal graph of the type described in Figure 1, where only some of the nodes are observable, it is now possible to determine by graphical techniques whether the causal effect of one variable on another can be computed from the (nonexperimental) distribution of the observables [19, 20]. If the answer is yes, then randomized experiments are not necessary and the effect of interventions can be predicted by symbolic manipulation of graphs and probabilities. If the answer is no, the method may suggest surrogate experiments that are either less objectionable or more economical than brute-force randomization of the putative cause. Alternatively, bounds similar to the ones discussed in this paper can be derived.

The graphical method described in [19, 20] has uncovered many new structures that permit the identification of causal effects from nonexperimental observations. For example, the structure of Figure 2 represents a class of observational studies in which the causal effect of  $X$  on  $Y$  can be determined by measuring a variable  $Z$  that

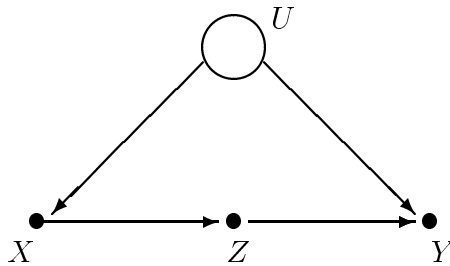


Figure 2: *Causal inference using an intermediate variable  $Z$ .*

*mediates* the interaction between the cause and its effect. This stands contrary to most of the literature on statistical experimentation, which considers the measurement of intermediate variables affected by the action to be useless for, if not harmful to, causal inference [7, 21]. The relevance of such structures in practical situations can be seen if, for instance, we identify  $X$  with smoking,  $Y$  with lung cancer,  $Z$  with the

amount of tar deposited in a subject's lungs, and  $U$  with an unobserved carcinogenic genotype that, according to the tobacco industry, also induces an inborn craving for nicotine. In this case, it is possible to determine graphically that measurement of  $Z$  renders randomization unnecessary, which provides us with the means to quantify, from nonexperimental data, the causal effect of smoking on cancer. (Assuming, of course, that the data  $P(x, y, z)$  is made available and that we believe that smoking does not have any direct causal effect on lung cancer except that mediated by tar deposits.) Moreover, the causal effect of  $X$  on  $Y$  can be written in close mathematical form as

$$E[P(y|x, u)] = \sum_z P(z|x) \sum_{x'} P(y|x', z')P(x') \quad (14)$$

The intuition behind Eq. (14) and an illustration of how it can be used to predict the effect of interventions from nonexperimental data are deferred to Appendix III, as this causal model (lacking an encouraging instrument) lies outside the main theme of this paper.

## Acknowledgments

This investigation owes much to the results of Alex Balke, as reported in [3, 4]; Alex's encounter with Dr. Pearson was not entirely fantastical. I value the encouragement of James Robins and Sander Greenland, who have recognized the usefulness of these results vis-à-vis intent-to-treat analysis. This research was partially supported by Air Force grant #F49620-94-1-0173, NSF grant #IRI-9200918, and Northrop-Rockwell Micro grant #93-124.

## References

- [1] J.D. Angrist and G.W. Imbens, Source of identifying information in evaluation models, Discussion Paper 1568, Department of Economics, Harvard University, Cambridge, MA, 1991.

- [2] J.D. Angrist, G.W. Imbens, and D.B. Rubin, Identification of causal effects using instrumental variables, Report No. 136 Department of Economics, Harvard University, Cambridge, MA, 1993. Submitted to *JASA*.
- [3] A.A. Balke and J. Pearl, Nonparametric bounds on causal effects from partial compliance data, Technical Report No. 199, Cognitive Systems Laboratory, UCLA Computer Science Department, Los Angeles, CA, September 1993. Submitted to *JASA*.
- [4] A.A. Balke and J. Pearl, Counterfactual probabilities: Computational methods, bounds, and applications, in: R. Lopez de Mantaras and D. Poole, eds., *Proceedings of the Tenth Conference on Uncertainty in Artificial Intelligence* (Morgan Kaufmann, San Mateo, CA, 1994) 46-54.
- [5] H.S. Bloom, Accounting for no-shows in experimental evaluation designs, *Evaluation Review*, Vol. 8, No. 2 (1984) 225–246.
- [6] R.J. Bowden and D.A. Turkington, *Instrumental Variables* (Cambridge University Press, Cambridge, UK, 1984).
- [7] D.R. Cox, *The Planning of Experiments* (John Wiley and Sons, New York, 1958).
- [8] J.T. Cushing and E. McMullin, eds., *Philosophical Consequences of Quantum Theory* (University of Notre Dame Press, South Bend, IA, 1989).
- [9] B. Efron and D. Feldman, Compliance as an explanatory variable in clinical trials, *Journal of the American Statistical Association*, Vol. 86, No. 413 (1991) 9–26.
- [10] A.S. Goldberger, Structural equation methods in the social sciences, *Econometrica*, Vol. 40 (1972) 979–1001.

- [11] J.J. Heckman, Randomization and social policy evaluation, in: C. Manski and I. Garfinkle, eds., *Evaluations of Welfare and Training Programs* (Harvard University Press, Cambridge, MA, 1992) 201–230.
- [12] G.W. Imbens and J.D. Angrist, Identification and estimation of local average treatment effects, *Econometrica*, Vol. 62 (1994) 467–476.
- [13] M.S. Kramer and S. Shapiro, Scientific challenges in the application of randomized trials, *Journal of the American Medical Association*, Vol. 252 (1984) 2739–2745.
- [14] The Lipid Research Clinics Coronary Primary Prevention Trial results, Parts I and II, *Journal of the American Medical Association*, Vol. 251, No. 3 (1984) 351–374.
- [15] C.F. Manski, Nonparametric bounds on treatment effects, *American Economic Review, Papers and Proceedings*, Vol. 80 (1990) 319–323.
- [16] J. Palca, AIDS drug trials enter new age, *Science Magazine*, October 6 (1989) 19–21.
- [17] J. Pearl, From Bayesian networks to causal networks, in: A. Gammerman, ed., *Bayesian Networks and Probabilistic Reasoning* (Alfred Walter Ltd., London, 1994) 1–31.
- [18] J. Pearl, A note on testing exogeneity of instrumental variables, Technical Report No. R-211-S, Cognitive Systems Laboratory, UCLA Computer Science Department, Los Angeles, CA, December 1993. Expanded version submitted to UAI-95.
- [19] J. Pearl, Causal diagrams for experimental research, Technical Report No. R-218-L, Cognitive Systems Laboratory, UCLA Computer Science Department, Los Angeles, CA, May 1993. To appear in *Biometrika*.



- [20] J. Pearl, A probabilistic calculus of actions, in: R. Lopez de Mantaras and D. Poole, eds., *Proceedings of the Tenth Conference on Uncertainty in Artificial Intelligence* (Morgan Kaufmann, San Mateo, CA, 1994) 454–462.
- [21] J.W. Pratt and R. Schlaifer, On the interpretation and observation of laws, *Journal of Econometrics*, Vol. 39 (1988) 23–52.
- [22] J.M. Robins, The analysis of randomized and non-randomized AIDS treatment trials using a new approach to causal inference in longitudinal studies, in: L. Sechrest, H. Freeman, and A. Mulley, eds., *Health Service Research Methodology: A Focus on AIDS* (NCHSR, U.S. Public Health Service, 1989) 113–159.
- [23] P. Spirtes, C. Glymour, and R. Schienes, *Causation, Prediction, and Search* (Springer-Verlag, New York, 1993).
- [24] P. Suppes, Probabilistic causality in space and time, in: B. Skyrms and W.L. Harper, eds., *Causation, Chance, and Credence* (Kluwer Academic Publishers, Dordrecht, The Netherlands, 1988) 135–151.
- [25] S. Wright, Appendix, in: P.G. Wright, *The Tariff on Animal and Vegetable Oils* (Macmillan, New York, 1928).

## Appendix I. Derivation of bounds for $\alpha$ and $\alpha^*$

To prove Eq. (5), we write

$$P(y, d|z) = \sum_u P(y|d, u) P(d|z, u) P(u) \quad (15)$$

and define the following four functions:

$$f_0(u) = P(y_1|d_0, u) \quad g_0(u) = P(d_1|u, z_0) \quad (16)$$

$$f_1(u) = P(y_1|d_1, u) \quad g_1(u) = P(d_1|u, z_1) \quad (17)$$

This permits us to express six independent components of  $P(y, d|z)$  as expectations of these functions:

$$\begin{aligned} P(y_1, d_0|z_0) &= E[f_0(1 - g_0)] = a \\ P(y_1, d_0|z_1) &= E[f_0(1 - g_1)] = b \\ P(d_1|z_0) &= E(g_0) = c \\ P(d_1|z_1) &= E(g_1) = d \\ P(y_1, d_1|z_0) &= E[f_1 \cdot g_0] = e \\ P(y_1, d_1|z_1) &= E[f_1 \cdot g_1] = h \end{aligned} \quad (18)$$

For any two random variables  $X$  and  $Y$  such that  $0 \leq X \leq 1, 0 \leq Y \leq 1$  we have

$$1 + E(XY) - E(Y) \geq E(X) \geq E(XY)$$

since  $E[(1 - X)(1 - Y)] \geq 0$ . This inequality holds for any pair of  $f, g$  functions (since they lie between 0 and 1) and we can write:

$$\begin{aligned} 1 + E(f_1 g_0) - E(g_0) &\geq E(f_1) \geq E(f_1 g_0) \\ 1 + E(f_1 g_1) - E(g_1) &\geq E(f_1) \geq E(f_1 g_1) \\ 1 + E[f_0(1 - g_0)] - E(1 - g_0) &\geq E(f_0) \geq E[f_0(1 - g_0)] \\ 1 + E[f_0(1 - g_1)] - E(1 - g_1) &\geq E(f_0) \geq E[f_0(1 - g_1)] \end{aligned}$$

or

$$\begin{aligned} \max[h; e] &\leq E(f_1) \leq \min[(1 + e - c); (1 + h - d)] \\ \max[a; b] &\leq E(f_0) \leq \min[(a + c); (b + d)] \end{aligned} \quad (19)$$

Lower bounding  $E(f_1)$  and upper bounding  $E(f_0)$  provides a lower bound for their difference:

$$\begin{aligned} E(f_1) - E(f_0) &\geq \max[e; h] - \min[(a + c); (b + d)] \\ &\geq h - (a + c) \end{aligned} \tag{20}$$

Substituting back the  $P(y, d|z)$  expressions from Eqs. (16) through (18) yields the lower bound of Eq. (5). Similarly, the difference can be upper bounded by

$$\begin{aligned} E(f_1) - E(f_0) &\leq \min[(1 + e - c); (1 + h - d)] - \max[a; b] \\ &\leq 1 + h - d - a \end{aligned}$$

thus proving Eq. (5).

To evaluate

$$\alpha^* = \sum_u [P(y_1|d_1, u) - P(y_1|d_0, u)]P(u|d_1)$$

we define

$$\begin{aligned}\Delta(u) &\equiv P(y_1|d_1, u) - P(y_1|d_0, u) = f_1(u) - f_0(u) \\ q &\equiv P(z_1)\end{aligned}$$

and write

$$\begin{aligned}\alpha^* &= \sum_u \Delta(u)P(u|d_1) \\ &= \frac{1}{P(d_1)} \sum_u \Delta(u)P(d_1|u)P(u) \\ &= \frac{1}{P(d_1)} \sum_u \sum_z \Delta(u)P(d_1|u, z)P(z)P(u) \\ &= \frac{1}{P(d_1)} \sum_u \Delta(u)P(u)[P(z_1)g_1(u) + P(z_0)g_0(u)] \\ &= \frac{1}{P(d_1)} E\{[f_1(u) - f_0(u)][qg_1(u) + (1-q)g_0(u)]\} \\ &= \frac{1}{P(d_1)} E[qf_1g_1 + (1-q)f_1g_0 - qf_0g_1 - (1-q)f_0g_0] \\ &= \frac{1}{P(d_1)} [qh + (1-q)e - qE(f_0g_1) - (1-q)E(f_0g_0)] \\ &= \frac{1}{P(d_1)} \{qh + [1-q]e - q[E(f_0) - b] - (1-q)[E(f_0) - a]\} \\ &= \frac{1}{P(d_1)} [q(h+b) + (1-q)(e+a) - E(f_0)]\end{aligned}\tag{21}$$

Substituting the expressions for  $(h+b)$  and  $(e+a)$  from Eq. (18) and using

$$a \leq E(f_0) < a + c$$

from Eq. (19), we obtain the follow upper and lower bounds on  $\alpha^*$ :

$$\begin{aligned}\frac{1}{P(d_1)} [P(y_1) - P(d_1|z_0) - P(y_1, d_0|z_0)] &\leq \alpha^* \\ \frac{1}{P(d_1)} [P(y_1) - P(y_1, d_0|z_0)] &\geq \alpha^*\end{aligned}\tag{22}$$

Alternatively, collecting common terms in both expressions of Eq. (22), we get

$$-\frac{P(y_0, d_1|z_0)}{P(d_1)} \leq \alpha^* - \frac{P(y_1|z_1) - P(y_1|z_0)}{P(d_1)/P(z_1)} \leq \frac{P(y_1, d_1|z_0)}{P(d_1)}$$

Thus,

$$\alpha^* = \frac{P(y_1|z_1) - P(y_1|z_0)}{P(d_1|z_1)}$$

if  $P(d_1|z_0) = 0$ , which proves Eqs. (8) and (9).

## Appendix II. Proof of the Instrumental Inequality

**Definition** (instrument): A variable  $Z$  is said to be an instrument relative to an ordered pair of variables  $(X, Y)$  if  $X$  and  $Y$  are generated by the following process:

$$\begin{aligned} x &= g(z, u) \\ y &= h(x, u) \end{aligned} \tag{23}$$

where  $g$  and  $h$  are arbitrary deterministic functions, and  $U$  is an arbitrary, unobserved random variable, independent of  $Z$ .

Eq. (23) represents the two-step process described in Figure 1, with  $x$  replacing  $d$ . Our problem is to determine, from the observed probability distribution  $P(x, y, z)$ , whether  $Z$  can be exogenous relative to  $(X, Y)$ , that is, whether there exist two functions  $g$  and  $h$  and a probability distribution on  $U$  and  $Z$  (with  $Z$  and  $U$  independent) such that the distribution generated by the two equations corresponds precisely to the observed distribution  $P(x, y, z)$ .

**Theorem:** A necessary condition for discrete variables  $X, Y$ , and  $Z$  to be generated by an instrumental process as defined in Eq. (23) is that the conditional distribution  $P(x, y|z)$  satisfies

$$\max_x \sum_y [\max_z P(x, y|z)] \leq 1 \tag{24}$$

**Proof:** If the probability distribution  $P(x, y, z)$  is generated by the instrumental process defined in Eq. (23), then it can be expressed in the form

$$P(x, y, z) = \sum_u P(y|x, u)P(x|z, u)P(u)P(z)$$

This can be seen by decomposing  $P(x, y, z, u, v)$  into product form along the order  $(y, x, v, u, z)$  and using the independence relations imposed by the model of Eq. (23),

as displayed in the graph of Figure 1. Therefore,

$$\begin{aligned} P(x, y|z) &= \sum_u P_1(y|x, u)P_2(x|z, u)P(u) \\ &= E_u[P(y|x, u)P(x|z, u)] \end{aligned} \tag{25}$$

If Eq. (25) holds for every triplet  $(x, y, z)$ , it certainly holds for a select set of triplets  $(x, y, z(x, y))$ , where  $z(x, y)$  is chosen so as to maximize  $P(x, y|z)$ . Thus, summing Eq. (25) over  $y$ , gives

$$\sum_y P[x, y|z(x, y)] = E_u \sum_y P(y|x, u)P[x|z(x, y), u] \tag{26}$$

For any fixed  $x$  and  $u$ , the term  $P[x|z(x, y), u]$  can be considered a function of  $y$  which is bounded from above by unity. The summation on the r.h.s. of Eq. (26) represents a convex sum of such  $P$  terms and, hence, it must also be bounded by unity which, after substituting out  $z(x, y)$ , gives

$$\sum_y \max_z P(x, y|z) \leq 1 \tag{27}$$

Moreover, since this inequality must hold for every  $x$ , we can write

$$\max_x \sum_y [\max_z P(x, y|z)] \leq 1 \tag{28}$$

which proves the theorem.  $\square$

Extending the instrumental inequality to the case where  $Z$  or  $Y$  is continuous presents no special difficulty. If  $f(y|x, z)$  is the conditional density function of  $Y$  given  $X$  and  $Z$ , then tracing the proof above gives a condition similar to Eq. (24):

$$\int_y \max_z [f(y|x, z)P(x|z)] dy \leq 1 \quad \forall x \tag{29}$$

However, the transition to continuous  $X$  signals a drastic change in behavior, and it seems that Eq. (23) induces no constraint whatsoever on the observed density [18].

## Appendix III. Smoking and the Genotype Theory: An Illustration

To illustrate how the measurement of mediating variables can enable us to predict causal effects, consider again the smoking-cancer example cited at the end of Section 7. According to many, the tobacco industry has managed to stay anti-smoking legislation by arguing that the observed correlation between smoking and lung cancer could be explained by some sort of carcinogenic genotype ( $U$ ) which involves inborn craving for nicotine.<sup>9</sup>

The amount of tar ( $Z$ ) deposited in a person's lungs is a variable that promises to meet the conditions specified by the structure of Figure 2. To justify the missing link between  $X$  and  $Y$ , we must assume that smoking cigarettes ( $X$ ) has no effect on the production of lung cancer ( $Y$ ) except that mediated through tar deposits. To justify the missing link between  $U$  and  $Z$ , we must assume that, even if a genotype is aggravating the production of lung cancer, it nevertheless has no effect on the amount of tar in the lungs except indirectly, through cigarette smoking.

To demonstrate how we can assess the degree to which cigarette smoking increases (or decreases) lung cancer risk, we will construct a hypothetical study in which the three variables,  $X$ ,  $Y$ , and  $Z$ , were measured simultaneously on a large, randomly selected sample from the population. To simplify the exposition, we will further assume that all three variables are binary, taking on true (1) or false (0) values. A hypothetical data set from a study on the relations among tar, cancer, and cigarette smoking is presented in Table 1.

---

<sup>9</sup>For an excellent historical account of this debate, see [23, pp. 291–302].

	Group Type	$P(x, z)$ Group Size (% of Population)	$P(Y = 1 x, z)$ % of Cancer Cases in Group
$X = 0, Z = 0$	Non-smokers, No tar	47.5	10
$X = 1, Z = 0$	Smokers, No tar	2.5	90
$X = 0, Z = 1$	Non-smokers, Tar	2.5	5
$X = 1, Z = 1$	Smokers, Tar	47.5	85

Table 1

The table shows that 95% of smokers and 5% of non-smokers have developed high levels of tar in their lungs. Moreover, 81.51% of subjects with tar deposits have developed lung cancer, compared to only 14% among those with no tar deposits. Finally, within each of the two groups, tar and no tar, smokers show a much higher percentage of cancer than non-smokers do.

These results seem to prove that smoking is a major contributor to lung cancer. However, the tobacco industry might argue that the table tells a different story—that smoking actually decreases, not increases, one’s risk of lung cancer. Their argument goes as follows. If you decide to smoke, then your chances of building up tar deposits are 95%, compared to 5% if you decide not to smoke. To evaluate the effect of tar deposits, we look separately at two groups, smokers and non-smokers. The table shows that tar deposits have a protective effect in both groups: in smokers, tar deposits lower cancer rates from 90% to 85%; in non-smokers, they lower cancer rates from 10% to 5%. Thus, regardless of whether I have a natural craving for nicotine, I should be seeking the protective effect of tar deposits in my lungs, and smoking offers a very effective means of acquiring them.

To settle the dispute between the two interpretations, we note that, while both arguments are based on stratification, the anti-smoking argument invokes an illegal stratification over a variable ( $Z$ ) that is affected by the treatment ( $X$ ). The tobacco industry’s argument, on the the hand, is made up of two steps, neither of which



involves stratification over treatment-affected variables: stratify over smoking to find the effect of tar deposit on lung cancer, then average (not stratify) over tar deposits when we consider each of the decision alternatives, smoking vs. non-smoking. This is indeed the intuition behind the formula in Eq. (14) and, given the causal assumptions of Figure 2, the tobacco industry’s argument is the correct one (see [19, 20] for formal derivation).

To illustrate the use of Eq. (14), let us use the data in Table 1 to calculate the probability that a randomly selected person will develop cancer ( $y_1 : Y = 1$ ) under each of the following two actions: smoking ( $x_1 : X = 1$ ) or not smoking ( $x_0 : X = 0$ ).

Substituting the appropriate values of  $P(y|x)$ ,  $P(y|x, z)$ , and  $P(x)$  gives

$$\begin{aligned}
 E[P(y_1|x_1, u)] &= .05(.10 \times .50 + .90 \times .50) + .95(.05 \times .50 + .85 \times .50) \\
 &= .05 \times .50 + .95 \times .45 = .4525 \\
 E[P(y_1|x_0, u)] &= .95(.10 \times .50 + .90 \times .50) + .05(.05 \times .50 + .85 \times .50) \\
 &= .95 \times .50 + .05 \times .45 = .4975
 \end{aligned} \tag{30}$$

Thus, contrary to expectation, the data prove smoking to be somewhat beneficial to one’s health.

The data in Table 1 are obviously unrealistic and were deliberately crafted so as to support the genotype theory. However, this exercise was meant to demonstrate how reasonable qualitative assumptions about the workings of mechanisms can produce precise quantitative assessments of causal effects when coupled with nonexperimental data. In reality, we would expect observational studies involving mediating variables to refute the genotype theory by showing, for example, that the mediating consequences of smoking, such as tar deposits, tend to increase, not decrease, the risk of cancer in smokers and non-smokers alike. The estimand given in Eq. (14) could then be used for quantifying the causal effect of smoking on cancer.