

# An Economic Analysis of Exclusion Restrictions for Instrumental Variable Estimation

Gerard J. van den Berg \*

April 16, 2006

## Abstract

Instrumental variable estimation requires untestable exclusion restrictions. With policy effects on individual outcomes, there is typically a time interval between the moment the agent realizes that he may be exposed to the policy and the actual exposure. In such cases there is an incentive for the agent to acquire information on the value of the IV. This leads to violation of the exclusion restriction. We analyze this in a dynamic economic model framework. This provides a foundation of exclusion restrictions in terms of economic behavior. The results are used to describe policy evaluation settings in which instrumental variables are likely or unlikely to make sense. For the latter cases we analyze the asymptotic bias. The exclusion restriction is more likely to be violated if the outcome of interest strongly depends on interactions between the agent's effort before the outcome is realized and the actual treatment status. The bias has the same sign as this interaction effect. Violation does not causally depend on the weakness of the candidate instrument or the size of the average treatment effect. With social experiments, violation is more likely if the treatment and control groups are to be of similar size. We also address side-effects of the treatment.

---

\*Princeton University, Free University Amsterdam, IFAU-Uppsala, IZA, IFS, CREST, and CEPR.

Keywords: treatment, policy evaluation, information, selection effects, randomization.

Thanks to Jim Heckman, Geert Ridder, Ulrich Müller, Angus Deaton, and participants in the 2005 Econometric Society World Congress (London), the 2005 IZA Prize Scientific Workshop, the 2005 Conference on the Econometric Evaluation of Public Policies (Paris), and seminars at Princeton and Columbia/NYU for useful suggestions. The work in this paper is part of the ESRC and IFAU Research Programs "Advancing Programme Evaluation Methods".

# 1 Introduction

Instrumental variable estimation has since long been a standard econometric technique for dealing with endogeneity and selection issues in general, and for non-experimental policy evaluation in particular (see e.g. Angrist, Imbens and Rubin, 1996, Heckman, LaLonde and Smith, 1999, and Blundell and MaCurdy, 1999, for surveys). Basically, if one is interested in the effect of a “treatment variable” on an outcome variable, and the treatment is not exogenously assigned, then one may perform causal inference by exploiting the presence of variables that causally affect the treatment status but do not have a direct causal effect on the outcome. The latter restriction is called an exclusion restriction. Exclusion restrictions are identifying restrictions, so they can not be tested. This means that empirical results critically depend on the validity of the exclusion restriction, and that this restriction needs to be justified on a priori grounds.

With policy effects on individual outcomes, there is typically a time interval between the moment the agent realizes that he may be exposed to the policy and the actual exposure. For example, unemployed workers are aware of the existence of policies leading to treatments at some point of time in the future. As long as the instrumental variable affecting the treatment does not have a causal effect on the individual’s behavior, the exclusion restriction is not violated. Often, a sufficient condition for this is that the agent does not observe the value of the instrumental variable. However, there is an incentive for the agent to acquire information on this value. After all, the probability of exposure to treatment is a determinant of the optimal strategy, and the more the agent knows about it, the better he can fine-tune his behavior in response to this, and the higher his expected present value will be. The agent’s strategy affects the outcome of interest. Thus, the acquisition of the value of the variable that is used by the econometrician as instrumental variable leads to violation of the exclusion restriction and to incorrect empirical inference.<sup>1</sup>

As an example, consider participation in a job search assistance program for unemployed individuals, where the policy intensity differs across two otherwise identical geographical regions. For example, in one region, the budget for the

---

<sup>1</sup>Earlier studies mentioning similar arguments, tacitly assuming that acquisition is free, include Abbring and Van den Berg (2003). For a very recent exposition, see Heckman and Navarro (2005). Modelling that agents use available information on determinants of future (policy) events that affect the outcome of interest goes back to at least the rational expectations literature; see e.g. Hansen and Sargent (1980). Note that we are not concerned with mechanical program lock-in effects that may affect the outcomes of participants before the end of the actual treatment participation.

program per potential participant may be larger, so that the individual probability of being treated is larger, holding everything else constant. An individual may be aware of the distribution of policy intensities but not know his personally relevant intensity, in which cases a regional dummy indicator may be a valid instrumental variable. If the individual finds out his relevant intensity then he will typically use this information before the treatment is realized. For example, it may be optimal to reduce the job search effort more if the treatment probability is large, because it is cheaper to provide effort after the treatment. This may lead to an under-estimation of the program effect on the employment rate say one year after inflow into unemployment.

In this paper we investigate, in the context of a dynamic economic framework, *under which conditions* it is optimal for the agent to acquire the value of the intended instrumental variable. This provides a foundation of exclusion restrictions in terms of economic behavior that takes costs and benefits into account. The results are used to describe the policy evaluation settings in which instrumental variables are likely or unlikely to make sense. This is especially useful since by definition no empirical evidence is available on the validity of exclusion restrictions. We also analyze the asymptotic bias of the instrumental variable estimator in case of violation of the exclusion restriction.

At first sight one may think that information acquisition does not take place if and only if the acquisition costs are high, and that therefore the conclusion is simply going to be that instrumental variables estimation is particularly useful to study policy effects for agents with scarce resources. For active labor market policy analysis this would mean that it is particularly useful for agents at the bottom of the labor market, which coincides with the target group of most of these policy measures. However, this line of reasoning ignores the role of the value of the information that is acquired. We show that this leads to a different set of conclusions. The results point at the importance of the extent to which the treatment status and the agent's effort interact in the outcome.

The literature on instrumental variable estimation has recently been concerned with the use of so-called weak instruments, i.e. instrumental variables that are only weakly related to the treatment status (see e.g. Stock, Wright and Yogo, 2002, for a survey). It is sometimes argued that weak instruments have the advantage that they are less likely to be used by agents as direct causal inputs into the outcome of interest. We argue that in certain cases this line of reasoning is incorrect, and therefore this advantage of weak instruments may have been over-estimated.

Our results can be straightforwardly applied to situations in which the treat-

ment variable is not a policy variable. Also, the costs of information may cover not only monetary costs but also other types of effort.<sup>2</sup>

We extend our framework and results in two directions. First, we allow agents to selectively affect their treatment status. In particular, we consider when they selectively choose to become non-compliers, if they can acquire information on determinants of the treatment assignment process.

To explain the second extension, notice that the outcome is a function of the agent's efforts and the treatment status. In the baseline analyses, we assume that the expected utility that an agent ultimately derives from a certain combination of effort and treatment is the same as the expected value of the corresponding outcome. For example, an unemployed individual may only be concerned about the probability of finding work, which is also the outcome in which the researcher is interested. This assumption is in line with dynamic economic theories such as job search theory. However, an agent may also be concerned about side-effects of the treatment that are not reflected in the measured outcome. We therefore extend the model to allow the agent's ultimate utility to be systematically different from the corresponding outcome. This extension also covers cases in which the agent does not know the treatment effect and can only make a best guess on it, as is the case in medical trials and pilot experimental evaluations of novel labor market programs. In all of these cases, the decision whether to acquire the value of the candidate instrument is driven by the agent's utility function, whereas the magnitude of the asymptotic bias also depends on the actual effects of effort and treatment on the outcome.

Throughout the paper we use the above example about participation in a job search assistance program as the leading example.<sup>3</sup> We also examine randomized experiments, in particular double-blind experiments of say a medication to treat a disease. The randomized intention to treat then equals the treatment status, and this is unobserved until the outcome is observed. However, in the case of life-

---

<sup>2</sup>Of course, exclusion restrictions for instrumental variable estimation may be violated for other reasons than those considered in this paper. Notably, the agent's value of the candidate instrument may be affected by unobserved characteristics of the agent that have a direct causal effect on the outcome variable. (This is prevented if the candidate instrument is the result of a deliberate randomization, like a deliberately randomized intention-to-treat variable in a social experiment.) Rosenzweig and Wolpin (2000) consider violations in a dynamic setup that are due to events occurring between the treatment and the outcome, whereas we are concerned with behavior before realization of the treatment.

<sup>3</sup>This is essentially equivalent to the example of using geographical distance to college as an instrument to study the returns to education. In this case, parents of children who live far from a college may provide substitute educational support.

threatening diseases, an individual has an incentive to find out whether he receives medication or a placebo, for example by sending one tablet to a laboratory. If he discovers that he receives a placebo then he may choose a different lifestyle, which in turn affects the outcome. Alternatively, he may drop out and apply for participation in another medical experiment. Yet another option is to share the tablets among participating individuals. Epstein (1996) and Schuklenk (2003) provide examples concerning experiments of AIDS medication. We show that this example can be translated into our model framework.

The paper is organized as follows: Section 2 presents the model framework, Section 3 derives the results concerning information acquisition, Section 4 discusses the implications for instrumental variable estimation, notably the asymptotic bias of the estimator in case of violation of the exclusion restriction. We also discuss testing the null hypothesis of no causal treatment effect. Section 5 concludes.

## 2 The model framework

The main insights and results can be derived in a simple model framework with three time periods.

Consider an assignment process leading to the actual treatment status of an agent or decision maker.<sup>4</sup> We assume that this process depends on a variable  $Z$ . At the onset of the first period,  $Z$  is realized separately for each agent. We take the distribution of  $Z$  as exogenously given and assume that  $Z$  is dispersed, so  $\text{var}(Z) > 0$ .  $Z$  is the candidate instrumental variable. Each agent knows the distribution of  $Z$  across all agents. However, an agent does not necessarily know his own realization of  $Z$ . More precisely, in Period 1, each agent decides whether to acquire information on  $Z$  at cost  $\gamma > 0$  or not.

At the onset of Period 2, the agent determines his optimal strategy or effort  $s$ . The agent may or may not know his private value of  $Z$  when determining  $s$ , but we assume that in Period 2 the agent does not know yet his treatment status  $Y$ . Without this assumption the analysis would be irrelevant, because there would never be any incentive to acquire information on the policy intensity. As we shall see, the analysis can allow for additional time periods and for agents to modify their strategy upon learning their value of  $Y$ , as long as their behavior before learning  $Y$  has an effect on the outcome. The effort  $s$  involves costs  $c(s)$  to be paid in Period 2.

---

<sup>4</sup>The model can be expressed in terms of counterfactual outcomes; see Angrist, Imbens and Rubin (1996) for an exposition.

In Period 3, the agent’s actual treatment status  $Y$  is realized. Both  $s$  and  $Y$  affect the outcome  $U$ , which is also realized in Period 3, simultaneously with or after the realization of  $Y$ . We express the outcome  $U$  given  $Y = y$  and given  $s$  as  $W \cdot f(y, s) + \varepsilon$ , with  $0 \leq f(y, s)$  and  $0 < W < \infty$ . Here  $W$  is just a multiplicative constant in the outcomes, and we merely introduce it to facilitate the analysis of effects of multiplicative changes in the outcomes. In the first leading example,  $f$  may be the probability of making an income gain, and  $W$  may be the expected income gain. The term  $\varepsilon$  is an idiosyncratic outcome component. We assume that  $\mathbb{E}(\varepsilon|Z) = 0$ , but all other determinants of  $Y$  may be correlated to  $\varepsilon$ . The latter dependence captures the endogeneity of the actual treatment.

Summarizing, the sequence of events is as follows: the treatment assignment intensity  $Z$  is realized,  $Z$  is acquired or not, the effort  $s$  is chosen, the treatment  $Y$  is realized, and the outcome  $Wf(y, s) + \varepsilon$  is realized.

The above framework gives rise to a “reduced form” treatment evaluation model. First, as will become clear below, we may assume without loss of generality that  $\mathbb{E}(Y|Z = z) = z$ . This leads to a “treatment equation”,

$$Y = Z + \omega \tag{1}$$

where  $\mathbb{E}(\omega|Z) = 0$ . Secondly, we have an “outcome equation”,

$$U = Wf(Y, s) + \varepsilon \tag{2}$$

with  $\mathbb{E}(\varepsilon|Z) = 0$ . The analyst observes  $U, Y$ , and  $Z$ .

Suppose that, more general than equation (1), we would specify  $\mathbb{E}(Y|Z = z) = g(z)$  for some function  $g$ . Typically, it is not difficult to estimate  $g$ , and this is why we simply redefine  $g(z)$  as our  $z$ . This presupposes that  $g$  varies with  $z$ . In other words,  $Z$  as a candidate instrumental variable must be informative. Also, note that the specification  $Y = Z + \omega$  with  $\mathbb{E}(\omega|Z) = 0$  can also capture discrete  $Y$ . Notably, if  $Y$  is binary, one may define  $Y = 1 \iff Y^* > 0$  with  $Y^* = Z + \omega^*$  and  $\omega^*$  being uniformly distributed on the interval  $[-1, 0]$ . To facilitate the exposition we assume that  $\Pr(Z \geq 0, Y \geq 0) = 1$ .

Suppose that in Period 1 the agent does not acquire his personal realization of  $Z$ . Then  $Z$  only affects  $U$  by way of  $Y$ , so  $Z$  is a valid instrumental variable (IV) because the corresponding exclusion restriction (ER) is satisfied. Now suppose that in Period 1 the agent does acquire his personal realization of  $Z$ . Then  $Z$  may affect his value of  $s$ . In that case, from equation (2), there is a causal effect of  $Z$  on the outcome, resulting in a violation of the ER needed for instrumental variable estimation (IVE). Before we analyze this, we first derive in the next

subsection the agent's optimal behavior concerning  $s$  and concerning acquisition of  $Z$ .

### 3 Economic behavior

#### 3.1 Present values

An agent maximizes his expected present value. To focus on the main issue we consider risk neutral agents. If the agent does not know his value of  $Z$  then the expected present value  $R_0$  at the onset of Period 2 equals

$$R_0 = \max_{s \in \mathcal{S}} -c(s) + \frac{1}{1+r} \mathbb{E}_Z \mathbb{E}_{Y|Z} (W \cdot f(Y, s)) \quad (3)$$

$\mathcal{S}$  denotes the choice set of  $s$ , with  $\mathcal{S} \subset [0, \infty)$ . We denote the optimal  $s$  by  $s_0$ .

Now suppose that the agent knows that his value of  $Z$  is  $z$ . The expected present value  $R(z)$  at the onset of Period 2 is

$$R(z) = \max_{s \in \mathcal{S}} -c(s) + \frac{1}{1+r} \mathbb{E}_{Y|Z=z} (W \cdot f(Y, s)) \quad (i = 1, 2) \quad (4)$$

where  $r$  is the discount factor. The optimal  $s$  can be expressed as  $s(z)$ .

The value of information  $V$  in Period 1 equals

$$V = \frac{1}{1+r} (\mathbb{E}_Z R(Z) - R_0) \quad (5)$$

It is optimal to acquire  $Z$  in Period 1 if and only if  $V > \gamma$ . A central issue of the paper is under which conditions this occurs. For the moment, we simplify the above expressions by subsuming the parameter  $W/(1+r)$  into the function  $f$ , and  $1+r$  in (5) into  $V$ .

The first insight is that if  $f$  is additive in  $Y$  and  $s$  then the optimal  $s(z)$  in (4) does not depend on  $z$ , and it is equal to  $s_0$ . Consequently,  $V = 0$ , and the agent does not acquire  $Z$ . In sum,

**Proposition 1.** *If the outcome is additive in the treatment status and the effort of the agent then the exclusion restriction is satisfied.*

Throughout the paper we consider functions  $f$  that are positive and increasing. As will become clear below, the main practical distinction in the derivations will be whether  $\mathcal{S}$  is discrete or not. In the continuous case we often assume that  $c$  is quadratic, with  $c(s) = \frac{1}{2}c_0s^2$  and  $c_0 > 0$ . Also, many results will be derived for the following functional form for  $f$ ,

$$f(y, s) = \psi_0 + \psi_1 s + \psi_2 y + \rho y s \tag{6}$$

with suitable restrictions on  $\psi_1, \psi_2, \psi_3, \rho$  guaranteeing that  $f$  is positive and increasing in the relevant intervals for  $y$  and  $s$ . We do however also provide results for non-parametric specifications of  $f$ . The functional form in (6) is concise and allows for explicit expressions and results for the quantities of interest. The interaction parameter  $\rho$  captures the degree of complementarity ( $\rho > 0$ ) or substitutability ( $\rho < 0$ ) of treatment and effort, in the outcome. The functional form is less restrictive as may seem. First of all, with binary  $s$  and  $y$ , (6) is non-parametric because the four possible values of  $f$  (i.e.,  $f(1, 1), f(1, 0), f(0, 1)$ , and  $f(0, 0)$ ) are represented by  $\psi_0, \psi_1, \psi_2$ , and  $\rho$ . Secondly, as will become clear below, we may generalize the term  $\psi_0 + \psi_2 y$  at no cost to a general function  $k_2(y)$  (provided that the resulting  $f$  is positive and increasing). So all results based on (6) generalize in this respect. The same applies if we replace  $\psi_1 s$  by a function  $k_1(s)$ , provided that  $k_1(s) - c(s)$  is quadratic (and again  $f$  is positive and increasing), where quadraticness is merely needed to ensure explicit expressions for the optimal  $s$ . More in general, the right-hand side of (6) can be interpreted as the first part of an expansion of the underlying  $f$ . It is also useful to point out that in the related literature on decision making with a noisy signal about the unknown state of the world, the general effect of the shape of the outcome function ( $f$ ) on the value of information ( $V$ ) is typically too hard to analyze in terms of the model primitives, if no parametric assumptions are made on  $f$  (see Persico, 2000, and Athey and Levin, 2001). We return to this literature below.

### 3.2 Optimal behavior with continuous effort

Suppose that the choice set  $\mathcal{S}$  of effort  $s$  is an interval. We do not restrict  $Z$  or  $Y$  to be discrete or continuous, so the results below are valid in both cases.

We start with the model in which the functions  $f$  and  $c$  satisfy (6) and the quadratic specification  $c(s) = \frac{1}{2}c_0 s^2$ , respectively. Within this framework we first consider the case  $\rho > 0$ . This means that treatment and effort are complements in the outcome. The requirements that  $f$  is positive and increasing are then fulfilled by way of the restrictions that  $\psi_0 > 0, \psi_1 \geq 0, \psi_2 \geq 0$ . Also, we assume that  $c_0 > 0$ . The optimal  $s$  is always positive, and we assume that the lower and upper boundary of  $\mathcal{S}$  are not binding for the optimal  $s$ .

Let  $\bar{z} := \mathbb{E}(Z)$  denote the population mean of  $Z$ . It is easy to derive that the optimal effort equals

$$\begin{aligned}
s(z) &= \frac{\psi_1 + \rho z}{c_0} \\
s_0 &= \frac{\psi_1 + \rho \bar{z}}{c_0}
\end{aligned}
\tag{7}$$

Note that  $s(z)$  increases in  $z$ . This was to be expected. The complementarity of  $f$  implies that the marginal return of effort is higher if the expected (beneficial) treatment level is higher. Also, the optimal effort is higher if the cost of it is lower, if the marginal return of it is higher, and if the degree of complementarity is higher. The optimal  $s$  does not depend on the marginal effect  $\psi_2$  of the treatment  $Y$ .

By substituting (7) into equations (4) and (3), we obtain that

$$V = \frac{1}{2c_0} \rho^2 \text{var}(Z) \tag{8}$$

Consequently, the ER is violated iff  $\frac{1}{2c_0} \rho^2 \text{var}(Z) > \gamma$ .

Before we interpret this result we first analyze the case  $\rho < 0$ . This means that treatment and effort are substitutes in the outcome. In the case of participation in active labor market programs, this case may be more realistic than the case  $\rho > 0$ . For example, there may be an upper bound on the outcome, and the effort has to compete with efforts for other activities outside of the model.<sup>5</sup> We also assume again that  $c_0 > 0$ . The requirement that  $f$  is positive and increasing now leads to more complex restrictions on the parameters of  $f$ . Notably, the parameters need to ensure that  $Y$  and  $s$  are bounded from above. For a start, we take  $\psi_0 > 0, \psi_1 > 0, \psi_2 > 0$ . Next, we ensure that  $f$  increases in  $s$  for every  $y$  in the support of  $Y$ . Sufficient for this is that  $\Pr(Y < \frac{\psi_1}{-\rho} | Z = z) = 1$  for every  $z$  in the support of  $Z$ , because this implies that  $\Pr(Y < \frac{\psi_1}{-\rho}) = 1$ , which implies the desired property of  $f$ . Note that  $\Pr(Y < \frac{\psi_1}{-\rho} | Z = z) = 1$  together with  $\mathbb{E}(Y | Z = z) = z$  also implies that  $\Pr(Z < \frac{\psi_1}{-\rho}) = 1$ , which in turn implies that the optimal  $s$  is positive. Finally, we ensure that  $f$  increases in  $Y$  for every  $s$ . This is satisfied if  $s < \frac{\psi_2}{-\rho}$ . We want to have a sufficient condition for this in terms of the model parameters. The optimal  $s$  decreases in  $z$ , so the largest possible value of  $s$  as a function of  $z$  is achieved at  $z = 0$ . This value equals  $\psi_1/c_0$ . Thus,  $f$  has the desired property if  $\psi_1/c_0 < \psi_2/(-\rho)$  or, equivalently,  $\psi_2 c_0 + \rho \psi_1 > 0$ . In sum, we require

---

<sup>5</sup>This does not rule out that the treatment and the efforts *after* the realization of the treatment are complements in their effect on outcomes in subsequent time periods. These are not in the present model but one can extend it to include them.

$$\psi_0 > 0, \quad \psi_1 > 0, \quad \psi_2 > 0, \quad \psi_2 c_0 + \rho \psi_1 > 0, \quad \Pr(Y < \frac{\psi_1}{-\rho} | Z = z) = 1, \quad (9)$$

where the last requirement applies for all  $z$  in the support of  $Z$ . The requirements also imply that the optimal  $s$  satisfies  $0 < s \leq \psi_1/c_0$ .

We again assume that the lower and upper boundary of  $\mathcal{S}$  are not binding for the optimal  $s$ . It is not difficult to see that the expressions for the latter are the same as (7), with  $s(z)$  now decreasing in  $z$ . The resulting expression for  $V$  is also the same as in (8). We thus obtain,

**Proposition 2.** *Consider the model with continuous effort, quadratic costs of effort, and the outcome function (6) with the conditions that ensure that it increases in effort and the treatment status. Then the exclusion restriction is violated iff  $\rho^2 \text{var}(Z)/(2c_0) > \gamma$ .*

This means, first of all, that the ER is likely to be violated if the treatment status and the effort before the treatment strongly interact in their effect on the outcome. This is because in that case the optimal effort  $s(z)$  is very responsive to the agent's value  $z$ , and the loss of choosing the wrong amount of effort is larger. For example, if  $\rho > 0$  then knowing that  $z$  is large leads to an optimal effort  $s(z)$  that is also very large, while knowing that  $z$  is small leads to a small  $s(z)$ , and in both of these cases the alternative choice of an intermediate effort level  $s_0$  entails a substantial loss.

Violation is also more likely if  $Z$  has a large variance. In that case, the candidate instrument generates a large range of mean policy intensities  $Y|Z$ , and it is more likely that not acquiring  $Z$  leads to a large loss. If the effort cost parameter  $c_0$  is large then violation is less likely, because then the optimal effort is small whether one acquires  $z$  or not. Violation of the ER is more likely if  $\gamma$  is small, which is trivial to understand.

It is also useful to discuss which model parameters do *not* affect the likelihood that the ER is violated. First, consider the parameter  $\psi_2$ . Until now we have not defined summary treatment effects yet. However, it is clear that any such measure depends on  $\psi_2$ . But this parameter does not affect the value of information. Therefore, the size of the (average) treatment effect does not affect the validity of the ER.

Secondly, consider the strength of the candidate instrument. This is usually defined as the strength of the association between  $Z$  and  $Y$ , for example as measured by the correlation coefficient  $R^2(Y, Z)$ , which in our model reduces to  $\text{var}(Z)/\text{var}(Y)$ . This quantity does not have a direct effect on the validity of the

ER. If the residual variance  $\text{var}(\omega)$  in the “treatment equation” (1) is of a higher order of magnitude than  $\text{var}(Z)$  then the candidate instrument is weak but it may nevertheless have a large variance by itself, and the latter makes it likely that the ER is violated. The underlying reason is that agents are only concerned with the *mean* of the treatment status  $Y$  given  $z$ , when they decide on their optimal effort  $s$ .<sup>6</sup>

Thirdly, the validity of the ER does not depend on  $\psi_0$  and  $\psi_1$ . The fact that it does not depend on the  $\psi_i$  parameters reflects the fact that additive effects of treatment status and effort do not affect the ER (recall Proposition 1).

We now return to the related literature on decision making with a noisy signal about the state of the world. In this literature, agents receive a signal (say,  $Z$ ) about the unknown state of the world (say,  $Y$ ) and they have to decide on which action (say  $s$ ) to take. The outcome (say,  $f$ ) depends on  $Y$  and  $s$  (see Gollier, 2001, for a recent overview of models with signals, effort, and outcomes). The main difference with our setup is that in this literature the focus is on the strength of the causal effect from the state of the world  $Y$  on the signal  $Z$  (or, equivalently, on the quality of the signal), whereas in our setup  $Z$  causally affects  $Y$ . Also, this literature restricts attention to outcome functions  $f$  that satisfy generalized notions of complementarity in the state of the world and the action of the individual (like supermodularity).

Nevertheless, some of the results from this literature are directly applicable to our context. Athey and Levin (2001) present a generalized version of the following. Consider our model with continuous effort, where the treatment status  $Y$  increases in  $Z$  in the sense of first-order stochastic dominance. The outcome function increases in effort and the treatment status, and the cross-derivative is positive. Then there is a monotone positive relation between  $z$  and the optimal effort  $s(z)$ .

If no parametric assumptions are made on  $f$  then the general effect of the

---

<sup>6</sup>More general functional forms of  $f(y, s)$  may have somewhat different implications. If  $f(y, s) = \rho k(y)s$  for some increasing function  $k$  then it can be shown that

$$V = \frac{1}{2c_0} \rho^2 \text{var}_Z(\mathbb{E}(k(Y)|Z))$$

If, for example,  $k(y) = y^3$  and  $\omega \equiv Y - Z$  is symmetric, then  $\text{var}_Z(\mathbb{E}(k(Y)|Z))$  can be shown to equal  $\text{var}(Z^3) + 9(\text{var}(\omega))^2 \text{var}(Z)$ , which in turn equals

$$\text{var}(Z^3) + 9(\text{var}(Z))^3 \left[ \frac{1}{R^2(Y, Z)} - 1 \right]^2$$

Therefore, the value of information  $V$  increases in  $\text{var}(\omega)$ , and increases in the weakness of the candidate IV as conventionally defined.

shape of the outcome function ( $f$ ) on the value of information ( $V$ ) is typically too hard to analyze in terms of the model primitives (see Persico, 2000, and Athey and Levin, 2001). Results for *given* (i.e., not optimally determined) functions  $s(z)$  emphasize the importance of the degree of complementarity of the outcome function on the value of information.

### 3.3 Optimal behavior with discrete effort

Now let  $\mathcal{S}$  be discrete. For expositional reasons, we simplify the analysis by assuming that  $s$  is binary, i.e. is taken from the set  $\{0, 1\}$ . We adopt specification (6) for  $f$ . The cost of effort function  $c(s)$  is now represented by its two possible values  $c(0)$  and  $c(1)$ . By analogy to the previous subsection, we denote  $c(1) - c(0)$  by  $c_0$ , which may be called *the* cost of effort. The expressions for the present values can now be simplified to

$$\begin{aligned} R(z) &= \psi_0 - c(0) + \psi_2 z + \max\{0, \psi_1 + \rho z - c_0\} \\ R_0 &= R(\bar{z}) \end{aligned} \tag{10}$$

and the agent chooses  $s = 1$  iff the second term in the maximum exceeds 0, so we replace equation (7) by

$$\begin{aligned} s(z) &= \mathbb{I}\left(\frac{\psi_1 + \rho z}{c_0} > 1\right) \\ s_0 &= \mathbb{I}\left(\frac{\psi_1 + \rho \bar{z}}{c_0} > 1\right) \end{aligned} \tag{11}$$

where  $\mathbb{I}(\cdot)$  denotes the indicator function. The intuition behind these expressions is exactly as for equation (7). The difference is that  $s(z)$  is now discrete instead of continuous, and the optimal choice of  $s$  is now insensitive to small changes in the parameter values.

To facilitate the exposition, we make the additional assumption that  $Z$  has two possible values in the population of agents:  $\Pr(Z = z_1) = p = 1 - \Pr(Z = z_2)$ , with  $z_1 \neq z_2$  and  $0 < p < 1$  and normalization  $z_1 > z_2$ . It is again useful to distinguish between the two cases in which treatment and effort are substitutes ( $\rho < 0$ ) or complements ( $\rho > 0$ ) for the outcome. This time we start with the former case, which requires again some restrictions on the range of values of the model parameters. By analogy to the previous subsection, we impose

$$\psi_0 > 0, \quad \psi_1 > 0, \quad \psi_2 > 0, \quad \psi_2 + \rho > 0, \quad \Pr(Y < \frac{\psi_1}{-\rho} | Z = z_i) = 1$$

where the last requirement applies for both  $z_i$ . In fact, the results below also apply if  $\psi_2 + \rho = 0$ , meaning that with effort  $s = 1$  it is irrelevant whether the treatment is realized or not. This describes situations in which the outcome of interest is the transition from unemployment to work if this is achieved with certainty by the effort  $s = 1$  but the treatment by itself cannot achieve this.

With  $\rho < 0$ , equation (11) implies that

$$0 \leq s(z_1) \leq s_0 \leq s(z_2) \leq 1$$

and we can distinguish between the following four “regimes”,

- Regime 1.  $\psi_1 + \rho z_1 > c_0$ . Then  $s(z_1) = s_0 = s(z_2) = 1$ .
- Regime 2.  $\psi_1 + \rho z_1 \leq c_0 < \psi_1 + \rho \bar{z}$ . Then  $0 = s(z_1) < s_0 = s(z_2) = 1$ .
- Regime 3.  $\psi_1 + \rho \bar{z} \leq c_0 < \psi_1 + \rho z_2$ . Then  $0 = s(z_1) = s_0 < s(z_2) = 1$ .
- Regime 4.  $\psi_1 + \rho z_2 \leq c_0$ . Then  $0 = s(z_1) = s_0 = s(z_2)$ .

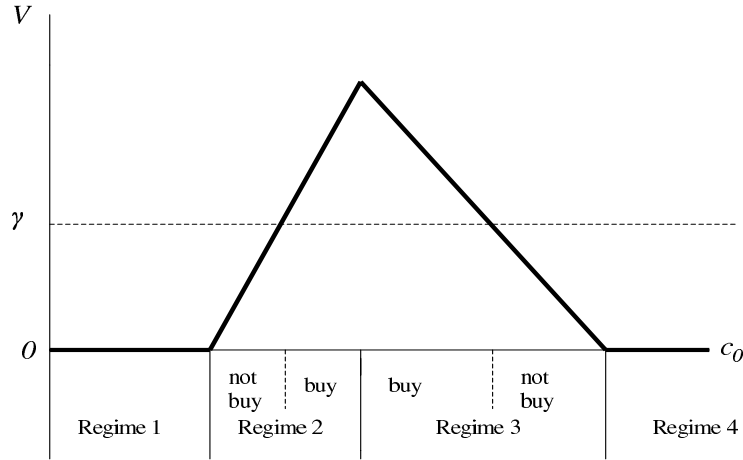
Next, we consider  $V$ . In Regime 1,  $V = 0$ , implying that the information is not bought. This is because the agent knows that he will always provide effort, under every policy, whether the policy is known or not, so the information is irrelevant for the optimal behavior. The same line of reasoning applies to Regime 4. So, for certain extreme parameter values, the agent does not acquire his value of  $Z$ . Now consider Regimes 2 and 3. We obtain that in Regime 2,

$$V = p(c_0 - \psi_1 - \rho z_1)$$

(recall that  $p := \Pr(Z = z_1)$ ). This result is particularly easy to interpret. In Regime 2, the information on  $Z$  is valuable if and only if in truth he has policy intensity  $z_1$ , because he *does* provide effort if he has no information. Therefore the value  $V$  equals minus the expected loss of making such a wrong<sup>7</sup> decision. More precisely,  $V$  equals the product of [ the probability that, when  $Z$  is not acquired, an effort  $s_0$  is chosen that is not optimal in the light of the actual  $Z$  ] and [ minus the loss of choosing  $s_0$  given that the actual  $Z$  would lead to another

---

<sup>7</sup>In this section, “wrong” is used in the sense of “wrong if the agent knows his actual value of  $Z$ ”.



**Figure 1.** The value of information  $V$  as a function (indicated by the fat solid curve) of the cost of effort  $c_0 := c(I) - c(0)$ .

Note: the parameter  $\gamma$  denotes the information acquisition cost.

choice of  $s$  ]. In Regime 2, the statement in the first square brackets equals  $\Pr(Z = z_1)$ , which equals  $p$ . The statement in the second square brackets equals [ the additional cost  $c_0$  of choosing  $s_0 = 1$  compared to the costs one would have made if one would know that  $Z = z_1$  ] minus [ the additional expected return of choosing  $s_0 = 1$  compared to the expected return if one would know that  $Z = z_1$  ].

Similarly, in Regime 3,

$$V = (1 - p)(\psi_1 + \rho z_2 - c_0)$$

Note that  $V > 0$  if and only if the optimal effort when knowing that  $Z = z_1$  differs from the optimal effort when knowing that  $Z = z_2$ . This is intuitively clear: only in these cases does the knowledge of  $Z$  potentially have an effect on the effort provided.

Consider Figure 1, plotting  $V$  against  $c_0$ . The maximum of  $V$  as a function of  $c_0$  is attained at the boundary between Regimes 2 and 3, which means  $c_0 = \psi_1 + \rho \bar{z}$ . This value  $V_{max}$  equals

$$V_{max} = |\rho|p(1 - p)(z_1 - z_2)$$

$V_{max} > \gamma$  iff there are values  $c(1), c(0)$  of the cost-of-effort function for which it is optimal to buy the information.

The expression for  $V_{max}$  is similar to the expression (8) for  $V$  in the continuous case. In particular, note that  $\text{var}(Z) = p(1-p)(z_1 - z_2)^2$ . The comparative statics results concerning the effects of  $\rho, \gamma$  and  $\text{var}(Z)$  for the continuous case therefore carry over to the present discrete case. For example, concerning  $p$  we find that acquisition is more likely if  $p$  is not too large and not too small. To understand this, note that  $p = 1/2$  maximizes the variance of  $Z$  in the population, for given values of  $z_1, z_2$ . The a priori uncertainty concerning which value of  $Z$  applies is largest for such intermediate values of  $p$ , and so is therefore the a priori probability of providing the wrong amount of effort.

What is particular about the discrete case is that for certain parameter values (namely in Regimes 1 and 4) the optimal  $s$  does not depend on  $Z$ . Due to the continuity of the value functions, the ensuing low value of information also applies for parameter values such that one is in one of the other regimes but close to Regimes 1 and 4. In those cases the loss of making a wrong decision on  $s$  is too small to justify the acquisition of one's value of  $Z$ . To analyze the associated comparative statics effects, one needs to examine for which parameter values one is likely to end up in Regimes 1 or 4. For example, acquisition is more likely if  $c_0$  is not too large and not too small. To understand this, note that for high or low values of  $c_0$  the information does not lead to behavioral changes and so is useless.

We may re-interpret the effects of the costs parameters  $c_0$  and  $\gamma$  in terms of marginal utility losses for a risk-averse agent who cannot transfer resources between time periods and who has a per-period utility function that displays decreasing absolute risk aversion. If the agent has a low per-period baseline income then the marginal utility losses of costs  $c_0$  in Period 2 and  $\gamma$  in Period 1 are relatively high. The acquisition of information is unlikely for two reasons: it is expensive and it is useless because it does not affect optimal effort. This makes violation of the exclusion restriction unlikely. If the agent has a very high per-period baseline income then the marginal utility loss of costs  $c_0$  in Period 2 is small, and acquisition is again unlikely, but now only for the second of the above two reasons.<sup>8</sup> Again, the exclusion restriction is then relatively easy to justify. For agents with per-period baseline incomes in between these extremes it is more likely that information on  $Z$  is acquired, so then the exclusion restriction

---

<sup>8</sup>Unless of course the expected marginal return of providing effort is always very small as well. The problem with characterizing comparative statics effects on optimal behavior if the latter is a highly non-linear function of a large number of parameters is that there are often joint limiting values of a subset of parameters that “push” the result in any desired direction.

is more easily violated. From this point of view, instrumental variable analyses that restrict attention to agents with low or high resources (e.g. income) are more likely to be valid than analyses that include agents with intermediate resource levels.

The results for the case  $\rho > 0$  are a mirror image of those for  $\rho < 0$ . The optimal efforts satisfy  $0 \leq s(z_2) \leq s(\bar{z}) \leq s(z_1) \leq 1$ , and these can be characterized in terms of the model parameters by applying equation (11). For sake of brevity we do not present the other results for this case.

One may combine the discrete case of this subsection with the continuous case of the previous subsection, e.g. by allowing the effort  $s$  to attain all values in a fixed interval, e.g.  $[0, 1]$ , with bounds that are binding for certain model parameter values. Sufficiently large parameter value changes then typically lead to results that correspond to those in this subsection.

### 3.4 Examples

Recall that in the first leading example, the treatment is participation in a job search assistance program for unemployed individuals, and the policy intensity differs across regions. An unemployed individual knows that there is a chance that he may enroll in a job search assistance program, and this affects his job search strategy before actual enrollment. If the program is attractive and if the applicable regional probability of enrollment into the program is known, then this probability affects the optimal private job search effort before enrollment into the program. Typically, the effort is lower if the probability of enrollment is higher, leading to a lower exit probability to work before enrollment. We now apply the results in order to inquire the conditions under which  $Z$  is a valid instrumental variable, that is, the conditions under which the agent does not acquire his value of  $Z$ .

First, it is relatively easy for agents to learn the specific situation in their own region and the effect of this on the rate at which they may expect to be treated. This questions the relatively common approach to use regional variation in the budget for (or, more generally, an indicator of the geographical availability of) active labor market programs as an instrumental variable to study causal effects of the program.

The agents' incentive to acquire information is even larger if the policy intensity varies strongly across the regions or if the agents take the a priori probability of whether one lives in the high-budget region to be equal to around 1/2 which means that about half of the agents are exposed to the high-budget situation. In

the limiting case where the policy intensity  $Z$  is binary, we have that  $Y \equiv Z$ , so the candidate instrument has maximum strength. It follows that in this case, strong instruments are more likely to be invalid instruments. In fact, the case where  $Z$  is binary is equivalent to a deliberate randomization of the binary intention to treat  $Z$  with full compliance. Typically, in experiments, the intention to treat is randomized with probability equal to  $1/2$ . Our results suggest that (provided that the individual randomization outcome can only be acquired at a positive cost when one has to decide the effort level  $s$ ) it may be better to use a smaller or larger probability, because this reduces the likelihood that the agent has an incentive to acquire and use the value of the candidate instrumental variable.

So far we have only examined a single candidate instrumental variable. Clearly, if there are many such variables, each giving only limited information on the treatment assignment process, and if the cost of information acquisition is linear in the number of variables on which information is acquired, then this reduces the likelihood that the exclusion restriction is violated. Other obvious results concern the timing of events (i.e., the relative lengths of the time periods in the model). For example, with a small amount of time between the moment at which the policy intensity is determined and the moment at which the treatment is realized, the scope for information acquisition is reduced. If treatment and outcome are realized close in time then acquisition is attractive from the point of view that its expected future returns are discounted less heavily.

### 3.5 An economic model for selective non-compliance

Suppose that agents can manipulate the probability distribution of their treatment status  $Y$ , by way of choosing an action  $s$  before  $Y$  is realized. As an extreme example, if the treatment status  $Y$  is binary, they may switch treatment status by choosing an appropriate  $s$ . Agents' optimal  $s$  may depend on  $Z$ , which may be acquired at a cost.

This model framework can be reformulated in terms of our framework, by allowing  $s$  to affect the distribution of  $Y|Z$ , which we denote by  $G(Y|Z; s)$ . If the agent does not acquire  $Z$  then he chooses action or effort  $s_0$  leading to  $G(Y|Z; s_0)$  which we denote by  $H(Y|Z)$ . We start off by assuming that  $s$  does not affect the expected outcome function  $f$ , so we write  $f := f(Y)$  instead of  $f(Y, s)$ . Note that this implies that the exclusion restriction is satisfied.

To see the connection to our framework, note that determining the optimal  $s$  involves calculation of  $\mathbb{E}_{Y|Z=z;s}f(Y)$ . This equals  $\int f(y)dG(y|z; s)$ , which can be

rewritten as<sup>9</sup>

$$\int f(y)dG(y|z; s) = \int \left[ \frac{f(y)dG(y|z; s)}{dH(y|z)} \right] dH(y|z)$$

The term in square brackets can now be interpreted as a new outcome function  $\tilde{f}$  and we can use the results derived earlier in the section to study the extent to which agents acquire their value of  $z$ .

In general, IV estimation can handle a certain amount of selective non-compliance, where it is (often tacitly) assumed that the compliance decision is made after the realization of the treatment status  $Y$ . In the setup of this subsection, the compliance decision is made *before* the realization of  $Y$ , but in this setup the ER is not violated either, because the actual outcome  $f$  only depends on  $Y$ .

Obviously, the setup of this subsection and the framework earlier in the paper can be combined by allowing the original  $f$  to depend on  $s$  as well. In either case the actual outcome function  $f$  differs from the transformed outcome function  $\tilde{f}$ . We address the bias for such cases in Subsection 4.5.

## 4 The magnitude of the bias of the instrumental variable estimator if the exclusion restriction is violated

### 4.1 The parameter of interest and the estimator

In this section we address the bias of the instrumental variable estimator due to the exclusion restriction violations that we consider. It is useful to return to the “reduced form” model representation from Section 2. IVE involves the estimation of the effect of  $Y$  on  $U$ , holding all other determinants of  $U$  constant. In the model, this is the partial derivative or first difference of  $f$  with respect to its first argument. The classical IV regression estimator, if applied to data on  $U, Y$ , and  $Z$ , estimates  $\text{cov}(U, Z)/\text{cov}(Y, Z)$ . More precisely, estimation involves that these two covariances are replaced by their sample equivalents, and then the probability limit of the estimator, which we denote by  $\hat{\beta}_{IV}$ , satisfies

---

<sup>9</sup>If  $z$  is a location parameter of the distribution of  $[Y|Z = z; s]$ , so that  $G(Y|z; s)$  can be expressed as  $G_0(Y - z|s)$  with  $G_0$  functionally independent of  $z$ , then this can be further simplified to  $\int f(y + z)dG_0(y|s)$ .

$$\widehat{\beta}_{IV} \rightarrow \frac{\text{cov}(U, Z)}{\text{cov}(Y, Z)}$$

In our model framework, if the ER is valid, this equals

$$\frac{\text{cov}(f(Y, s_0), Z)}{\text{cov}(Y, Z)} \tag{12}$$

We simply define this to be the parameter of interest  $\beta$ . It captures the mean slope of the outcome as a function of treatment status for a given fixed effort.<sup>10</sup> Note that this definition of  $\beta$  is particularly sensible if  $f$  is linear in  $Y$ , as is the case in specification (6), because then the slope of the outcome as a function of treatment status for a given fixed effort is a constant equal to  $\beta$ . Moreover, as we shall see,  $\beta$  is a local average treatment effect if  $Z$  is discrete with two points of support. Note that we assume population homogeneity of all model parameters, functions, and distributions.

Throughout the upcoming subsections we assume that the function  $f$  governing the observed outcomes is the same as the function  $f$  that agents use in Periods 1 and 2 to assess their expected utility in Period 3. This entails that agents know the average treatment effect that the researcher aims to estimate. In Subsection 4.5 we relax this assumption and show how the results change.

Expression (12) can be simplified by noting that  $\text{cov}(Y, Z)$  equals  $\text{var}(Z)$ . Also, for any two random variables  $X_1, X_2$ , there holds that  $\text{cov}(X_1, X_2) = \text{cov}(\mathbb{E}(X_1|X_2), X_2)$ . This results in

$$\beta = \frac{\text{cov}(\mathbb{E}(f(Y, s_0)|Z), Z)}{\text{var}(Z)} \tag{13}$$

If the ER does not apply, then

$$\widehat{\beta}_{IV} \rightarrow \frac{\text{cov}(f(Y, s(Z)), Z)}{\text{cov}(Y, Z)} = \frac{\text{cov}(\mathbb{E}(f(Y, s(Z))|Z), Z)}{\text{var}(Z)} \tag{14}$$

Note that  $\widehat{\beta}_{IV}$  captures the over-all effect of  $Z$  on the outcome. If the ER is violated then the over-all effect does not equal the causal treatment effect but also includes the causal chain that runs by way of the effort  $s$ . From equations (13) and (14) it follows that

$$\widehat{\beta}_{IV} - \beta \rightarrow \frac{\text{cov}(\mathbb{E}(f(Y, s(Z)) - f(Y, s_0)|Z), Z)}{\text{var}(Z)} \tag{15}$$

---

<sup>10</sup> This is even clearer for the approximation of  $\beta$  obtained by applying Stein's Lemma. For two random variables  $X_1, X_2$  and a function  $g$ , this Lemma states that  $\text{cov}(g(X_1), X_2) \approx \mathbb{E}(g'(X_1))\text{cov}(X_1, X_2)$ . Application leads to  $\beta \approx \mathbb{E}(\frac{\partial f}{\partial y}(Y, s_0))$ .

Somewhat loosely, this is an average of the effect on the outcome of the dependence of  $s$  on  $Z$ .<sup>11</sup> If treatment and effort are complements in the outcome then one may expect that the optimal effort is an increasing function of  $z$ , and consequently one may expect the asymptotic bias to be positive. In the next subsections we examine this more formally.

If  $Z$  describes the assigned treatment (as opposed to the actual treatment  $Y$ ) then the over-all effect is usually called the intention-to-treat (ITT) effect on the outcome. This can be decomposed into the actual treatment effect  $\beta$  and the announcement or ex ante effect of the treatment (see e.g. Abbring and Van den Berg, 2003, 2005, for this terminology). The latter thus equals  $\widehat{\beta}_{IV} - \beta$ , which is the asymptotic bias of the IV estimator  $\widehat{\beta}_{IV}$  of  $\beta$ .

From an econometric regression point of view, one may state that the asymptotic bias of the IV estimator results from the fact that the size of the causal effect of treatment on outcome depends on the candidate IV. An alternative way to look at the asymptotic bias is to write the outcome equation as  $U = f(Y, s_0) + (f(Y, s(Z)) - f(Y, s_0)) + \varepsilon$ . By ignoring the dependence of  $s$  on  $Z$ , an IV regression analysis takes the sum of the second and third terms in the right-hand side as the residual term in the outcome equation. Consequently, the candidate IV  $Z$  is correlated to the error term in the outcome equation.

Note that instrumental variable estimators are typically biased in case of ER violations even if there is no selectivity in the treatment.

## 4.2 Continuous effort

Suppose that the ER is violated. We are first going to examine the asymptotic bias in the model with continuous effort, quadratic costs of effort, and the outcome function (6), with the conditions that ensure that it increases in effort and the

---

<sup>11</sup> This is again more transparent for approximations of the asymptotic bias. Application of Stein's Lemma (see footnote 10) leads to  $\widehat{\beta}_{IV} - \beta \approx \mathbb{E}_Z(\frac{d}{dz}\mathbb{E}(f(Y, s(Z)) - f(Y, s_0)|Z))$ . Application of the Delta method approximation (which, for a random variable  $X$  and a function  $g$ , states that  $\text{cov}(g(X), X) \approx g'(\mathbb{E}(X))\text{var}(X)$ ) leads to  $\widehat{\beta}_{IV} - \beta \approx \frac{d}{dz}\mathbb{E}(f(Y, s(z)) - f(Y, s_0)|z)|_{z=\bar{z}}$ . It can be shown that if  $z$  is a location parameter of the distribution of  $[Y|Z = z]$  and if  $s(\bar{z}) = s_0$  and  $s(z)$  is differentiable then this simplifies further to

$$\widehat{\beta}_{IV} - \beta \approx s'(\bar{z})\mathbb{E}_\omega\left(\frac{\partial f}{\partial s}(\omega + \bar{z}, s_0)\right)$$

where  $\omega = Y - Z$  (see equation (1)). Clearly,  $s'(\bar{z})$  captures the responsiveness of effort to the value of  $Z$ , and the second term in the right hand side captures the effect of effort on the outcome. We return in the next subsection to the quality of these approximations in special cases.

treatment status. Subsequently we generalize the cost-of-effort function and the outcome function.

With the outcome function (6),  $\beta$  as defined above is the average treatment effect in the population  $\partial f(y, s)/\partial y$  which in this specific case does not depend on  $y$ , and which is evaluated at  $s = s_0$ . By substituting  $f$  into (13), and substituting  $s_0$ , we obtain,

$$\beta = \psi_2 + \rho s_0 = \psi_2 + \rho \frac{\psi_1 + \rho \bar{z}}{c_0} \quad (16)$$

which is always positive, by virtue of the conditions ensuring that  $f$  is increasing.

Similarly, by substituting (6) and  $s(z)$  from (7) into (14), we obtain,

$$\widehat{\beta}_{IV} \rightarrow \psi_2 + \frac{2\rho\psi_1}{c_0} + \frac{\rho^2 \text{cov}(Z, Z^2)}{c_0 \text{var}(Z)} \quad (17)$$

As a result,

$$\widehat{\beta}_{IV} - \beta \rightarrow \frac{\rho}{c_0} \left[ \psi_1 - \rho \bar{z} + \rho \frac{\text{cov}(Z, Z^2)}{\text{var}(Z)} \right] \quad (18)$$

In Appendix 1 we prove the following:

**Proposition 3.** *Consider the model with continuous effort, quadratic costs of effort, and the outcome function (6) with the conditions that ensure that it increases in effort and the treatment status. If the exclusion restriction is violated then the asymptotic bias of the IV estimator  $\widehat{\beta}_{IV}$  has the same sign as  $\rho$ .*

This result reinforces earlier results on the importance of the interaction parameter. For  $\rho > 0$  the result is very intuitive. In that case, if the agent acquires his value  $z$  of  $Z$ , then effort is increasing in  $z$ ; the estimated treatment effect is then boosted, and the causal effect of the treatment is over-estimated, so asymptotically  $\widehat{\beta}_{IV} > \beta$ . In Appendix 1 we demonstrate that if  $\rho > 0$  then  $\widehat{\beta}_{IV}$  and  $\widehat{\beta}_{IV} - \beta$  increase in  $\rho$ . If  $\rho$  increases then it is optimal to make the effort  $s(z)$  more responsive to  $z$ , so that the announcement effect and therefore the asymptotic bias increase. For  $\rho < 0$  the additive and the interaction effects of  $s(z)$  have opposite signs, but the former dominates. Note also that the asymptotic bias decreases in  $c_0$ . If costs of effort are low then it pays off to supply a large amount of effort for certain  $z$ , thus making effort more responsive to  $z$ .

In Appendix 1 we also demonstrate that (18) can be rewritten as

$$\widehat{\beta}_{IV} \rightarrow 2\beta - \psi_2 + \frac{\rho^2 \bar{z}}{c_0} \xi \sigma \quad (19)$$

where  $\xi$  and  $\sigma$  are the skewness and the coefficient of variation of  $Z$ , so  $\xi := \kappa_3/\kappa_2^{\frac{3}{2}}$  and  $\sigma := \kappa_2^{\frac{1}{2}}/\bar{z}$ , with  $\kappa_i := \mathbb{E}(Z - \bar{z})^i$ . From this equation some additional comparative statics results follow. As  $\beta$  does not depend on moments of  $Z$  higher than the first, it follows that  $\widehat{\beta}_{IV}$  and  $\widehat{\beta}_{IV} - \beta$  increase in the skewness  $\xi$  of  $Z$ . This makes sense: e.g. if  $Z$  is highly skewed to the right then relatively many agents have a very high treatment value and a very high effort level, leading to a very high correlation between outcome and candidate instrument.

Notice that if  $Z$  is symmetric (i.e.  $\xi = 0$ ) then the variance of  $Z$  does not affect the asymptotic bias of  $\widehat{\beta}_{IV}$ . In this case, increasing  $\text{var}(Z)$  leads to an increasing value of information, and therefore an increasing likelihood of ER violation, but not to an increasing asymptotic bias of the IV estimator.<sup>12</sup> So in this sense the asymptotic bias is even less sensitive to the strength of the candidate instrument than the ER.

Equation (19) can also be used to bound  $\beta$  further. We consider using the sign of the skewness of  $Z$ , which is always observable. In particular, in case of a balanced experiment,  $\xi = 0$ . We also consider cases where one has a priori knowledge on the sign of  $\rho$  and on whether the additive effect  $\psi_2$  of treatment on outcome is zero or positive.

**Corollary 1.** *Under the conditions of Proposition 3,*

- (i). *If  $Z$  is symmetric or skewed to the left then, asymptotically,  $\beta \geq \frac{1}{2}\widehat{\beta}_{IV}$ .*
- (ii). *If  $Z$  is symmetric and  $\psi_2 = 0$  then, asymptotically,  $\beta = \frac{1}{2}\widehat{\beta}_{IV}$ .*
- (iii). *If  $Z$  is skewed to the right and  $\psi_2 = 0$  then, asymptotically,  $\beta < \frac{1}{2}\widehat{\beta}_{IV}$ .*
- (iv). *If  $Z$  is symmetric or skewed to the left and  $\rho > 0$  then, asymptotically,  $\frac{1}{2}\widehat{\beta}_{IV} \leq \beta < \widehat{\beta}_{IV}$ .*

Extension to other cases is straightforward. Notice that  $\psi_2 = 0$  implies that  $\rho > 0$ .

We now consider more general models. First, if we replace the additive term  $\psi_0 + \psi_2 y$  in the outcome function by a general function  $k_2(y)$  then all asymptotic bias results remain valid. Next, we relax the assumption that costs of effort are quadratic.

**Proposition 4.** *Consider the model with continuous effort and the outcome function (6) with  $\psi_0 > 0, \psi_1, \psi_2 \geq 0$ , and  $\rho > 0$ . Let the cost of effort be increasing,*

---

<sup>12</sup>With our outcome and cost functions, the case  $\xi = 0$  is also the case for which the asymptotic bias approximations by Stein's Lemma and the Delta method, in footnote 11, are exact.

convex, and differentiable and lead to an interior solution for the optimal effort. If the exclusion restriction is violated then the asymptotic bias of the IV estimator  $\widehat{\beta}_{IV}$  is positive.

See Appendix 2 for the proof.

It is more difficult to extend the result to the case where treatment and effort are substitutes in the outcome (i.e.,  $\rho < 0$ ). This is because with  $\rho < 0$  the additive treatment effect and the interaction effect of treatment and effort have opposite signs.<sup>13</sup>

Next, we relax the assumption that  $f$  satisfies (6), so we do not make any parametric assumptions on the cost function and the outcome function. We use the monotonicity result of Athey and Levin (2001) (see Subsection 3.2) as an input. As noted above, in the literature on decision making with signals on the unknown state of the world, comparative statics are sometimes too hard to analyze in terms of the model primitives, and one can only derive results for given (i.e., not optimally determined) effort functions. To a certain extent our next result shares this feature, in that we assume that the optimal effort function  $s(z)$  satisfies  $s(\bar{z}) = s_0$  without translating this into model primitives. The assumption means that the agent's optimal effort if he knows that he has the average treatment intensity is equal to the agent's optimal effort if he does not know his intensity. In the cases considered in the previous propositions, this assumption actually follows from the assumptions on the model primitives.

**Proposition 5.** *Let the outcome function  $f(y, s)$  be non-negative, increasing and supermodular in  $y$  and  $s$ . Let the treatment status  $Y$  increase in  $Z$  in the sense of first-order stochastic dominance. Let the optimal effort function  $s(z)$  satisfy  $s(\bar{z}) = s_0$ . If the exclusion restriction is violated then the asymptotic bias of the IV estimator  $\widehat{\beta}_{IV}$  is non-negative.*

See Appendix 3 for the proof. Supermodularity captures complementarity and is satisfied if the cross-derivative of  $f(y, s)$  with respect to  $y$  and  $s$  is positive. The

---

<sup>13</sup>See e.g. equation (33) in Appendix 2. Specifically, with less parametric assumptions, it is cumbersome to formulate and exploit restrictions ensuring that the outcome increases in treatment and effort everywhere if treatment and effort are substitutes. However, using the Delta method approximation (see footnote 11) the desired result follows. We substitute the result from Appendix 2 that  $s(\bar{z}) = s_0$  into the approximation, to obtain:

$$\widehat{\beta}_{IV} - \beta \approx (\psi_1 + \rho\bar{z})s'(\bar{z})$$

With  $\rho < 0$ , the function  $s(z)$  is decreasing (see Appendix 2), and the requirement that  $f$  is increasing implies that  $\bar{z} < \psi_1/(-\rho)$ . The right-hand side of the above approximation is therefore negative.

assumption that  $Y$  increases in  $Z$  in the sense of first-order stochastic dominance is not nested with the assumption that  $\mathbb{E}(Y|Z = z) = z$ , although the former is weaker as a condition for  $\mathbb{E}(Y|Z = z)$  to increase in  $z$ .

Note that Proposition 5 does not require continuity of effort  $s$ . In the next subsection we consider the discrete case in more detail.

### 4.3 Discrete effort

Let  $s$  be binary and let  $Z$  have a discrete distribution with points of support  $z_1 > z_2$ , like in Subsection 3.3. It can be shown that  $\beta$  in (13) then simplifies to

$$\beta = \frac{\mathbb{E}(f(Y, s_0)|z_1) - \mathbb{E}(f(Y, s_0)|z_2)}{z_1 - z_2}$$

which is a local average treatment effect (compare Imbens and Angrist, 1994; note that the denominator equals  $\mathbb{E}(Y|z_1) - \mathbb{E}(Y|z_2)$ ). If the ER is valid then the IV estimator  $\hat{\beta}_{IV}$  converges to this number. Accordingly, we can use a Wald estimator as IV estimator.

If the ER is violated then we can simplify (14) to

$$\hat{\beta}_{IV} \rightarrow \frac{\mathbb{E}(f(Y, s(z_1))|z_1) - \mathbb{E}(f(Y, s(z_2))|z_2)}{z_1 - z_2} \quad (20)$$

so that again  $\hat{\beta}_{IV}$  captures the over-all effect of  $Z$  on the outcome.

Now let us proceed by taking the outcome function  $f$  to satisfy (6). From the above (as well as from the previous subsection) it immediately follows that  $\beta = \psi_2 + \rho s_0$ . If the ER is violated then necessarily  $s(z_1) \neq s(z_2)$ . With  $\rho < 0$ , violation implies that  $0 = s(z_1) < s(z_2) = 1$ , whereas with  $\rho > 0$ , this is reversed. By elaborating on equation (14) or on equation (20), we obtain

$$\begin{aligned} \hat{\beta}_{IV} &\rightarrow \psi_2 - \frac{\psi_1 + \rho z_2}{z_1 - z_2} && \text{if } \rho < 0 \\ \hat{\beta}_{IV} &\rightarrow \psi_2 + \frac{\psi_1 + \rho z_1}{z_1 - z_2} && \text{if } \rho > 0 \end{aligned}$$

Depending on the sign of  $\rho$  and the value of  $s_0$ , we have four different expressions for the asymptotic bias,

$$\begin{aligned} \hat{\beta}_{IV} - \beta &\rightarrow -\frac{\psi_1 + \rho z_2}{z_1 - z_2} && \text{if } \rho < 0 \text{ and } s_0 = 0 \\ \hat{\beta}_{IV} - \beta &\rightarrow -\frac{\psi_1 + \rho z_1}{z_1 - z_2} && \text{if } \rho < 0 \text{ and } s_0 = 1 \end{aligned}$$

$$\widehat{\beta}_{IV} - \beta \rightarrow \frac{\psi_1 + \rho z_1}{z_1 - z_2} \quad \text{if } \rho > 0 \text{ and } s_0 = 0$$

$$\widehat{\beta}_{IV} - \beta \rightarrow \frac{\psi_1 + \rho z_2}{z_1 - z_2} \quad \text{if } \rho > 0 \text{ and } s_0 = 1$$

Thus,

**Proposition 6.** *Consider the model with binary effort, a candidate IV with a discrete distribution with two points of support, and the outcome function (6) with the conditions that ensure that it increases in effort and the treatment status. If the exclusion restriction is violated then the asymptotic bias of the IV estimator  $\widehat{\beta}_{IV}$  has the same sign as  $\rho$ .*

The bias terms are larger if  $z_1$  and  $z_2$  are close. However, recall that this result is derived under the simplifying assumptions that all agents are in Regime 2 and acquire information on  $Z$ , so that  $z_1$  and  $z_2$  can not be too close.

From the interpretation of the bias as capturing the announcement effect or ex ante effect of the treatment, it follows that a large bias term is equivalent to a large ex ante effect. This means that a large bias is often associated to a high value of information.

#### 4.4 Testing for a causal treatment effect in absence of an exclusion restriction

Despite the fact that the IV estimator is asymptotically biased if the ER is violated, we can still use standard tests to inquire whether the treatment has a causal effect on the outcome. To see this, notice that in absence of a treatment effect the value of information is zero, so there is no acquisition of information, implying that the ER is satisfied and the IV treatment effect estimator is asymptotically equal to zero. The standard asymptotic tests of the null hypothesis of a zero treatment effect therefore have the correct size. This exploits the insight that in our model any violation of the ER is behaviorally triggered by a non-zero treatment effect.

One may wonder what the power is of such tests if the null hypothesis of no treatment effect is incorrect in reality and, in addition, the ER is violated. We shed some light on this by examining the (asymptotic) signs of  $\beta$ ,  $\widehat{\beta}_{IV}$ , and  $\widehat{\beta}_{IV} - \beta$ . According to expression (13), the treatment effect parameter  $\beta$  has the same sign as  $\text{cov}(\mathbb{E}(f(Y, s_0)|Z), Z)$ . If the treatment status  $Y$  increases in  $Z$  in the sense of first-order stochastic dominance, and  $f(y, s)$  strictly increases in  $y$ , then  $\mathbb{E}(f(Y, s_0)|Z)$  increases in  $Z$ , and consequently  $\beta > 0$ .

To proceed, it matters again whether treatment and effort are complements or substitutes in the outcome. For the former case we saw earlier in this section that asymptotically  $\widehat{\beta}_{IV} \geq \beta$ . Therefore one may expect a higher power of the standard asymptotic tests of the null hypothesis of a zero treatment effect, compared to when the ER is not violated. This makes sense: the true treatment effect is magnified by the agents' efforts.

For the substitution case with our parametric functional forms of  $f(y, s)$  and  $c(s)$ , we saw that asymptotically  $\widehat{\beta}_{IV} < \beta$ . The next example shows that it is even possible that asymptotically  $\widehat{\beta}_{IV} = 0$ .

*Example 1.* Let  $\psi_1 = \psi_2 = c_0 = 1$  and  $\rho = -3/4$ . Note that these values satisfy the parameter restriction  $\psi_2 c_0 + \rho \psi_1 > 0$  (see (9)). In turn, they imply that  $\Pr(Z < 4/3) = 1$ . By substituting our  $f$  and  $s(z)$  into (14) we obtain that asymptotically  $\widehat{\beta}_{IV} \geq 0$  iff  $\text{cov}(Z, -\frac{1}{2}Z + \frac{9}{16}Z^2) \geq 0$ . The function  $-\frac{1}{2}Z + \frac{9}{16}Z^2$  strictly decreases in  $Z$  on the interval  $(0, \frac{4}{9})$ , so if this interval includes the support of  $Z$  then asymptotically  $\widehat{\beta}_{IV} < 0$ . If  $Z$  has a discrete distribution with as only two points of support  $1/3$  and  $5/9$  then  $\widehat{\beta}_{IV} \rightarrow 0$ .

In this pathological case the asymptotic power is zero. In general, in the substitution case, one may expect the power to be lower if the ER is violated than if the ER is satisfied.

One way to improve the power is to look for evidence of higher-order dependencies between the candidate instrument  $Z$  and the outcome  $U$ . Under the null hypothesis, these must be absent as well.

## 4.5 Systematic difference between the outcome and the agent's utility in the outcome period

We now allow the function  $f^*$  governing the observed outcomes to differ from the function  $f$  that agents use in Periods 1 and 2 to evaluate their expected utility in Period 3. This is relevant in a number of cases. The function  $f^*$  may not take account of side-effects that make the treatment unattractive to the agent. Also, agents may not know the value of the average treatment effect that the researcher aims to estimate, and the assumptions they make about this in Periods 1 and 2 may be deviate systematically from the actual  $f^*$ .

The analysis in Section 3 describes how the choice of the effort level  $s$  and the validity of the ER depend on  $f$ . These decisions do not depend on  $f^*$ , so all results from that section also apply to the current framework.

To proceed, we start assuming that  $f$  satisfies (6), that effort is continuous, and that the actual outcomes  $U$  are generated by  $U = f^*(Y, s) + \varepsilon$  with

$$f^*(y, s) = \psi_0^* + \psi_1^*s + \psi_2^*y + \rho^*ys \quad (21)$$

The case  $\psi_2 < \psi_2^*$  is particularly interesting as it captures situations where the agent perceives a disutility of the treatment that is not revealed in the outcome. Also, with  $\rho < \rho^* < 0$ , the agent dislikes the combination of high treatment and high effort more than is warranted by the corresponding actual outcome.

Equations (16) and (17), expressing the parameter of interest  $\beta$  and the probability limit of its estimator  $\widehat{\beta}_{IV}$  in terms of the model parameters, are now replaced by

$$\beta = \psi_2^* + \rho^*s_0 = \psi_2^* + \rho^* \frac{\psi_1 + \rho\bar{z}}{c_0} \quad (22)$$

$$\widehat{\beta}_{IV} \rightarrow \psi_2^* + \frac{\rho\psi_1^* + \rho^*\psi_1}{c_0} + \frac{\rho\rho^* \text{cov}(Z, Z^2)}{c_0 \text{var}(Z)} \quad (23)$$

In Appendix 4 we prove the following extension of Proposition 3,

**Proposition 7.** *Consider the model with continuous effort, quadratic costs of effort, the agents' outcome utility function (6) with the conditions that ensure that it increases in effort and the treatment status, and the outcome function (21) with the conditions that (i) if  $\rho^* \geq 0$  then  $\psi_1^* \geq 0$ , and (ii) if  $\rho^* < 0$  then  $\psi_1^* > 0$  and  $\Pr(Z < -\psi_1^*/\rho^*) = 1$ . If the exclusion restriction is violated then the asymptotic bias of the IV estimator  $\widehat{\beta}_{IV}$  has the same sign as  $\rho$ , with the single exception that if  $\rho^* = \psi_1^* = 0$  then this bias is zero.*

The conditions (i) and (ii) on the range of  $\psi_1^*$  and the support of  $Z$  serve to ensure that  $f^*(y, s)$  is increasing in  $s$ . This implies that we assume that both  $f$  and  $f^*$  are increasing in  $s$ . Note that to some extent any disutility of  $s$  can be accommodated for by the cost function  $c(s)$  in Period 2. Note that we do not assume that the actual average treatment effect has the same sign as the perceived expected treatment effect, since we do not make assumptions on (the sign of)  $\psi_2^*, \rho^*$ , and  $\psi_0^*$ .

Proposition 7 implies that of the two interaction parameters, it is the interaction in the agent's objective function that drives the asymptotic bias. The underlying reason is that the sign of this interaction determines whether the agent's effort is increasing or decreasing in  $z$ . If it is increasing then this boosts the estimated average treatment effect regardless of how the actual outcomes

generated. Notice also that the asymptotic bias is completely independent of the actual and the perceived additive treatment effect parameters  $\psi_2^*$  and  $\psi_2$  and of the strength of the candidate instrument. The other propositions of Subsections 4.2 and 4.3 can be generalized accordingly.

It is sometimes plausible that the agent can have access to more information on the determinants of the treatment assignment process than the analyst can observe. As this is somewhat related to the topic of this subsection, we briefly discuss its main implications for the asymptotic bias here, using a simple framework where the assignment process can be captured by the following “treatment equation”

$$Y = Z_1 + Z_2 + \omega \tag{24}$$

with  $\mathbb{E}(\omega|Z) = 0$  and also  $\mathbb{E}(\varepsilon|Z) = 0$ , where  $Z := (Z_1, Z_2)$ . We define  $Z_1$  and  $Z_2$  such that  $\mathbb{E}(Z_1) = 0$ , and, consequently,  $\mathbb{E}(Y|Z_1 = z_1, Z_2 = z_2) = z_2$ . The analyst observes  $U, Y$ , and  $Z_2$  but not  $Z_1$ . The agent is able to acquire his values of  $Z_1$  and of  $Z_2$ . We adopt the usual functional forms for  $f(y, s)$  and  $c(s)$  and we start off with the assumption that  $Z_1 \perp\!\!\!\perp Z_2$ . Notice that the parameter of interest  $\beta$  is now defined using  $Z_2$ , as follows:  $\beta := \text{cov}(f(Y, s_0), Z_2) / \text{cov}(Y, Z_2)$ .

We can distinguish between four cases, depending on which  $Z_i$  are acquired by the agent. Firstly, suppose the ER applies to both  $Z_1$  and  $Z_2$ . Then either (or both) can be used as instrumental variables. Secondly, suppose the ER only applies to  $Z_1$ . Then  $Z_1$  is a valid instrument but it is unobserved to the analyst, so it cannot be used for inference. In this case,  $Z_1$  can be subsumed into  $\omega$ , which leads to the main model framework of this paper. Thirdly, suppose the ER only applies to  $Z_2$ . Then  $Z_2$  can be used for IV estimation.

Fourthly, and most interestingly, suppose the ER is violated for both  $Z_1$  and  $Z_2$ . We redefine  $Z_1$  and  $Z_2$  such that  $\mathbb{E}(Z_1) = 0$ . By analogy to Subsection 4.1,  $\beta$  can now be shown to equal  $\psi_2 + \rho s_0$ , where  $s_0 = (\psi_1 + \rho \bar{z}_2) / c_0$ . The IV estimator  $\widehat{\beta}_{IV}$  converges to  $\text{cov}(f(Y, s(Z)), Z_2) / \text{cov}(Y, Z_2)$ , where  $s(z) = (\psi_1 + \rho(z_1 + z_2)) / c_0$ . This can be shown to lead to exactly the same expression for  $\widehat{\beta}_{IV}$  as (17) in Subsection 4.2, with  $Z_2$  now replacing  $Z$ . Consequently, the results of the previous subsections all apply here. In sum, the results of the paper are robust with respect to whether the agent has more information on the treatment assignment process than the analyst. Note though that the maintained assumption that the private information is orthogonal to the shared information (i.e.  $Z_1 \perp\!\!\!\perp Z_2$ ) is crucial. Without this assumption it is more difficult to assess the asymptotic bias.

## 4.6 Changing the variance of the candidate instrument or the randomization probability

By definition, it is difficult to empirically test the predictions of the paper, because exclusion restrictions are untestable. As an alternative approach, one may look for exogenous variation in the model parameters that affect the value of information. For a certain range of parameter values the ER is satisfied, whereas for another range it is not. With sufficient variation of these parameters in the data, it can be verified whether the estimated treatment effect behaves as predicted as a function of the parameter. Note that this does not involve a comparison between different values of a candidate IV. Rather, the comparison is at a deeper level, namely between different evaluation settings.

In fact, of all the determinants of  $V$  in the continuous case of Subsection 3.2, only the variance of  $Z$  is more or less directly observable. We therefore compare settings with different distributions of the candidate IV, notably with different variances. To proceed, let effort be continuous, and let  $Z$  have a distribution with separate parameters capturing the mean  $\bar{z}$  and the variance  $\sigma_z^2$ . The data should now include settings with different values of  $\sigma_z^2$ .

From the results in the paper it follows that the ER is satisfied iff  $\sigma_z^2 < 2c_0\gamma/\rho^2$ . For these (smaller) values of  $\sigma_z^2$ , the IV estimator  $\widehat{\beta}_{IV}$  estimates the policy effect parameter  $\beta$ . This parameter depends on  $\bar{z}$  but not on  $\sigma_z^2$  or any other feature of the distribution of  $Z$ . For the values of  $\sigma_z^2$  larger than  $2c_0\gamma/\rho^2$ , the ER is violated, and  $\widehat{\beta}_{IV}$  estimates a number different from  $\beta$ , where the sign of the difference is determined by  $\rho$ . In general, this number varies itself with  $\sigma_z^2$ . Moreover,  $\widehat{\beta}_{IV}$  as a function of  $\sigma_z^2$  is discontinuous at  $2c_0\gamma/\rho^2$ .

In sum, the model predicts that as  $\sigma_z^2$  increases, there is a point at which  $\widehat{\beta}_{IV}$  makes a discontinuous jump. This can be verified empirically. To the left of the discontinuity point,  $\widehat{\beta}_{IV}$  is constant, but this depends crucially on the usual functional forms for  $f(y, s)$  and  $c(s)$ . If one has a priori knowledge on the sign of the interaction of  $y$  and  $s$  in  $f(y, s)$  then one can verify empirically that the jump has the same sign.

Now suppose that  $Z \equiv Y$  is binary, with  $p := \Pr(Z = 1)$ . This captures experiments in which participants can only observe their treatment status at a cost. The parameter  $p$  is then the randomization probability or treatment assignment probability. This is the only parameter of the distribution of  $Z$ , so it is not possible to vary  $\text{var}(Z)$  while keeping  $\bar{z}$  constant. Specifically,  $\text{var}(Z) = p(1 - p)$  while  $\mathbb{E}(Z) = p$ . We therefore aim to compare settings with different values of  $p$ , which is directly observable.

From the results in the paper it follows that the ER is violated iff  $p(1-p) > 2c_0\gamma/\rho^2$ . This translates into a symmetric non-empty interval around  $p = 1/2$  provided that  $\rho^2 > 8c_0\gamma$ , which we assume here. The interval is strictly embedded in  $(0, 1)$ . For values of  $p$  close to 0 or 1, the IV estimator  $\widehat{\beta}_{IV}$  estimates the policy effect parameter  $\beta$ , which equals

$$\beta = \psi_2 + \rho \frac{\psi_1 + \rho p}{c_0}$$

This parameter depends on  $p$ , meaning that the average treatment effect differs across experiments with different randomization probabilities, even if the ER is satisfied, like in a medical double-blind experiment. It is interesting to examine this in some detail. If  $p$  is high then the agents know that it is likely that they have been assigned to the treatment group. With  $\rho > 0$  ( $\rho < 0$ ), the marginal expected return of effort is higher (lower) if the expected treatment is higher, so agents then have an incentive to provide more (less) effort  $s_0$ . This boosts the treatment effect.

This provides an alternative explanation for the evidence gathered by Malani (2006b) that the estimated treatment effect in double-blind medical trials increases with the announced treatment probability. Malani (2006b) attributes this fact to placebo effects. (Malani, 2006a, uses similar data to study how enrollment among heterogeneous individuals in a medical trial depends on the treatment probability.) In our framework, agents rationally adjust their efforts in response to the probability of being assigned to the treatment group.

For the values of  $p$  close to  $1/2$ , the ER is violated, and  $\widehat{\beta}_{IV}$  estimates a number different from  $\beta$ , where the sign of the difference is determined by  $\rho$ . Specifically,

$$\widehat{\beta}_{IV} \rightarrow \psi_2 + \frac{2\rho\psi_1}{c_0} + \frac{\rho^2}{c_0}$$

This does not depend on  $p$ , but that result depends crucially on the functional forms for  $f(y, s)$  and  $c(s)$ . More importantly,  $\widehat{\beta}_{IV}$  as a function of  $p$  is discontinuous at  $1/2 \pm 1/2\sqrt{1 - 8c_0\gamma/\rho^2}$ .

In sum, the model predicts that as  $p$  increases from 0 to 1, there are two points at which  $\widehat{\beta}_{IV}$  makes a discontinuous jump, and these jumps have opposite signs. This can be verified empirically. If one has a priori knowledge on the sign of the interaction of  $y$  and  $s$  in  $f(y, s)$  then one can verify empirically that the sequence of the signs of the jumps is correct. For example, if the interaction sign is negative, so treatment and effort are substitutes, then one would expect the

estimated treatment effect to be relatively small for treatment probabilities close to  $1/2$ .

## 5 Conclusions

Exclusion restrictions for instrumental variable estimation are untestable and therefore need to be justified externally. We consider situations in which there is a time interval between the moment the agent realizes that he may be exposed to the policy and the actual exposure. We economically analyze the decision whether to acquire information concerning the value of the candidate instrumental variable.

The results suggest, first, that the exclusion restriction is more likely to be violated if the candidate instrument covers a large shift in policy intensity or if it divides the population into groups of similar size. We also find that the exclusion restriction is more likely to be violated if the outcome of interest strongly depends on interactions between the agent’s effort before the outcome is realized on the one hand, and the agent’s treatment status on the other.

Deliberate randomization of the intention to treat, like in the case of social experiments, does not help. The randomization outcome is typically available to the agent at low cost, and the typical randomization probability of  $1/2$  corresponds to a high incentive to acquire one’s realization. In fact, it may be better to use a smaller or larger probability, because this reduces the incentive to acquire and use the value of the candidate instrumental variable.

Having a weak instrument does not help either. Weakness of the candidate instrument, as defined or measured in ways proposed in the literature, is not directly related to the likelihood that the exclusion restriction is violated.

With discrete effort, instrumental variable analyses that restrict attention to agents with low or high resources (e.g. income) are more likely to be valid than analyses that include agents with intermediate resource levels. The reason is that for the former groups, the information is useless because it does not affect optimal effort. In addition, for the low resource agents, it may be too expensive.

Finally, concerning the bias in case of violation of the exclusion restriction, we find that typically, it is large if the value of information is large.

Suppose the circumstances are such that it is plausible that the exclusion restriction is violated, so IV cannot be applied. One way to proceed is to estimate a structural economic model. Alternatively, with sufficient variation in the timing of treatment and the outcome of interest, then one may follow the so-called “timing-of-events” approach (Abbring and Van den Berg, 2003), that

is, impose some semi-parametric structure and exploit the variation in the timing of events for identification of a causal treatment effect. Note that the value of  $Z$ , if observed, does not have any influence on optimal behavior *after* the actual treatment. Abbring and Van den Berg (2005) exploit this so-called “ex post exclusion restriction” for identification of selection effects. They also demonstrate that, with sufficient semi-parametric structure, the information in the timing of events enables identification of the ex ante effect and the treatment effect.

# Appendix

## Appendix 1 Proof of Proposition 3 and some implications

First of all, recall that  $\Pr(Z \geq 0) = 1$  and  $\text{var}(Z) > 0$ . From (18), if  $\rho = 0$  then the asymptotic bias is zero.

Now consider the case  $\rho > 0$ . From Subsection 3.2, the conditions on the outcome function that are listed in the proposition amount to  $\psi_0 > 0, \psi_1 \geq 0, \psi_2 \geq 0$ . Also,  $c_0 > 0$ . By substituting  $\text{cov}(Z, Z^2) = \mathbb{E}(Z^3) - \bar{z}\mathbb{E}(Z^2)$  into (18), it follows that asymptotically  $\hat{\beta}_{IV} - \beta$  has the same sign as

$$\psi_1 + \rho\bar{z} \left[ \frac{\tilde{\mu}_3 - 2\tilde{\mu}_2 + 1}{\tilde{\mu}_2 - 1} \right] \quad (25)$$

where  $\tilde{\mu}_i := \mathbb{E}Z^i/\bar{z}^i$  for  $i = 2, 3$ . Note that the denominator of the above term in square brackets is proportional to  $\text{var}(Z)$  which is positive. Consequently, the term in square brackets is positive iff the numerator is positive. The latter is equivalent to  $\tilde{\mu}_3 > \tilde{\mu}_2^2 - (\tilde{\mu}_2 - 1)^2$ . From Shohat and Tamarkin (1943)'s results for the so-called Stieltjes Moment Problem it follows that a random variable  $Z$  with  $\Pr(Z \geq 0) = 1$  and  $\text{var}(Z) > 0$  necessarily satisfies  $\mathbb{E}(Z)\mathbb{E}(Z^3) > (\mathbb{E}(Z^2))^2$ , with the exception of the special case in which the support of  $Z$  consists of two mass points, one of which is zero. The inequality is equivalent to  $\tilde{\mu}_3 > \tilde{\mu}_2^2$ , which in turn implies that  $\tilde{\mu}_3 > \tilde{\mu}_2^2 - (\tilde{\mu}_2 - 1)^2$ . For the special case in which the support of  $Z$  consists of two mass points, one of which is zero, there holds that  $\tilde{\mu}_3 = \tilde{\mu}_2^2$ , but also that  $\tilde{\mu}_2 - 1 > 0$ , so that again  $\tilde{\mu}_3 > \tilde{\mu}_2^2 - (\tilde{\mu}_2 - 1)^2$ . Consequently, the asymptotic bias is always positive.

Now consider the case  $\rho < 0$ . From Subsection 3.2, the conditions on the outcome function that are listed in the proposition amount to  $\psi_0 > 0, \psi_1 > 0, \psi_2 > 0, \Pr(Y < \frac{\psi_1}{-\rho} | Z = z) = 1$  for all possible realizations  $z$ , and  $\psi_2 c_0 + \rho\psi_1 > 0$ , with  $c_0 > 0$ . From equation (18), and by analogy to equation (25), we obtain that asymptotically  $\hat{\beta}_{IV} - \beta$  has the same sign as

$$\frac{\psi_1}{\rho} + \frac{\mu_3 - 2\bar{z}\mu_2 + \bar{z}^3}{\mu_2 - \bar{z}^2} \quad (26)$$

where  $\mu_i := \mathbb{E}(Z^i)$ . In Subsection 3.2 we saw that the parameter inequalities for this case imply that  $\Pr(Z < \frac{\psi_1}{-\rho}) = 1$ , so  $Z$  satisfies  $\Pr(Z \in [0, -\psi_1/\rho]) = 1$ . We denote  $-\psi_1/\rho$  by  $z_u$ . The term (26) is negative iff

$$-z_u(\mu_2 - \bar{z}^2) + \mu_3 - 2\bar{z}\mu_2 + \bar{z}^3 \quad (27)$$

is negative. For convenience, we rewrite this in terms of central moments of  $Z$ . Let  $\kappa_i := \mathbb{E}(Z - \bar{z})^i$ . Then (27) equals

$$\kappa_3 - (z_u - \bar{z})\kappa_2 \quad (28)$$

To sign this, we apply results for the so-called Hausdorff Moment Problem (see Shohat and Tamarkin, 1943, and Frontini and Tagliani, 1995). In particular, a random variable  $Z_0$  with  $\Pr(Z_0 \in [0, 1]) = 1$  and  $\text{var}(Z_0) > 0$  necessarily satisfies

$$(1 - \mathbb{E}(Z_0))(\mathbb{E}(Z_0^2) - \mathbb{E}(Z_0^3)) > (\mathbb{E}(Z_0) - \mathbb{E}(Z_0^2))$$

We take  $Z = z_u \cdot Z_0$ , so that  $Z$  satisfies  $\Pr(Z \in [0, z_u]) = 1$  and  $\text{var}(Z) > 0$ , as required, and we use the notation  $\bar{z}$  and  $\mu_i$  to denote its moments. Clearly,  $\bar{z} = z_u \mathbb{E}(Z_0)$  and  $\mu_i = z_u^i \mathbb{E}(Z_0^i)$ . The above moment inequality can now be written as

$$z_u^2(\mu_2 - \bar{z}) - z_u(\mu_3 - \bar{z}\mu_2) + \bar{z}\mu_3 - \mu_2^2 > 0$$

and in terms of central moments of  $Z$  this simplifies to

$$0 > \kappa_2^2 + (z_u - \bar{z})[\kappa_3 - (z_u - \bar{z})\kappa_2]$$

Since  $z_u > \bar{z}$ , this implies that  $\kappa_3 - (z_u - \bar{z})\kappa_2 < 0$ . Consequently, the asymptotic bias is always negative. This completes the proof of Proposition 3.  $\square$

Notice that the result for  $\rho < 0$  also applies if we allow for  $Y$  and  $Z$  to be able to attain the value  $z_u$ . Also, notice that by substituting  $\text{cov}(Z, Z^2) = \mathbb{E}(Z^3) - \bar{z}\mathbb{E}(Z^2)$  into (18) and by rewriting this in terms of central moments  $\kappa_i$  of  $Z$ , we obtain

$$\widehat{\beta}_{IV} - \beta \rightarrow \frac{\rho\psi_1}{c_0} + \frac{\rho^2\bar{z}}{c_0} \left[ \frac{\kappa_3}{\bar{z}\kappa_2} + 1 \right] \quad (29)$$

By substituting  $\beta$  into the right-hand side, and by noting that  $\kappa_3/(\bar{z}\kappa_2)$  equals the skewness of  $Z$  times the coefficient of variation of  $Z$ , we obtain equation (19) in the main text. Also, by combining equation (29) with the proof for the case  $\rho > 0$ , it becomes clear that in that case  $\widehat{\beta}_{IV}$  and  $\widehat{\beta}_{IV} - \beta$  asymptotically increase in  $\rho$ .

## Appendix 2 Proof of Proposition 4

Because  $c(s)$  is not parameterized, there are no parameterized expressions for  $s(z)$  and  $s_0$  either. However, from equations (3) and (4) it follows that with

the outcome function (6) and with the assumptions on  $f$  and on  $c(s)$  listed in the proposition, the optimal effort level  $s(z)$  in absence of the ER increases in  $z$  (whereas with  $\rho < 0$  it decreases in  $z$ ). Moreover,

$$c'(s_0) = \psi_1 + \rho\bar{z} = c'(s(\bar{z}))$$

implying that  $s(\bar{z}) = s_0$ . The latter means that the agent's optimal effort if he knows that he has the average treatment intensity is equal to the agent's optimal effort if he does not know his intensity. We will use this below. Note that the function  $s(z)$  cannot be constant over the support of  $Z$ .

If we substitute (6) into (13) then we obtain that

$$\beta = \psi_2 + \rho s_0 \tag{30}$$

Similarly, if we substitute (6) into (14) then we obtain that

$$\widehat{\beta}_{IV} \rightarrow \psi_2 + \frac{\psi_1 \text{cov}(Z, s(Z))}{\text{var}(Z)} + \frac{\rho \text{cov}(Z, Zs(Z))}{\text{var}(Z)} \tag{31}$$

Since  $s(z)$  is increasing, the second term on the right-hand side is positive. By comparing (31) to (30), it then follows that if the third term on the right-hand side of (31) is larger than or equal to  $\rho s_0$  then asymptotically  $\widehat{\beta}_{IV} > \beta$ . Since  $\rho > 0$ , this requires that

$$\frac{\text{cov}(Z, Zs(Z))}{\text{var}(Z)} \geq s_0$$

which is equivalent to  $\text{cov}(Z, Z[s(Z) - s_0]) \geq 0$ . The latter can be written as

$$\mathbb{E}(Z^2 s(Z)) - \bar{z}\mathbb{E}(Zs(Z)) - s_0\mathbb{E}(Z^2) + s_0\bar{z}^2 \geq 0$$

which is equivalent to

$$\mathbb{E}[Z(Z - \bar{z})(s(Z) - s_0)] \geq 0 \tag{32}$$

Earlier in this proof we derived that  $s(z)$  is increasing and that  $s(\bar{z}) = s_0$ . These two facts imply that  $s(z) \geq s_0 \iff z > \bar{z}$ . Consequently, (32) is true, and this completes the proof.  $\square$

Note that we can rewrite equation (31) as

$$\widehat{\beta}_{IV} \rightarrow \beta + \frac{\psi_1 \text{cov}(Z, s(Z))}{\text{var}(Z)} + \frac{\rho \mathbb{E}[Z(Z - \bar{z})(s(Z) - s_0)]}{\text{var}(Z)} \tag{33}$$

### Appendix 3 Proof of Proposition 5

In terms of our model and notation, Athey and Levin (2001) prove that the assumptions on the function  $f$  imply that the optimal effort level  $s(z)$  in absence of the ER increases in  $z$ .

The numerator on the right-hand side of (15),  $\text{cov}(\mathbb{E}(f(Y, s(Z)) - f(Y, s_0)|Z), Z)$  determines the sign of the asymptotic bias. Following the line of reasoning that leads up to (32) in Appendix 2, it is easy to see that this numerator can be written as

$$\mathbb{E}[(Z - \bar{z}) \cdot \mathbb{E}(f(Y, s(Z)) - f(Y, s_0)|Z)] \quad (34)$$

Since  $s$  increases,  $s(z) \geq s(\bar{z}) \iff z \geq \bar{z}$ . Now  $f$  is non-negative, increasing and supermodular, so for any given  $y$  there holds that  $f(y, s(z)) \geq f(y, s(\bar{z}))$  if and only if  $z \geq \bar{z}$ . By assumption,  $s(\bar{z}) = s_0$ , so for any given  $y$  there holds that  $f(y, s(z)) \geq f(y, s_0)$  if and only if  $z \geq \bar{z}$ . The expectation over  $Y|Z = z$  of  $f(Y, s(z)) - f(Y, s_0)$  then also has the property that it is non-negative if and only if  $z \geq \bar{z}$ . Consequently, (34) is non-negative, and this completes the proof.  $\square$

### Appendix 4 Proof of Proposition 7

Equations (22) and (23) imply that

$$\widehat{\beta}_{IV} - \beta \rightarrow \frac{\rho}{c_0} \left[ \psi_1^* - \rho^* \bar{z} + \rho^* \frac{\text{cov}(Z, Z^2)}{\text{var}(Z)} \right] \quad (35)$$

Clearly, if  $\rho = 0$  or if  $\rho^* = \psi_1^* = 0$  then the asymptotic bias is zero.

Now consider the other cases, so  $\rho \neq 0$ , and, in addition,  $\rho^* \neq 0$  and/or  $\psi_1^* > 0$ . The proof in Appendix 1 of Proposition 3 demonstrates that if  $\rho^* > 0$  then

$$\psi_1^* - \rho^* \bar{z} + \rho^* \frac{\text{cov}(Z, Z^2)}{\text{var}(Z)} \quad (36)$$

is positive. If  $\rho^* = 0$  then  $\psi_1^* > 0$  and again (36) is positive. Similarly, the proof in Appendix 1 demonstrates that if  $\rho^* < 0$  then the expression in (36) divided by  $\rho^*$  is negative, taking into account the assumptions that then  $\psi_1^* > 0$  and  $\Pr(Z < -\psi_1^*/\rho^*) = 1$ . So if  $\rho^* < 0$  then again (36) is positive. By comparing (35) to (36) it follows that the asymptotic bias has the same sign as  $\rho$ . Note that we do not require assumptions on the range of values of  $\psi_0^*$  and  $\psi_2^*$  and the support of  $Y|Z$  in terms of  $-\psi_1^*/\rho^*$  that correspond to those made concerning the parameters of  $f$ .  $\square$

## References

- Abbring, J.H. and G.J. van den Berg (2003), “The non-parametric identification of treatment effects in duration models”, *Econometrica*, 71, 1491–1517.
- Abbring, J.H. and G.J. van den Berg (2005), “Social experiments and instrumental variables with duration outcomes”, Working paper, Free University Amsterdam, Amsterdam.
- Angrist, J.D., G.W. Imbens, and D.B. Rubin (1996), “Identification of causal effects using instrumental variables”, *Journal of the American Statistical Association*, 91, 444–455.
- Athey, S. and J. Levin (2001), “The value of information in monotone decision problems”, Working paper, Stanford University, Stanford.
- Blundell, R. and T. MaCurdy (1999), “Labor supply”, in O. Ashenfelter and D. Card, editors, *Handbook of Labor Economics, Volume III*, North-Holland, Amsterdam.
- Epstein, S. (1996), *Impure Science: AIDS, Activism, and the Politics of Knowledge*, University of California Press, Berkeley.
- Frontini, M. and A. Tagliani (1995), “Maximum entropy in the finite Hausdorff moment problem”, Working paper, Politecnico di Milano, Milan.
- Gollier, C. (2001), *The Economics of Risk and Time*, MIT Press, Cambridge.
- Hansen, L.P. and T.J. Sargent (1980), “Dynamic linear rational expectations models”, *Journal of Economic Dynamics and Control*, 2, 7–46.
- Heckman, J.J., R.J. LaLonde, and J.A. Smith (1999), “The economics and econometrics of active labor market programs”, in O. Ashenfelter and D. Card, editors, *Handbook of Labor Economics, Volume III*, North-Holland, Amsterdam.
- Heckman, J.J. and S. Navarro (2005), “Dynamic discrete choice and dynamic treatment effects”, Working paper, University of Chicago, Chicago.
- Imbens, G.W. and J.D. Angrist (1994), “Identification and estimation of local average treatment effects”, *Econometrica*, 62, 467–475.
- Malani, A. (2006a), “Patient enrollment in medical trials: selection bias in a randomized experiment”, Working paper, University of Virginia, Charlottesville.
- Malani, A. (2006b), “Identifying placebo effects with data from clinical trials”, *Journal of Political Economy*, 114, forthcoming.
- Persico, N. (2000), “Information acquisition in auctions”, *Econometrica*, 68, 135–148.

- Rosenzweig, M.R. and K.I. Wolpin (2000), “Natural “natural experiments” in economics”, *Journal of Economic Literature*, 38, 827–874.
- Schuklenk, U. (2003), “AIDS: bioethics and public policy”, *New Review of Bioethics*, 1, 127–144.
- Shohat, J.A. and J.D. Tamarkin (1943), *The Problem of Moments*, American Mathematical Society, New York.
- Stock, J.H., J.H. Wright, and M. Yogo (2002), “A survey of weak instruments and weak identification in generalized method of moments”, *Journal of Business and Economic Statistics*, 20, 518–529.