

Instrumental Variables

Application and Limitations

Edwin P. Martens,*† Wiebe R. Pestman,† Anthonius de Boer,* Svetlana V. Belitser,*
and Olaf H. Klungel*

Abstract: To correct for confounding, the method of instrumental variables (IV) has been proposed. Its use in medical literature is still rather limited because of unfamiliarity or inapplicability. By introducing the method in a nontechnical way, we show that IV in a linear model is quite easy to understand and easy to apply once an appropriate instrumental variable has been identified. We also point out some limitations of the IV estimator when the instrumental variable is only weakly correlated with the exposure. The IV estimator will be imprecise (large standard error), biased when sample size is small, and biased in large samples when one of the assumptions is only slightly violated. For these reasons, it is advised to use an IV that is strongly correlated with exposure. However, we further show that under the assumptions required for the validity of the method, this correlation between IV and exposure is limited. Its maximum is low when confounding is strong, such as in case of confounding by indication. Finally, we show that in a study in which strong confounding is to be expected and an IV has been used that is moderately or strongly related to exposure, it is likely that the assumptions of IV are violated, resulting in a biased effect estimate. We conclude that instrumental variables can be useful in case of moderate confounding but are less useful when strong confounding exists, because strong instruments cannot be found and assumptions will be easily violated.

(*Epidemiology* 2006;17: 260–267)

In medical research, randomized, controlled trials (RCTs) remain the gold standard in assessing the effect of one variable of interest, often a specified treatment. Nevertheless, observational studies are often used in estimating such an effect.¹ In epidemiologic as well as sociologic and economic

research, observational studies are the standard for exploring causal relationships between an exposure and an outcome variable. The main problem of estimating the effect in such studies is the potential bias resulting from confounding between the variable of interest and alternative explanations for the outcome (confounders). Traditionally, standard methods such as stratification, matching, and multiple regression techniques have been used to deal with confounding. In the epidemiologic literature, some other methods have been proposed^{2,3} of which the method of propensity scores is best known.⁴ In most of these methods, adjustment can be made only for observed confounders.

A method that has the potential to adjust for all confounders, whether observed or not, is the method of instrumental variables (IV). This method is well known in economics and econometrics as the estimation of simultaneous regression equations⁵ and is also referred to as structural equations and two-stage least squares. This method has a long tradition in economic literature, but has entered more recently into the medical research literature with increased focus on the validity of the instruments. Introductory texts on instrumental variables can be found in Greenland⁶ and Zohoori and Savitz.⁷

One of the earliest applications of IV in the medical field is probably the research of Permutt and Hebel,⁸ who estimated the effect of smoking of pregnant women on their child's birth weight, using an encouragement to stop smoking as the instrumental variable. More recent examples can be found in Beck et al,⁹ Brooks et al,¹⁰ Earle et al,¹¹ Hadley et al,¹² Leigh and Schembri,¹³ McClellan,¹⁴ and McIntosh.¹⁵ However, it has been argued that the application of this method is limited because of its strong assumptions, making it difficult in practice to find a suitable instrumental variable.¹⁶

The objectives of this article are first to introduce the application of the method of IV in epidemiology in a nontechnical way and second, to show the limitations of this method, from which it follows that IV is less useful for solving large confounding problems such as confounding by indication.

A SIMPLE LINEAR INSTRUMENTAL VARIABLES MODEL

In an RCT, the main purpose is to estimate the effect of one explanatory factor (the treatment) on an outcome variable. Because treatments have been randomly assigned to individuals, the treatment variable is in general independent

Submitted 8 February 2005; accepted 16 November 2005.

From the *Department of Pharmacoepidemiology and Pharmacotherapy, Utrecht Institute of Pharmaceutical Sciences (UIPS), Utrecht University, Utrecht, The Netherlands; and the †Centre for Biostatistics, Utrecht University, Utrecht, The Netherlands.

Supported by the Utrecht Institute of Pharmaceutical Sciences (UIPS), Utrecht University, Utrecht, The Netherlands.

Correspondence: Olaf H. Klungel, Department of Pharmacoepidemiology and Pharmacotherapy, Utrecht Institute of Pharmaceutical Sciences (UIPS), Utrecht University, Sorbonnelaan 16, 3584 CA Utrecht, The Netherlands. E-mail: o.h.klungel@pharm.uu.nl.

Copyright © 2006 by Lippincott Williams & Wilkins

ISSN: 1044-3983/06/1703-0260

DOI: 10.1097/01.ede.0000215160.88317.cb

of other explanatory factors. In case of a continuous outcome and a linear model, this randomization procedure allows one to estimate the treatment effect by means of ordinary least squares with a well-known unbiased estimator (see, for instance, Pestman¹⁷). In observational studies, on the other hand, one has no control over this explanatory factor (further denoted as *exposure*) so that ordinary least squares as an estimation method will generally be biased because of the existence of unmeasured *confounders*. For example, one cannot directly estimate the effect of cigarette smoking on health without considering confounding factors such as age and socioeconomic position.

One way to adjust for all possible confounding factors, whether observed or not, is to make use of an instrumental variable. The idea is that the causal effect of exposure on outcome can be captured by using the relationship between the exposure and another variable, the instrumental variable. How this variable can be selected and which conditions have to be fulfilled is discussed subsequently. First, we illustrate the model and its estimator.

The Instrumental Variables Model and Its Estimator

A simple linear model for IV estimation consists of 2 equations:

$$Y = \alpha + \beta X + E \tag{1}$$

$$X = \gamma + \delta Z + F \tag{2}$$

where Y is the outcome variable, X is the exposure, Z is the instrumental variable, and E and F are errors. In this set of structural equations, the variable X is *endogenous*, which means that it is explained by other variables in the model, in this case the instrumental variable Z . Z is supposed to be linearly related to X and *exogenous*, ie, explained by variables outside the model. For simplicity, we restrict ourselves to one instrumental variable, 2 equations, and no other explaining variables. Under conditions further outlined in the next section, it can be proved that equation (3) presents an asymptotically unbiased estimate of the effect of X on Y ¹⁸:

$$\hat{\beta}_{iv} = \frac{\frac{1}{n-1} \sum_{i=1}^n (z_i - \bar{z})(y_i - \bar{y})}{\frac{1}{n-1} \sum_{i=1}^n (z_i - \bar{z})(x_i - \bar{x})} = \frac{\hat{\sigma}_{Z,Y}}{\hat{\sigma}_{Z,X}} \tag{3}$$

where $\hat{\sigma}_{Z,Y}$ is the sample covariance of Z and Y and $\hat{\sigma}_{Z,X}$ is the sample covariance of Z and X . It is more convenient to express the IV estimator in terms of 2 ordinary least squares estimators:

$$\hat{\beta}_{iv} = \frac{\hat{\sigma}_{Z,Y}}{\hat{\sigma}_{Z,X}} = \frac{\hat{\sigma}_{Z,Y}/\hat{\sigma}_Z^2}{\hat{\sigma}_{Z,X}/\hat{\sigma}_Z^2} = \frac{\hat{\beta}_{ols(Z \rightarrow Y)}}{\hat{\beta}_{ols(Z \rightarrow X)}} \tag{4}$$

The numerator equals the effect of the instrumental variable on the outcome, whereas in the denominator, the effect of the IV on the exposure is given. In case of a dichotomous IV, the numerator equals simply the difference in mean outcome between $Z = 0$ and $Z = 1$ and the denominator equals the difference in mean exposure. When the outcome and exposure variable are also dichotomous and linearity is still assumed, this model is known as a linear probability model. In that case, the IV estimator presented here can be simply expressed as probabilities¹⁸:

$$\hat{\beta}_{iv} = \frac{P(Y = 1|Z = 1) - P(Y = 1|Z = 0)}{P(X = 1|Z = 1) - P(X = 1|Z = 0)} \tag{5}$$

where $P(Y = 1|Z = 1) - P(Y = 1|Z = 0)$ equals the risk difference of an event between $Z = 1$ and $Z = 0$.

How to Obtain a Valid Instrumental Variable

One can imagine that a method that claims to adjust for all possible confounders without randomization of treatments puts high requirements on the IV to be used for estimation. When this method is applied, 3 important assumptions have been made. The first assumption is the existence of at least some correlation between the IV and the exposure, because otherwise, equation (2) would be useless and the denominator of equation (4) would be equal to zero. In addition to this formal condition, it is important that this correlation should not be too small (see “Implications of Weak Instruments”).

The second assumption is that the relationship between the instrumental variable and the exposure is not confounded by other variables so that equation (2) is estimated without bias. This is the same as saying that the correlation between the IV and the error F must be equal to zero. One way to achieve this is to use as IV a variable that is *controlled by the researcher*. An example can be found in Permutt and Hebel,⁸ in which a randomized encouragement to stop smoking was used as the IV to estimate the effect of smoking by pregnant women on child’s birth weight. The researchers used 2 encouragement regimes, an encouragement to stop smoking versus no encouragement, randomly assigned to pregnant smoking women. Alternatively, in some situations, a *natural randomization process* can be used as the IV. An example, also known as Mendelian randomization, can be found in genetics in which alleles are considered to be allocated at random in offspring with the same parents.^{19,20} In a study on the causality between low serum cholesterol and cancer, a genetic determinant of serum cholesterol was used as the instrumental variable.^{21,22} When neither an active randomization nor a natural randomization is feasible to obtain an IV, the only possibility is to select an IV on *theoretical grounds*, assuming and reasoning that the relationship between the IV and the exposure can be estimated without bias. Such an example can be found in Leigh and Schembri¹³ in which the observed cigarette price per region was used as the IV in a study on the relationship between smoking and health. The authors argued that there was no bias in estimating the relationship between cigarette price and smoking because the price elasticities in their study (the percentage change in

number of cigarettes smoked related to the percentage change in cigarette price) matched the price elasticities mentioned in the literature.

The third assumption for an IV is most crucial and states that there should be no correlation between the IV and the error E (further referred to as *the main assumption*). This means that the instrumental variable should influence the outcome neither directly nor indirectly by its relationship with other variables. Whether this assumption is valid can be argued only theoretically, and cannot be tested empirically.

These 3 assumptions can be summarized as follows:

1. $\rho_{Z,X} \neq 0$, no zero-correlation between IV and exposure;
2. $\rho_{Z,F} = 0$, no correlation between IV and other factors explaining X (error F); and
3. $\rho_{Z,E} = 0$, no correlation between IV and other factors explaining Y (error E), main assumption.

It should be noted that confounders of the X - Y relation are not explicitly mentioned in these assumptions and that these confounders are part of both errors E and F . In the special case that $\rho_{E,F} = 1$, the assumption could be formulated by referring to confounders only.⁶

Numeric Example of Instrumental Variable Application

As an example of IV estimation, we use the research of Permutt and Hebel.⁸ Here the effect of smoking (X) by pregnant women on child's birth weight (Y) was studied. The instrumental variable (Z) was the randomization procedure used to assign women to an encouragement program to stop smoking, which fulfills the second assumption. To apply IV estimation, first the intention-to-treat estimator $\beta_{ols(Z \rightarrow Y)}$ needs to be calculated. In case of a dichotomous IV, this simply equals the difference in mean birth weight between women who were encouraged to stop smoking and women who were not ($\beta_{ols(Z \rightarrow Y)} = 98$ g). Next, we calculate the difference between encouragement groups in the fraction of women who stopped smoking ($\beta_{ols(Z \rightarrow X)} = 0.43 - 0.20 = 0.23$). The ratio equals the IV estimator $= 98 / (0.43 - 0.20) = 430$ g, indicating that stopping smoking raises average birth weight by 430 g. Figure 1 illustrates this calculation, in which "actually stopped smoking" is denoted as $X = 1$ and "continued to smoke" as $X = 0$.

The encouragement-smoking relationship and the encouragement-birth weight relationship are represented by the solid lines in the lower and upper panel, respectively. Under the assumptions of IV estimation, the effect of smoking on birth weight is known only when smoking is changed from 0.43 to 0.20, in which in fact interest is in a change from $X = 0$ to $X = 1$. Extending this difference to a difference from 0 to 1, indicated by the dotted line in the lower panel, and using the relationship between Z and Y in the upper panel, the intention-to-treat estimator of 98 g is "extended" to become the IV estimator of 430 g. Reminding that our second assumption has been fulfilled by randomization, the possible bias of the IV estimator mainly depends on the assumption that there should be no effect from encouragement on child's birth weight other than by means of changing smoking behavior. Such an effect cannot be ruled out completely, for

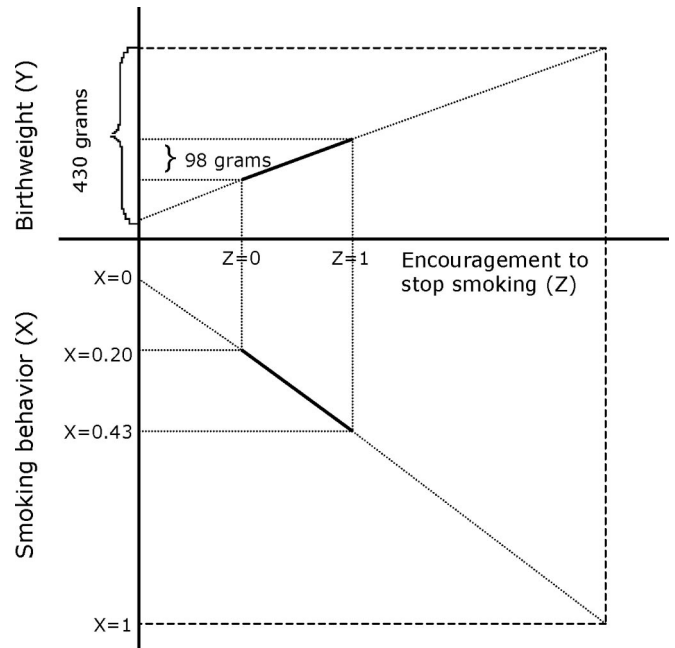


FIGURE 1. The instrumental variable estimator in the study of Permutt and Hebel.⁸

instance, because women who were encouraged to stop smoking could become also more motivated to change other health-related behavior as well (for instance, nutrition). Birth weight will then be influenced by encouragement independently of smoking, which will lead to an overestimation of the effect of stopping smoking.

IMPLICATIONS OF WEAK INSTRUMENTS

In the previous sections, the method and application of instrumental variables in a linear model were introduced in a nontechnical way. Here we focus on the implications when the correlation between the instrumental variable and the exposure is small or when the instrument is weak. We refer to this correlation as $\rho_{Z,X}$.

Large Standard Error

A weak instrument means that the denominator in equation (4) is small. The smaller this covariance, the more sensitive the IV estimate will be to small changes. This sensitivity is mentioned by various authors^{16,23} and can be deduced from the formula for the standard error:

$$\hat{\sigma}\beta_{iv} = \frac{\sigma_Z \sigma_E}{\sigma_{Z,X}} \quad (6)$$

where σ_Z is the standard deviation of Z , σ_E is the standard deviation of E , and $\sigma_{Z,X}$ is the covariance of Z and X . This covariance in the denominator behaves as a multiplier, which means that a small covariance (and hence a small correlation) will lead to a large standard error. In Figure 1, this sensitivity is reflected by the fact that the slope estimate in the lower

panel becomes less reliable when the difference in X between $Z = 0$ and $Z = 1$ becomes smaller.

Bias When Sample Size Is Small

An important characteristic of an estimator is that it should equal on average the true value (*unbiasedness*). Assuming that the assumptions of IV are not violated, the IV estimator is only *asymptotically* unbiased, meaning that on average bias will exist when the estimator $\hat{\beta}_{iv}$ is used in smaller samples. This bias appears because the relationship between the instrumental variable and the exposure is in general unknown and has to be estimated by equation (2). As is usual in regression, overfitting generates a bias that depends on both the sample size and the correlation between the IV and the exposure. With moderate sample size and a weak instrument, this bias can become substantial.²⁴ It can be shown that this bias will be in the direction of the ordinary least squares estimator $\hat{\beta}_{ols}$ calculated in the simple linear regression of outcome on exposure.^{23,25} Information on the magnitude of the small sample bias is contained in the F -statistic of the regression in equation (2), which can be expressed as

$$F = \frac{\hat{\rho}_{Z,X}^2 (n - 2)}{1 - \hat{\rho}_{Z,X}^2} \quad (7)$$

An F -value not far from 1 indicates a large small sample bias, whereas a value of 10 seems to be sufficient for the bias to be negligible.¹⁶ For example, in a sample of 250 independent observations, the correlation between Z and X should be at least 0.20 to reach an F -value of 10. Another solution to deal with possible small sample bias is to use other IV estimators.^{16,26}

Bias When the Main Assumption Is Only Slightly Violated

Every violation of the main assumption of IV will naturally result in a biased estimator. More interesting is that only a small violation of this assumption will result in a large bias in case of a weak instrument because of its multiplicative effect in the estimator. Bound et al²³ expressed this bias in infinitely large samples (inconsistency) as a relative measure compared with the bias in the ordinary least squares estimator

$$\frac{\lim \hat{\beta}_{iv} - \beta}{\lim \hat{\beta}_{ols} - \beta} = \frac{\rho_{Z,E}/\rho_{X,E}}{\rho_{Z,X}} \quad (8)$$

where *lim* is the limit as sample size increases. From this formula, it can be seen that even a small correlation between the instrumental variable and the error ($\rho_{Z,E}$ in the numerator) will produce a large inconsistency in the IV estimate relative to the ordinary least squares estimate when the instrument is weak, ie, when $\rho_{Z,X}$ is small. Thus, when Z has some small direct effect on Y , or an indirect effect other than through X , the IV estimate will be increasingly biased when the instrument becomes weaker, even in very large samples.

It can be concluded that a small correlation between the IV and the exposure can be a threat for the validity of the IV method, mainly in combination with a small sample or a

possible violation of the main assumption. Although known from the literature, this aspect is often overlooked.

A LIMIT ON THE STRENGTH OF INSTRUMENTS

From the last section, it follows that the correlation between a possible instrumental variable and exposure (the strength of the IV $\rho_{Z,X}$) has to be as strong as possible, which also intuitively makes sense. However, in practice, it is often difficult to obtain an IV that is strongly related to exposure. One reason can be found in the existence of an upper bound on this correlation, which depends on the amount of confounding (indicated by $\rho_{X,E}$), the correlation between the errors in the model ($\rho_{E,F}$), and the degree of violation of the main assumption ($\rho_{Z,E}$). We further explore the relationship between these correlations and distinguish between a situation in which the main assumption is fulfilled and one in which it is not.

When the Main Assumption Has Been Fulfilled

In case the main assumption of IV has been fulfilled, which means that the IV changes the outcome only through its relationship with the exposure, it can be shown that

$$|\rho_{Z,X}| = \sqrt{1 - \frac{\rho_{X,E}^2}{\rho_{E,F}^2}} \quad (9)$$

of which the proof is given in Appendix A. Equation (9) indicates that there is a maximum on the strength of the instrumental variable and that this maximum decreases when the amount of confounding increases. In case of considerable confounding, the maximum correlation between IV and exposure will be quite low. This relationship between the correlations is illustrated in Figure 2.

The relation between the strength of the IV $\rho_{Z,X}$ and the amount of confounding $\rho_{X,E}$ is illustrated by curves representing various levels of the correlation between the errors $\rho_{E,F}$. It can be seen that the maximum correlation between the potential instrumental variable and exposure becomes smaller when the amount of confounding becomes larger. When, for example, there is considerable confounding by indication ($\rho_{X,E} = 0.8$), the maximum strength of the IV is 0.6. Probably this maximum will be even lower because the correlation between the errors will generally be less than 1.0. When, for instance, $\rho_{E,F} = 0.85$, this maximum drops to only 0.34. Of the 3 correlations presented in equation (9) and Figure 2, the correlation between the errors is most difficult to understand. For the main message, however, its existence is not essential, as is illustrated in Figure 3 using vectors.

In Figure 3A, the angle between X and E is close to 90°, meaning that their correlation is small (small confounding). Because Z has to be uncorrelated with E according to the third IV assumption (perpendicular), the angle between X and Z will be automatically small, indicating a strong IV. In contrast, Figure 3B shows that a large confounding problem (small angle between X and E) implies a weak instrument (large angle and small correlation between X and Z). The tradeoff between these correlations is an important characteristic of IV estimation. (Note that we simplified the figure by

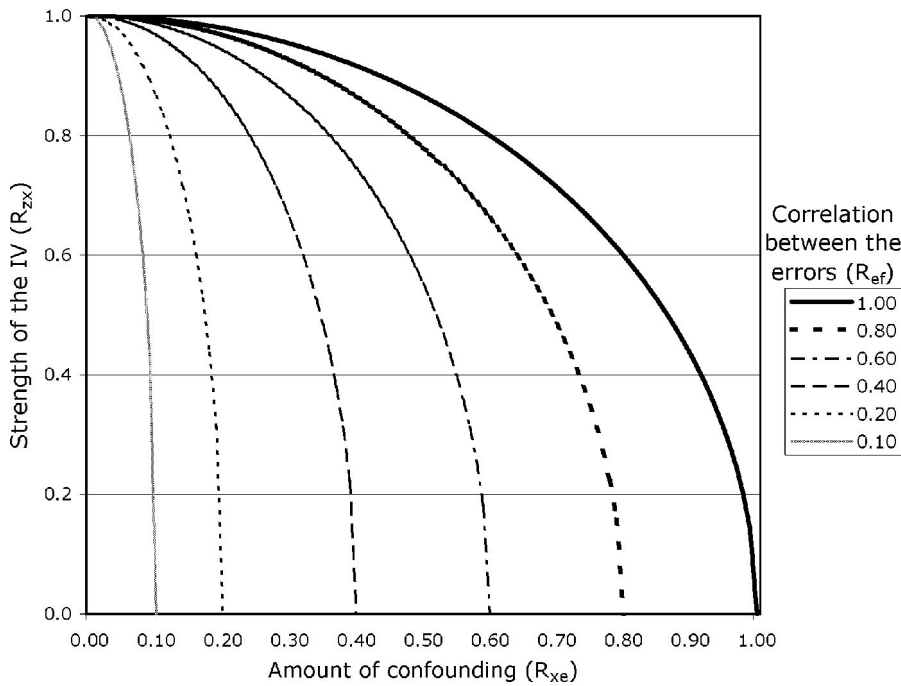


FIGURE 2. Relationship between strength of an instrumental variable ($\rho_{Z,X}$) and amount of confounding ($\rho_{X,E}$) for different error correlation levels ($\rho_{E,F}$) when main assumption has been fulfilled ($\rho_{Z,E} = 0$).

choosing Z in the same plane as X and Y to remove $\rho_{E,F}$ from the figure because it equals its maximum of 1.0. See Appendix B for the situation in which Z is not in this plane.)

As has been said, the correlation between the errors $\rho_{E,F}$ also plays a role. To better understand its meaning, we give 2 examples. In Permutt and Hebel,⁸ it is likely that this correlation will be small. Other reasons for birth weight variation besides smoking include socioeconomic conditions, inadequate nutrition, abuse, genetic factors, ethnic factors, physical work conditions, and chronic diseases. Because these explanatory factors for birth weight will be only partly overlapping with the reasons for *noncompliance*, ie, to continue smoking while encouraged to stop, $\rho_{E,F}$ is expected to be small. When, on the other hand, this correlation approaches 1, it means that the set of variables accounting for the unexplained variation in the outcome Y (error E) is strongly correlated with the unexplained instrumental variance (error F). An example of such a large correlation is a case of strong confounding by indication, in which unob-

served health problems are the main reason for getting an illness and also for receiving preventive treatment. That causes variables E and F to be strongly correlated and the maximum strength of the IV to be relatively small (see the right side of Fig. 2).

When the Main Assumption Has Not Been Fulfilled

When the main assumption has not been (completely) fulfilled, the correlation between Z and E is not equal to 0. Because the correlation between the errors plays a minor role, this correlation has been set to its maximum value of 1. In that case, the next inequality holds:

$$\rho_{Z,X} \leq |\rho_{Z,E}| |\rho_{X,E}| + \sqrt{1 - \rho_{Z,E}^2} \sqrt{1 - \rho_{X,E}^2} \quad (10)$$

Like equation (9), this expression states that in case of considerable confounding, the strength of the instrumental variable is bound to a relatively small value. It further states that a tradeoff exists between $\rho_{Z,X}$ and $\rho_{Z,E}$: given a certain degree of confounding, the strength of the IV can be enlarged by relaxing the main assumption. In practice, this means that when IV is applied to a situation in which a considerable amount of confounding is to be expected and a very strong instrument has been found, it is very likely that the main assumption has been violated.

The Effect on Bias

The limit of the correlation between exposure and instrumental variable has an indirect effect on the bias, because the correlation to be found in practice will be low. This has several disadvantages that can be illustrated using some previous numeric examples. Suppose we deal with

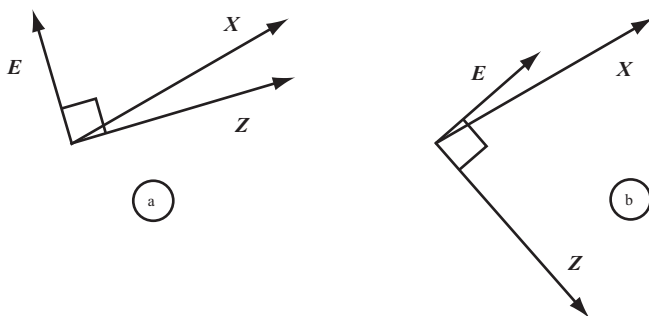


FIGURE 3. Relationship among X , Z , and E expressed in vectors.

strong confounding by indication, say $\rho_{X,E} = 0.80$. As has been argued before, this will naturally imply a strong but imperfect correlation between the errors, say $\rho_{E,F} = 0.85$. In that case, the limit of the correlation between exposure and IV will be $\rho_{Z,X} = 0.34$. Restricting ourselves to instrumental variables that fulfill the main assumption ($\rho_{Z,E} = 0$), it will be practically impossible to find an IV that possesses the characteristic of being maximally correlated with exposure, which implies that this correlation will be lower than 0.34, for instance 0.20. With such a small correlation, the effect on the bias will be substantial when sample size falls below 250 observations. Because we cannot be sure that the main assumption has been fulfilled, care must be taken even with larger samples sizes.

DISCUSSION

We have focused on the method of instrumental variables for its ability to adjust for confounding in nonrandomized studies. We have explained the method and its application in a linear model and focused on the correlation between the IV and the exposure. When this correlation is very small, this method will lead to an increased standard error of the estimate, a considerable bias when sample size is small, and a bias even in large samples when the main assumption is only slightly violated. Furthermore, we demonstrated the existence of an upper bound on the correlation between the IV and the exposure. This upper bound is not a practical limitation when confounding is small or moderate because the maximum strength of the IV is still very high. When, on the other hand, considerable confounding by indication exists, the maximum correlation between any potential IV and the exposure will be quite low, resulting possibly in a fairly weak instrument to fulfill the main assumption. Because of a tradeoff between violation of this main assumption and the strength of the IV, the presence of considerable confounding and a strong instrument will probably indicate a violation of the main assumption and thus a biased estimate.

This article serves as an introduction on the method of instrumental variables demonstrating its merits and limitations. Complexities such as more equations, more instruments, the inclusion of covariates, and nonlinearity of the model have been left out. More equations could be added with more than 2 endogenous variables, although it is unlikely to be useful in epidemiology when estimating an exposure (treatment) effect. In equation (2), multiple instruments could be used; this extension does not change the basic ideas behind this method.²⁷ An advantage of more than one instrumental variable is that a test on the exogeneity of the instruments is possible.¹⁶ Another extension is the inclusion of measured covariates in both equations.²⁷

We limited the model to linear regression, assuming that the outcome and the exposure are both continuous variables, while in medical research, dichotomous outcomes or exposures are more common. The main reason for this choice is simplicity: the application and implications can be more easily presented in a linear framework. A dichotomous outcome or dichotomous exposure can easily fit into this model when linearity is assumed using a *linear probability model*.

Although less known, the results from this model are practically indistinguishable from logistic and probit regression analyses as long as the estimated probabilities range between 0.2 and 0.8.^{28,29} When risk ratios or log odds are to be analyzed, like in logistic regression analysis, the presented IV estimator cannot be used and more complex IV estimators are required. We refer to the literature for IV estimation in such cases or in nonlinear models in general.^{6,30,31} The limitations when instruments are weak, and the impossibility of finding strong instruments in the presence of strong confounding, apply in a similar way.

When assessing the validity of study results, investigators should report both the correlation between IV and exposure (or difference in means) and the F-value resulting from equation (2) and given in equation (7). When either of these is small, instrumental variables will not produce unbiased and reasonably precise estimates of exposure effect. Furthermore, it should be made clear whether the IV is randomized by the researcher, randomized by nature, or is simply an observed variable. In the latter case, evidence should be given that the various categories of the instrumental variable have similar distributions on important characteristics. Additionally, the assumption that the IV determines outcome only by means of exposure is crucial. Because this cannot be checked, it should be argued theoretically that a direct or indirect relationship between the IV and the outcome is negligible. Finally, in a study in which considerable confounding can be expected (eg, strong confounding by indication), one should be aware that the existence of a very strong instrument within the IV assumptions is impossible. Whether the instrument is sufficiently correlated with exposure depends on the number of observations and the plausibility of the main assumption.

We conclude that the method of IV can be useful in case of moderate confounding but is less useful when strong confounding (by indication) exists, because strong instruments cannot be found and assumptions will be easily violated.

REFERENCES

1. Concato J, Shah N, Horwitz RI. Randomized, controlled trials, observational studies, and the hierarchy of research designs. *N Engl J Med*. 2000;342:1887–1892.
2. McMahon AD. Approaches to combat with confounding by indication in observational studies of intended drug effects. *Pharmacoepidemiol Drug Saf*. 2003;12:551–558.
3. Klungel OH, Martens EP, Psaty BM, et al. Methods to assess intended effects of drug treatment in observational studies are reviewed. *J Clin Epidemiol*. 2004;57:1223–1231.
4. Rosenbaum PR, Rubin DB. The central role of the propensity score in observational studies for causal effects. *Biometrika*. 1983;70:41–55.
5. Theil H. *Principles of Econometrics*. New York: Wiley; 1971.
6. Greenland S. An introduction to instrumental variables for epidemiologists. *Int J Epidemiol*. 2000;29:722–729.
7. Zohoori N, Savitz DA. Econometric approaches to epidemiologic data: relating endogeneity and unobserved heterogeneity to confounding. *Ann Epidemiol*. 1997;7:251–257. Erratum in *Ann Epidemiol*. 1997;7:431.
8. Permutt TH, Hebel JR. Simultaneous-equation estimation in a clinical trial of the effect of smoking on birth weight. *Biometrics*. 1989;45:619–622.
9. Beck CA, Penrod J, Gyorkos TW, et al. Does aggressive care following acute myocardial infarction reduce mortality? Analysis with instrumental variables to compare effectiveness in Canadian and United States patient populations. *Health Serv Res*. 2003;38:1423–1440.

10. Brooks JM, Chrischilles EA, Scott SD, et al. Was breast conserving surgery underutilized for early stage breast cancer? Instrumental variables evidence for stage II patients from Iowa. *Health Serv Res.* 2003; 38:1385–1402. Erratum in *Health Serv Res.* 2004;39:693.
11. Earle CC, Tsai JS, Gelber RD, et al. Effectiveness of chemotherapy for advanced lung cancer in the elderly: instrumental variable and propensity analysis. *J Clin Oncol.* 2001;19:1064–1070.
12. Hadley J, Polsky D, Mandelblatt JS, et al. An exploratory instrumental variable analysis of the outcomes of localized breast cancer treatments in a medicare population. *Health Econ.* 2003;12:171–186.
13. Leigh JP, Schembri M. Instrumental variables technique: cigarette price provided better estimate of effects of smoking on SF-12. *J Clin Epidemiol.* 2004;57:284–293.
14. McClellan M, McNeil BJ, Newhouse JP. Does more intensive treatment of acute myocardial infarction in the elderly reduce mortality? Analysis using instrumental variables. *JAMA.* 1994;272:859–866.
15. McIntosh MW. Instrumental variables when evaluating screening trials: estimating the benefit of detecting cancer by screening. *Stat Med.* 1999;18:2775–2794.
16. Staiger D, Stock JH. Instrumental variables regression with weak instruments. *Econometrica.* 1997;65:557–586.
17. Pestman WR. *Mathematical Statistics.* Walter de Gruyter, 1998.
18. Angrist JD, Imbens GW, Rubin DB. Identification of causal effects using instrumental variables. *JASA.* 1996;91:444–455.
19. Thomas DC, Conti DV. Commentary: the concept of ‘mendelian randomization.’ *Int J Epidemiol.* 2004;33:21–25.
20. Minelli C, Thompson JR, Tobin MD, et al. An integrated approach to the meta-analysis of genetic association studies using mendelian randomization. *Am J Epidemiol.* 2004;160:445–452.
21. Katan MB. Apolipoprotein E isoforms, serum cholesterol, and cancer. *Lancet.* 1986;1:507–508.
22. Smith GD, Ebrahim S. Mendelian randomization: prospects, potentials, and limitations. *Int J Epidemiol.* 2004;33:30–42.
23. Bound J, Jaeger DA, Baker RM. Problems with instrumental variables estimation when the correlation between the instruments and the endogenous explanatory variable is weak. *JASA.* 1995;90:443–450.
24. Sawa T. The exact sampling distribution of ordinary least squares and two-stage least squares estimators. *J Am Stat Assoc.* 1969;64:923–937.
25. Nelson CR, Startz R. Some further results on the exact small sample properties of the instrumental variable estimator. *Econometrica.* 1990;58:967–976.
26. Angrist JD, Krueger AB. Split sample instrumental variables. *J Bus Econ Stat.* 1995;13:225–235.
27. Angrist JD, Imbens GW. Two-stage least squares estimation of average causal effects in models with variable treatment intensity. *JASA.* 1995;90:431–442.
28. Cox DR, Snell EJ. *Analysis of Binary Data.* Chapman and Hall, 1989.
29. Cox DR, Wermuth N. A comment on the coefficient of determination for binary responses. *American Statistician.* 1992;46:1–4.
30. Bowden RJ, Turkington DA. A comparative study of instrumental variables estimators for nonlinear simultaneous models. *J Am Stat Assoc.* 1981;76:988–995.
31. Amemiya T. The nonlinear two-stage least-squares estimator. *Journal of Econometrics.* 1974;2:105–110.

APPENDIX A

Theorem 1

The correlation between Z and X, $\rho_{Z,X}$ is bound to obey the equality

$$|\rho_{Z,X}| = \sqrt{1 - \frac{\rho_{X,E}^2}{\rho_{E,F}^2}} \tag{11}$$

Proof: According to the model, one has

$$\begin{cases} Y = \alpha + \beta X + E \\ X = \gamma + \delta Z + F \end{cases}$$

with $\sigma_{Z,E} = 0$ and $\sigma_{Z,F} = 0$

It follows from this that $\sigma_{X,E} = \sigma\gamma_E + \delta\sigma_{Z,E} + \sigma_{F,E} = 0 + 0 + \sigma_{E,F} = \sigma_{E,F}$. Using this expression for $\sigma_{X,E}$, one derives that

$$\rho_{X,E} = \frac{\sigma_{X,E}}{\sigma_X \sigma_E} = \frac{\sigma_{E,F}}{\sigma_X \sigma_E} \frac{\sigma_F}{\sigma_F} = \frac{\sigma_F}{\rho_{E,F} \sigma_X} = \pm \sqrt{\frac{\sigma_F^2}{\rho_{E,F}^2 \sigma_X^2}} = \pm \sqrt{\rho_{E,F}^2 (1 - \rho_{Z,X}^2)}$$

Squaring, rearranging terms, and taking square roots will give

$$|\rho_{Z,X}| = \sqrt{1 - \frac{\rho_{X,E}^2}{\rho_{E,F}^2}}$$

which proves the theorem.

APPENDIX B

The condition $\rho_{E,F} = 1$ is equivalent to the condition that Z is in the same plane as X and E as can be seen in Figure 4. For simplicity, we assume that the expectation values of the variables X, Y, and Z are all equal to zero.

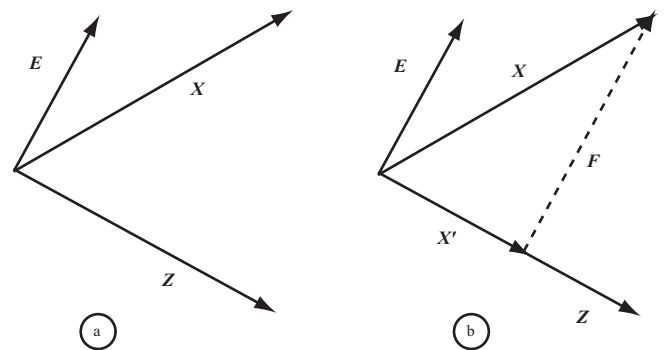


FIGURE 4. Relationship among X, Z, E, and F expressed in vectors.

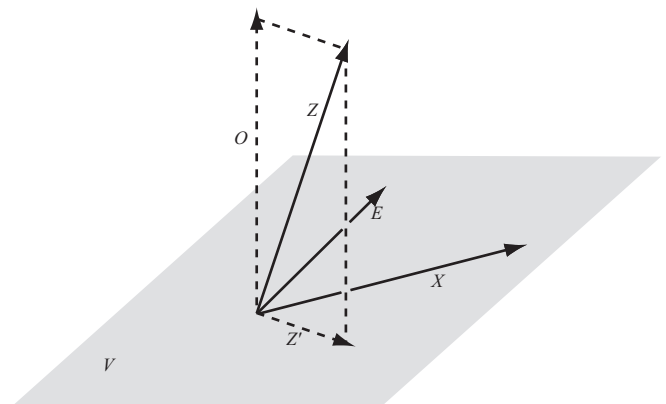


FIGURE 5. Three-dimensional picture of X, Z, E, and noise O expressed in vectors.

According to the IV condition that $\rho_{Z,E} = 0$ (these are perpendicular in panel a) and the condition that $\rho_{Z,F} = 0$, it follows from panel b that E and F necessarily point in the same or opposite direction, implying $\rho_{E,F} = 1$. In this situation, there is (up to scalar multiples) only one instrumental variable Z possible in the plane spanned by E and X . As has been argued in the text, it is not likely that this correlation equals 1. This is visualized in Figure 5 in which Z is not in the plane spanned by X and E , meaning that F , which is in the

plane spanned by X and Z and perpendicular to Z , can impossibly point in the same direction as E . Consequently, one then has $\rho_{E,F} < 1$. Here Z' is the projection of Z on the plane spanned by E and X . The vector Z can now be decomposed as $Z = Z' + O$ where Z' is in the plane spanned by E and X and where O is perpendicular to this plane. The vector O can be referred to as *noise* because it is uncorrelated to both X and Y . Note that the variable Z' is an instrumental variable itself.

The Editors of EPIDEMIOLOGY
present

***The Changing Face
of Epidemiology***

A series of symposia and commentaries
on rapidly-evolving aspects of epidemiologic research

Coming in 2006:
**How is 'big' epidemiology
changing epidemiology?**

The growth of multi-center,
multi-disciplinary, multi-investigator studies

Organized by
Jonathan Samet and Sholom Wacholder, Editors
with talks by
Bob Hoover, Muin Khoury, and Roberta Ness

2006 Congress of Epidemiology
21-24 June 2006
Seattle