

QSEP

**RESEARCH INSTITUTE FOR QUANTITATIVE
STUDIES IN ECONOMICS AND POPULATION**

**EXPLORING THE USE OF A NONPARAMETRICALLY
GENERATED INSTRUMENTAL VARIABLE IN THE
ESTIMATION OF A LINEAR PARAMETRIC EQUATION**

FRANK T. DENTON

QSEP Research Report No. 390

**EXPLORING THE USE OF A NONPARAMETRICALLY
GENERATED INSTRUMENTAL VARIABLE IN THE ESTIMATION OF A
LINEAR PARAMETRIC EQUATION**

FRANK T. DENTON

QSEP Research Report No. 390

November 2004

F.T. Denton is a QSEP Research Associate and a faculty member in the McMaster University Department of Economics.

The Research Institute for Quantitative Studies in Economics and Population (QSEP) is an interdisciplinary institute established at McMaster University to encourage and facilitate theoretical and empirical studies in economics, population, and related fields. For further information about QSEP visit our web site <http://socserv2.mcmaster.ca/qsep> or contact Secretary, QSEP Research Institute, Kenneth Taylor Hall, Room 426, McMaster University, Hamilton, Ontario, Canada, L8S 4M4, FAX: 905 521 8232, Email: qsep@mcmaster.ca. The Research Report series provides a vehicle for distributing the results of studies undertaken by QSEP associates. The author takes full responsibility for all expressions of opinion.

November 2004

EXPLORING THE USE OF A NONPARAMETRICALLY GENERATED INSTRUMENTAL
VARIABLE IN THE ESTIMATION OF A LINEAR PARAMETRIC EQUATION

Frank T. Denton

McMaster University

ABSTRACT

The use of a nonparametrically generated instrumental variable in estimating a single-equation linear parametric model is explored, using kernel and other smoothing functions. The method, termed IVOS (Instrumental Variables Obtained by Smoothing), is applied in the estimation of measurement error and endogenous regressor models. Asymptotic and small-sample properties are investigated by simulation, using artificial data sets. IVOS is easy to apply and the simulation results exhibit good statistical properties. It can be used in situations in which standard IV cannot because suitable instruments are not available.

JEL Classification: C13; C14; C21

Keywords: single equation models; nonparametric; instrumental variables

EXPLORING THE USE OF A NONPARAMETRICALLY GENERATED INSTRUMENTAL VARIABLE IN THE ESTIMATION OF A LINEAR PARAMETRIC EQUATION

Frank T. Denton¹

McMaster University

1. INTRODUCTION

The method of instrumental variables (IV) has a long history of use in the estimation of linear regression models. There is interesting uncertainty about who had the original idea (Stock and Trebbi, 2003) but no uncertainty about its importance in applied econometrics, whether as IV *per se* or the equivalent but more structured two stage least squares (2SLS) method developed by Theil (1953) and Basmann (1957). The standard IV/2SLS procedure (henceforth just IV) is basic textbook material and the associated statistical theory well established. However, the application of standard IV to a single equation model requires the existence of one or more variables external to the model and uncorrelated with the model's error term, and that is often a problem. For example, analysis of the effect of income on an index of health status must contend with the fact that health status may also affect income. Income must thus be viewed as an endogenous regressor in a health-on-income regression model, and one looks for variables that could serve as instruments for income. Age, sex, and education are often available from surveys of health status and would be excellent IV candidates. But age, sex, and education have their own effects on health, properly belong in the regression model as additional regressors, and hence are ruled out as instruments in the standard method. This problem is well known in the quantitative literature on income-health effects (Buckley et al., 2004, for example), and in many other contexts as well. A similar type of problem arises in the case of a model in which one of the regressors is subject to random measurement error, and hence correlated with the model's error term: standard IV cannot be applied in the absence of a suitable instrumental variable external to the model.

¹My thanks to Lonnie Magee for helpful comments on an earlier version of this paper.

The difficulty in situations of this kind lies often with the standard IV requirement that an instrumental variable must not be one of the model's regressors (more generally, must be linearly independent of them). Earlier suggestions have included the use of powers or (with time series) lagged values of the exogenous regressors as instruments. A promising approach today, though, seems to be the use of nonparametric functions, and that is the approach explored in this paper. The paper is in no way intended as a contribution to the literature on the estimation of nonparametric or semiparametric models with endogenous regressors (Newey, 1990, Pagan and Ullah, Ch. 6, 1999, Newey and Powell, forthcoming, Hall and Horowitz, 2004, and other publications cited therein). Rather it takes from that literature the idea of generating an instrumental variable nonparametrically and applies it to the estimation of a *parametric* model. Its sharper focus is the estimation of a linear parametric model but the procedure can be used in estimating a nonlinear parametric model in the same way that standard IV can be adapted to the estimation of a nonlinear one. The procedure simply provides a way of generating an instrumental variable, which can then be used in a familiar way.

If it is reasonable to generate an instrumental variable nonparametrically one might ask why the equation to be estimated is restricted to being parametric; why should it too not be nonparametric? That may be a telling question in some circumstances. However, the point of view adopted here is that of an investigator who believes there is good reason to estimate an equation with parametric structure but has no reason for assuming a particular functional form for the IV generating equation.

The use of a nonparametric method to generate an instrumental variable is essentially a smoothing device that "averages out" random components of a regressor that are causing the estimation problem. Early averaging out procedures include Wald's (1940) method, which involves sorting observations into two groups and fitting a straight line to the group averages, and Bartlett's (1949) method, a modification of Wald's. (See also Neyman and Scott, 1957, Madansky, 1959, Ware, 1972, and Pakes, 1982, for discussion of such methods.) Given modern nonparametric methods and software, their application in smoothing or averaging out errors seems a natural extension of the earlier approach.

2. MODIFYING THE STANDARD IV METHOD

Consider the single-equation model $Y = X\theta + Z\lambda + u$ where Y and u are $nx1$, X is nxk_1 , Z is nxk_2 (k_1 and $k_2 > 0$), and the parameter vectors are dimensioned accordingly.

Assume $E(u) = 0$; that for every column Z_i of Z , $E(Z_i'u) = 0$, and that for every column X_i of X , $E(X_i'u) \neq 0$. The variables on the right side of the equation are thus divided into two groups, those that are correlated with u and those that are not. Ordinary least squares (OLS) is known to be biased and inconsistent for this model. To apply the standard IV method one looks then for an nxk_3 matrix Q of observations on a set of instrumental variables with $k_3 \geq k_1$, subject to the conditions (1) $\text{plim } n^{-1}Q'u = 0$, (2) at least k_1 of the columns of Q are linearly independent of the columns of Z , (3) $\text{plim } n^{-1}Q'X$ is finite and of full rank, and (4) the variables in Q are as highly linearly correlated as possible with the variables in X .

The standard IV procedure is equivalent to replacing X with $\hat{X} = Q(Q'Q)^{-1}Q'X$ in the model to be estimated and applying OLS. Suppose though that restriction (2) cannot be satisfied – that instrumental variables linearly independent of the Z variables are not available, or not available in sufficient number. The standard IV method cannot then be applied.

The inapplicability of the standard IV method in this case may be related to the requirement that the regression of X on Q be linear. In the absence of model specifications that impose such linearity, and are regarded as binding (see next section), one may be able to construct a matrix $\tilde{X} = G(Q, X)$, where G is some nonlinear regression function, and then proceed as before, substituting \tilde{X} for \hat{X} in the IV procedure. Indeed, if $k_3 = 0$ and $k_2 \geq k_1$,

one may be able to calculate \tilde{X} as $\tilde{X} = G(Z, X)$. (The nonlinearity of G then serves to identify the equation for Y .) Alternatively, if $0 < k_3 < k_1$, it may be possible to combine the set of available Q variables with the set of Z variables in calculating \tilde{X} . Even if $k_3 \geq k_1$ it may be desirable to use a nonlinearly generated \tilde{X} rather than \hat{X} on grounds of efficiency: the X variables may be more highly correlated with the \tilde{X} variables than with the \hat{X} variables.

The function G may be parametric or nonparametric; all that is necessary is that it produce results satisfying conditions (1) to (4). Modern software availability and computational speeds make the use of nonparametric procedures convenient and attractive, and those are the prime focus in this paper. Linear models are used for demonstration purposes but similar procedures could be used to generate instrumental variables for estimating nonlinear parametric models.

It is convenient to have a short label for the modified IV method. I shall refer to it as IVOS, standing for Instrumental Variables Obtained by Smoothing. Most of the applications of IVOS in this paper involve the use of a kernel smoother but some results are reported also for other smoothers.

3. IVOS AS A LIMITED INFORMATION ESTIMATOR

IVOS can be viewed as a method of estimating a single equation, without regard for any others. Alternatively, it can be viewed as a limited information method of estimating one equation in a parametric system of equations, where the interpretation of “limited” is an extension of the customary one.

Consider a system in which each equation is linear. Assuming proper identification, the system can be estimated by a full information method (full information maximum likelihood,

three stage least squares) or a limited information method (limited information maximum likelihood, two stage least squares), one equation at a time. A shorthand way of describing the limited information approach to the estimation of the j^{th} equation is to say that it uses information about the list of exogenous variables in the system but ignores the structural specifications of all equations except the j^{th} . However, that is not quite correct. The standard 2SLS procedure, for example, makes use of the fact that the equations of the system are linear, and to that extent takes account of the system's structure. The IVOS method takes the definition of limited information one step further: it takes account of the list of exogenous variables but ignores all aspects of the structure of the system outside the j^{th} equation, including its linearity. At least, that is one way of thinking about the method.

If the system is correctly specified, and is linear, there is some sacrifice of efficiency in ignoring the linearity of the reduced form in generating an instrumental variable, although asymptotically that should not matter. On the other hand, a method that ignores the assumed linearity of other equations in the system may be more robust to misspecification of the functional form of those equations, just as a parametric limited information method may be more robust than a full information method to errors in parameter restrictions.

4. DEMONSTRATION MODELS

Two models are specified for experiments with IVOS. The first is a measurement error model; the second is an endogenous regressor model.

Measurement error model: A response variable y is a function of two variables. The i^{th} observed value of the first is x_i , the true value is \bar{x}_i , and $x_i = \bar{x}_i + v_i$, where v_i is a random measurement error. The generating equation is thus

$$(1) y_i = \beta_0 + \beta_1 \bar{x}_i + \beta_2 z_i + u_i$$

where u_i is a random equation error. Since \bar{x}_i is not observed the equation must be estimated as

$$(2) y_i = \beta_0 + \beta_1 x_i + \beta_2 z_i + w_i$$

where $w_i = u_i - \beta_1 v_i$. Restrictions are $E(u_i) = E(v_i) = E(u_i v_i) = 0$, $E(u_i^2) = \sigma_u^2$,

and $E(v_i^2) = \sigma_v^2$, $\forall i$. In consequence, $E(u_i w_i) = \sigma_u^2$, $E(v_i w_i) = -\beta_1 \sigma_v^2$,

and $E(w_i^2) = \sigma_u^2 + \beta_1^2 \sigma_v^2$.

Endogenous regressor model: There are two forms of this model. The first has two regressors, x and z . Both are measured without error but x is endogenous within some larger but unspecified simultaneous system. Again, $x_i = \bar{x}_i + v_i$ but v_i is now interpreted as a component of x_i that is correlated with u_i , by virtue of the endogeneity. The other component, \bar{x}_i , is uncorrelated with u_i . The variable y is generated by

$$(3) y_i = \beta_0 + \beta_1 x_i + \beta_2 z_i + u_i$$

where z is exogenous. This is also the form in which the model is estimated. u and v are subject to the same restrictions as before. The second form of the endogenous regressor model is the same as the first except that there are two exogenous regressors, z_1 and z_2 :

$$(4) y_i = \beta_0 + \beta_1 x_i + \beta_2 z_{1i} + \beta_3 z_{2i} + u_i$$

5. DATA SETS

Two sets of data are used. Both are artificial, although one has a basis in real economic

series. They are created in such a way as to hold constant the distribution of the nonstochastic variables as the sample size increases.

Data set 1: This set is used in experiments with the measurement error model defined by equations (1) and (2) and the endogenous regressor model defined by equation (3). The \bar{x} variable for the first 100 observations is generated in both cases as the sequence 1,2,...,100. The sample size n is set as a multiple of 100 and the sequence is repeated every 100 observations: $\bar{x}_i = \bar{x}_{i-100}, i = 101, 102, \dots, n$. The exogenous, error-free z variable is generated as a sine function with period 50: $z_i = \sin(2\pi i / 50), i = 1, 2, \dots, n$. For convenience in subsequent calculations, the \bar{x} and z variables are standardized so as to have mean zero and unit variance. The x series is then calculated as $x_i = \bar{x}_i + u_i$ and the y series is generated either from equation (1) (for measurement error experiments) or equation (3) (for endogenous regressor experiments). For convenience in interpreting the results of the experiments, all coefficients in the models are set to unity: $\beta_0 = \beta_1 = \beta_2 = 1$. u and v are generated as zero-mean random normal numbers with $\sigma_u^2 = \sigma_v^2 = 0.5$. For the measurement error experiments they are generated independently; for the endogenous regressor experiments they are from a bivariate distribution with correlation coefficient $\rho = 0.9$. There are no variables in data set 1 external to the equation to be estimated, and hence no instrumental variables for application of the standard IV method. For application of IVOS, though, the internal instrumental variable z can be used.

Data set 2: Series consisting of 72 observations relating to the U.S. commercial loan market are used as a starting point for the construction of this data set. The series are from Table 9.3 of Maddala's (1992) textbook. The data set is used in endogenous regressor experiments based on equation (4) so that observations corresponding to the variables \bar{x} , x , z_1 , and z_2 are required.

Maddala's R series (average prime rate charged by banks) is used as a basis for the \bar{x} variable, his RD series (AAA corporate bond rate) as a basis for z_1 , and his X series (industrial production index) as a basis for z_2 . The \bar{x} , z_1 , and z_2 variables are standardized. The sample size is set as a multiple of 72 and the initial sequences of those variables are repeated for observations 73 to 144, 145 to 216, etc. The x variable is then obtained from $x_i = \bar{x}_i + v_i$, as before, and the y variable is generated from equation (4). As before, all coefficients are set to unity. The means, variances, and correlation coefficient of u and v are the same as in data set 1. Two external variables are included in data set 2 in addition to the internal ones, z_1 and z_2 . The external ones, z_3 and z_4 , are taken also from the Maddala table: RS (3-month treasury bill rate) is used for z_3 , Y (total bank deposits) for z_4 . Both variables are converted to standardized form. The existence of external instruments means that the standard IV method can be used to estimate equation (4), and the results compared with those of IVOS, both when the external instruments are ignored in applying that method, and when they are included..

6. APPROXIMATE ASYMPTOTICS WITH DATA SET 1

The first set of experiments assumes a sample size large enough to generate results that can be viewed as approximately asymptotic. Table 1 shows estimates of the coefficients in the measurement error model of equation (1) and the endogenous regressor model of equation (3), using data set 1. Estimates obtained by IVOS are shown together with OLS estimates, for comparison. The sample size is set at $n = 20,000$.

The instrumental variable used for the IVOS estimates in Table 1 is \tilde{x}_i , the Nadaraya-Watson univariate kernel-smoothed value of x_i , based on normal distribution weights (Härdle,

1990). The IVOS method is implemented (here and subsequently) using SHAZAM, Version 9 (SHAZAM, 2001). The bandwidth is set at the SHAZAM default value, which is an approximately optimal value for a normal kernel (SHAZAM, 2001, Silverman, 1986). (Experiments with alternative bandwidths and smoothing functions are reported below.) To facilitate comparisons of results for different estimators the same random numbers are generated in each case from a common seed. (This procedure is used also in the subsequent experiments.)

The OLS results are as expected. OLS underestimates by a wide margin the slope coefficients in the measurement error model and overestimates by a wide margin the slope coefficients in the endogenous regressor model. Applying (incorrectly) the standard t test to any OLS slope coefficient in either model rejects the null hypothesis that the coefficient is equal to 1 at virtually any significance level that one might choose. The IVOS estimates, on the other hand, are quite close to the true values – within less than 1 percent in all cases and within half of 1 percent in three of the four. These large-sample results based on data set 1 thus suggest good asymptotic properties for an estimator that performs well using only internal instrumental variables.

7. APPROXIMATE ASYMPTOTICS WITH ALTERNATIVE SMOOTHING FUNCTIONS

It is of interest to see the extent to which the results of Table 1 are affected by choosing a smoothing function other than the normal kernel function. Focusing now (and in all subsequent experiments) on the endogenous regressor model, Table 2 shows what happens when five other functions are employed to calculate \tilde{x} for the purpose of estimating equation (3):

Epanechnikov, quartic, triangular, and uniform kernel functions, and the weighted local least squares function loess. (For the weighting patterns implicit in these functions, see Yatchew, 2003.) The bandwidth for the kernel functions is the same as in Table 1.

The choice of smoothing function in calculating \tilde{x} has hardly any effect on the estimated coefficients. The kernel-based estimates are virtually identical in all cases and the loess-based estimates differ only slightly from the others. This suggests that the choice of a smoothing function can be based on convenience, as far as asymptotic properties are concerned.

8. APPROXIMATE ASYMPTOTICS WITH ALTERNATIVE BANDWIDTHS

A somewhat similar finding applies to the choice of bandwidth in using a kernel smoother. The normal kernel is used to generate the results in Table 3, combined with eight choices of bandwidth. The first three are the bandwidths that minimize the cross-validation mean square error statistic (CV), the generalized cross-validation statistic (GCV), and the Akaike information criterion (AIC). (Descriptions can be found in Ruppert, Wand, and Carroll, 2003; other criteria for choosing an optimum bandwidth are available but these three are representative.) The remaining five choices are the default bandwidth used previously and proportionate decreases or increases therefrom: 0.50 and 0.25 of the default bandwidth on the downside, 1.50 and 2.00 on the up side.

The coefficient estimates are quite insensitive to bandwidth selection. This is in fact a quite reasonable finding for a large sample. A larger bandwidth implies more smoothing, a smaller one less smoothing. One might define an optimum bandwidth and think of anything larger as implying oversmoothing. But asymptotically, oversmoothing should not matter. As long as the smoothing eliminates the correlation between \tilde{x} and u , and falls short of eliminating all correlation between \tilde{x} and \bar{x} , the degree of smoothing should be irrelevant as far as asymptotic properties are concerned. This is analogous to the fact that in the standard IV method weak correlation between an instrumental variable and the variable for which it serves as an instrument can still produce an estimator with good asymptotic properties. With a large enough sample size,

coefficient estimates should be similar over a wide range of bandwidth choices. As a practical procedure, one could establish a range over which coefficient estimates are approximately stable by trying a few alternatives and then simply choose a bandwidth from that range.

9. APPROXIMATE ASYMPTOTICS WITH DATA SET 2

The endogenous regressor model of equation (4) differs from that of equation (3) (which has been used up to this point) in that it has two internal exogenous variables (z_1, z_2), rather than one. Data set 2 is used now to estimate this equation, with sample size 21,600 (the core size 72 multiplied by 300). Unlike data set 1, data set 2 also provides external exogenous variables (z_3, z_4). This makes possible the application of standard IV as well as IVOS. Table 4 shows estimates of the coefficients in equation (4) based on OLS, standard IV, and four versions of IVOS, labelled IVOS-NP1, IVOS-NP2, IVOS-NP3, and IVOS-SPL. Definitions of the IVOS estimators are as follows:

IVOS-NP1: \tilde{x} is generated by a multivariate normal kernel function with the SHAZAM default bandwidth (SHAZAM, 2001), using only internal exogenous variables. (The SHAZAM multivariate procedure uses a method due to Rust (1988) which allows for covariance among regressors.)

IVOS-NP2: Same as IVOS-NP1 except that both internal and external exogenous variables are used in generating \tilde{x} ; z_3 and z_4 are then combined with \tilde{x} to form a set of three instrumental variables for estimating equation (4) by the usual IV procedure; z_3 and z_4 are thus treated as if they were instrumental variables independent of \tilde{x} .

IVOS-NP3: Same as IVOS-NP2 except that z_3 and z_4 are not treated as separate instrumental variables (\tilde{x} is assumed to capture all effects associated with z_3 and z_4 , as well as z_1 and z_2).

IVOS-SPL: This version of IVOS employs a spline function for smoothing purposes, rather than a kernel function; \tilde{x} is calculated by treating x as a cubic spline regression function of the internal exogenous variables, with six knots for each of z_1 and z_2 .

The IVOS-NP1 procedure ignores the external exogenous variables and thus simulates a situation in which only internal ones are available, and in which standard IV is therefore not possible. IVOS-NP2 adds \tilde{x} to the list of instrumental variables used in the standard IV procedure, and thus shows the “value added” by it to the standard procedure. IVOS-NP3 represents a more “natural” application of the IVOS method by using the \tilde{x} variable calculated in IVOS-NP2 *instead of* z_3 and z_4 .

The results in Table 4 are generally similar to the corresponding results in Table 1. Treating them as approximations to asymptotic results, OLS is seen to be markedly inconsistent, overstating β_1 badly and understating β_2 and β_3 . Standard IV performs well and so do all of the IVOS estimators. The spline version of IVOS performs about as well as the kernel-based versions. However, it is more cumbersome to set up (to establish the number and positioning of knots); the kernel-based versions are more flexible, and seem more appealing as practical tools. In any event, the multivariate IVOS procedures all measure up well against standard IV and the experiments suggest good asymptotic properties. The fact that IVOS-NP1 performs well suggests again considerable promise for the IVOS method in large-sample situations in which

lack of external instrumental variables precludes the use of standard IV.

10. SMALL-SAMPLE PROPERTIES

The previous discussion has been concerned with large-sample or approximate asymptotic properties of the IVOS estimator, based on $n = 20,000$ (data set 1) and $n = 21,600$ (data set 2). The present section and the next discuss results for small samples. IVOS is used again to estimate equation (4) by applying it to data set 2, but this time with sample sizes 72, 144, 360, and 720 in a series of Monte Carlo experiments. Results for the kernel-based NP1, NP2, and NP3 versions of IVOS (as defined in the previous section) are reported in Tables 5 and 6 along with OLS and standard IV results. The Monte Carlo experiments are repeated 20,000 times for each sample size. Differences of mean coefficient estimates from true values (1 in all cases) are reported in Table 5; they are interpreted as (estimated) biases. Root mean square errors (RMSE) are reported also.

The large biases associated with OLS are evident in Table 5, as expected. Standard IV, which uses z_3 and z_4 as instrumental variables, performs well: its biases are negligible and its RMSE values are notably lower than those of OLS for all slope coefficient estimates, with all sample sizes. The IVOS slope coefficient estimates exhibit only small biases with $n = 72$, and those diminish as the sample size increases; at $n = 720$ the biases for all three IVOS estimators are well under 1 percent. The RMSE values are roughly similar to those of standard IV and far below the OLS RMSE values. IVOS-NP1, which ignores the availability of external instrumental variables, produces slope coefficient RMSE values that are only a little higher than those produced by standard IV; for $n = 720$ the two sets of values are particularly close.

Table 6 shows the empirical probabilities of rejecting the null hypothesis $\beta_i = 1$ for each of $i = 0,1,2,3,4$ in asymptotic t tests using nominal Type I error probabilities of 0.05 and 0.01.

(An empirical probability is calculated as the proportion of cases out of the 20,000 Monte Carlo replications in which the t value exceeds the nominal critical value.) As expected, the calculated OLS probabilities of rejection are far in excess of what they would be if the theoretical test assumptions were valid. The standard IV probabilities are reasonably close to the nominal ones, even for $n = 72$. The IVOS probabilities are higher than the nominal ones but the differences decline with increasing sample size. (Example: at nominal probability 0.01, the proportion of rejections of null hypothesis $\beta_1 = 1$ for IVOS-NP1 is 0.0382 at $n = 72$, 0.0168 at $n = 720$.)

The results suggest the advisability of playing safe by choosing somewhat lower nominal Type I error probabilities in carrying out tests on coefficients with smaller sample sizes.

11. CONCLUSION

Creating instrumental variables by nonparametric smoothing is a useful procedure. Although the theory of the procedure is implicit in the nonparametric regression estimation literature the advantages of its application to parametric models seem not to have been exploited in applied econometrics. The procedure is especially useful in situations in which standard IV cannot be employed because external instrumental variables are not available, or not available in sufficient number, but there are internal variables that can be used. Such situations are common. The focus in this paper has been on the estimation of linear parametric models but the method can be adapted to apply to nonlinear parametric models in the same way that standard IV can be adapted. The method simply provides another way to generate instrumental variables, which can then be handled in a familiar manner.

The IVOS method can be implemented in various ways. Attention has been given mostly to kernel-based smoothing techniques in this paper. Those techniques are flexible and easy to apply with generally available software. Other smoothing techniques can be used though,

including nearest neighbour, locally weighted regression, and spline function regression.

IVOS performed well in the experiments reported here, including both large-sample and small-sample experiments. Kernel smoothing with large samples is particularly convenient in that oversmoothing is not a significant problem for asymptotic properties: the asymptotics are likely to be the same for alternative kernel functions and a wide range of bandwidth choices.

Now for the necessary cautionary remarks. The results reported in this paper are based on particular artificial data sets. As with all such experiments, different data may produce different results. The evidence in favour of IVOS is suggestive of good properties but there are no guarantees in particular practical situations (any more than there are guarantees for standard IV). An obvious requirement is a nonlinear relationship that can be exploited nonparametrically between a variable for which an instrument is needed and the candidate instrumental variable or variables that are available, internal or external. If internal instrumental variables are to be used (in the absence of external ones) a further requirement is that the model include such variables in sufficient number. The method will not work in every situation but it does appear to be a useful tool to have at hand in applied parametric econometrics.

REFERENCES

- Bartlett, M. S. (1949), "Fitting a Straight Line When Both Variables are Subject to Error," Biometrics, 5, 207-212.
- Basman, R. L. (1957), "A Generalized Classical Method of Estimation of Coefficients in a Structural Equation," Econometrica, 25, 77-83.
- Buckley, N. J., F. T. Denton, A. L. Robb, and B. G. Spencer (2004), "The Transition from Good to Poor Health: An Econometric Study of the Older Population," Journal of Health Economics, 23, 1013-1034.
- Hall, P. and J. L. Horowitz (2004), "Nonparametric Methods for Inference in the Presence of Instrumental Variables," unpublished.
- Härdle, W. (1990), Applied Nonparametric Regression, Cambridge University Press.
- Madansky, A. (1959), "The Fitting of Straight Lines When Both Variables are Subject to Error," Journal of the American Statistical Association, 54, 173-205. (Also "Corrigenda," Journal of the American Statistical Association, 54, 812).
- Maddala, G. S. (1992), Introduction to Econometrics, Second Edition, Macmillan Publishing Company.
- Newey, W. (1990), "Efficient Instrumental Variables Estimation of Nonlinear Models," Econometrica, 58, 757-837.
- Newey, W. and J. L. Powell (forthcoming), "Instrumental Variable Estimation of Nonparametric Models," Econometrica.
- Neyman, J. and E. L. Scott (1951), "On Certain Methods of Estimating the Linear Structural Relation," Annals of Mathematical Statistics, 22, 352-361.
- Pagan, A. and A. Ullah (1999), Nonparametric Econometrics, Cambridge University Press.
- Pakes, A. (1982), "On the Asymptotic Bias of Wald-Type Estimators of a Straight Line When

- Both Variables are Subject to Error,” International Economic Review, 23, 491-497.
- Ruppert, D., M. P. Wand, and R. J. Carroll (2003), Semiparametric Regression, Cambridge University Press.
- Rust, R. T. (1988), “Flexible Regression,” Journal of Marketing Research, 25, 10-24.
- SHAZAM (2001), SHAZAM User’s Reference Manual, Version 9, Northwest Econometrics Limited.
- Silverman, B. W. (1986), Density Estimation, Chapman and Hall.
- Stock, J. H. and F. Trebbi (2003), “Who Invented Instrumental Variable Regression?,” Journal of Economic Perspectives, 17, 177-194.
- Theil, H. (1953), “Repeated Least Squares Applied to Complete Equation Systems,” The Hague: Central Planning Bureau, unpublished.
- Wald, A. (1940), “The Fitting of Straight Lines if Both Variables are Subject to Error,” Annals of Mathematical Statistics, 11, 284-300.
- Ware, J. H. (1972), “The Fitting of Straight Lines When Both Variables are Subject to Error and the Ranks of the Means are Known,” Journal of the American Statistical Association, 67, 891-897.
- Yatchew, A. (2003), Semiparametric Regression for the Applied Econometrician, Cambridge University Press.

TABLE 1: LARGE-SAMPLE REGRESSION RESULTS: MEASUREMENT ERROR AND ENDOGENOUS REGRESSOR MODELS, DATA SET 1 (n=20,000)

(True values: $\beta_0 = \beta_1 = \beta_2 = 1$)

		Estimated coefficients (standard errors)		
		β_0	β_1	β_2
Measurement error model				
	OLS	1.0009 (0.0048)	0.7771 (0.0046)	0.9134 (0.0051)
	IVOS	1.0020 (0.0051)	1.0093 (0.0055)	1.0040 (0.0055)
Endogenous regressor model				
	OLS	0.9981 (0.0032)	1.2085 (0.0030)	1.0811 (0.0034)
	IVOS	0.9971 (0.0035)	1.0050 (0.0038)	1.0018 (0.0038)

TABLE 2: LARGE-SAMPLE REGRESSION RESULTS USING INSTRUMENTAL VARIABLE GENERATED BY ALTERNATIVE NONPARAMETRIC SMOOTHING FUNCTIONS: ENDOGENOUS REGRESSOR MODEL, DATA SET 1 (n=20,000)
(True values: $\beta_0 = \beta_1 = \beta_2 = 1$)

	Estimated coefficients (standard errors)		
	β_0	β_1	β_2
<u>Smoothing function</u>			
Normal kernel	0.9971 (0.0035)	1.0050 (0.0038)	1.0018 (0.0038)
Epanechnikov kernel	0.9971 (0.0035)	1.0053 (0.0038)	1.0019 (0.0038)
Quartic kernel	0.9971 (0.0035)	1.0054 (0.0038)	1.0020 (0.0038)
Triangular kernel	0.9971 (0.0035)	1.0055 (0.0038)	1.0020 (0.0038)
Uniform kernel	0.9971 (0.0035)	1.0052 (0.0038)	1.0019 (0.0038)
Local weights (loess)	0.9971 (0.0035)	1.0070 (0.0038)	1.0026 (0.0038)

TABLE 3: LARGE-SAMPLE REGRESSION RESULTS USING INSTRUMENTAL VARIABLE GENERATED BY NORMAL KERNEL FUNCTION WITH ALTERNATIVE BANDWIDTHS: ENDOGENOUS REGRESSOR MODEL, DATA SET 1 (n=20,000)
 (True values: $\beta_0 = \beta_1 = \beta_2 = 1$)

	Estimated coefficients (standard errors)		
	β_0	β_1	β_2
<u>Bandwidth selection criterion</u>			
Minimum CV	0.9971 (0.0035)	1.0062 (0.0038)	1.0023 (0.0038)
Minimum GCV	0.9971 (0.0035)	1.0060 (0.0038)	1.0022 (0.0038)
Minimum AIC	0.9971 (0.0035)	1.0060 (0.0038)	1.0022 (0.0038)
Default	0.9971 (0.0035)	1.0050 (0.0038)	1.0018 (0.0038)
Default x 0.25	0.9971 (0.0035)	1.0058 (0.0038)	1.0021 (0.0038)
Default x 0.50	0.9971 (0.0035)	1.0053 (0.0038)	1.0019 (0.0038)
Default x 1.50	0.9971 (0.0035)	1.0047 (0.0039)	1.0017 (0.0038)
Default x 2.00	0.9971 (0.0035)	1.0044 (0.0040)	1.0016 (0.0039)

TABLE 4: LARGE SAMPLE REGRESSION RESULTS: ENDOGENOUS REGRESSOR
 MODEL, DATA SET 2 (n=21,600)
 (True values: $\beta_0 = \beta_1 = \beta_2 = \beta_3 = 1$)

<u>Estimator</u>	Estimated coefficients (standard errors)			
	β_0	β_1	β_2	β_3
OLS	1.0030 (0.0030)	1.2494 (0.0032)	0.8596 (0.0038)	0.8676 (0.0037)
Standard IV	1.0033 (0.0034)	1.0023 (0.0052)	0.9993 (0.0048)	1.0013 (0.0047)
IVOS-NP1	1.0033 (0.0034)	1.0029 (0.0048)	0.9990 (0.0046)	1.0010 (0.0046)
IVOS-NP2	1.0033 (0.0034)	0.9999 (0.0043)	1.0007 (0.0045)	1.0027 (0.0044)
IVOS-NP3	1.0033 (0.0034)	1.0000 (0.0043)	1.0006 (0.0045)	1.0026 (0.0044)
IVOS-SPL	1.0033 (0.0034)	1.0065 (0.0051)	0.9970 (0.0047)	0.9991 (0.0047)

TABLE 5: SIMULATED SMALL-SAMPLE BIASES AND ROOT MEAN SQUARE ERRORS:
 ENDOGENOUS REGRESSOR MODEL, DATA SET 2

(True values: $\beta_0 = \beta_1 = \beta_2 = \beta_3 = 1$)

	Bias				RMSE			
	β_0	β_1	β_2	β_3	β_0	β_1	β_2	β_3
<u>Sample size and estimator</u>								
n=72: OLS	-0.0002	0.2390	-0.1357	-0.1288	0.0468	0.2436	0.1474	0.1414
Standard IV	-0.0003	0.0003	0.0006	-0.0001	0.0592	0.0924	0.0845	0.0829
IVOS-NP1	-0.0003	0.0437	-0.0241	-0.0234	0.0570	0.1109	0.0894	0.0875
IVOS-NP2	-0.0002	0.0600	-0.0334	-0.0323	0.0560	0.0900	0.0806	0.0792
IVOSNP-3	-0.0003	0.0502	-0.0278	-0.0270	0.0565	0.0848	0.0792	0.0780
n=144:OLS	0.0003	0.2447	-0.1388	-0.1318	0.0326	0.2468	0.1447	0.1379
Standard IV	0.0003	-0.0007	0.0007	0.0006	0.0415	0.0655	0.0602	0.0587
IVOS-NP1	0.0003	0.0243	-0.0135	-0.0129	0.0406	0.0755	0.0630	0.0616
IVOS-NP2	0.0003	0.0334	-0.0186	-0.0178	0.0402	0.0600	0.0569	0.0559
IVOS-NP3	0.0003	0.0290	-0.0162	-0.0154	0.0404	0.0582	0.0564	0.0555
n=360:OLS	-0.0001	0.2482	-0.1406	-0.1335	0.0208	0.2490	0.1429	0.1360
Standard IV	-0.0001	0.0001	0.0001	0.0000	0.0264	0.0410	0.0375	0.0370
IVOS-NP1	-0.0001	0.0120	-0.0067	-0.0064	0.0262	0.0453	0.0388	0.0382
IVOS-NP2	-0.0001	0.0159	-0.0089	-0.0085	0.0261	0.0363	0.0356	0.0352
IVOS-NP3	-0.0001	0.0145	-0.0081	-0.0078	0.0261	0.0358	0.0354	0.0350
n=720:OLS	0.0002	0.2493	-0.1413	-0.1339	0.0145	0.2497	0.1424	0.1351
Standard IV	0.0002	-0.0001	0.0000	0.0003	0.0185	0.0291	0.0264	0.0261
IVOS-NP1	0.0002	0.0067	-0.0038	-0.0034	0.0184	0.0308	0.0269	0.0265
IVOS-NP2	0.0002	0.0089	-0.0051	-0.0046	0.0183	0.0251	0.0249	0.0246
IVOS-NP3	0.0002	0.0083	-0.0047	-0.0042	0.0183	0.0249	0.0249	0.0246

TABLE 6: SIMULATED SMALL-SAMPLE TYPE I ERROR PROBABILITIES:
 ENDOGENOUS REGRESSOR MODEL, DATA SET 2
 (True values: $\beta_0 = \beta_1 = \beta_2 = \beta_3 = 1$)

	Nominal probability 0.05				Nominal probability 0.01			
	β_0	β_1	β_2	β_3	β_0	β_1	β_2	β_3
<u>Sample size and estimator</u>								
n=72: OLS	0.0138	0.9904	0.5398	0.5088	0.0052	0.9580	0.2918	0.2664
Standard IV	0.0510	0.0540	0.0526	0.0538	0.0094	0.0181	0.0134	0.0136
IVOS-NP1	0.0474	0.1013	0.0738	0.0720	0.0087	0.0382	0.0236	0.0226
IVOS-NP2	0.0513	0.1828	0.0870	0.0857	0.0100	0.0784	0.0286	0.0260
IVOSNP-3	0.0517	0.1507	0.0788	0.0762	0.0100	0.0608	0.0250	0.0234
n=144:OLS	0.0280	1.0000	0.8632	0.8326	0.0041	0.9998	0.6690	0.6189
Standard IV	0.0486	0.0550	0.0517	0.0512	0.0093	0.0134	0.0109	0.0108
IVOS-NP1	0.0470	0.0803	0.0640	0.0650	0.0091	0.0274	0.0171	0.0175
IVOS-NP2	0.0486	0.1254	0.0716	0.0724	0.0097	0.0476	0.0192	0.0196
IVOS-NP3	0.0488	0.1098	0.0682	0.0684	0.0096	0.0405	0.0172	0.0176
n=360:OLS	0.0272	1.0000	0.9985	0.9968	0.0036	1.0000	0.9917	0.9851
Standard IV	0.0484	0.0512	0.0504	0.0522	0.0100	0.0119	0.0104	0.0116
IVOS-NP1	0.0480	0.0690	0.0575	0.0591	0.0099	0.0203	0.0140	0.0153
IVOS-NP2	0.0488	0.0899	0.0600	0.0636	0.0100	0.0272	0.0162	0.0160
IVOS-NP3	0.0486	0.0836	0.0580	0.0621	0.0100	0.0249	0.0158	0.0157
n=720:OLS	0.0256	1.0000	1.0000	1.0000	0.0033	1.0000	1.0000	1.0000
Standard IV	0.0480	0.0525	0.0488	0.0487	0.0088	0.0112	0.0100	0.0099
IVOS-NP1	0.0478	0.0627	0.0513	0.0550	0.0088	0.0168	0.0120	0.0118
IVOS-NP2	0.0479	0.0736	0.0566	0.0554	0.0090	0.0220	0.0129	0.0116
IVOS-NP3	0.0480	0.0708	0.0552	0.0547	0.0090	0.0206	0.0122	0.0115

QSEP RESEARCH REPORTS - Recent Releases

Number	Title	Author(s)
No. 351:	Describing Disability among High and Low Income Status Older Adults in Canada	P. Raina M. Wong L.W. Chambers M. Denton A. Gafni
No. 352:	Some Demographic Consequences of Revising the Definition of Old to Reflect Future Changes in Life Table Probabilities	F.T. Denton B.G. Spencer
No. 353:	The Correlation Between Husband's and Wife's Education: Canada, 1971-1996	L. Magee J. Burbidge L. Robb
No. 354:	The Effect of Marginal Tax Rates on Taxable Income: A Panel Study of the 1988 Tax Flattening in Canada	M.-A. Sillamaa M.R. Veall
No. 355:	Population Change and the Requirements for Physicians: The Case of Ontario	F.T. Denton A. Gafni B.G. Spencer
No. 356:	2 ½ Proposals to Save Social Security	D. Fretz M.R. Veall
No. 357:	The Consequences of Caregiving: Does Employment Make A Difference?	C.L. Kemp C.J. Rosenthal
No. 358:	Exploring the Effects of Population Change on the Costs of Physician Services	F.T. Denton A. Gafni B.G. Spencer
No. 359:	Reflexive Planning for Later Life: A Conceptual Model and Evidence from Canada	M.A. Denton S. French A. Gafni A. Joshi C. Rosenthal S. Webb
No. 360:	Time Series Properties and Stochastic Forecasts: Some Econometrics of Mortality from The Canadian Laboratory	F.T. Denton C.H. Feaver B.G. Spencer
No. 361:	Linear Public Goods Experiments: A Meta-Analysis	J. Zelmer
No. 362:	The Timing and Duration of Women's Life Course Events: A Study of Mothers With At Least Two Children	K.M. Kobayashi A. Martin-Matthews C.J. Rosenthal S. Matthews

QSEP RESEARCH REPORTS - Recent Releases

Number	Title	Author(s)
No. 363:	Age-Gapped and Age-Condensed Lineages: Patterns of Intergenerational Age Structure among Canadian Families	A. Martin-Matthews K.M. Kobayashi C.J. Rosenthal S.H. Matthews
No. 364:	The Education Premium in Canada and the United States	J.B. Burbidge L. Magee A.L. Robb
No. 365:	Student Enrolment and Faculty Recruitment in Ontario: The Double Cohort, the Baby Boom Echo, and the Aging of University Faculty	B.G. Spencer
No. 366:	The Economic Well-Being of Older Women Who Become Divorced or Separated in Mid and Later Life	S. Davies M. Denton
No. 367:	Alternative Pasts, Possible Futures: A “What If” Study of the Effects of Fertility on the Canadian Population and Labour Force	F.T. Denton C.H. Feaver B.G. Spencer
No. 368:	Baby-Boom Aging and Average Living Standards	W. Scarth M. Souare
No. 369:	The Impact of Reference Pricing of Cardiovascular Drugs on Health Care Costs and Health Outcomes: Evidence from British Columbia – Volume I: Summary	P.V. Grootendorst L.R. Dolovich A.M. Holbrooke A.R. Levy B.J. O'Brien
No. 370:	The Impact of Reference Pricing of Cardiovascular Drugs on Health Care Costs and Health Outcomes: Evidence from British Columbia – Volume II: Technical Report	P.V. Grootendorst L.R. Dolovich A.M. Holbrooke A.R. Levy B.J. O'Brien
No. 371:	The Impact of Reference Pricing of Cardiovascular Drugs on Health Care Costs and Health Outcomes: Evidence from British Columbia – Volume III: ACE and CCB Literature Review	L.R. Dolovich A.M. Holbrook M. Woodruff
No. 372:	Do Drug Plans Matter? Effects of Drug Plan Eligibility on Drug Use Among the Elderly, Social Assistance Recipients and the General Population	P. Grootendorst M. Levine

QSEP RESEARCH REPORTS - Recent Releases

Number	Title	Author(s)
No. 373:	Student Enrolment and Faculty Recruitment in Ontario: The Double Cohort, the Baby Boom Echo, and the Aging of University Faculty	B.G. Spencer
No. 374:	Aggregation Effects on Price and Expenditure Elasticities in a Quadratic Almost Ideal Demand System	F.T. Denton D.C. Mountain
No. 375:	Age, Retirement and Expenditure Patterns: An Econometric Study of Older Canadian Households	F.T. Denton D.C. Mountain B.G. Spencer
No. 376:	Location of Adult Children as an Attraction for Black and White Elderly <i>Return</i> and <i>Onward</i> Migrants in the United States: Application of a Three-level Nested Logit Model with Census Data	K-L. Liaw W.H. Frey
No. 377:	The Dynamics of Food Deprivation and Overall Health: Evidence from the Canadian National Population Health Survey	L. McLeod M.R. Veall
No. 378:	Quebec's Lackluster Performance in Interprovincial Migration and Immigration: How, Why, and What Can Be Done?	K-L. Liaw L. Xu M. Qi
No. 379:	Out-of-Pocket Prescription Drug Expenditures and Public Prescription Drug Programs	S. Alan T.F. Crossley P. Grootendorst M.R. Veall
No. 380:	Population Aging, Productivity, and Growth in Living Standards	W. Scarth
No. 381:	The Transition from Good to Poor Health: An Econometric Study of the Older Population	N.J. Buckley F.T. Denton A.L. Robb B.G. Spencer
No. 382:	The Evolution of High Incomes In Canada, 1920-2000	E. Saez M.R. Veall
No. 383:	Population Change and Economic Growth: The Long-Term Outlook	F.T. Denton B.G. Spencer
No. 384:	The Economic Legacy of Divorced and Separated Women in Old Age	L. McDonald A.L. Robb

QSEP RESEARCH REPORTS - Recent Releases

Number	Title	Author(s)
No. 385:	National Catastrophic Drug Insurance Revisited: Who Would Benefit from Senator Kirby's Recommendations?	T.F. Crossley P.V. Grootendorst M.R. Veall
No. 386:	Wages in Canada: SCF, SLID, LFS and the Skill Premium	A.L. Robb L. Magee J.B. Burbidge
No. 387:	Socioeconomic Influence on the Health of Older People: Estimates Based on Two Longitudinal Surveys	N.J. Buckley F.T. Denton A.L. Robb B.G. Spencer
No. 388:	An Invitation to Multivariate Analysis: An Example About the Effect of Educational Attainment on Migration Propensities in Japan	A. Otomo K-L. Liaw
No. 389:	Financial Planning for Later Life: Subjective Understandings of Catalysts and Constraints	C.L. Kemp C.J. Rosenthal M. Denton
No. 390:	Exploring the Use of a Nonparametrically Generated Instrumental Variable in the Estimation of a Linear Parametric Equation	F.T. Denton