



A Note on an 'Errors in Variables' Model

Tony Lancaster

Journal of the American Statistical Association, Volume 61, Issue 313 (Mar., 1966),
128-135.

Stable URL:

<http://links.jstor.org/sici?sici=0162-1459%28196603%2961%3A313%3C128%3AANO%27I%3E2.0.CO%3B2-Q>

Your use of the JSTOR archive indicates your acceptance of JSTOR's Terms and Conditions of Use, available at <http://www.jstor.org/about/terms.html>. JSTOR's Terms and Conditions of Use provides, in part, that unless you have obtained prior permission, you may not download an entire issue of a journal or multiple copies of articles, and you may use content in the JSTOR archive only for your personal, non-commercial use.

Each copy of any part of a JSTOR transmission must contain the same copyright notice that appears on the screen or printed page of such transmission.

Journal of the American Statistical Association is published by American Statistical Association. Please contact the publisher for further permissions regarding the use of this work. Publisher contact information may be obtained at <http://www.jstor.org/journals/astata.html>.

Journal of the American Statistical Association
©1966 American Statistical Association

JSTOR and the JSTOR logo are trademarks of JSTOR, and are Registered in the U.S. Patent and Trademark Office. For more information on JSTOR contact jstor-info@umich.edu.

©2002 JSTOR

A NOTE ON AN 'ERRORS IN VARIABLES' MODEL

TONY LANCASTER

University of Birmingham

We consider an errors in variables model in which the 'true' part of the determining variable is generated by a simple forecasting mechanism. It is shown that the Least Squares errors in variables bias can be interpreted in terms of the parameters of the forecasting mechanism; and that the 'standard' result for this bias may no longer hold in this situation. An empirical illustration of the results is given.

SUPPOSE that economic theory leads us to the following behavioural equation:

$$C_{it} = \alpha + k \cdot Y_{it}^* + U_{it}. \quad (1)$$

C is an observable economic variable; α and k are constants; and U is an unobservable random disturbance in the equation distributed independently of Y^* . Y_{it}^* is a variable which represents the forecast, made by the i^{th} individual in the t^{th} period, of the value of a second variable, Y_i , one period hence. Y_{it}^* is a forecast of Y_{it+1} . To fix ideas we could interpret C_{it} as the current consumption of the i^{th} household, Y_{it} as its income and Y_{it}^* as its current forecast of its *next* period's income.^{1,2}

By its subjective character, Y^* is a variable on which it may be difficult to secure observations. Indeed, it is only one example of a whole class of variables which have recently appeared in econometric equations and which raise the same difficulty for the same reason. Though this class of variable does not appear to have a name, examples of others of similar type are: 'desired' inventories, 'planned' output, 'normal' price, 'permanent' income, 'expected' sales, and so on.

A problem that arises with models of this type is that of examining the relation between α and k and the coefficients in the regression of C_t on some readily observable variable related to Y^* . Suppose in particular that one computes the Least Squares (LS) estimator of C_{it} on Y_{it} for n sample observations, i.e. one uses actual income and not next period's forecast income as the determining variable. The analogous procedure with other economic variables of the same class would be the use of, say, actual inventories instead of the desired level, actual price instead of normal price, actual output instead of planned output, etc.

Analysis of the LS estimator has sometimes been based on the assumption that the differences between actual and forecast Y could be assumed to behave like random measurement error, and then the classical theory of errors in variables has been employed.³ The conceptual differences between Y^* and Y have

¹ We do not, of course, suppose this to be a good consumption-income model.

² Note that since we are going to consider samples of individuals at the same point of time the operative sampling subscript is i . The t subscript enters in because events of several periods are relevant to the behaviour of each individual at any given time. The actual time we are considering here is the t^{th} period also referred to as the current period.

³ cf. Liviatan (3).

been treated as though Y was merely an imperfect empirical approximation to Y^* .

An outline of this errors in variables analysis may be briefly given as follows. The measurement error in Y_i is denoted by Y_i^{**} , hence $Y_i = Y_i^* + Y_i^{**}$. This error is assumed to be uncorrelated with U_i and Y_i^* in the sense that their sampling covariances tend to zero as n becomes large. Then assuming that the variances of Y^* and Y^{**} tend similarly to constants denoted by σ^2 with an appropriate subscript we have the well known result on the LS estimator,

$$\text{plim } \hat{k}/k = \frac{\sigma_{y^*}^2}{\sigma_{y^*}^2 + \sigma_{y^{**}}^2}, \tag{2}$$

where $\hat{k} = \Sigma cy / \Sigma y^2$.⁴ The LS estimator is asymptotically biased in the direction of zero, in general.

Suppose now that each household forecasts its next period's income in the way described by the model of adaptive forecasting,⁵ the simplest version of which may be described by,

$$Y_{it}^* = Y_{it-1}^* + \beta(Y_{it} - Y_{it-1}^*) + V_{it} = Y_{it-1}^* + \beta E_{it} + V_{it}. \tag{3}$$

β is a constant, which has been called a coefficient of adaptive expectations; V_t is a random disturbance in the equation allowing for unsystematic influences on the forecast and assumed distributed independently of Y_t and Y_{t-1}^* . $Y_{it} - Y_{it-1}^*$ ($= E_{it}$) is the error made by the i^{th} household in its last period's forecast of its current income. It will be referred to as the forecast error to distinguish it from the 'measurement error', Y_i^{**} . This equation tells us that, apart from the disturbance, the forecast is adjusted from period to period by a proportion of the forecast error.

Let us consider the previous result on the large sample errors in variables bias in the light of this further specification of Y^* . A necessary condition for (2) to hold is that the sampling covariance of Y_i^* and Y_i^{**} tends to zero as the sample size increases, that the true and error parts of Y be uncorrelated. But combining (3) with the definition, $Y_{it}^{**} = Y_{it} - Y_{it}^*$ tells us that,

$$Y_{it}^{**} = (1 - \beta) \cdot E_{it} - V_{it} \tag{4}$$

The 'measurement error' is equal to a proportion of the forecast error less the disturbance term in the forecasting model. Using (4) and (3) we then find,

$$E y_{it} y_{it}^{**} = \beta(1 - \beta) \sigma_e^2 - \sigma_v^2 + (1 - \beta) \sigma_{e \cdot y_{it-1}^*}, \tag{5}$$

where the σ^2 's denote the population variance of the forecast errors and of V .⁶ $\sigma_{e \cdot y_{it-1}^*} = E(y_{it} - y_{it-1}^*) \cdot y_{it-1}^*$, and on the special assumption that forecasts and

⁴ Small letters will, throughout, denote variables measured from their means.

⁵ cf. Nerlove, (4), among many other descriptions of the model.

⁶ At the given time, t .

forecast errors are uncorrelated⁷ we may set this term at zero, and (5) reduces to

$$E y_i^* y_i^{**} = \beta(1 - \beta)\sigma_e^2 - \sigma_v^2 \quad (6)$$

Since there is no reason to suppose that (6) is equal to zero we see that there is no reason to suppose the sampling covariance of Y^* and Y^{**} tends to zero, and if it does not the result (2) is no longer correct.^{8,9}

If we continue to assume that forecasts and forecast errors are uncorrelated, a very simple alternative expression for the large sample bias in LS of C on Y may be derived. The assumptions in full are,

$$E y_{i-1}^* u_{it} = E e_{it} u_{it} = E e_{it} v_{it} = E y_{i-1}^* v_{it} = E y_{i-1}^* e_{it} = 0. \quad i = 1, \dots, n \quad (7)$$

These assumptions state that the disturbances in the two equations (1) and (3) are uncorrelated with the right hand determining variables in equation (3), and that the forecasts and forecast errors are uncorrelated. Using assumptions (7), together with (1) and (3), the numerator in the LS expression is, in expectation,

$$\begin{aligned} E \sum C_i y_{it} &= k E \sum y_{it}^* y_{it} + E \sum u_{it} y_{it}, \\ &= k E \sum (y_{i-1}^* + \beta e_{it} + v_{it}) \cdot y_{it}, \\ &= nk(\beta \sigma_e^2 + \sigma_v^2) \cdot \quad (8) \end{aligned}$$

The denominator is similarly,

$$\begin{aligned} E \sum y_{it}^2 &= E \sum (y_{it} - y_{i-1}^* + y_{i-1}^*)^2 \\ &= E \sum (y_{it} - y_{i-1}^*)^2 + E \sum y_{i-1}^{*2} \\ &= n(\sigma_e^2 + \sigma_v^2) \cdot \quad (9) \end{aligned}$$

Assuming that numerator and denominator tend in probability to (8) and (9) respectively an expression for the probability limit of the LS estimator is,

$$\text{plim } \hat{k}/k = \frac{\beta \sigma_e^2 + \sigma_v^2}{\sigma_e^2 + \sigma_v^2} \quad (10)$$

⁷ This does not seem unreasonable. In an economic context we could think of the forecasts as being generally positive and the forecast errors as being equally likely to be positive or negative for any level of the forecast.

⁸ Positive correlation between Y^* and Y^{**} could be interpreted as follows. Where current income exceeds the value forecast for the next period, that forecast is relatively high because, owing to current income being high, the forecast has itself been increased. Where current income is less than the value forecast for the next period that forecast is itself relatively low owing to its being adapted to the low current income.

⁹ A comment on Friedman's "Theory of the Consumption Function" (1) may be helpful here. In equation (1), interpreting C as consumption, Y^* as permanent income (very broadly average expected income) and setting $\alpha = 0$, with i , the operative subscript, denoting a household, we have Friedman's cross-sectional consumption function model. Interpreting Y as current ('measured') income then Y^{**} is the transitory component of income. Now if we make the assumption that each household forms its permanent income in the way described by equation (3)—an assumption Friedman did not make—then our results tell us that the transitory and permanent components of income can be expected to be correlated. It follows that the result (2), which was the basis of much of Friedman's empirical work and which he wrote as $b = kPy = k \cdot \sigma_y^* / \sigma_y \sigma_e^*$, is false. But, to repeat, Friedman did not make the assumption that each household formed its permanent income according to (3).

But by the assumption that forecasts and forecast errors are uncorrelated, the variance of the outcomes, Y_t , is equal to the sum of the variances of the forecasts Y_{t-1}^* and the forecast errors E_t . That is,

$$\sigma_y^2 = \sigma_{y^*_{t-1}}^2 + \sigma_e^2, \tag{11}$$

and dividing through by σ_y^2 we have,

$$1 = \frac{\sigma_{y^*_{t-1}}^2}{\sigma_y^2} + \frac{\sigma_e^2}{\sigma_y^2} = \rho^2 + (1 - \rho^2),$$

where ρ^2 may be interpreted as the proportion of the variance of the outcomes 'explained', in the population, by the forecasts. Our result may then be more interestingly written,

$$\text{plim } \hat{k}/k = \beta(1 - \rho^2) + \rho^2 = \rho^2(1 - \beta) + \beta. \tag{12}$$

That is to say, the proportionate LS bias is equal to one minus a weighted average of unity and the coefficient of adaptive expectations, or alternatively, to one minus a weighted average of unity and the square of the population correlation between forecasts and outcomes. The proportionate bias is, as in the errors in variables expression (2), in the direction of zero, given that $0 \leq \beta < 1$ and $\rho^2 < 1$. The size of the bias depends upon the accuracy of the forecasts, as measured by ρ^2 , and upon the proximity of β to unity. Some particular cases are,

- 1. $\rho^2 = 1$; Proportionate bias = 0,
- 2. $\rho^2 = 0$; Proportionate bias = $(1 - \beta)$,
- 3. $\beta = 1$; Proportionate bias = 0,
- 4. $\beta = 0$; Proportionate bias = $(1 - \rho^2)$.

In particular, the bias vanishes only when forecast and actual incomes are perfectly correlated or when the coefficient of expectations is unity, so that, apart from a random disturbance, the forecast of next period's income is the actual value of income in the current period.

Since (3) is only the most simple version of the model of adaptive forecasting it is worthwhile presenting the analogous results for the case in which the individuals are forecasting from a slightly more elaborate model. The obvious modification is to include a trend in the forecasts, for example, by writing,

$$Y_{it}^* = (1 + \alpha) \cdot Y_{it-1}^* + \bar{\beta}(1 + \alpha) \cdot (Y_{it} - Y_{it-1}^*) + V'_{it}. \tag{13}$$

Here the forecast is equal to $(1 + \alpha)$ times the forecast that was made for the current period, plus a proportion of the current forecast error also multiplied by the trend constant. α is assumed constant over i . For $\alpha = 0$ the model reverts to the previous one. If we make assumptions analogous to (7) except that V is replaced by V' , and in particular continue to suppose that forecasts and forecast errors are uncorrelated, we easily derive,

$$\text{plim } \hat{k}/k = (1 + \alpha)(\bar{\beta}(1 - \rho^2) + \rho^2), \tag{14}$$

which is analogous to (12). One difference between (14) and (12) is that even though $\bar{\beta}$ and ρ^2 are less than unity, the LS bias may, on the two parameter forecasting model, be away from zero if α is sufficiently great. The reason for this is most easily seen if in equation (13) we set $\bar{\beta} = 1$, which shows that in this case, on average, $Y_{it}^* = (1 + \alpha)Y_{it}$, by virtue of the multiplicative trend in the forecasts.

As investigators may attempt to analyse the problem of 'subjective' determining variables as one of 'errors in variables', they may also attempt to resolve the problem of estimating the relation between C and Y^* by the method of instrumental variables, (IV), which is a standard solution for the errors in variables problem [2, p. 165]. The question arises, on our interpretation of Y^* is it still true that instrumental variable estimation of the relation between the observables, C and Y , will provide a consistent estimator of k , as it does on the standard assumptions about the measurement error? Consider the estimator with instrument denoted by Z ,

$$\tilde{k} = \frac{\sum c_{it}z_{it}}{\sum y_{it}z_{it}} \quad (15)$$

The standard assumptions under which \tilde{k} is a consistent estimator of k are that Z be uncorrelated with Y^{**} and with U , and correlated with Y . Using (3) and (1) we may expand the estimator to find,

$$\tilde{k} = \frac{k(\sum y_{t-1z_{it}}^* + \beta \sum e_{it}z_{it} + \sum v_{it}z_{it}) + \sum u_{it}z_{it}}{\sum y_{t-1z_{it}}^* + \sum e_{it}z_{it}} \quad (16)$$

Evidently and in general it is necessary to assume that the instrument is uncorrelated with U , the random term in the relation between C and Y^* ; with V , the disturbance in the forecasting model; and with E , the forecast error. It must also be correlated with the current forecasts, Y_{t-1}^* . Where these conditions are satisfied the IV estimator will be a consistent estimator of k , under general conditions. The assumption that the instrument is uncorrelated with the forecast errors is analogous to the assumption that it is uncorrelated with the measurement errors in the usual model.

We now give some calculations which may be interpreted as estimates of the quantities β , ρ^2 and k . The data is 12 consecutive cross-sections of observations on the ordinary dividends and gross income¹⁰ of 23 British public companies. We shall work with the model described by equations (1) and (3), where C is dividends, Y income and the i subscript denotes the firm and runs from 1 to 23.¹¹ The variable Y^* , however, will not be interpreted as a firm's forecast of its next year's gross income for this would be a somewhat idiosyncratic model. Instead Y^* will be interpreted as a firm's estimate of its average anticipated gross income over some short future horizon and we shall refer to it as 'expected income' for brevity. But we shall still suppose that firms form their estimate of

¹⁰ i.e. profits before deduction of a depreciation allowance, but after tax.

¹¹ A detailed account of this application of the model will be published in the near future.

TABLE A. FORECAST ACCURACY AND LEAST SQUARES BIAS

Year		$\tilde{\beta}$	\hat{k}/\bar{k}	$\hat{\rho}^2$	C
(1950)	2	.760	.904	.60	.43
	3	.690	.909	.71	.18
	4	.559	.832	.62	1.46
	5	.310	.829	.75	.18
	6	.429	.729	.53	1.11
	7	.457	.666	.38	2.05
	8	.784	.912	.59	.87
	9	.851	.951	.67	.14
	10	.722	.896	.62	1.45
	11	1.091	1.034	.63	.44
	(1960)	12	.949	.964	.29
Mean		.691	.875	.58	

Notes: $\tilde{\beta}$ and \hat{k} are estimates of β and k formed from the coefficients of income and lagged dividends in equation (17). \hat{k} is the LS slope of dividends on current income alone in each year. C is the standardised proportionate change in mean group income in each year, taken as positive. The sample is 23 British public companies in the Electrical Goods Manufacturing Industry. Dividends is total ordinary dividends distributed in the accounting year, net of tax. Income is Trading Profit plus Non-Trading Income after tax and payment of debenture interest but before deduction of depreciation.

'expected income' according to equation (3). Hence $Y - Y_{t-1}^*$ is not now strictly a forecast error but is a deviation of current income from the level that had been expected to hold, on average, in the near future, and ρ^2 is now a measure of the correlation between current incomes and the previous expected average level.

Combining (1) and (3) we derive the relation between observable variables,

$$C_{it} = \beta\alpha + \beta k \cdot Y_{it} + (1 - \beta) \cdot C_{it-1} + W_{it} \quad (17)$$

which we shall estimate in the 11 consecutive years for which observations on the three variables involved are available. β is only assumed constant over firms, it may vary from year to year.¹² We also estimate in each of these 11 years the LS regression of dividends on current income, C on Y , giving us \hat{k} . Hence if we consider equation (12), we have computed \hat{k} , what we take as consistent estimates of β and k may be derived from the coefficients of current income and lagged dividends in (17), and so we may form an estimate of ρ^2 in each year. These estimates are listed in Table A.

The first column lists the year to year estimates of β , the second lists our

¹² In the derivation of (17) from (1), and (1) lagged one period, it is assumed that k is the same in each successive pair of years and hence constant over all 12 years. If, in fact, it changed from year to year a term in k_t/k_{t-1} enters multiplicatively into the coefficient of lagged dividends in (17). If the changes in k from year to year were relatively small, as evidence suggest they were, we may conveniently ignore this term.

Furthermore, it is possible to argue that, since the income of firms was on balance rising over these 12 years, 1949-60, the two parameter forecasting model described by (13) is more appropriate to the data. If this is so, the coefficients of income and of lagged dividends in (17) become $\tilde{\beta}(1+\alpha)k$ and $(1-\tilde{\beta})(1+\alpha)$ respectively. Where α is small relative to β and to the standard errors of the coefficient estimates, say of the order of .05 on average, there is justification for ignoring it in the determination of β and k from the estimated coefficients.

A final point is that the assumption that β , and also α if this term is included, is the same for all 23 firms in any year is a very strong one and (17) must be considered a rather crude model of the dividend-income relation.

estimates of the ratio of the LS estimator of k to k itself—one minus these values gives an estimate of the large sample LS bias in each year. Column 3 lists the derived estimates of ρ^2 , the square of the correlation between current and expected average gross income. The estimated LS bias varies from about zero to -33% with a mean of about -12% , and the estimated coefficients of adaptive expectations vary from .3 to about unity. The implied estimates of ρ^2 vary from .29 to .75 with a mean of .58.

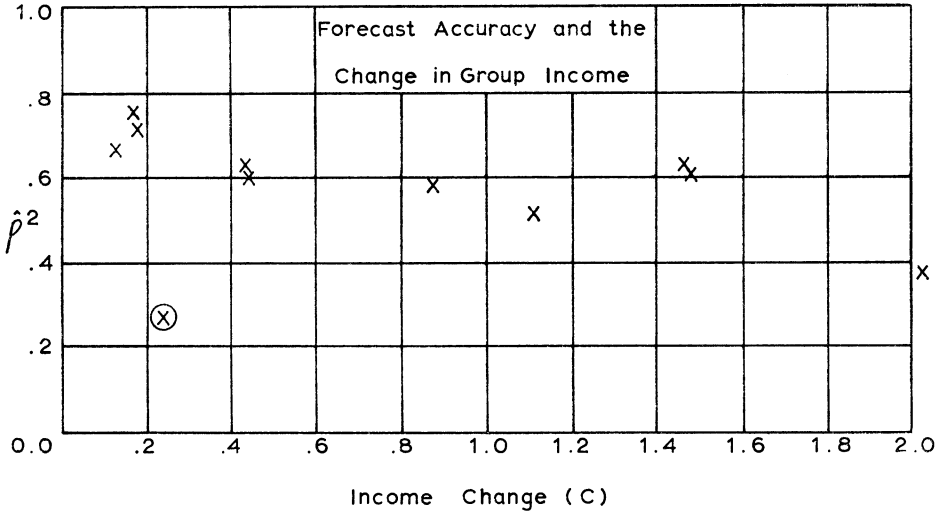
One test we might run in order to check whether our interpretation of the LS bias is at all accurate is to see whether the calculated values for ρ^2 vary from year to year as we might expect them to. ρ^2 measures the correlation between actual incomes and the levels that had been previously expected to obtain on average in the near future, according to our interpretation of the model. It seems reasonable to suppose that this correlation would be smaller in years when the income of firms, on average, has changed unusually over its level in the previous years, than when it has changed by a more normal amount. Let us suppose, in fact, that ρ^2 is a decreasing function of the extent to which the change in income over the previous year has been extreme. We measure this latter variable by taking the proportionate change in the mean income of the 23 firms in each of the 11 years, subtracting from it the mean proportionate change, and dividing through by the standard deviation of these changes. This gives us the standardized proportionate change in mean income and, taken without regard to sign, it is listed in column 4 as C . The larger is C the more 'unusual' was the change in mean income of the group relative to its experience over these 12 years. We now expect to find a negative relation between C and ρ^2 , and the data is plotted in the scatter diagram.

There appears to be a negative association between C and $\hat{\rho}^2$. If we include all 11 observations the correlation is $-.34$. If we omit the outlying observation for year 12, circled, the correlation is raised to $-.81$. The latter is significant at the 5% level by a z test, although the former is not.^{13,14} This suggests that our model is, in fact, relevant to the relation between business dividends and income, and that our interpretation of the LS bias in the regression of dividends on current income is not inaccurate.

The problem has been to examine what happens when one uses actual rather than forecast values as an explanatory variable in a Least Squares regression. This is a specification error that may arise by a mistake on the part of the investigator or by his supposing that actual levels could be used as an empirical proxy for the correct, forecast, values. On the assumption that the forecasts are generated by a very simple forecasting model, we have derived an expression for the LS bias in terms of two parameters of this model. It has been pointed out that analysis of the consequences of using actual rather than forecast values as determining variable in terms of a measurement "error in the vari-

¹³ There are, in fact, grounds for treating all results for year 12 (1960) as suspect. These have to do with the impact on the data and model of the abolition of the discriminatory tax on distributed profits in 1959.

¹⁴ It is of interest that there is quite a strong direct correlation between our estimate of the LS bias, \hat{k}/\bar{k} , and the change in income, C . Excluding year 12, the correlation is $-.66$. If we include year 12 the correlation is, in fact, better than that between C and ρ^2 . This suggests that our estimate of β is badly out in the final year.



ables" model may be inaccurate in this case, and is anyway less easy to interpret economically than the preceding result.¹⁵

Our analysis may also be carried over to situations in which the correct determining variable Y^* , is not strictly a forecast of next period's value of Y , but is something looser, for example an anticipated average value of Y , or an estimate of the 'normal' value of Y . All that is required is that this average or normal value be generated by the adaptive model (3) or some variant on this, though of course β and ρ^2 require to be reinterpreted according to the particular context.

We have also pointed out that while the analysis of the problem as one of errors in variables is likely to be incorrect, an errors in variables solution to the estimation problem, the method of instrumental variables, may still provide a consistent estimator. The basic requirement for consistency is that the instrument is uncorrelated with the current forecast errors.

Finally, an empirical illustration of the results has been given which offers support for the analysis and suggests that the model is relevant to the relation between company dividends and income.

REFERENCES

- [1] Friedman, M.: "A Theory of the Consumption Function," National Bureau of Economic Research, 1957.
- [2] Johnston, J.: "Econometric Methods," McGraw-Hill, 1963.
- [3] Liviatan, N.: "Errors in Variables and Engel Curve Analysis," *Econometrica*, 1961.
- [4] Nerlove, M.: "Distributed Lags and Demand Analysis," U. S. Dept. of Agriculture, 1958.

¹⁵ Recent textbooks of Econometrics, particularly Johnston (2), have included accounts of the classical errors in variables problem in estimation. To the extent that this inclusion is motivated by the fact that the general problem alluded to in the second paragraph of this paper can be set up, formally, as one of errors in the variables, our results suggest the emphasis is misplaced. The problem of 'subjective' determining variables in econometric equations appears to require an analysis of its own.