# The Simple Geometry of Correlated Regressors and IV Corrections

Richard J. Butler

*Brigham Young University, USA*

**Abstract:** In medical research, frequently other important determinants, correlated with the key treatment variable, are omitted from the analysis. This omission yields biased and inconsistent estimates. For example, leaving out correlated (with, say, $X_1$) determinants of $Y$ from regressions yield biased estimates of key parameters (say $\hat{\beta}_1$). Instrumental variable estimation solves this problem by constructing similar triangles to retrieve consistent estimates. This article illustrates the geometry of correlated regressor bias, and the simple IV geometric solution.

**Keywords:** Omitted variable bias (OVB), classical measurement error (CME), simultaneous equation models (SEM), instrumental variables, orthogonal projections.

## SAMPLE REGRESSION ONTO THE X1, X2 PLANE

Consider the costs of treating low-back back pain ($Y$) upon admission into a disability claim, as it varies with (the time-of-admission) self-reported level of back pain ($X_1$) and another factor ($X_2$), so that the sample regression function expressed in terms of vectors $(Y, \overline{1}, X_1, X_2, \hat{\mu})$ and the estimated sample parameters $(\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2)$ is as follows (note that $\overline{1}$ is the vector of ones associated with the intercept term, $\hat{\beta}_0$):

$$Y = \hat{\beta}_0 \overline{1} + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \hat{\mu}$$

Averaging both sides of this equation, and subtracting that resulting averaged equation from above (to eliminate the constant term), we arrive at the demeaned specification as follows:

$$Y^* = \hat{\beta}_1 X_1^* + \hat{\beta}_2 X_2^* + \hat{\mu}$$

so that $Y^* = Y - \overline{Y}\overline{1}$, $X_1^* = X_1 - \overline{X}_1 \overline{1}$, and $X_2^* = X_2 - \overline{X}_2 \overline{1}$, where (again) $\overline{1}$ is a vector of ones, and a bar above the other variables indicates its sample mean value.

When variables are deviated from their means, then their (dot) product is a covariance, and uncorrelated variables have zero covariance. Statistical "independence" for these vectors means that their covariance is zero, that their dot product is zero, and equivalently in vector spaces, that the are orthogonal or at right angles to each other. Hence, any of the three right hand side vectors $(X_1^*, X_2^*, \hat{\mu})$ orthogonal to the other two (say, $\hat{\mu}$), does not affect the projection of— that is, the regression of— $Y^*$ onto the space spanned by the remaining two vectors (say, $X_1^*, X_2^*$). When one variable is at right angles (uncorrelated) to the other variables, we can safely ignore it when doing our analysis.

With no simultaneous causality in the population regression, E($\mu$|X)=0, sample regression is fitted by taking the residuals to be orthogonal to $X_1^*, X_2^*$ this yields sample orthogonality (right angle, or uncorrelated) conditions that minimizes the sum of squared residuals:

orthogonality $X_1^* : (X_1^*)' \hat{\mu} = 0$ or
$$\sum_i X_{1i}^* \left( Y_1^* - \left( \hat{\beta}_1 X_{1i}^* + \hat{\beta}_2 X_{2i}^* \right) \right) = 0$$

orthogonality $X_2^* : (X_2^*)' \hat{\mu} = 0$ or
$$\sum_i X_{2i}^* \left( Y_i^* - \left( \hat{\beta}_1 X_{1i}^* + \hat{\beta}_2 X_{2i}^* \right) \right) = 0$$

Solving these two equations yields the predicted value of the demeaned dependent variable, namely $\hat{Y}^*$ in the $(X_1^*, X_1^*)$ plane, and simultaneously yields the estimated coefficient values $\hat{\beta}_1, \hat{\beta}_2$. Omitted Variable Bias (OVB) and Classical Measurement Error (CME) in the estimation of the effect of $X_1^*$ on output $\hat{Y}^*$ are biases arising from the omission of a correlated regressor (say the omission of $X_2^*$ from the regression in the OVB case). These biases are graphically illustrated in sections 2 and 3.

With OVB or CME, the sample error is assumed to be generated by a process that makes it orthogonal to the regression plane (i.e., generated under the assumption that E($\mu$|X)=0), and so $\hat{\mu}$ plays no role in generating either OVB or CME biases. In section 4, this assumption is dropped and E($\mu$|X)≠0 because a simultaneous equation model (SEM) process is assumed to be generating $Y^*$.

However, in terms of the geometric decomposition of $Y^*$, SEM bias is graphically similar to OVB bias and

*Address correspondence to this author at the Brigham Young University, USA; Tel: 801 422 1372; Fax: 801 422 2859; E-mail: richard_butler@byu.edu

CME bias, with $Y^*$ projected onto the space generated by the endogenous $X_1^*$ variable and the correlated (with $X_1^*$) error $\hat{\mu}$. In all three cases (OVB, CME, SEM), the instrumental variable (IV) solution is to find the appropriate instrumental variable estimator that yields a consistent estimator for $\hat{\beta}_1$ by means of similar triangles. This is illustrated in the fifth section. The final section offers an empirical example of CME and its IV solution.

## OMITTED VARIABLE BIAS, E($\hat{\mu}$ |X)=0

If the covariance between $X_{1i}^*, X_{2i}^*$ is not zero, then omitting $X_{2i}^*$ from the specification will bias the estimate of $\hat{\beta}_1$. This is the essence of both omitted variable bias and classical measurement error bias, which share a common geometric structure.

Figure **1** indicates the bias in the OLS estimate of $\hat{\beta}_1$ when we leave out $X_2^*$ from the regression specification (ignore the influence of $X_2^*$). The intuition is simple: when we properly include both $X_1^*$ and $X_2^*$, the predicted value of $Y^*$ is the parallelogram formed by $\hat{\beta}_1 X_1^*$ and $\hat{\beta}_2 X_2^*$. Leaving out $X_2^*$ means the $Y^*$ is projected directly onto the $X_1^*$ vector, rather than the $X_1^*, X_2^*$ plane.

The bias is readily seen looking directly down onto the regression plane ($Y^*$ and $\hat{Y}^*$ are perfectly aligned from this angle), as indicated in Figure **1**. Let $P_{X_1^*}$ be the projection operator that casts variable vectors into the $X_1^*$-space (here, a line) and so $P_{X_1^*} X_2^*$ is the regression of $X_2^*$ (the dependent variable) onto $X_1^*$ (the independent variable).
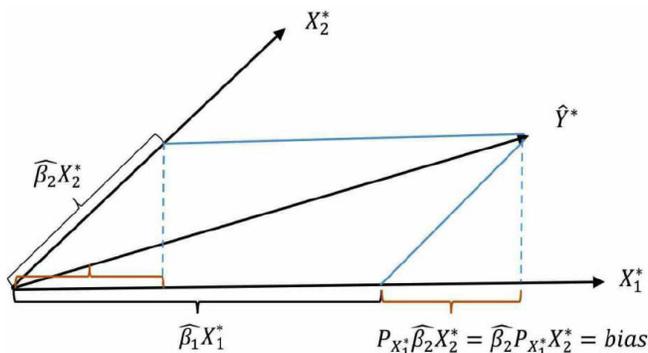


**Figure 1:** The Regression Plane View of OVB.

When $X_2^*$ is omitted from the specification, $Y^*$ is projected directly onto the $X_1^*$ line instead of the

$(X_1^*, X_2^*)$ plane. In the plane, the parallelogram provides the best estimated parameter variables, $\hat{\beta}_1, \hat{\beta}_2$ give the sample. Omitting $X_2^*$ from the specification biases upward the estimated $\hat{\beta}_1$ parameter in this example. The added OVB bias equals $\hat{\beta}_2$ times the coefficient resulting from the regression of $X_2^*$ on $X_1^*$, which is given as $P_{X_1^*} X_2^*$ in Figure **1**.

There will be no bias from omitting $X_2^*$ either when $\hat{\beta}_2$ =0, or when the $X_1^*, X_2^*$ vectors are orthogonal, and hence, uncorrelated (with covariance of zero). A formal proof of OVB in terms of projection operators is straightforward. Multiply through the model $Y^* = \hat{\beta}_1 X_1^* + \hat{\beta}_2 X_2^* + \hat{\mu}$

by $P_{X_1^*}$ to get

$$P_{X_1^*} Y^* = \hat{\beta}_1 P_{X_1^*} X_1^* + \hat{\beta}_2 P_{X_1^*} X_2^* + P_{X_1^*} \hat{\mu}$$

where the far right hand term is zero (from the orthogonality conditions), and the far left hand term is the estimated effect when omitting $X_2^*$ from the model. The first term on right hand side of the equality sign, $\hat{\beta}_1 P_{X_1^*} X_1^* = \hat{\beta}_1 X_1^*$, is the true effect; while the second hand term is the bias effect from omitting $X_2^*$.

### Recapping OVB without the Demeaning

To summarize the findings in terms of regressions with the constants reinserted into the sample estimates (before the demeaning removed the constants):

"true" for the sample: $Y = \hat{\beta}_0 \bar{1} + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \hat{\mu}$

actual sample estimates: $Y = \hat{\alpha}_0 \bar{1} + \hat{\alpha}_1 X_1 + \in$

auxiliary (partial correlation) regression:
$X_1 = \hat{\gamma}_0 \bar{1} + \hat{\gamma}_1 X_2 + \xi$

The omitted variable bias formula derived via Figure **1** above is

$$\hat{\alpha}_1 X_1^* = \hat{\beta}_1 X_1^* + \hat{\beta}_2 \hat{\gamma}_1 X_1^*$$

or as usually written,

$\hat{\alpha}_1 = \hat{\beta}_1 + \hat{\beta}_2 \hat{\gamma}_1$ (the far right hand term being the bias from omitting $X_2^*$).

## CLASSICAL MEASUREMENT ERROR (CME), $E(\hat{\mu}|X)=0$

This is a special case of omitted variable bias. Consider a simple regression model in vector format, where the reported independent variable is subject to measurement error ($\in$).

$$X_1 = X_1^t + \in$$

where $\in$ is uncorrelated with the true value of the independent variable, $X_1^t$, but is positively correlated with the reported or observed value of the independent variable, $X_1$. In the empirical example below, $X_1^t$ is the "true" value from the Roland-Morris back pain index, a self-reported measure of back-pain generally reflecting true back-pain with some error ($\in$). Substituting into the simple regression for the true value of the independent variable, $X_1^t = X_1 - \in$, into the simple regression model $Y = \hat{\beta}_0 \bar{1} + \hat{\beta}_1 X_1^t + \hat{\mu}$ and demeaning the model to get rid of the constant term, we get the following sample regression function:

$$Y^* = \hat{\beta}_1 X_1^* - \hat{\beta}_2 \in + \hat{\mu}$$

where the sample errors have zero means, so $\hat{\mu}^* = \hat{\mu}$ and $\in^* = \in$. Classical Measurement Error (CME) is a special case of OVB where we know that $X_1^*$ and $\in$ are positively correlated (so that $X_1^*$ and $-\in$ are negatively correlated and there is no need for an auxiliary regression to sign the bias), and the $\hat{\beta}_1$ is the coefficient both for the observed $X_1^*$ and the omitted variable $\in$. Regardless of the sign of $\hat{\beta}_1$, it is reorienting both $X_1^*$ and $-\in$ the same way in the regression plane. Suppose that $\hat{\beta}_1 > 0$. Then analogous to the OVB diagram, we have CME in Figure **2**.

While $\hat{\beta}_1 > 0$ in Figure **2**, if $\hat{\beta}_1 < 0$, then the diagram is symmetrically reflected on each axis, so the conclusion
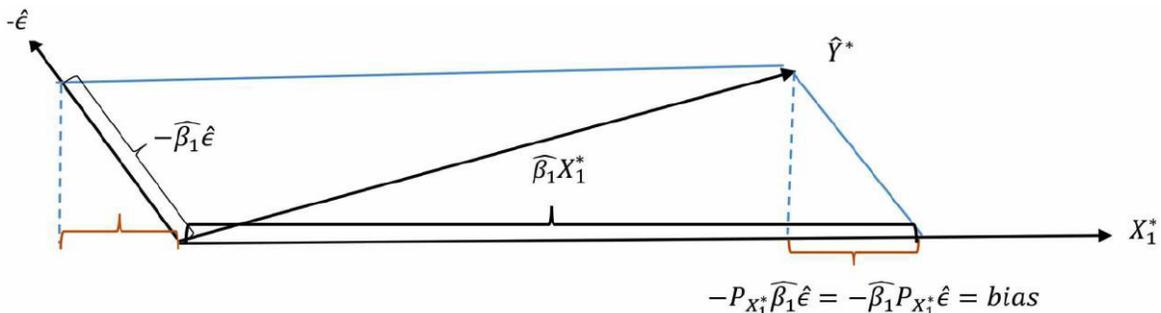
still holds: *for CME, the estimated $\hat{\beta}_1$ is biased towards zero*. CME results in "attenuation" bias.

The proof is also a simple extension of the omitted variable bias proof above. Multiply through the model $Y^* = \hat{\beta}_1 X_1^* - \hat{\beta}_1 \in + \hat{\mu}$ by $P_{X_1^*}$ to get

$$P_{X_1^*} Y^* = \hat{\beta}_1 P_{X_1^*} X_1^* - \hat{\beta}_1 P_{X_1^*} \in + P_{X_1^*} \hat{\mu}$$

where the far right hand term is zero again from the orthogonality of $X_1^*$ and the residual $\hat{\mu}$, and the far left hand term is the estimated effect when omitting $\in$ from the model (that is, from having Classical Measurement Error that is not accounted for in the regression). The first term on right hand side of the equality sign, $\hat{\beta}_1 P_{X_1^*} X_1^* = \hat{\beta}_1 X_1^*$, is the true effect; while the second hand term is the bias effect from omitting $\in$, which always causes $\hat{\beta}_1$ to be biased towards zero. (Again, the mutual scaling of $X_1^*$ and $\in$ by the same scalar $\hat{\beta}_1$ guarantees that the bias is towards zero regardless of the sign of $\hat{\beta}_1$).

## SIMULTANEOUS EQUATION MODEL (SEM), $E(\hat{\mu}|X)\neq0$

While OVB and CME obviously generate bias yielding the same geometric diagram (through the omission of an important correlated variable), it may not be obvious that the SEM shares this same geometric structure. Suppose we want to estimate the demand response for insurance purchases ($\hat{\beta}_1$) in the following simplified supply and demand setting:

Demand: $P = \hat{\beta}_0 + \hat{\beta}_1 Q + \hat{\mu}$

where *P, price* and *Q, quantity* are the key endogenous variables.

Supply: $Q = \hat{\alpha}_0 + \hat{\alpha}_1 P + \hat{\alpha}_2 \, roi + \in$



**Figure 2:** The Regression Plane View of CME.

We are only interested in estimating the demand equation (in particular, just $\hat{\beta}_1$). *roi*, return on investment, serves as the Instrumental Variable or *Z* (IV-type variable) for getting a consistent estimate of the demand equation, as it only enters into this market as an exogenous supply factor. OLS estimation of $\hat{\beta}_1$ will be inherently inconsistent because $E(\mu|X) \neq 0$, and so the (estimated) error will be correlated with *Q* the right hand side endogenous variable in the demand equation by construction: $\hat{\mu}$ determines *P* (positive covariance) from the demand equation, and *P* has a positive impact on *Q* from the supply equation (another positive covariance), so the system suggests an inherently positive correlation between $\hat{\mu}$ and *Q* in the demand equation.

The geometrical problem (bias) here is the same as that pictured in Figure **1**, with $\hat{\mu}$ replacing $X_2^*$ and with $\hat{\beta}_2 = 1$.

## SIMILAR TRIANGLES: IV GEOMETRY FIXES OBV, CME, AND SEM

To get consistent estimates of our regression coefficients when we have OVB, CME, or endogenity due to SEM, we need instrumental variables. These are one or more variables Z (our generic indicator of an identifying instrumental variable) which are uncorrelated with the sample model error, $\hat{\mu}$, correlated with the appropriate right hand side variable, *X*, and in the case of OVB and CME, *Z* is also uncorrelated with the omitted "regressor" ( $X_2^*$ and $\in$, respectively).

These *Z*s allow us to obtain consistent parameter estimates because they form similar triangles, with proportionate sides, where the factors of

proportionately will be our consistently estimated coefficients. Consider classical measurement error in Figure **1**, and the right triangles formed with the instrumental *Z*:

As in the previous CME geometric view (Figure **2** above), we employ sample values, deviate all variable values from their means to get rid of the constant term, letting $\hat{\beta}_1$ represent the ``appropriate'' value for our sample. An orthogonal projection of $Y^*$ (recall in this "overhead" view that $Y^*$ and $\hat{Y}^*$ are perfectly aligned) onto $X_1^*$ yields an estimated effect that is biased towards zero (attenuation bias), because we have failed to account for the negative covariance between $X_1^*$ and $-\hat{\in}$.

To fix the attenuation bias, and get consistent parameter values using IV (instrumental variable) estimators, we first regress the correlated regressor $X_1^*$ on the IV *Z* to get a predicted value of $X_1^*$, namely $\hat{X}_1^*$. This predicted value of $X_1^*$ is orthogonal to the measurement error, because *Z* is orthogonal to the measurement error. Then we regress $Y^*$ on this predicted value to get $\hat{\beta}_1 \hat{X}_1^*$.

This process yields similar right triangles, which recovers the appropriate $\hat{\beta}_1$. That is, the ratio of $\hat{\gamma} Z^*$ to $\hat{\beta}_1 \hat{X}_1^*$. is the same as the ratio of to $X_1^*$ to $\hat{\beta}_1 X_1^*$. From the first stage of our IV estimation, we have $Z^* = \hat{X}^* / \left[ \dfrac{\text{cov}(X^*, Z^*)}{\text{var}(Z^*)} \right] Z^*$, and in the second stage we have $\hat{Y}^* = \left[ \dfrac{\text{cov}(Y^*, Z^*)}{\text{var}(Z^*)} \right] Z^*$. Hence, the effect of $X^*$ *on* $Y^*$ on as mediated by $Z^*$ is the IV estimator:



**Figure 3:** Instrumental Variables to Correct Classical Measurement Error (CME) Bias.

$$\hat{\beta}_1^{IV} = \text{cov}(Y^*, Z^*) / \text{cov}(X^*, Z^*)$$

Since $Z$ is orthogonal to the sample residual vector $\hat{\mu}$ and in the case of OVB and CME, $Z$ is also orthogonal to the omitted factors (respectively, $X_1^*$ and $\hat{\in}$), it essentially only affects $Y^*$ through its correlation with $X_1^*$. The geometry of IV estimation makes this clearer. In the same manner, IV estimation for the SEM model, yields consistent estimates for the demand response by use of similar triangles.
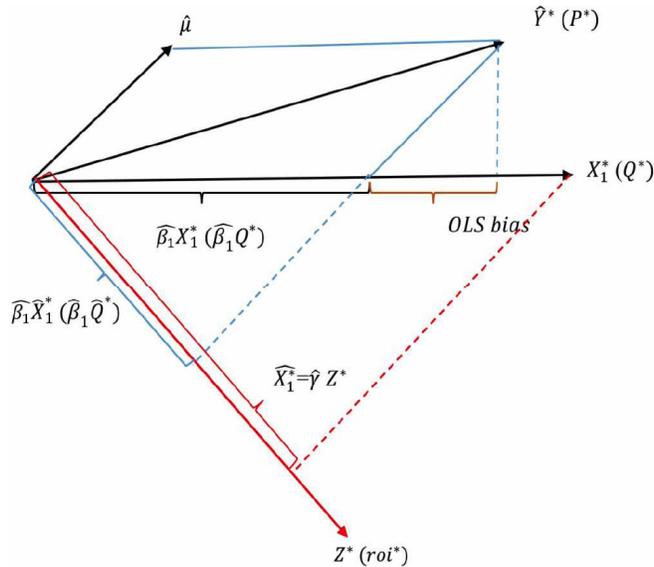


**Figure 4:** Instrumental Variables to Correct for SEM Bias.

The SEM system of the previous section suggests *roi* (the return on investment) as a natural IV or $Z$ variable. In the first stage of IV estimation, we regress $Q$ on *roi* (our $Z$ in this example). In the second stage, we regress $P$ on the predicted value of $Q$ (predicted with *roi*). This results again in similar triangles, yielding a consistent estimate of $\hat{\beta}_1$.

These IV diagrams indicate why the instrumental variable must be orthogonal to the omitted, correlated regressor: that orthogonality is used to construct similar triangles to retrieve consistent estimates. To get similar triangles, the IV $Z$ variable must be correlated with respective $X$ but uncorrelated with the correlated omitted factor (in Figure **4**, this correlated omitted factor is $\hat{\mu}$).

**EXAMPLE FROM A PROSPECTIVE LOW BACK PAIN STUDY OF CLAIM COSTS**

Low back pain (LBP) is the most costly claim type in the US workers' compensation system. In a prospective survey of LBP costs and initial patient satisfaction with their medical care, Butler and Johnson (2008) [1] collect three measures of back pain disability at the onset of LBP claims: Roland-Morris scale of back disability, physical SF12 scale of activity limitations, and a self-assessed measure of backpain (0 to 100%). The reported Roland-Morris index at claim onset is a noisy measure of back limitations, subject to measurement error. Hence, estimated effect of the Roland-Morris on subsequent workers' compensation costs likely suffers from attenuation bias. This possibility is explored in Table **1**, where total claim costs are regressed on the Roland-Morris scale and other control variables.

Consistent with research in workers' compensation, the socio-demographic variables in Table **1** have the expected signs: older workers and males have higher workers' compensation costs while those reporting "good experience" with their initial treatment have lower costs. Workers in states allowing employee choice of initial treating physician (*ee choice*=1) also have lower costs, although the results are not statistically significant.

In this sample, the Roland-Morris scale has a (normalized) mean of 48 (out of a 100), and a standard deviation of 32, with higher values indicating more disability limitations. The OLS results in the far left hand column indicates that a one standard deviation increase in scale values (32) increases costs by about 48 percent (32 *.015). Figure **2** suggests that this is biased towards zero.

The use of *SF12* and *backpain* variables as instrumental variables (IV) for the Roland-Morris index in separate analyses going from left to right, mitigates the attenuation bias as indicated graphically in Figure **3**. In the second column from the left, the *SF12* IV estimates increases costs from 48 percent to 80 percent, and for the *backpain* IV estimate from 48 percent to 58 percent for a one standard deviation increase in Roland-Morris scale values.

The right hand column in Table **1** employs both instrumental variables to predict the Roland-Morris value in the first stage (instead of one at a time, as in the middle columns). Figure **3** indicates what happens when the IV assumption of orthogonality with the measurement error fails: if the IV is even slightly negatively correlated with the error, the estimated Roland-Morris coefficient will tend to be too small in this example, and if the IV is even slightly positively correlated with the measurement error, the estimated Roland-Morris coefficient will tend to be too large. The

**Table 1:   Attenuation Bias in Roland Morris Scale for Back Disability (probsignif)**

| Variables | OLS | 2SLS:SF12 IV | 2SLS: backpain IV | 2SLS: both |
|---|---|---|---|---|
| Intercept | 5.834 | 5.4413 | 5.7401 | 5.5696 |
|  | (<.0001) | (<.0001) | (<.0001) | (<.0001) |
| Age | 0.014 | 0.019 | 0.0146 | 0.0142 |
|  | (0.0964) | (0.1269) | (0.1020) | (0.1146) |
| Male | 0.275 | 0.285 | 0.278 | 0.282 |
|  | (0.1630) | (0.1585) | (0.1602) | (0.1597) |
| Good experience | -0.495 | -0.233 | -0.432 | -0.318 |
|  | (0.0323) | (0.3525) | (0.0832) | (0.1899) |
| EE choice | -0.184 | -0.203 | -0.188 | -0.197 |
|  | (.4480) | (.4132) | (.4376) | (.4219) |
| Roland-Morris | 0.015 | 0.025 | 0.018 | 0.022 |
|  | (<.0001) | (<.0001) | (<.0001) | (<.0001) |
| R-squared | .3146 | .3208 | .2925 | .3208 |

Note: N=243. Prospective survey of completed indemnity claims for workers compensation LBP only, with IV estimates for measures of back disability (Roland-Morris, SF12, backpain rating). All specifications also included dummy variables for each firm--whose coefficients are not reported here. These firm-specific controls hold sick leave and human resource disability policies between firms constant. The firm specific effects were jointly significant at better than the one percent level. The probability significance level for the over-identification test for the specification in the far right hand specification is .077.

*over-identification test* is a test of the coherency of these alternative IV estimates: do they satisfy the appropriate orthogonality condition and estimate the same effect [2, 3]? In this case, the over-identification test for the right hand specification has a probability significance level of .077, suggesting the instruments did not both wholly satisfy the orthogonality conditions necessary as indicated in Figures **3** and **4**. (Coherent IVs will generate estimated differences that are statistically insignificant from each other in over-identification tests).

Leaving out correlated (with $X_1^*$) determinants of $Y^*$ (say correlated vectors $X_2^*$, $\in$, or $\hat{\mu}$ respectively for OVB, CME, or SEM) from our regressions yield biased estimates of key parameters ($\hat{\beta}_1$). Instrumental variable estimation solves this problem by constructing similar triangles to retrieve consistent estimates. When we have more instrumental variables than correlated regressors, the over-identification test indicates whether those IVs provide coherent estimates in the sense of identifying the same coefficient [4-7].

**Stata and SAS computer code [4-7] for this model (right hand column)**

STATA:

ivreg cost age male good_expee_choice (RM=SF12 backpain)

SAS:

procsyslin 2sls;

endogeneous cost RM;

instruments   age   male   good_expee_choice   SF12 backpain;

model cost =age male good_expee_choice RM;

run;

## REFERENCES

[1]     Butler RJ, William GJ. Satisfaction with Low Back Pain Care. The Spine Journal 2008; 8(3): 510-21. http://dx.doi.org/10.1016/j.spinee.2007.04.006

[2]     Angrist JD, Jorn-Steffen P. Mostly Harmless Econometrics. (Princeton University Press, 2009).

[3]     Parentea P, Santos S. A cautionary note on tests of overidentifying restrictions. Economics Letters, 115 (2), 314-317. http://dx.doi.org/10.1016/j.econlet.2011.12.047

[4]     Butler RJ. The FIRE Project: Econometric Applications for Finance, Insurance and Risk Management online book, with taped lectures and auxiliary material 2016: https://sites.google.com/site/fireeconometrics/

[5]     Chen Y, Becky AB. Use of instrumental variable in prescription drug research with observational data: a systematic review. Journal of Clinical Epidemiology 2011; 4(6): 687-700. http://dx.doi.org/10.1016/j.jclinepi.2010.09.006

[6]    Didelez V, Nuala S. Mendelian randomization as an instrumental variable approach to causal inference. Stat Methods Med Res 2007; 16(4): 309-330.
http://dx.doi.org/10.1177/0962280206077743

[7]    Martens EP, Pestman WR, de Boer A, Belitser SV, Klungel OH. Instrumental Variables: Application and Limitations. Epidemiology 2006; 17(3): 260-267.
http://dx.doi.org/10.1097/01.ede.0000215160.88317.cb