

Grootendorst (2007) and Deaton (1997) recount what appears to be the earliest application of the method of instrumental variables:

## The First IV Study

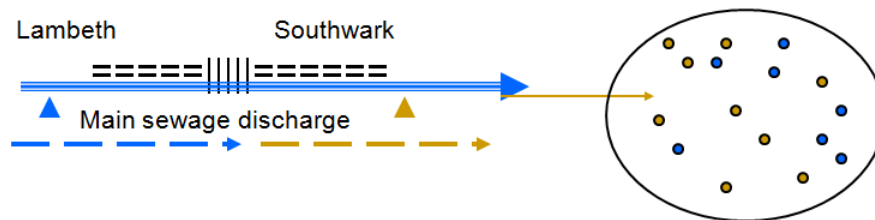
(Snow, J., *On the Mode of Communication of Cholera*, 1855)

London Cholera epidemic, ca 1853-4

Cholera =  $f(\text{Water Purity}, u) + \varepsilon$ .

- Effect of water purity on cholera?
- Purity =  $f(\text{cholera prone environment (poor, garbage in streets, rodents, etc.)})$ . Regression does not work.

Two London water companies



Although IV theory has been developed primarily by economists, the method originated in epidemiology. IV was used to investigate the route of cholera transmission during the London cholera epidemic of 1853–54. A scientist from that era, John Snow, hypothesized that cholera was waterborne. To test this, he could have tested whether those who drank purer water had lower risk of contracting cholera. In other words, he could have assessed the correlation between water purity ( $x$ ) and cholera incidence ( $y$ ). Yet, as Deaton (1997) notes, this would not have been convincing: “The people who drank impure water were also more likely to be poor, and to live in an environment contaminated in many ways, not least by the ‘poison miasmas’ that were then thought to be the cause of cholera.” Snow instead identified an instrument that was strongly correlated with water purity yet uncorrelated with other determinants of cholera incidence, both observed and unobserved. This instrument was the identity of the company supplying households with drinking water. At the time, Londoners received drinking water directly from the Thames River. One company, the Lambeth Water Company, drew water at a point in the Thames above the main sewage discharge; another, the Southwark and Vauxhall Company, took water below the discharge. Hence the instrument  $z$  was strongly correlated with water purity  $x$ . The instrument was also uncorrelated with the unobserved determinants of cholera incidence ( $y$ ). According to Snow (1855, pp. 74–75), the households served by the two companies were quite similar; indeed: “the mixing of the supply is of the most intimate kind. The pipes of each Company go down all the streets, and into nearly all the courts and alleys. ... The experiment, too, is on the grandest scale. No fewer than three hundred thousand people of both sexes, of every age and occupation, and of every rank and station, from gentlefolks down to the very poor, were divided into two groups

without their choice, and in most cases, without their knowledge; one group supplied with water containing the sewage of London, and amongst it, whatever might have come from the cholera patients, the other group having water quite free from such impurity.”

A stylized sketch of Snow’s experiment is useful for suggesting how the instrumental variable estimator works. The theory states that

$$\text{Cholera Occurrence} = f(\text{Impure Water, Other Factors}).$$

For simplicity, denote the occurrence of cholera in household  $i$  with

$$c_i = \alpha + \delta w_i + \varepsilon_i,$$

where

$c_i$  represents the presence of cholera,

$w_i = 1$  if the household has (measurably) impure water, 0 if not,

and  $\delta$  is the sought after causal effect of the water impurity on the prevalence of cholera. It would seem that one could simply compute  $d = (\bar{c} | w=1) - (\bar{c} | w=0)$ , which would be the result of a regression of  $c$  on  $w$ , to assess the effect of impure water on the prevalence of cholera. The flaw in this strategy is that a cholera prone environment,  $u$ , affects both the water quality,  $w$ , and the other factors,  $\varepsilon$ . Interpret this to say that both  $\text{Cov}(w, u)$  and  $\text{Cov}(\varepsilon, u)$  are nonzero and therefore,  $\text{Cov}(w, \varepsilon)$  is nonzero. The endogeneity of  $w$  in the equation invalidates the regression estimator of  $\delta$ . The pernicious effect of the common influence,  $u$ , works through the unobserved factors,  $\varepsilon$ . The implication is that  $E[c|w] \neq \alpha + \delta w$  because  $E[\varepsilon|w] \neq 0$ . Rather,

$$E[c|w=1] = \alpha + \delta + E[\varepsilon|w=1]$$

$$E[c|w=0] = \alpha + E[\varepsilon|w=0]$$

so,

$$E[c|w=1] - E[c|w=0] = \delta + \{E[\varepsilon|w=1] - E[\varepsilon|w=0]\}.$$

It follows that comparing the cholera rates of households with bad water to those with good water,  $E[c|w=1] - E[c|w=0]$  does not reveal only the impact of the bad water on the prevalence of cholera. It partly reveals the impact of bad water on some other factor in  $\varepsilon$  that, in turn impacts the cholera prevalence. Snow’s IV approach based on the water supplying company works as follows: Define

$$l = \begin{cases} 1 & \text{if water is supplied by Lambeth,} \\ 0 & \text{if Southwark \& Vauxhall.} \end{cases}$$

To establish the *relevance* of this instrument, Snow argued that

$$E[w|l=1] \neq E[w|l=0].$$

Snow’s theory was that water supply was the culprit, and Lambeth supplied purer water than Southwark. This can be verified observationally. The instrument is *exogenous* if

$$E[\varepsilon|l=1] = E[\varepsilon|l=0].$$

This is the theory of the instrument. Water is supplied randomly to houses. Homeowners do not even know who supplies their water. The assumption is not that the unobserved factor  $\varepsilon$  is unaffected

by the water quality. It is that the other factors, not the water quality, are present in equal measure in households supplied by the two different water suppliers. This is Snow's argument that the households supplied by the two water companies are otherwise similar. The assignment is random.

To use the instrument, we note  $E[c|l] = \delta E[w|l] + E[\varepsilon|l]$ , so

$$E[c|l=1] = \alpha + \delta E[w|l=1] + E[\varepsilon|l=1]$$

$$E[c|l=0] = \alpha + \delta E[w|l=0] + E[\varepsilon|l=0]$$

This produces an estimating equation,

$$E[c|l=1] - E[c|l=0] = \delta\{E[w|l=1] - E[w|l=0]\} + \{E[\varepsilon|l=1] - E[\varepsilon|l=0]\}$$

The second term in braces is zero if  $l$  is exogenous, which was assumed. The IV estimator is then

$$\hat{\delta} = \frac{E[c|l=1] - E[c|l=0]}{E[w|l=1] - E[w|l=0]}.$$

Note that the nonzero denominator results from the relevance condition. We can see that  $\delta$  is analogous to  $\text{Cov}(c,l)/\text{Cov}(w,l)$ , which is (8-6).

To operationalize the estimator, we will use

$$P(c|l=1) = \hat{E}(c|l=1) = \bar{c}_1 = \text{proportion of households supplied by Lambeth that have cholera,}$$

$$P(w|l=1) = \hat{E}(w|l=1) = \bar{w}_1 = \text{proportion of households supplied by Lambeth that have bad water,}$$

$$P(c|l=0) = \hat{E}(c|l=0) = \bar{c}_0 = \text{proportion of households supplied by Vauxhall that have cholera,}$$

$$P(w|l=0) = \hat{E}(w|l=0) = \bar{w}_0 = \text{proportion of households supplied by Vauxhall that have bad water.}$$

To complete this development of Snow's experiment, we can show that the estimator  $\hat{\delta}$  is an application of (8-6). Define three dummy variables,  $c_i = 1$  if household  $i$  suffers from cholera and 0 if not,  $w_i = 1$  if household  $i$  receives impure water and 0 if not, and  $l_i = 1$  if household  $i$  receives its water from Lambeth and 0 if from Vauxhall; let  $\mathbf{c}$ ,  $\mathbf{w}$  and  $\mathbf{l}$  denote the column vectors of  $n$  observations on the three variables and let  $\mathbf{i}$  denote a column of ones. For the model  $c_i = \alpha + \delta w_i + \varepsilon_i$ , we have  $\mathbf{Z} = [\mathbf{i}, \mathbf{l}]$ ,  $\mathbf{X} = [\mathbf{i}, \mathbf{w}]$  and  $\mathbf{y} = \mathbf{c}$ . The estimator is

$$\begin{pmatrix} a \\ d \end{pmatrix} = [\mathbf{Z}'\mathbf{X}]^{-1}\mathbf{Z}'\mathbf{y} = \begin{bmatrix} \mathbf{i}'\mathbf{i} & \mathbf{i}'\mathbf{w} \\ \mathbf{l}'\mathbf{i} & \mathbf{l}'\mathbf{w} \end{bmatrix}^{-1} \begin{pmatrix} \mathbf{i}'\mathbf{c} \\ \mathbf{l}'\mathbf{c} \end{pmatrix} = \begin{bmatrix} n & n\bar{w} \\ n_1 & n_1\bar{w}_1 \end{bmatrix}^{-1} \begin{pmatrix} n\bar{c} \\ n_1\bar{c}_1 \end{pmatrix} = \frac{1}{nn_1(\bar{w}_1 - \bar{w})} \begin{bmatrix} n_1\bar{w}_1 & -n\bar{w} \\ -n_1 & n \end{bmatrix} \begin{pmatrix} n\bar{c} \\ n_1\bar{c}_1 \end{pmatrix}.$$

Collecting terms,  $d = (\bar{c}_1 - \bar{c}) / (\bar{w}_1 - \bar{w})$ . Since  $n = n_0 + n_1$ ,  $\bar{c}_1 = (n_0\bar{c}_0 + n_1\bar{c}_1)/n$  and

$$\bar{c} = (n_0\bar{c}_0 + n_1\bar{c}_1)/n, \text{ so } \bar{c}_1 - \bar{c} = (n_0/n)(\bar{c}_1 - \bar{c}_0). \quad \text{Likewise, } \bar{w}_1 - \bar{w} = (n_0/n)(\bar{w}_1 - \bar{w}_0) \text{ so}$$

$d = (\bar{c}_1 - \bar{c}_0) / (\bar{w}_1 - \bar{w}_0) = \hat{\delta}$ . This estimator based on the difference in means is the Wald (1940) estimator.