

Omitted Variables¹

STA442/2101 Fall 2016

¹See last slide for copyright information.

Overview

- 1 Omitted Variables
- 2 Instrumental Variables

A Practical Data Analysis Problem

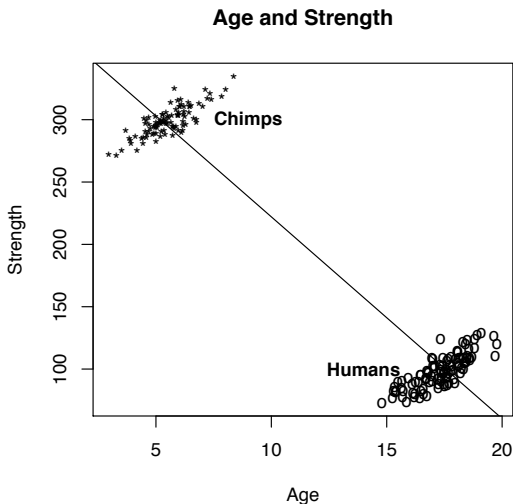
When more explanatory variables are added to a regression model and these additional explanatory variables are correlated with explanatory variables already in the model (as they usually are in an observational study),

- Statistical significance can appear when it was not present originally.
- Statistical significance that was originally present can disappear.
- Even the signs of the $\hat{\beta}$ s can change, reversing the interpretation of how their variables are related to the response variable.

An extreme, artificial example

To make a point

Suppose that in a certain population, the correlation between age and strength is $r = -0.93$.



The fixed x regression model

$$Y_i = \beta_0 + \beta_1 x_{i,1} + \cdots + \beta_k x_{i,p-1} + \epsilon_i, \text{ with } \epsilon_i \sim N(0, \sigma^2)$$

- If viewed as conditional on $\mathbf{X}_i = \mathbf{x}_i$, this model implies independence of ϵ_i and \mathbf{X}_i , because the conditional distribution of ϵ_i given $\mathbf{X}_i = \mathbf{x}_i$ does not depend on \mathbf{x}_i .
- What is ϵ_i ? *Everything else* that affects Y_i .
- So the usual model says that if the explanatory variables are random, they have *zero covariance* with all other variables that are related to Y_i , but are not included in the model.
- For observational data, this assumption is almost always violated.
- Does it matter?

Example

The explanatory variables are random.

Suppose that the variables X_2 and X_3 affect Y and are correlated with X_1 , but they are not part of the data set. The values of the response variable are generated as follows:

$$Y_i = \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \beta_3 X_{i,3} + \epsilon_i,$$

independently for $i = 1, \dots, n$, where $\epsilon_i \sim N(0, \sigma^2)$. The explanatory variables are random, with expected value and variance-covariance matrix

$$E \begin{pmatrix} X_{i,1} \\ X_{i,2} \\ X_{i,3} \end{pmatrix} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{pmatrix} \quad \text{and} \quad V \begin{pmatrix} X_{i,1} \\ X_{i,2} \\ X_{i,3} \end{pmatrix} = \begin{pmatrix} \phi_{11} & \phi_{12} & \phi_{13} \\ & \phi_{22} & \phi_{23} \\ & & \phi_{33} \end{pmatrix},$$

where ϵ_i is independent of $X_{i,1}$, $X_{i,2}$ and $X_{i,3}$.

Absorb X_2 and X_3

Since X_2 and X_3 are not observed, they are absorbed by the intercept and error term.

$$\begin{aligned} Y_i &= \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \beta_3 X_{i,3} + \epsilon_i \\ &= (\beta_0 + \beta_2 \mu_2 + \beta_3 \mu_3) + \beta_1 X_{i,1} + (\beta_2 X_{i,2} + \beta_3 X_{i,3} - \beta_2 \mu_2 - \beta_3 \mu_3 + \epsilon_i) \\ &= \beta'_0 + \beta_1 X_{i,1} + \epsilon'_i. \end{aligned}$$

And,

$$\text{Cov}(X_{i,1}, \epsilon'_i) = \beta_2 \phi_{12} + \beta_3 \phi_{13} \neq 0$$

The “True” Model

Almost always closer to the truth than the usual model, for observational data

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i,$$

where $E(X_i) = \mu_x$, $Var(X_i) = \sigma_x^2$, $E(\epsilon_i) = 0$, $Var(\epsilon_i) = \sigma_\epsilon^2$, and $Cov(X_i, \epsilon_i) = c$.

Under this model,

$$\sigma_{xy} = Cov(X_i, Y_i) = Cov(X_i, \beta_0 + \beta_1 X_i + \epsilon_i) = \beta_1 \sigma_x^2 + c$$

Estimate β_1 as usual with least squares

$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} \\ &= \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2} \\ &= \frac{\hat{\sigma}_{xy}}{\hat{\sigma}_x^2} \\ &\xrightarrow{\text{a.s.}} \frac{\sigma_{xy}}{\sigma_x^2} \\ &= \frac{\beta_1 \sigma_x^2 + c}{\sigma_x^2} \\ &= \beta_1 + \frac{c}{\sigma_x^2}\end{aligned}$$

$$\widehat{\beta}_1 \xrightarrow{a.s.} \beta_1 + \frac{c}{\sigma_x^2}$$

It converges to the wrong thing.

- $\widehat{\beta}_1$ is inconsistent.
- For large samples it could be almost anything, depending on the value of c , the covariance between X_i and ϵ_i .
- Small sample estimates could be accurate, but only by chance.
- The only time $\widehat{\beta}_1$ behaves properly is when $c = 0$.
- Test $H_0 : \beta_1 = 0$: Probability of Type I error goes almost surely to one.

All this applies to multiple regression

Of course

When a regression model fails to include all the explanatory variables that contribute to the response variable, and those omitted explanatory variables have non-zero covariance with variables that are in the model, the regression coefficients are inconsistent. Estimation and inference are almost guaranteed to be misleading, especially for large samples.

Correlation-Causation

- The problem of omitted variables is the technical version of the correlation-causation issue.
- The omitted variables are “confounding” variables.
- With random assignment and good procedure, x and ϵ have zero covariance.
- But random assignment is not always possible.
- Most applications of regression to observational data provide very poor information about the regression coefficients.
- Is bad information better than no information at all?

How about another estimation method?

Other than ordinary least squares

- Can *any* other method be successful?
- This is a very practical question, because almost all regressions with observational data have the disease.

For simplicity, assume normality

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

- Assume (X_i, ϵ_i) are bivariate normal.
- This makes (X_i, Y_i) bivariate normal.
- $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{i.i.d.}{\sim} N_2(\mathbf{m}, \mathbf{V})$, where

$$\mathbf{m} = \begin{pmatrix} m_1 \\ m_2 \end{pmatrix} = \begin{pmatrix} \mu_x \\ \beta_0 + \beta_1 \mu_x \end{pmatrix}$$

and

$$\mathbf{V} = \begin{pmatrix} v_{11} & v_{12} \\ & v_{22} \end{pmatrix} = \begin{pmatrix} \sigma_x^2 & \beta_1 \sigma_x^2 + c \\ & \beta_1^2 \sigma_x^2 + 2\beta_1 c + \sigma_\epsilon^2 \end{pmatrix}.$$

- All you can ever learn from the data are the approximate values of \mathbf{m} and \mathbf{V} .
- Even if you knew \mathbf{m} and \mathbf{V} exactly, could you know β_1 ?

Five equations in six unknowns

The parameter is $\theta = (\mu_x, \sigma_x^2, \sigma_\epsilon^2, c, \beta_0, \beta_1)$. The distribution of the data is determined by

$$\begin{pmatrix} m_1 \\ m_2 \end{pmatrix} = \begin{pmatrix} \mu_x \\ \beta_0 + \beta_1 \mu_x \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} v_{11} & v_{12} \\ & v_{22} \end{pmatrix} = \begin{pmatrix} \sigma_x^2 & \beta_1 \sigma_x^2 + c \\ & \beta_1^2 \sigma_x^2 + 2\beta_1 c + \sigma_\epsilon^2 \end{pmatrix}$$

- $\mu_x = m_1$ and $\sigma_x^2 = v_{11}$.
- The remaining 3 equations in 4 unknowns have infinitely many solutions.
- So infinitely many sets of parameter values yield the *same distribution of the sample data*.
- This is serious trouble – lack of parameter identifiability.
- *Definition:* If a parameter is a function of the distribution of the observable data, it is said to be *identifiable*.

Skipping the High School algebra

$$\theta = (\mu_x, \sigma_x^2, \sigma_\epsilon^2, c, \beta_0, \beta_1)$$

- For *any* given \mathbf{m} and \mathbf{V} , all the points in a one-dimensional subset of the 6-dimensional parameter space yield \mathbf{m} and \mathbf{V} , and hence the same distribution of the sample data.
- In that subset, values of β_1 range from $-\infty$ to $-\infty$, so \mathbf{m} and \mathbf{V} could have been produced by *any* value of β_1 .
- There is no way to distinguish between the possible values of β_1 based on sample data.
- The problem is fatal, if all you can observe is a single X and a single Y .

Details for the record

$$\theta = (\mu_x, \sigma_x^2, \sigma_\epsilon^2, c, \beta_0, \beta_1)$$

For *any* given \mathbf{m} and \mathbf{V} , all the points in a one-dimensional subset of the 6-dimensional parameter space yield \mathbf{m} and \mathbf{V} , and hence the same distribution of the sample data.

- $\mu_x = m_1$ and $\sigma_x^2 = v_{11}$ remain fixed.
- $\sigma_\epsilon^2 \geq |\mathbf{V}|/v_{11}$
- When $\sigma_\epsilon^2 = |\mathbf{V}|/v_{11}$, $\beta_1 = v_{12}/v_{11}$
- For $\sigma_\epsilon^2 > |\mathbf{V}|/v_{11}$, two values of β_1 are compatible with \mathbf{m} and \mathbf{V} .
- As σ_ϵ^2 increases, the lower β_1 goes to $-\infty$ and the upper β_1 goes to $-\infty$.
- β_0 and c are linear functions of β_1 :
 - $\beta_0 = m_2 - \beta_1 m_1$
 - $c = v_{12} - \beta_1 v_{11}$
- This set of parameter values is geometrically interesting.

Instrumental Variables (Wright, 1928)

A partial solution

- An instrumental variable is a variable that is correlated with an explanatory variable, but is not correlated with any error terms and has no direct effect on the response variable.
- Usually, the instrumental variable *influences* the explanatory variable.
- An instrumental variable is often not the main focus of attention; it's just a tool.

A Simple Example

What is the contribution of income to credit card debt?

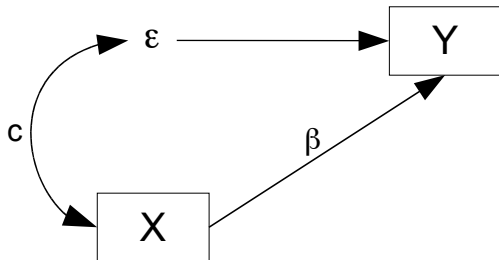
$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i,$$

where $E(X_i) = \mu_x$, $Var(X_i) = \sigma_x^2$, $E(\epsilon_i) = 0$, $Var(\epsilon_i) = \sigma_\epsilon^2$, and $Cov(X_i, \epsilon_i) = c$.

A path diagram

$$Y_i = \alpha + \beta X_i + \epsilon_i,$$

where $E(X_i) = \mu$, $Var(X_i) = \sigma_x^2$, $E(\epsilon_i) = 0$, $Var(\epsilon_i) = \sigma_\epsilon^2$, and $Cov(X_i, \epsilon_i) = c$.



Least squares estimate of β is inconsistent, and so is every other possible estimate. If the data are normal.

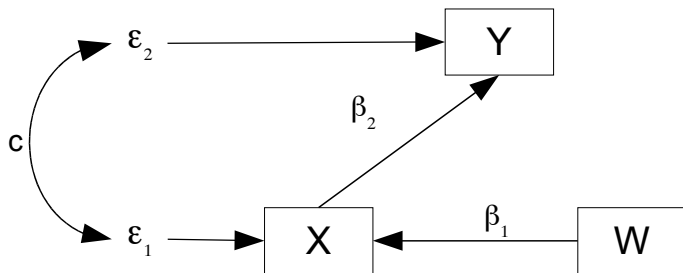
Add an instrumental variable

X is income, Y is credit card debt.

Focus the study on real estate agents in many cities. Include median price of resale home W_i .

$$X_i = \alpha_1 + \beta_1 W_i + \epsilon_{i1}$$

$$Y_i = \alpha_2 + \beta_2 X_i + \epsilon_{i2}$$



Main interest is in β_2 .

Base estimation and inference on the covariance matrix of (W_i, X_i, Y_i) : Call it $V = [v_{ij}]$

Based on $X_i = \alpha_1 + \beta_1 W_i + \epsilon_{i1}$ and $Y_i = \alpha_2 + \beta_2 X_i + \epsilon_{i2}$,

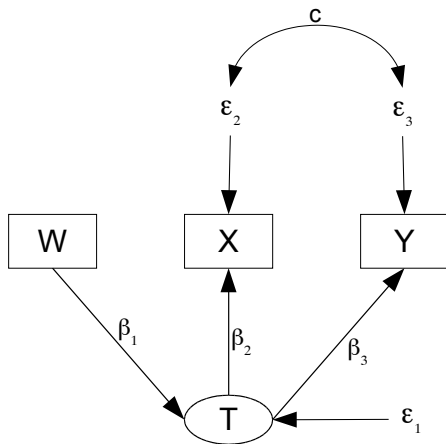
$$V = \begin{array}{c|ccc} & W & X & Y \\ \hline W & \sigma_w^2 & \beta_1 \sigma_w^2 & \beta_1 \beta_2 \sigma_w^2 \\ X & & \beta_1^2 \sigma_w^2 + \sigma_1^2 & \beta_2 (\beta_1^2 \sigma_w^2 + \sigma_1^2) + c \\ Y & & & \beta_1^2 \beta_2^2 \sigma_w^2 + \beta_2^2 \sigma_1^2 + 2\beta_2 c + \sigma_2^2 \end{array}$$

$$\beta_2 = \frac{v_{13}}{v_{12}}$$

And all the other parameters are identifiable too.

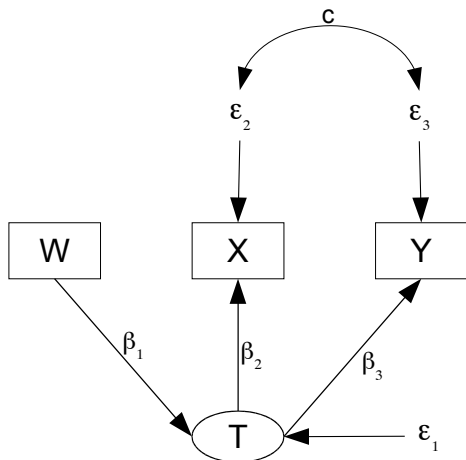
But of course there is measurement error

W is housing prices, T is true income, X is reported income, Y is credit card debt.



Main interest is in β_3 .

Read the model equations from the path diagram



$$T_i = \alpha_1 + \beta_1 W_i + \epsilon_{i,1}$$

$$X_i = \alpha_2 + \beta_2 T_i + \epsilon_{i,2}$$

$$Y_i = \alpha_3 + \beta_3 T_i + \epsilon_{i,3}$$

Covariance matrix

W is housing price, X is reported income, Y is credit card debt.

Main interest is in β_3 , which links true income to credit card debt.

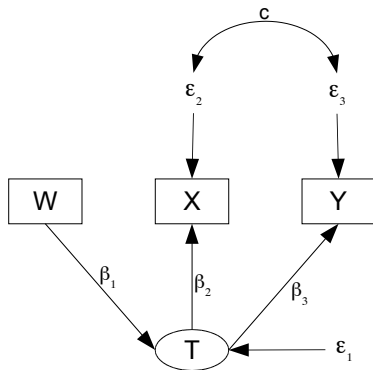
$$\begin{array}{ccc}
 w & x & y \\
 \left(\begin{array}{ccc}
 \sigma_w^2 & \beta_1\beta_2\sigma_w^2 & \beta_1\beta_3\sigma_w^2 \\
 \beta_1\beta_2\sigma_w^2 & \beta_1^2\beta_2^2\sigma_w^2 + \beta_2^2\sigma_{11} + \sigma_{22} & \beta_1^2\beta_2\beta_3\sigma_w^2 + \beta_2\beta_3\sigma_{11} + c \\
 \beta_1\beta_3\sigma_w^2 & \beta_1^2\beta_2\beta_3\sigma_w^2 + \beta_2\beta_3\sigma_{11} + c & \beta_1^2\beta_3^2\sigma_w^2 + \beta_3^2\sigma_{11} + \sigma_{33}
 \end{array} \right)
 \end{array}$$

- β_1 is the link between housing price and true income. It should be positive.
- So we can tell whether β_3 is positive, negative or zero.
- $H_0 : v_{13} = 0$ is equivalent to $H_0 : \beta_3 = 0$.
- Confidence interval?

We can do better

W is housing prices, T is true income, X is reported income, Y is credit card debt.

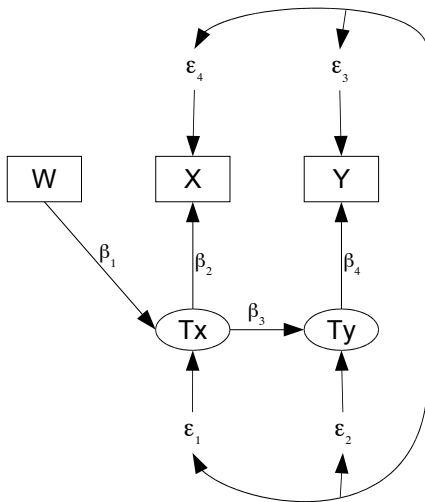
Here's the model again.



- Credit card debt is measured with error.
- True income and credit card debt have other common influences, so their error terms should be correlated.

A better model

W is housing prices, X is reported income, Y is credit card debt.



$cov(\epsilon) = [\sigma_{ij}]$. Main interest is in β_3 .

Covariance matrix of (W_i, X_i, Y_i)

W is housing price, X is reported income, Y is credit card debt.

Main interest is in β_3 , which links true income to true credit card debt.

	w	x	y
w	σ_w^2	$\beta_1\beta_2\sigma_w^2$	$\beta_1\beta_3\beta_4\sigma_w^2$
x		$\beta_1^2\beta_2^2\sigma_w^2 + \beta_2^2\sigma_{11} + 2\beta_2\sigma_{14} + \sigma_{44}$	$\beta_1^2\beta_2\beta_3\beta_4\sigma_w^2 + \beta_2\beta_3\beta_4\sigma_{11} + \beta_2\beta_4\sigma_{12} + \beta_3\beta_4\sigma_{14} + \beta_2\sigma_{13} + \beta_4\sigma_{24} + \sigma_{34}$
y			$\beta_1^2\beta_3^2\beta_4^2\sigma_w^2 + \beta_3^2\beta_4^2\sigma_{11} + 2\beta_3\beta_4^2\sigma_{12} + 2\beta_3\beta_4\sigma_{13} + \beta_4^2\sigma_{22} + 2\beta_4\sigma_{23} + \sigma_{33}$

β_1 and β_4 should be positive.

Comments

- Instrumental variables help with measurement error as well as with omitted variables.
- If there is measurement error, regression coefficients of interest are not identifiable and cannot be estimated consistently, but their signs can.
- Good instrumental variables are not easy to find.
- They will not just happen to be in the data set, except by a miracle.
- They really have to come from another universe, but still have a strong and clear effect.
- Wright's original example was tax policy for cooking oil.
- Time series applications are common.

Copyright Information

This slide show was prepared by **Jerry Brunner**, Department of Statistics, University of Toronto. It is licensed under a **Creative Commons Attribution - ShareAlike 3.0 Unported License**. Use any part of it as you like and share the result freely. The \LaTeX source code is available from the course website:
<http://www.utstat.toronto.edu/~brunner/oldclass/appliedf16>