

Multiple Regression Analysis

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + u$$

1. Estimation

Parallels with Simple Regression

- ◆ β_0 is still the intercept
- ◆ β_1 to β_k all called slope parameters
- ◆ u is still the error term (or disturbance)
- ◆ Still need to make a zero conditional mean assumption, so now assume that
- ◆ $E(u/x_1, x_2, \dots, x_k) = 0$
- ◆ Still minimizing the sum of squared residuals, so have $k+1$ first order conditions

Interpreting Multiple Regression

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_k x_k, \text{ so}$$

$$\Delta \hat{y} = \Delta \hat{\beta}_1 x_1 + \Delta \hat{\beta}_2 x_2 + \dots + \Delta \hat{\beta}_k x_k,$$

so holding x_2, \dots, x_k fixed implies that

$$\Delta \hat{y} = \Delta \hat{\beta}_1 x_1, \text{ that is each } \beta \text{ has}$$

a *ceteris paribus* interpretation

A “Partialling Out” Interpretation

Consider the case where $k = 2$, i.e.

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2, \text{ then}$$

$$\hat{\beta}_1 = \left(\sum \hat{r}_{i1} y_i \right) / \sum \hat{r}_{i1}^2, \text{ where } \hat{r}_{i1} \text{ are}$$

the residuals from the estimated

$$\text{regression } \hat{x}_1 = \hat{\gamma}_0 + \hat{\gamma}_2 \hat{x}_2$$

“Partialling Out” continued

- ◆ Previous equation implies that regressing y on x_1 and x_2 gives same effect of x_1 as regressing y on residuals from a regression of x_1 on x_2
- ◆ This means only the part of x_{i1} that is uncorrelated with x_{i2} is being related to y_i so we’re estimating the effect of x_1 on y after x_2 has been “partialled out”

Simple vs Multiple Reg Estimate

Compare the simple regression $\tilde{y} = \tilde{\beta}_0 + \tilde{\beta}_1 x_1$

with the multiple regression $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2$

Generally, $\tilde{\beta}_1 \neq \hat{\beta}_1$ unless :

$\hat{\beta}_2 = 0$ (i.e. no partial effect of x_2) OR

x_1 and x_2 are uncorrelated in the sample

Goodness-of-Fit

We can think of each observation as being made up of an explained part, and an unexplained part,

$y_i = \hat{y}_i + \hat{u}_i$ We then define the following :

$\sum (y_i - \bar{y})^2$ is the total sum of squares (SST)

$\sum (\hat{y}_i - \bar{y})^2$ is the explained sum of squares (SSE)

$\sum \hat{u}_i^2$ is the residual sum of squares (SSR)

Then $SST = SSE + SSR$

Goodness-of-Fit (continued)

- ◆ How do we think about how well our sample regression line fits our sample data?
- ◆ Can compute the fraction of the total sum of squares (SST) that is explained by the model, call this the R-squared of regression
- ◆ $R^2 = SSE/SST = 1 - SSR/SST$

Goodness-of-Fit (continued)

We can also think of R^2 as being equal to the squared correlation coefficient between the actual y_i and the values \hat{y}_i

$$R^2 = \frac{\left(\sum (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})\right)^2}{\left(\sum (y_i - \bar{y})^2\right)\left(\sum (\hat{y}_i - \bar{\hat{y}})^2\right)}$$

More about R^2 -squared

- ◆ R^2 can never decrease when another independent variable is added to a regression, and usually will increase
- ◆ Because R^2 will usually increase with the number of independent variables, it is not a good way to compare models

Assumptions for Unbiasedness

- ◆ Population model is linear in parameters:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + u$$

- ◆ We can use a random sample of size n , $\{(x_{i1}, x_{i2}, \dots, x_{ik}, y_i) : i=1, 2, \dots, n\}$, from the population model, so that the sample model is $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + u_i$

- ◆ $E(u/x_1, x_2, \dots, x_k) = 0$, implying that all of the explanatory variables are exogenous

- ◆ None of the x 's is constant, and there are no exact linear relationships among them

Too Many or Too Few Variables

- ◆ What happens if we include variables in our specification that don't belong?
- ◆ There is no effect on our parameter estimate, and OLS remains unbiased

- ◆ What if we exclude a variable from our specification that does belong?
- ◆ OLS will usually be biased

Omitted Variable Bias

Suppose the true model is given as

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u, \text{ but we}$$

estimate $\tilde{y} = \tilde{\beta}_0 + \tilde{\beta}_1 x_1 + u$, then

$$\tilde{\beta}_1 = \frac{\sum (x_{i1} - \bar{x}_1) y_i}{\sum (x_{i1} - \bar{x}_1)^2}$$

Omitted Variable Bias (cont)

Recall the true model, so that

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + u_i, \text{ so the}$$

numerator becomes

$$\sum (x_{i1} - \bar{x}_1)(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + u_i) =$$

$$\beta_1 \sum (x_{i1} - \bar{x}_1)^2 + \beta_2 \sum (x_{i1} - \bar{x}_1)x_{i2} + \sum (x_{i1} - \bar{x}_1)u_i$$

Omitted Variable Bias (cont)

$$\tilde{\beta} = \beta_1 + \beta_2 \frac{\sum (x_{i1} - \bar{x}_1)x_{i2}}{\sum ((x_{i1} - \bar{x}_1)^2)} + \frac{\sum (x_{i1} - \bar{x}_1)u_i}{\sum ((x_{i1} - \bar{x}_1)^2)}$$

since $E(u_i) = 0$, taking expectations we have

$$E(\tilde{\beta}_1) = \beta_1 + \beta_2 \frac{\sum (x_{i1} - \bar{x}_1)x_{i2}}{\sum ((x_{i1} - \bar{x}_1)^2)}$$

Omitted Variable Bias (cont)

Consider the regression of x_2 on x_1

$$\tilde{x}_2 = \tilde{\delta}_0 + \tilde{\delta}_1 x_1 \text{ then } \tilde{\delta}_1 = \frac{\sum (x_{i1} - \bar{x}_1)x_{i2}}{\sum ((x_{i1} - \bar{x}_1)^2)}$$

$$\text{so } E(\tilde{\beta}_1) = \beta_1 + \beta_2 \tilde{\delta}_1$$

Summary of Direction of Bias

	$\text{Corr}(x_1, x_2) > 0$	$\text{Corr}(x_1, x_2) < 0$
$\beta_2 > 0$	Positive bias	Negative bias
$\beta_2 < 0$	Negative bias	Positive bias

Omitted Variable Bias Summary

- ◆ Two cases where bias is equal to zero
 - $\beta_2 = 0$, that is x_2 doesn't really belong in model
 - x_1 and x_2 are uncorrelated in the sample
- ◆ If correlation between x_2 , x_1 and x_2 , y is the same direction, bias will be positive
- ◆ If correlation between x_2 , x_1 and x_2 , y is the opposite direction, bias will be negative

The More General Case

- ◆ Technically, can only sign the bias for the more general case if all of the included x 's are uncorrelated
- ◆ Typically, then, we work through the bias assuming the x 's are uncorrelated, as a useful guide even if this assumption is not strictly true

Variance of the OLS Estimators

- ◆ Now we know that the sampling distribution of our estimate is centered around the true parameter
- ◆ Want to think about how spread out this distribution is
- ◆ Much easier to think about this variance under an additional assumption, so
- ◆ Assume $\text{Var}(u/x_1, x_2, \dots, x_k) = \sigma^2$
(Homoskedasticity)

Variance of OLS (cont)

- ◆ Let \mathbf{x} stand for (x_1, x_2, \dots, x_k)
- ◆ Assuming that $\text{Var}(u|\mathbf{x}) = \sigma^2$ also implies that $\text{Var}(y|\mathbf{x}) = \sigma^2$
- ◆ The 4 assumptions for unbiasedness, plus this homoskedasticity assumption are known as the Gauss-Markov assumptions

Variance of OLS (cont)

Given the Gauss - Markov Assumptions

$$\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{SST_j(1 - R_j^2)}, \text{ where}$$

$SST_j = \sum (x_{ij} - \bar{x}_j)^2$ and R_j^2 is the R^2 from regressing x_j on all other x 's

Components of OLS Variances

- ◆ The error variance: a larger σ^2 implies a larger variance for the OLS estimators
- ◆ The total sample variation: a larger SST_j implies a smaller variance for the estimators
- ◆ Linear relationships among the independent variables: a larger R_j^2 implies a larger variance for the estimators

Estimating the Error Variance

- ◆ We don't know what the error variance, σ^2 , is, because we don't observe the errors, u_i
- ◆ What we observe are the residuals, \hat{u}_i
- ◆ We can use the residuals to form an estimate of the error variance

Error Variance Estimate (cont)

$$\hat{\sigma}^2 = \left(\sum \hat{u}_i^2 \right) / (n - k - 1) \equiv SSR / df$$

$$\text{thus, } se(\hat{\beta}_j) = \hat{\sigma} \left[SST_j (1 - R_j^2) \right]^{1/2}$$

- ◆ $df = n - (k + 1)$, or $df = n - k - 1$
- ◆ df (i.e. degrees of freedom) is the (number of observations) – (number of estimated parameters)

The Gauss-Markov Theorem

- ◆ Given our 5 Gauss-Markov Assumptions it can be shown that OLS is “BLUE”
- ◆ Best
- ◆ Linear
- ◆ Unbiased
- ◆ Estimator
- ◆ Thus, if the assumptions hold, use OLS