



Multiple Linear Regression

CIVL 7012/8012





THE UNIVERSITY OF **MEMPHIS**.

Dreamers. Thinkers. Doers.

Multiple Regression Analysis (MLR)

- Allows us to explicitly control for many factors those simultaneously affect the dependent variable
- This is important for
 - examining theories
 - assessing various policies of independent variables
- MLR can accommodate many independent variables that may be correlated with the dependent variable we can infer causality.
 - In such instances simple regression analysis may be misleading or underestimate the model strength

THE UNIVERSITY OF

MLR Motivation

- Incorporate more explanatory factors into the model
- Explicitly hold fixed other factors that otherwise would be in *u*
- Allow for more flexible functional forms
- Can take as many as independent variables



MLR Notation

THE UNIVERSITY OF

• Explains "y" in terms of $x_1, x_2, ..., x_k$



4



MLR Example-1

Wage equation







THE UNIVERSITY OF **MEMPHIS**.

MLR Example-2

• Average test score, student spending, and income $avgscore = \beta_0 + \beta_1 expend + \beta_2 avginc + u$ Average standardized test score of school Per student spending at this school Average family income of students at this school









MLR Example-4

• Example: CEO salary, sales and CEO tenure

 $\begin{array}{c} \log(salary) = \beta_0 + \beta_1 \log(sales) + \beta_2 ceoten + \beta_3 ceoten^2 + u \\ \hline \\ \text{Log of CEO salary} \\ \end{array}$

- Model assumes a constant elasticity relationship between CEO salary and the sales of his or her firm
- Model assumes a quadratic relationship between CEO salary and his or her tenure with the firm
- Meaning of linear regression
 - The model has to be linear in the parameters (not in the variables)

THE UNIVERSITY OF

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
- $E(u | x_1, x_2, ..., x_k) = 0$
- Still minimizing the sum of squared residuals, so have k+1 first order conditions



Interpreting Multiple Regression (1)

THE UNIVERSITY OF MEMPHIS.

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

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

so holding x_2, \ldots, 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



THE UNIVERSITY OF

Interpreting Multiple Regression (2)

Interpretation of the multiple regression model

$$\beta_j = \frac{\partial y}{\partial x_j} \longleftarrow$$

By how much does the dependent variable change if the j-th independent variable is increased by one unit, <u>holding all</u> other independent variables and the error term constant

- The multiple linear regression model manages to hold the values of other explanatory variables fixed even if, in reality, they are correlated with the explanatory variable under consideration
- Ceteris paribus-interpretation
- It has still to be assumed that unobserved factors do not change if the explanatory variables are changed



MLR Estimation (1)

OLS Estimation of the multiple regression model

• Random sample

THE UNIVERSITY OF

1EMPHIS.

 $\{(x_{i1}, x_{i2}, \dots, x_{ik}, y_i) : i = 1, \dots n\}$

Regression residuals

$$\widehat{u}_i = y_i - \widehat{\beta}_0 - \widehat{\beta}_1 x_{i1} - \widehat{\beta}_2 x_{i2} - \ldots - \widehat{\beta}_k x_{ik}$$

Minimize sum of squared residuals

 \boldsymbol{n}

$$\min \sum_{i=1}^{n} \widehat{u}_{i}^{2} \rightarrow \widehat{\beta}_{0}, \widehat{\beta}_{1}, \widehat{\beta}_{2}, \dots, \widehat{\beta}_{k}$$

Minimization will be carried out by computer

THE UNIVERSITY OF **MEMPHIS**

MLR Estimation (2)

- Estimates can be derived from the first order conditions
 - Properties of OLS on any sample of data
 - Fitted values and residuals

$$\widehat{y}_i = \widehat{\beta}_0 + \widehat{\beta}_1 x_{i1} + \widehat{\beta}_2 x_{i2} + \ldots + \widehat{\beta}_k x_{ik} \qquad \widehat{u}_i = y_i - \widehat{y}_i$$
Fitted or predicted values Residuals

• Algebraic properties of OLS regression



Deviations from regression line sum up to zero

$$\sum_{i=1}^{n} x_{ij} \hat{u}_i = 0$$

Correlations between deviations and regressors are

$$\bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}_1 + \ldots + \hat{\beta}_k \bar{x}_k$$

Sample averages of y and of the regressors lie on regression



Goodness-of-fit (1)

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 : $\sum (y_i - \overline{y})^2$ is the total sum of squares(SST) $\sum_{i} (\hat{y}_{i} - \overline{y})^{2}$ is the explained sum of squares(SSE) $\sum \hat{u}_i^2$ is the residual sum of squares (SSR) Then SST = SSE + SSR



THE UNIVERSITY OF MEMPHIS.

Dreamers. Thinkers. Doers.

Goodness-of-fit (2)

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

$R^2 = SSE/SST = 1 - SSR/SST$



THE UNIVERSITY OF **MEMPHIS**

Goodness-of-Fit (3)

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

www.memphis.edu





More about R-squared

• *R*² can never decrease when another independent variable is added to a regression, and usually will increase

 Because R² will usually increase with the number of independent variables, it is not a good way to compare models



THE UNIVERSITY OF

Assumptions on MLR (1)

- <u>Standard assumptions for the multiple regression model</u>
- Assumption MLR.1 (Linear in parameters)
- $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_k x_k + u \checkmark$

In the population, the relationship between y and the explanatory variables is linear

Assumption MLR.2 (Random sampling)

 $\{(x_{i1}, x_{i2}, \dots, x_{ik}, y_i): i = 1, \dots n\}$ The data is a random sample drawn from the population

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \ldots + \beta_k x_{ik} + u_i$$

Each data point therefore follows the population equation



THE UNIVERSITY OF **MEMPHIS**.

Assumptions on MLR (2)

- Standard assumptions for the multiple regression model (cont.)
- Assumption MLR.3 (No perfect collinearity)

"In the sample (and therefore in the population), none of the independent variables is constant and there are no exact relationships among the independent variables"

Remarks on MLR.3

- The assumption only rules out <u>perfect</u> collinearity/correlation between explanatory variables; imperfect correlation is allowed
- If an explanatory variable is a perfect linear combination of other explanatory variables it is superfluous and may be eliminated
- Constant variables are also ruled out (collinear with intercept)



THE UNIVERSITY OF **MEMPHIS**

Assumptions on MLR (3)

- <u>Standard assumptions for the multiple regression model (cont.)</u>
- Assumption MLR.4 (Zero conditional mean)

 $E(u_i|x_{i1}, x_{i2}, \ldots, x_{ik}) = 0 \longleftarrow$

The value of the explanatory variables must contain no information about the mean of the unobserved factors

 In a multiple regression model, the zero conditional mean assumption is much more likely to hold because fewer things end up in the error



THE UNIVERSITY OF

Assumptions on MLR (4)

- Discussion of the zero mean conditional assumption
 - Explanatory variables that are correlated with the error term are called <u>endogenous</u>; endogeneity is a violation of assumption MLR.4
 - Explanatory variables that are uncorrelated with the error term are called <u>exogenous</u>; MLR.4 holds if all explanat. var. are exogenous
 - Exogeneity is the key assumption for a causal interpretation of the regression, and for unbiasedness of the OLS estimators

Theorem 3.1 (Unbiasedness of OLS)

 $MLR.1-MLR.4 \Rightarrow E(\hat{\beta}_j) = \beta_j, \quad j = 0, 1, \dots, k$

 Unbiasedness is an average property in repeated samples; in a given sample, the estimates may still be far away from the true values

THE UNIVERSITY OF **MEMPHIS**

MLR Unbiasedness

Population model is linear in parameters: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + ... + \beta_k x_k + u$ \diamondsuit 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 $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \ldots + \beta_k x_{ik} + u_i$ $(E(u/x_1, x_2, \dots, x_k) = 0, \text{ implying that all of the }$ explanatory variables are exogenous \diamond None of the x's is constant, and there are no exact linear relationships among them

THE UNIVERSITY OF

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:

 $\beta_2 = 0$ (i.e. no partial effect of x_2) OR x_1 and x_2 are uncorrelated in the sample



THE UNIVERSITY OF **MEMPHIS**

Simple vs Multiple Reg Estimate

$\widetilde{\beta_1} = \widehat{\beta_1} + \widehat{\beta_1} \widetilde{\delta_1}$

- $\widetilde{\delta_1}$ is the slope coefficient of the regression of x_{i2} on x_{i1}
- The above equation shows that how $\widetilde{\beta_1}$ differs from partial effect of x1 on \hat{y}
- The relationship shows two distinct cases
 - The partial effect of x2 on \hat{y} is zero, i.e. $\hat{\beta}_2 = 0$; or
 - x1 and x2 are uncorrelated in the sample, i.e. $\widetilde{\delta_1} = Q_4$



Including /Omitting Irrelevant Variables

Including irrelevant variables in a regression model

 $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + u$

No problem because $E(\hat{\beta}_3) = \beta_3 = 0$. = 0 in the population

However, including irrevelant variables may increase sampling variance.

Omitting relevant variables: the simple case

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u \longleftarrow$$
 True model (contains x₁ and x₂)

 $y = \alpha_0 + \alpha_1 x_1 + w$ Estimated model (x₂ is omitted)





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

THE UNIVERSITY OF **MEMPHIS**.

Goodness-of-fit: Adjusted R-square (1)

- More on goodness-of-fit and selection of regressors
- General remarks on R-squared
 - A high R-squared does not imply that there is a causal interpretation
 - A low R-squared does not preclude precise estimation of partial effects
- Adjusted R-squared

- What is the ordinary R-squared supposed to measure?

$$R^{2} = 1 - \frac{SSR}{SST} = 1 - \frac{(SSR/n)}{(SST/n)}$$

THE UNIVERSITY OF **MEMPHIS**

Dreamers. Thinkers. Doers.



• Adjusted R-squared (cont.)

Correct degrees of freedom of nominator and denominator

- A better estimate taking into account degrees of freedom would be

$$\bar{R}^2 = 1 - \frac{(SSR/(n-k-1))}{(SST/(n-1))^2} = adjusted \ R^2$$

- The adjusted R-squared imposes a penalty for adding new regressors
- The adjusted R-squared increases if, and only if, the t-statistic of a newly added regressor is greater than one in absolute value
- Relationship between R-squared and adjusted R-squared

$$\bar{R}^2 = 1 - (1 - R^2)(n - 1)/(n - k - 1)$$
 The adjusted R-squared may even get negative

THE UNIVERSITY OF **MEMPHIS**.

Dreamers. Thinkers. Doers.

Goodness-of-fit: Adjusted R-square (3)

- Using adjusted R-squared to choose between nonnested models
 - Models are nonnested if neither model is a special case of the other

$$rdintens = \beta_0 + \beta_1 \log(sales) + u \leftarrow R^2 = .061, \bar{R}^2 = .030$$
$$rdintens = \beta_0 + \beta_1 sales + \beta_2 sales^2 + u \leftarrow R^2 = .148, \bar{R}^2 = .090$$

- A comparison between the R-squared of both models would be unfair to the first model because the first model contains fewer parameters
- In the given example, even after adjusting for the difference in degrees of freedom, the quadratic model is preferred



Incorporating Non-linearities in SLR

- Incorporating nonlinearities: Semi-logarithmic form
- Regression of log wages on years of eduction

$$\log(wage) = \beta_0 + \beta_1 educ + u$$
Natural logarithm of wage

• This changes the interpretation of the regression coefficient:

$$\beta_{1} = \frac{\partial \log(wage)}{\partial educ} = \frac{1}{wage} \cdot \frac{\partial wage}{\partial educ} = \underbrace{\frac{\partial wage}{wage}}_{\substack{wage}} \leftarrow \underbrace{\frac{\partial educ}{\partial educ}}_{\substack{wage}} \leftarrow \underbrace{\frac{\partial educ}{\partial educ}}_{\substack{wage}} \leftarrow \underbrace{\frac{\partial wage}{\partial educ}}_{\substack{wage}} \leftarrow \underbrace{$$



Incorporating Non-linearities in SLR



38

educ



Incorporating Non-linearities in SLR

- Incorporating nonlinearities: Log-logarithmic form
- CEO salary and firm sales

$$\log(salary) = \beta_0 + \beta_1 \log(sales) + u$$

Natural logarithm of CEO salary

Natural logarithm of his/her firm's sales

• This changes the interpretation of the regression coefficient:





Incorporating Non-linearities in SLR

TABLE 2.3 Summary of Functional Forms Involving Logarithms			
Model	Dependent Variable	Independent Variable	Interpretation of β ₁
Level-level	у	Х	$\Delta y = \beta_1 \Delta x$
Level-log	У	log(x)	$\Delta y = (\beta_1 / 100) \% \Delta x$
Log-level	log(y)	х	$\Delta y = (100\beta_1)\Delta x$
Log-log	log(y)	log(x)	$\Delta y = \beta_1 \Delta x$

Source: Woolridge, Introductory Econometrics: A Modern Approach



Introducing Quadratic Forms

- $y = \beta_0 + \beta_1 x + \beta_2 x^2$
- When beta1 >0; and beta2<0, then

•
$$x^* = \frac{\beta_1}{-2\beta_2}$$

THE UNIVERSITY OF



Source: Woolridge, Introductory Econometrics: A Modern Approach