# Multiple Regression Analysis

1. Motivation
2. Estimation
3. Expected Value
4. Variance

# Motivation

- Examples
- Ceteris Paribus: Public vs. Private University
  - https://www.youtube.com/watch?v=iPBV3BlV7jk&list=PL-uRhZ_p-BM5ovNRg-G6hDib27OCvcyW8&index=2
- Interpretation

Estimation **Ass?s:**

$$E(u) = 0$$

$$E(u \mid x_1, x_2, \ldots, x_k) = E(u) = 0$$

$$E(x_1 u) = 0, \quad E(x_2 u) = 0$$

$$\ldots \quad E(x_k u) = 0$$

- Model with $k$ independent variables

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

- Objective: estimate $\beta_0, \beta_1, \ldots \beta_k$

$$u = y - \beta_0 - \beta_1 x_1 - \ldots - \beta_k x_k$$

$$E(y - \beta_0 - \beta_1 x_1 - \ldots - \beta_k x_k) = 0$$

$$E[x_1(y - \beta_0 - \beta_1 x_1 - \ldots - \beta_k x_k)] = 0$$

$$\vdots$$

$$E[x_k(y - \beta_0 - \beta_1 x_1 - \ldots - \beta_k x_k)] = 0$$

# Estimation (cont.)

e.g. $x_1 \longrightarrow$ educ

$exper_i \to x_{i2}$ $x_2 \to exper$

$educ_i \longrightarrow x_{i1}$

- Sample analogs

$$n^{-1}\sum_{i=1}^{n}\left(y_i - \widehat{\beta}_0 - \widehat{\beta}_1 x_{i1} - \ldots - \widehat{\beta}_k x_{ik}\right) = 0$$

$$n^{-1}\sum_{i=1}^{n} x_{i1}\left(y_i - \widehat{\beta}_0 - \widehat{\beta}_1 x_{i1} - \ldots - \widehat{\beta}_k x_{ik}\right) = 0$$

$$\vdots$$

$$n^{-1}\sum_{i=1}^{n} x_{ik}\left(y_i - \widehat{\beta}_0 - \widehat{\beta}_1 x_{i1} - \ldots - \widehat{\beta}_k x_{ik}\right) = 0$$

$x_{ij}$ : observation $i$ for variable $x_j$

# Estimation (cont.)

$$\hat{\beta}_0, \hat{\beta}_1, \ldots, \hat{\beta}_k$$

- OLS estimates:
- Fitted value:
- Residual:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \ldots + \hat{\beta}_k x_{ik}$$

$$\hat{u}_i = y_i - \hat{y}_i$$

# Estimation (cont.)

Example

| $y$ (wage) | $x_1$ (educ) | $x_2$ (exper) |
|------------|--------------|---------------|
| 3.1        | 11           | 2             |
| 3.2        | 12           | 22            |
| 3          | 11           | 2             |
| 6          | 8            | 44            |
| 5.3        | 12           | 7             |
| 8.8        | 16           | 9             |
| 11         | 18           | 15            |
| 5          | 12           | 5             |
| 3.6        | 12           | 26            |
| 18         | 17           | 22            |

## Estimation (cont.)

$$\overset{\checkmark}{\overline{y}} - \overset{\checkmark}{\widehat{\beta}_0} - \overset{\checkmark}{\widehat{\beta}_1}\overline{x_1} - \overset{\checkmark}{\widehat{\beta}_2}\overline{x_2} = 0$$

$$\overline{x_1 y} - \widehat{\beta}_0\overline{x_1} - \widehat{\beta}_1\overline{(x_1)^2} - \widehat{\beta}_2\overline{x_1 x_2} = 0$$

$$\overline{x_2 y} - \widehat{\beta}_0\overline{x_2} - \widehat{\beta}_1\overline{x_1 x_2} - \widehat{\beta}_2\overline{(x_2)^2} = 0$$

$$6.742 - \widehat{\beta}_0 - 12.9\widehat{\beta}_1 - 15.4\widehat{\beta}_2 = 0$$

$$97.234 - 12.9\widehat{\beta}_0 - 175.1\widehat{\beta}_1 - 190.4\widehat{\beta}_2 = 0$$

$$115.064 - 15.4\widehat{\beta}_0 - 190.4\widehat{\beta}_1 - 396.8\widehat{\beta}_2 = 0$$

$$\widehat{\beta}_0 = -12.317 \qquad \widehat{\beta}_1 = 1.312 \qquad \widehat{\beta}_2 = 0.138$$

Sum/avg. value of residuals $= 0$

Correl. b/w residuals &

each explanatory

variable $= 0$

Properties of OLS

correl. b/w residuals & fitted value $= 0$

If each $x_j = \overline{x_j}$, $\hat{y} = \overline{y}$.

Goodness-of-fit

- R-squared

$$R^2 = \frac{SSE}{SST}$$

$$= 1 - \frac{SSR}{SST}$$

Non-decreasing in the number of independent variables, $k$

- Adjusted R-squared

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

As $k \uparrow$, $SSR \downarrow$ but $(n - k - 1)$ also $\downarrow$. $\bar{R}^2$ may $\downarrow$ or $\uparrow$

# Expected Value

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

n price      sq. ft.      age      #BRs      quality

- Under certain assumptions, the OLS estimators are unbiased so that

$$wage = \beta_0 + \beta_1 \, age$$
$$+ \beta_2 \, educ \qquad E(\hat{\beta}_j) = \beta_j \qquad j = 0, 1, \ldots, k$$

- Assumptions $+ \beta_3 \, exper + \ldots + \beta_k x_k + u$
  - (MLR.1) Linear in Parameters
  - (MLR.2) Random Sampling
  - ✓ (MLR.3) No Perfect Collinearity    Variation in each regressor
  - ✓ (MLR.4) Zero Conditional Mean

$$age = 6 + educ + exper$$
$$E(u \mid x_1, x_2, \ldots, x_k) = 0$$

No linear relationship among explanatory vars.

$$n \geq k + 1$$

# Expected Value (cont.)

Omitted Variable Bias

- Will You Make More Going to a Private University?
  - https://www.youtube.com/watch?v=6YrIDhaUQOE
- Population model satisfying the assumptions (MLR.1) to (MLR.4)

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$$
$$\text{wage} = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{IQ} + u$$
$$\text{bweight} = \beta_0 + \beta_1 \text{smoking} + \beta_2 \text{alcohol} + u$$

- Omitting a relevant variable

OLS estimator
of $\beta_1$ : likely
biased

$$y = \beta_0 + \beta_1 x_1 + u$$
$$\text{wage} = \beta_0 + \beta_1 \text{educ} + u$$
$$\text{bweight} = \beta_0 + \beta_1 \text{smoking} + u$$

Bias depends
on $\beta_2$
and
correl. b/w
$x_1$ and $x_2$

# Expected Value (cont.)

- Fitted values from the regression where $x_2$ is omitted $\longrightarrow \tilde{\beta}_0, \tilde{\beta}_1$.

$$\tilde{y} = \tilde{\beta}_0 + \tilde{\beta}_1 x_1$$
$$\widetilde{\text{wage}} = \tilde{\beta}_0 + \tilde{\beta}_1 \, educ$$
$$\widetilde{\text{bweight}} = \tilde{\beta}_0 + \tilde{\beta}_1 \, smoking$$

- Fitted values from the regression if $x_2$ is included $\longrightarrow \hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2$$
$$\widehat{\text{wage}} = \hat{\beta}_0 + \hat{\beta}_1 \, educ + \hat{\beta}_2 \, IQ$$
$$\widehat{\text{bweight}} = \hat{\beta}_0 + \hat{\beta}_1 \, smok + \hat{\beta}_2 \, alcoh.$$

- Fitted values from the regression of $x_2$ on $x_1$

$$\widetilde{x_2} = \tilde{\delta}_0 + \tilde{\delta}_1 x_1$$
$$\widetilde{\text{IQ}} = \tilde{\delta}_0 + \tilde{\delta}_1 \, educ$$
$$\widetilde{\text{alcohol}} = \tilde{\delta}_0 + \tilde{\delta}_1 \, smoking$$

- Relationship between $\tilde{\beta}_1$ and $\hat{\beta}_1$

$$\tilde{\beta}_1 = \hat{\beta}_1 + \hat{\beta}_2 \tilde{\delta}_1$$

- If $\beta_1$ is estimated with $x_2$ omitted

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

- Bias depends on

$$\beta_2 \ \& \ \text{correl. b/w } x_1 \text{ and } x_2$$

$$\text{Bias} = 0$$
$$\text{if } \beta_2 = 0$$
$$\text{or } \tilde{\delta}_1 = 0$$

# Expected Value (cont.)

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + u \longrightarrow \text{political activism}$$

FDI ← $x_1$

env. reg. ← $x_2$

input prices → $x_3$

infrastructure → 

more complicated derivation of bias if $u$ corr. with $x_1$ but not $x_2$ and $x_3$

- Additional explanatory variables
- Other sources of bias
- Inclusion of irrelevant regressors
  - Exercise caution
  - May affect the variance of OLS estimators

OLS estimators of all $\beta$'s biased if $x_1$ corr. w/ $x_2$ and $x_3$.

$$\text{wage} = \beta_0 + \beta_1 \text{edu} + \beta_2 IQ + \beta_3 \text{exper}$$
$$+ \beta_4 \text{discrim.}$$
$$+ \beta_5 \text{occupation}$$
$$+ u$$

measurement error in $x$ or $y$ (eg. crime)

simultaneity

$x$ → police PC

$y$ → crime PC

trade

sample selection data observed if $y > $ threshold (e.g.

# Expected Value (cont.)

- $n = 500$, reps $= 500$, corr($x_1, x_2 = 0.4$), corr($x_1, u = 0$), corr($x_2, u = 0$)
- $y = 1 + 2x_1 + x_2 + u$

# Expected Value (cont.)

- $n = 500$, reps $= 500$, corr$(x_1, x_2 = 0.4)$, corr$(x_1, u = -0.6)$, corr$(x_2, u = 0.2)$
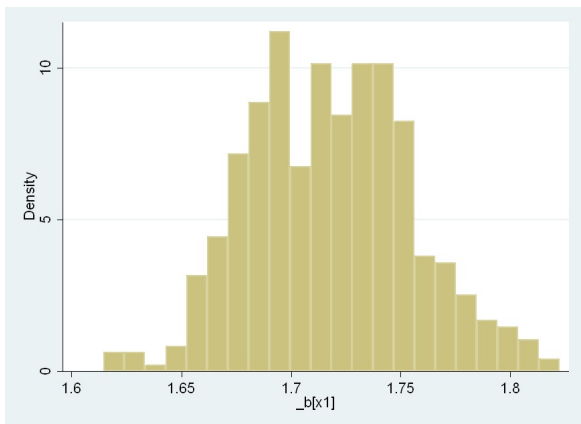- $y = 1 + 2x_1 + x_2 + u$

# Expected Value (cont.)

- $n = 500$, reps $= 500$, corr$(x_1, x_2 = 0.4)$, corr$(x_1, u = 0)$, corr$(x_2, u = 0.6)$

- $y = 1 + 2x_1 + x_2 + u$

# Variance

- Model
$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_k x_k + u$$

- (MLR.5) Homoskedasticity
$$Var(u|x_1, x_2, ..., x_k) = \sigma^2$$

- Alternatively
$$Var(y|x_1, x_2, ..., x_k) = \sigma^2$$

- If $Var(u|x_1, x_2, ..., x_k)$ depends on $x_j$ $\longrightarrow$ heteroskedasticity

- Under Assumptions MLR.1 to MLR.5 (Gauss Markov assumptions)

$$Var(\hat{\beta}_j) = \frac{\sigma^2}{SST_j(1 - R_j^2)} \text{ except } j = 0$$

- $SST_j = \sum_{i=1}^{n} (x_{ij} - \overline{x}_j)^2$
- $R_j^2$: $\longrightarrow$ $R^2$ from regression of $x_j$ on other $x$'s

- Example

$$wage = \beta_0 + \beta_1 educ + \beta_2 IQ + \beta_3 exper + u$$

$$Var(\hat{\beta}_1) = \frac{\sigma^2}{SST_{educ}(1 - R_{educ}^2)}$$

educ : avg. educ.

- $SST_{educ} = \sum_{i=1}^{n} (educ_i - \overline{educ})^2$
- $R_{educ}^2$:

$\longrightarrow$ $R^2$ from reg. of educ. on IQ and exper

- $Var(\hat{\beta}_j) \uparrow$ with $\sigma^2$ and thus *may* $\downarrow$ with additional regressors
- $Var(\hat{\beta}_j) \downarrow$ with $SST_j$ and thus likely to $\downarrow$ with $n$
- $Var(\hat{\beta}_j) \uparrow$ with $R_j^2$
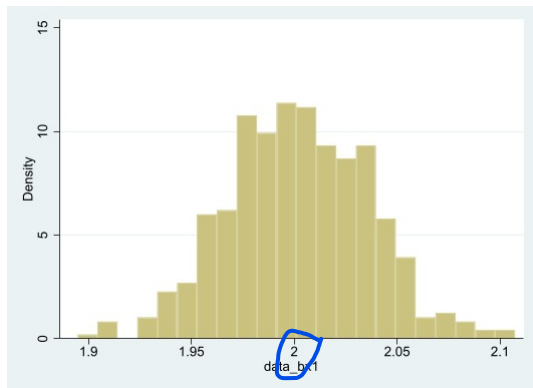  - $R_j^2$ "close" to one - the "problem" of *multicollinearity;*
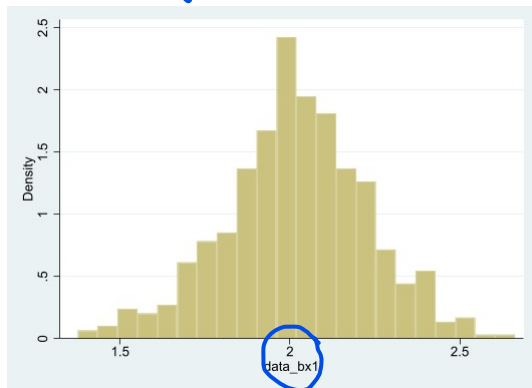  - $R_j^2 = 1$ ruled out by *does not violate*
    MLR.3 *ass^n s reqd.*

    *for unbiasedness*

# Variance (cont.)

- $n = 1000$, reps $= 500$, corr$(x_1, x_2 = 0.4)$, corr$(x_1, u = 0)$, corr$(x_2, u = 0)$
- $y = 1 + 2x_1 + x_2 + u$

# Variance (cont.)

- $n = 1000$, reps $= 500$, corr$(x_1, x_2 = 0.99)$, corr$(x_1, u = 0)$, corr$(x_2, u = 0)$
- $y = 1 + 2x_1 + x_2 + u$

# Variance (cont.)

- Additional thoughts on the inclusion of irrelevant regressors
  - May $\uparrow Var(\hat{\beta}_j)$    if $R_j^2$ high
  - Likely to $\downarrow Var(\hat{\beta}_j)$    "   "   low

# Variance (cont.)

$k$ : # of regressors

$n - (k+1)$ : $n$ − # of $\beta$'s

e.g.
SLR :
$k = 1$

- Under MLR.1 to MLR.5, $\hat{\sigma}^2$ is an unbiased estimator of $\sigma^2$

$$\hat{\sigma}^2 = (n - k - 1)^{-1} \sum_{i=1}^{n} \hat{u}_i^2$$

- $\hat{\sigma}$: standard error of the regression
- Standard error of each $\hat{\beta}_j$

$$se(\hat{\beta}_j) = \frac{\hat{\sigma}}{\sqrt{SST_j(1 - R_j^2)}} \text{ except } j = 0$$