

CPP 523: Foundations of Eval I

Regression Specification

Prof. Jesse Lecy

PRACTICE EXAM

NAME_____

Instructions: You have four hours to complete the exam once it is started. You can use notes, calculators, and statistical software. You are NOT allowed to work with anyone else, or share questions and solutions with others. Good luck!

Please give non-mathematical definitions to the following statistical concepts:

(1) The Standard Error:

On average, how far is the "best guess" from the "truth". The best guess is the statistic that we calculate from a sample - the sample mean or the slope of a regression. The truth is the population mean or slope.

(2) The 95% confidence interval of a slope:

The range over which we expect the true value of the statistic (the population statistic) to fall.

(3) R-Squared:

The proportion or percentage of the variance of the dependent variable we were able to explain with the regression model.

(4) Name three things that will reduce the standard error of a regression slope.

1. Increase the sample size.

2. Decrease the residual of Y (add control variables uncorrelated with the policy variable of interest).

3. Increase the variance of X (only possible if designing a sampling framework or an experiment where you have control over X – instead of caffeine dosages from 0 to 100mg, give dosages from 0 to 500mg – variance of X will increase as a result).

(5) Name two sins of the Seven Sins that will always increase the standard error of a regression slope.

1. Measurement error in the dependent variable.

2. High multicollinearity (adding a control variable that is highly-correlated with X).

(6) Control variables that are uncorrelated with our policy variable will not cause omitted variable bias if we do not include them in a regression. Why do we include them in the model? For example, Teacher Quality in the Classroom Size example from class.

They explain more of Y, or stated differently account for more of the residual of Y. As a result, including them will decrease the standard errors, which makes confidence intervals smaller, improving our likelihood of achieving statistically significant results and increasing the efficiency of our estimates.

- (7) Calculate the slope and the intercept for a simple bivariate regression model
 ($Y = b_0 + b_1X + e$) from the following information:

$$\begin{array}{ll} \bar{x}: & 4 \\ \bar{y}: & 11 \\ \text{var}(x): & 3 \\ \text{var}(y): & 7 \\ \text{cov}(x,y): & 6 \end{array}$$

$$b_1 = \text{cov}(x,y)/\text{var}(x) = 6 / 3 = 2$$

$$b_0 = 11 - 2(4) = 3$$

- (8) Now using the slope and intercept that you calculated above, calculate the residual (prediction error – column e) for the following three cases.

<u>X</u>	<u>Y</u>	<u>\hat{Y}</u>	<u>e</u>
5	15	$3 + (2)(5) = 13$	2
6	11	$3 + (2)(6) = 15$	-4
8	21	$3 + (2)(8) = 19$	2

(9) Consider the following regression:

$$Y = \beta_0 + \beta_1 X_1 + \varepsilon$$

Consider the case where $B_1 = 6$, and $SE_{B_1} = 2.61$.

Using $t=2.58$, calculate the 99% confidence interval for B_1 . Is the slope statistically significant at this level? How do we know?

$$6 \pm (2.58)(2.61): \\ -0.7 < B_1 < 12.7$$

It is **NOT** significant because it contains zero, meaning we can't tell for certain whether the intervention has a positive or negative impact on the outcome (at a 99% level of confidence).

(10) Consider the following regression results.

Table VI Multiple regression results on channel satisfaction

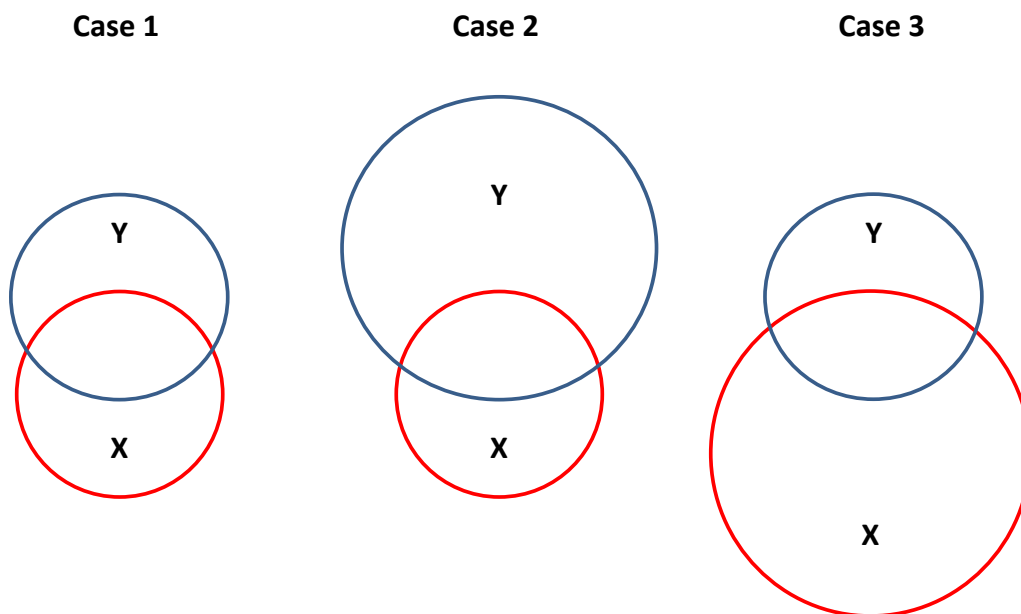
Independent variables	Beta	T-statistic	Significance
Constant	0	-0.76	0.4540
Relative performance	0.75	5.34	0.0001*
Experience	0.15	1.04	0.3063
Control	0.44	4.98	0.0560**
Changeability	0.14	0.80	0.4319
Uncertainty	0.04	0.31	0.7592
Monitoring	0.12	0.78	0.4406
Intermediate mode	0.21	2.58	0.5638
Hierarchy mode	-0.02	-0.06	0.9511

Notes: $R^2 = 0.52$; $F = 18.23$; $n = 45$; * = $P \leq 0.01$ (1 tailed t-test); ** = $P \leq 0.10$ (1 tailed t-test)

What is the largest level of confidence you can choose for the confidence interval around the slope estimate for the **Experience** coefficient before it crosses zero?

$$1 - 0.3063 = 0.6937 \text{ confidence interval}$$

(11) Consider the following cases:



- a. Holding $cov(x,y)$ constant across all cases, which case(s) will have the largest standard error?

$$SE \sim \text{residual}(y) / \text{var}(x)$$

$$\text{Case 1: } SE \sim a / b$$

$$\text{Case 2: } SE \sim A / b \ll \text{CASE 2}$$

$$\text{Case 3: } SE \sim a / B$$

- b. Holding $cov(x,y)$ constant across all cases, which case(s) will have the smallest slope?

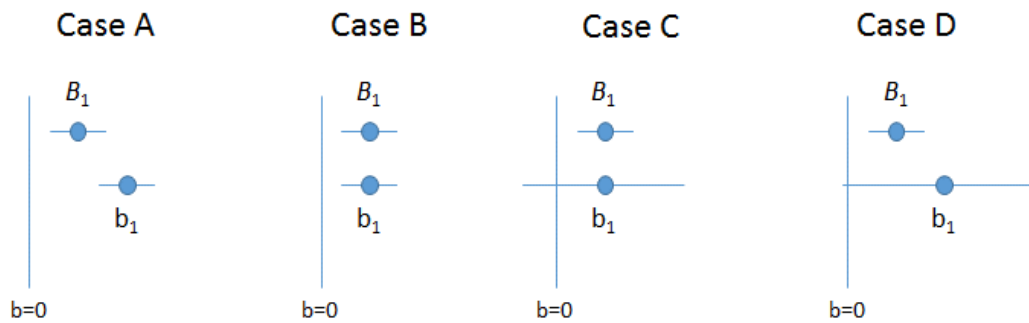
$$\text{slope} \sim \text{cov}(x,y) / \text{var}(x)$$

Covariance is the same across all cases, so Case 3 has the biggest $\text{var}(x)$ and thus the smallest slope.

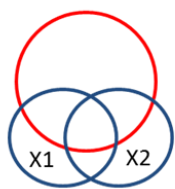
- (12) Consider four cases below. The full regression in this case is with X1 being the policy variable:

$$Y = \beta_0 + \beta_1 X1 + \beta_2 X2 + \beta_3 X3 + e$$

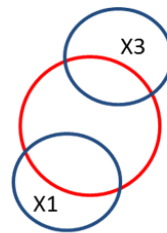
Write the correct case letter under each Venn diagram.



Unbiased but inefficient



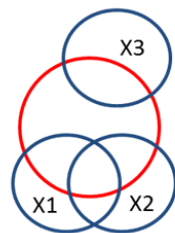
Case C



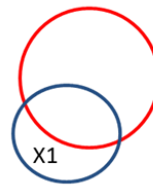
Case A

Biased but efficient

Unbiased and efficient



Case B



Case D

Biased and inefficient

BONUS (3pts): Go back to the model that attempts to discern the effects of class size on test scores:

DV: Test-Scores	Model 1	Model 2	Model 3	Model 4	Model 5
Constant	141.675*** (4.571)	-179.058*** (1.609)	141.479*** (4.601)	141.553*** (4.596)	-179.725*** (1.575)
Class Size	-0.433*** (0.021)	-0.468*** (0.003)		-0.377 (0.207)	-0.267*** (0.029)
Quality of Instruction		62.125*** (0.285)			62.169*** (0.278)
Socio-Economic Status			43.420*** (2.076)	5.649 (20.853)	20.344*** (2.918)
R-squared	0.307	0.986	0.305	0.307	0.986
N	1000	1000	1000	1000	1000

Now think about another model:

$$SES = \pi_0 + \pi_1 ClassSize + e$$

What is the exact slope for the regression of **SES** on **Class Size**, π_1 ? Show your math.

BONUS (4 pts): Think back to the model that we have studied looking at the relationship between classroom size and test scores:

$$TestScore = \beta_0 + \beta_1 ClassSize + \beta_2 SES + \beta_3 TeacherQuality + \varepsilon \quad (1)$$

Now think about a different way to run the regression model. What if we constructed it in the following way:

$$ClassSize = \pi_0 + \pi_1 SES + e_1 \quad (2)$$

$$TestScore = b_0 + b_1 e_1 + b_2 SES + b_3 TeacherQuality + \gamma \quad (3)$$

In this case the e_1 in model (3) is the residual term from model (2). Using a Venn diagram to justify your response, answer the following questions:

Does $b_1 = \beta_1$?

Does $b_2 = \beta_2$?

Does $b_3 = \beta_3$?

Does $\varepsilon = \gamma$?