

# Pset4

Gabriel Wang

December 8, 2025

## 1 Part 1: Reading

```
1  # -----
2  # 1. Create one simulated dataset and print summary
3  # -----
4
5  n0 <- 1000  # any moderate sample size
6
7  C <- rnorm(n0)
8  X <- 0.6*C + rnorm(n0)
9  Y <- 1.0*X + 0.4*C + rnorm(n0)
10
11 # Fit the true data-generating model
12 true_model <- lm(Y ~ X + C)
13 stargazer(true_model)
14 > print(summary(true_model))
15
16 Call:
17 lm(formula = Y ~ X + C)
18
19 Residuals:
20      Min       1Q   Median       3Q      Max
21 -3.2747 -0.6593 -0.0103  0.7197  3.1937
22
23 Coefficients:
24             Estimate Std. Error t value Pr(>|t|)
25 (Intercept)  0.002381   0.031652   0.075    0.94
26 X            1.009568   0.031083  32.480 <2e-16 ***
27 C            0.326033   0.037769   8.632 <2e-16 ***
28 ---
29 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
30
31 Residual standard error: 1.001 on 997 degrees of freedom
32 Multiple R-squared:  0.6554,    Adjusted R-squared:  0.6547
33 F-statistic: 947.9 on 2 and 997 DF,  p-value: < 2.2e-16
```

Table 1:

	Dependent variable:
	Y
X	1.004*** (0.011)
C	0.407*** (0.013)
Constant	0.009 (0.011)
Observations	8,000
R <sup>2</sup>	0.677
Adjusted R <sup>2</sup>	0.677
Residual Std. Error	0.986 (df = 7997)
F Statistic	8,380.917*** (df = 2; 7997)
Note:	*p<0.1; **p<0.05; ***p<0.01

### Question a

```

1 # Sample sizes we want to study
2 n_vals <- c(50, 200, 500, 1000)
3
4 # A data frame to store all simulation results
5 results <- data.frame()
6
7 # Number of simulations for each sample size
8 R <- 2000
9
10 for (n in n_vals) {
11
12   # A vector to store the 2000 estimated coefficients for this n
13   beta_hats <- numeric(R)
14
15   for (r in 1:R) {
16
17     # Confounder C
18     C <- rnorm(n)
19
20     # Treatment X depends on C (confounding relationship)
21     X <- 0.6 * C + rnorm(n)
22
23     # Outcome Y follows the true data-generating model:
24     # Y = 1.0*X + 0.4*C + error
25     Y <- 1.0 * X + 0.4 * C + rnorm(n)
26
27     # Estimate the linear model and store only the coefficient on X (1)
28     beta_hats[r] <- coef(lm(Y ~ X + C))[2]
29   }
30
31   # Combine results into a long data frame for plotting
32   results <- rbind(
33     results,
34     data.frame(
35       n = rep(n, R),
36       beta_hat = beta_hats
37     )
38   )
39 }
40

```

```

41 # -----
42 # Plot sampling distributions for different sample sizes
43 # -----
44
45 library(ggplot2)
46
47 ggplot(results, aes(x = beta_hat)) +
48   geom_density(fill = "skyblue", alpha = 0.5) +
49
50   # One panel per sample size
51   facet_wrap(~ n, scales = "free", ncol = 3) +
52
53   # True value of 1 (red dashed line)
54   geom_vline(xintercept = 1, color = "red", linetype = "dashed", size = 1) +
55
56   theme_minimal() +
57   labs(
58     title = "Sampling Distribution of 1 as Sample Size Increases",
59     subtitle = "Red dashed line = true 1 = 1.0",
60     x = "Estimated 1",
61     y = "Density"
62   )
63
64   ## proof of CLT
65 # Step 1: summarize simulated sampling distribution
66 results_summary <- results %>%
67   group_by(n) %>%
68   summarize(
69     mean_beta = mean(beta_hat),
70     sd_beta = sd(beta_hat)
71   )
72
73 # Step 2: compute theoretical standard deviation = 1 / sqrt(n)
74 sd_theoretical <- data.frame(
75   n = n_vals,
76   sd_theoretical = 1 / sqrt(n_vals)
77 )
78
79 # Step 3: merge simulated results with theoretical SD
80 summary_combined <- results_summary %>%
81   left_join(sd_theoretical, by = "n")
82
83 # Step 4: print with stargazer
84 stargazer(summary_combined,
85           summary = FALSE,
86
87           title = "Simulated vs. Theoretical Standard Deviations of Beta1",
88           rownames = FALSE)

```

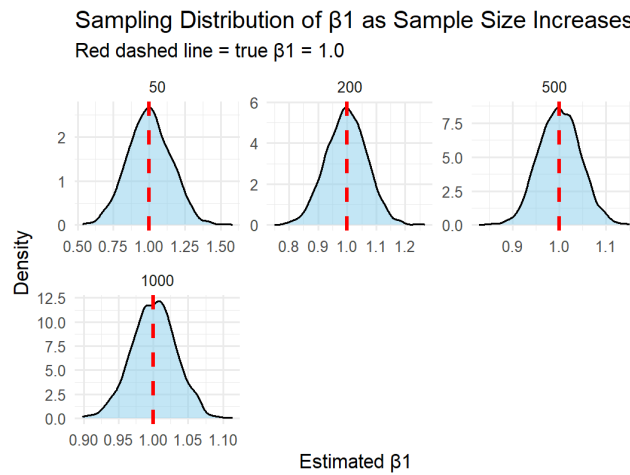


Figure 1: Distribution as N increases

Table 2: Simulated vs. Theoretical Standard Deviations of Beta1

n	mean_beta	sd_beta	sd_theoretical
50	0.999761980621645	0.144303258648783	0.14142135623731
200	0.999088603861918	0.0723473587392907	0.0707106781186548
500	1.00062695672353	0.0454422380872791	0.0447213595499958
1000	1.00023682102431	0.0314158881354247	0.0316227766016838

The simulated standard deviation  $sd\_beta$  closely matches the theoretical value  $\frac{1}{\sqrt{N}}$ . The mean of the coefficient is also converging to 1. This confirms that the estimator's variance decreases at the rate predicted by the Central Limit Theorem.

### Question b

```

1 # -----
2 # 2. Bootstrap standard error
3 # -----
4
5 B <- 2000 # number of bootstrap replications
6 boot <- numeric(B)
7
8 for (b in 1:B) {
9
10   # Resample rows with replacement
11   resample <- sample(1:n, size = n, replace = TRUE)
12
13   # Fit model on bootstrap sample
14   boot_model <- lm(Y[resample] ~ X[resample] + C[resample])
15
16   # Store coefficient on treatment variable X
17   boot[b] <- coef(boot_model)[2]
18 }
19
20 # Bootstrapped standard error of 1
21 boot_se <- sd(boot)
22 boot_se

```

The boot se is 0.03130765.

### Question c

```

1 for (n in n_vals) {

```

```

2
3 # A vector to store the 2000 estimated coefficients for this n
4 beta_hats_omitted <- numeric(R)
5
6 for (r in 1:R) {
7
8   # Confounder C
9   C <- rnorm(n)
10
11   # Treatment X depends on C (confounding relationship)
12   X <- 0.6 * C + rnorm(n)
13
14   # Outcome Y follows the true data-generating model:
15   # Y = 1.0*X + 0.4*C + error
16   Y <- 1.0 * X + 0.4 * C + rnorm(n)
17
18   # Estimate the linear model and store only the coefficient on X (1)
19   beta_hats_omitted [r] <- coef(lm(Y ~ X))[2]
20 }
21
22 # Combine results into a long data frame for plotting
23 results_omitted <- rbind(
24   results_omitted,
25   data.frame(
26     n = rep(n, R),
27     beta_hats_omitted = beta_hats_omitted
28   )
29 )
30 }
31
32 # -----
33 # Plot sampling distributions for different sample sizes
34 # -----
35
36 ggplot(results_omitted, aes(x = beta_hats_omitted)) +
37   geom_density(fill = "skyblue", alpha = 0.5) +
38
39   # One panel per sample size
40   facet_wrap(~ n, scales = "free", ncol = 3) +
41
42   # True value of 1 (red dashed line)
43   geom_vline(xintercept = 1, color = "red", linetype = "dashed", size = 1) +
44
45   theme_minimal() +
46   labs(
47     title = "Sampling Distribution of 1 with omitted confounding variable as Sample Size
48             Increases",
49     subtitle = "Red dashed line = true 1 = 1.0",
50     x = "Estimated 1 with omitted confounding variable",
51     y = "Density"
52   )

```

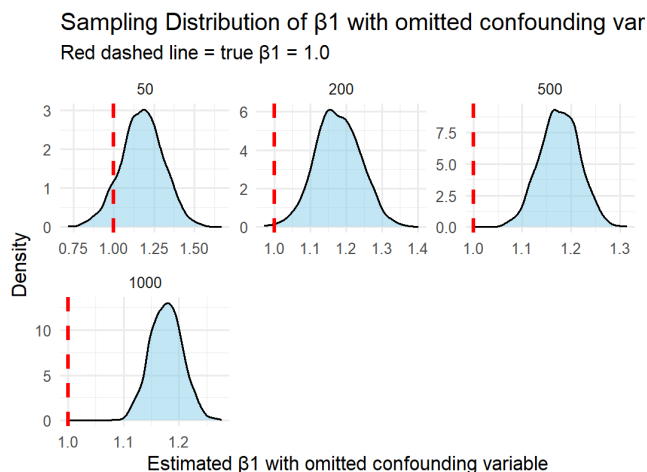


Figure 2: Omitting Confounders

When the confounding variable  $C$  is omitted, the sampling distribution of the treatment coefficient becomes centered around a biased value rather than the true parameter  $\beta_1 = 1$ . In the data-generating process,

$$X = 0.6C + u, \quad Y = 1.0X + 0.4C + \varepsilon,$$

the omitted-variable bias formula is

$$\hat{\beta}_1^{omit} = \beta_1 + \beta_2 \frac{\text{Cov}(X, C)}{\text{Var}(X)}.$$

Since  $\beta_1 = 1$ ,  $\beta_2 = 0.4$ , and  $\frac{\text{Cov}(X, C)}{\text{Var}(X)} = 0.6$ , the expected biased coefficient is

$$\hat{\beta}_1^{omit} = 1 + 0.4(0.6) = 1.24.$$

This matches the simulation results: the sampling distribution is centered around approximately 1.24 instead of 1.0. Although the distribution becomes narrower as  $n$  increases, it converges to the *wrong* value.

This implies that statistical tests based on this biased sampling distribution (t-tests, confidence intervals, p-values) will be misleading. The estimator becomes precise but incorrect, showing why the confounder must be included for valid inference.

## 2 Part 2: Simulation

### Question 1

```
1 data <- read_dta("C:/WindowsD/NYU/summer/quant1/lab/assignment/ 9/gss.dta")
2 View(data)
3 library(stargazer)
4 library(sandwich)
5 library(tidyverse)
6 data$educ
7 dim(data)
8 # if the individual's highest school year completed is larger than 16, I code him
   higheduc =1
9 data$higheduc <- ifelse(data$educ > 16, 1, 0)
10 range(data$educ)
11 difference_mean <- t.test(gun ~ higheduc, data)
12 print(difference_mean)
13 > print(difference_mean)
14
15      Welch Two Sample t-test
16
17 data:  gun by higheduc
18 t = -0.99138, df = 277.6, p-value =
19 0.3224
20 alternative hypothesis: true difference in means between group 0 and group 1 is not equal
   to 0
21 95 percent confidence interval:
```

```

22 -0.10560930 0.03486502
23 sample estimates:
24 mean in group 0 mean in group 1
25 0.3372694 0.3726415

```

Using the 2018 General Social Survey, I test whether gun ownership differs between highly educated individuals (16+ years of schooling) and others.

**Hypotheses:**

$$H_0 : \mu_{\text{high}} = \mu_{\text{low}}, \quad H_1 : \mu_{\text{high}} \neq \mu_{\text{low}}.$$

**Choice of test:** I use a two-sample t-test because the population variance is unknown. Although the sample size is large, the t-test is the standard approach and is asymptotically equivalent to the z-test.

**Significance level:**  $\alpha = 0.05$ .

**Results:** The test yields a p-value of 0.3224. Since  $p > 0.05$ , I fail to reject the null hypothesis.

**Interpretation:** Statistically, there is no evidence of a difference in mean gun-ownership rates across high and low education groups. Substantively, education does not appear to predict gun ownership in this sample.

## Question 2

```

1
2 m_lpm <- lm(data = data, gun ~ educ + income + educ:income)
3 # Fit a Linear Probability Model (LPM) predicting 'gun' using:
4 #   - educ: education level
5 #   - income: respondent income
6 #   - educ:income: interaction between education and income
7 # The model estimates how education, income, and their interaction affect gun policy
  support.
8
9 vm_lpm <- vcovHC(m_lpm, type = 'HC1')
10 # Compute heteroskedasticity-consistent (HC1) robust -variancecovariance matrix
11 # This generates robust standard errors to correct for heteroskedasticity in the LPM.
12
13 stargazer(m_lpm, se = list(sqrt(diag(vm_lpm))))
14 # Produce a formatted regression table using stargazer.
15 # The 'se' option replaces the default standard errors with the robust SEs
16 # (square roots of the diagonal elements of vm_lpm).
17
18 # calculate the pvalue
19 t_values <- coef(m_lpm) / sqrt(diag(vm_lpm))
20 print(t_values)
21 > print(t_values)
22 (Intercept)      educ      income
23 0.5260547    2.8729479    4.6404480
24 educ:income
25 -3.5376800

```

Table 3:

<i>Dependent variable:</i>	
	gun
educ	0.015*** (0.005)
income	0.005*** (0.001)
educ:income	−0.0003*** (0.0001)
Constant	0.037 (0.071)
Observations	1,443
R <sup>2</sup>	0.043
Adjusted R <sup>2</sup>	0.041
Residual Std. Error	0.467 (df = 1439)
F Statistic	21.564*** (df = 3; 1439)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

Using the same data, I fit a linear probability model predicting gun ownership from education, income, and their interaction.

#### Education

The coefficient for education is 0.015, meaning that each additional year of schooling is associated with a 1.5 percentage-point increase in the probability of owning a gun, holding income constant. The standard error is 0.005, indicating a very precise estimate. The t-value (approximately 2.87) implies that the coefficient is statistically different from zero. The p-value is below 0.01, so the effect is statistically significant.

#### Income

The coefficient for income is 0.005, meaning that higher income is associated with a higher probability of gun ownership, all else equal. The standard error is 0.001, also indicating high precision. The t-value (about 4.64) shows strong statistical significance, and the p-value is below 0.01.

#### Interaction: Education × Income

The interaction term has a coefficient of −0.0003. This negative sign indicates that the effect of education on gun ownership becomes weaker as income increases, and similarly, the effect of income becomes weaker at higher levels of education. The standard error is 0.0001, suggesting a precise estimate. The t-value ( $|t| = 3.54$ ) shows that the coefficient differs significantly from zero, and the p-value is below 0.01, confirming statistical significance.