

# Handout 7: Multiple Regression

EC 282: Introduction to Econometrics

Spring 2026

**Instructions:** Run the provided R code and answer the questions. Show your work for calculations. This handout introduces a **new dataset**—we move from cross-country comparisons to individual-level U.S. wage data.

## 1 Setup: The Current Population Survey (CPS) Wage Data

How much does an extra year of education increase your hourly wage? The answer seems simple—until you realize that education is correlated with experience, gender, and many other factors. We use data from the 1976 Current Population Survey (CPS) to explore these issues.

```
1 # Install the wooldridge package (run once)
2 # install.packages("wooldridge")
3
4 library(wooldridge)
5 library(ggplot2)
6
7 options(scipen = 999)
8
9 # Load the wage data
10 data(wage1)
11
12 cat("Number of workers:", nrow(wage1), "\n")
13 cat("Variables:", paste(names(wage1)[1:12], collapse = ", "), "\n")
14
15 # Key variables:
16 # wage      = average hourly earnings ($)
17 # educ      = years of education
18 # exper     = years of potential experience
19 # tenure    = years with current employer
20 # female    = 1 if female, 0 if male
21 # nonwhite  = 1 if nonwhite, 0 if white
22 # married   = 1 if married
23 # northcen, south, west = region dummies (northeast is omitted)
24
25 summary(wage1[, c("wage", "educ", "exper", "tenure",
26                  "female", "nonwhite", "married")])
```

**Question 1.1:** Examine the summary statistics. What is the average hourly wage? What is the average education level? What fraction of the sample is female?

**Question 1.2:** Create a scatterplot of wages against education. What pattern do you see?

```
1 ggplot(wage1, aes(x = educ, y = wage)) +  
2   geom_jitter(color = "steelblue", alpha = 0.5, width = 0.3) +  
3   geom_smooth(method = "lm", color = "red", se = FALSE) +  
4   labs(title = "Hourly Wage vs. Years of Education",  
5         x = "Years of Education",  
6         y = "Hourly Wage ($)") +  
7   theme_minimal()
```

## 2 Starting Point: Simple Regression

We begin with a simple regression of wages on education—the “naive” estimate of the return to schooling.

**Question 2.1:** Run the simple regression and examine the output.

```
1 reg1 <- lm(wage ~ educ, data = wage1)
2 summary(reg1)
```

- Write down the estimated equation:  $\widehat{\text{Wage}}_i = \hat{\beta}_0 + \hat{\beta}_1 \times \text{Educ}_i$
- Interpret  $\hat{\beta}_1$  in a sentence. What does it mean economically?
- Is  $\hat{\beta}_1$  statistically significant? At what level?
- What is the  $R^2$ ? What fraction of wage variation is explained by education alone?

**Question 2.2:** Think about what is in the error term  $u_i$ . Name at least two variables that are:

- Correlated with education, AND
- Likely affect wages

For each, predict the *direction* of the omitted variable bias. Does it make  $\hat{\beta}_1$  too large or too small?

## 3 Multiple Regression and Partial Effects

Now we add control variables to address omitted variable bias. The key idea:  $\hat{\beta}_1$  in a multiple regression measures the effect of education on wages **holding other variables constant**.

**Question 3.1:** Add experience and tenure to the regression.

```
1 reg2 <- lm(wage ~ educ + exper + tenure, data = wage1)
2 summary(reg2)
```

- (a) Write down the estimated equation.
- (b) Interpret the coefficient on **educ**. How does the interpretation differ from the simple regression?
- (c) Interpret the coefficient on **tenure**. What does “holding education and experience constant” mean in practice?
- (d) Compare the education coefficient to Model 1. Did it go up or down? Were you expecting this direction?

**Question 3.2:** The coefficient on education *increased* when we added controls (from  $\approx 0.54$  to  $\approx 0.60$ ). This is surprising—usually we expect omitted variable bias to make the simple regression coefficient too large. Verify the omitted variable bias formula for the two-variable case:

$$\tilde{\beta}_1 = \hat{\beta}_1 + \hat{\beta}_2 \times \hat{\alpha}_1$$

where  $\hat{\alpha}_1$  is the slope from regressing `exper` on `educ`.

```

1 # Step 1: Run the auxiliary regression
2 aux_reg <- lm(exper ~ educ, data = wage1)
3 alpha1_hat <- coef(aux_reg)["educ"]
4 cat("alpha1_hat (exper ~ educ):", round(alpha1_hat, 4), "\n")
5
6 # Step 2: Check the correlation
7 cat("Correlation(educ, exper):", round(cor(wage1$educ,
8     wage1$exper), 4), "\n")
9
10 # Step 3: Verify the OVB formula (for the 2-variable case)
11 # Simple regression: wage ~ educ
12 beta1_tilde <- coef(reg1)["educ"]
13
14 # Two-variable regression: wage ~ educ + exper
15 reg_two <- lm(wage ~ educ + exper, data = wage1)
16 beta1_hat <- coef(reg_two)["educ"]
17 beta2_hat <- coef(reg_two)["exper"]
18
19 cat("\nbeta1_tilde (simple):", round(beta1_tilde, 4), "\n")
20 cat("beta1_hat (multiple):", round(beta1_hat, 4), "\n")
21 cat("beta2_hat * alpha1_hat:", round(beta2_hat * alpha1_hat, 4),
22     "\n")
23 cat("beta1_hat + beta2_hat * alpha1_hat:",
24     round(beta1_hat + beta2_hat * alpha1_hat, 4), "\n")

```

- What is the sign of  $\hat{\alpha}_1$ ? What does this tell you about the relationship between education and experience?
- What is the sign of  $\hat{\beta}_2$  (the effect of experience on wages)?
- Using the OVB formula, explain *why* the simple regression *underestimates* the return to education.
- Does the formula hold exactly? Verify numerically.

**Question 3.3:** Now add demographic controls.

```
1 reg3 <- lm(wage ~ educ + exper + tenure + female + nonwhite,  
2           data = wage1)  
3 summary(reg3)
```

- (a) Interpret the coefficient on `female`. What does it mean in economic terms?
- (b) Is the coefficient on `nonwhite` statistically significant? What does this tell us?
- (c) How did the education coefficient change compared to Model 2? Why might this be?

## 4 Measures of Fit: SER, RMSE, and Adjusted $R^2$

As we add variables, we need better tools to evaluate model fit than  $R^2$  alone.

**Question 4.1:** Compute measures of fit for each model.

```

1 models <- list(reg1, reg2, reg3)
2 model_names <- c("(1) Educ only", "(2) + Exper, Tenure",
3                 "(3) + Female, Nonwhite")
4
5 cat("Model Comparison:\n")
6 cat(sprintf("%-25s %8s %8s %8s %8s\n",
7             "Model", "R2", "Adj R2", "SER", "RMSE"))
8 cat(paste(rep("-", 60), collapse = ""), "\n")
9
10 for (i in 1:3) {
11   s <- summary(models[[i]])
12   n <- nrow(models[[i]]$model)
13   k <- length(coef(models[[i]])) - 1
14   RSS <- sum(residuals(models[[i]])^2)
15   SER <- sqrt(RSS / (n - k - 1))
16   RMSE <- sqrt(RSS / n)
17   cat(sprintf("%-25s %8.4f %8.4f %8.4f %8.4f\n",
18             model_names[i], s$r.squared, s$adj.r.squared, SER, RMSE))
19 }

```

- As we add variables,  $R^2$  always increases. Verify this from the table.
- Does Adjusted  $R^2$  also always increase? Why is Adjusted  $R^2$  a better measure than  $R^2$  for comparing models with different numbers of regressors?
- The SER tells you the typical size of prediction errors, in the same units as  $Y$  (dollars per hour). How large is it in Model 3?

**Question 4.2:** Now add an irrelevant variable (`numdep` = number of dependents) that shouldn't affect wages directly.

```

1 reg_irrel <- lm(wage ~ educ + exper + tenure + numdep,
2                data = wage1)
3
4 cat("Model 2:      R2 =", round(summary(reg2)$r.squared, 4),
5     " Adj R2 =", round(summary(reg2)$adj.r.squared, 4), "\n")
6 cat("+ numdep:    R2 =",
7     round(summary(reg_irrel)$r.squared, 4),
8     " Adj R2 =",
9     round(summary(reg_irrel)$adj.r.squared, 4), "\n")

```

What happens to  $R^2$  vs. Adjusted  $R^2$  when you add a variable that doesn't belong? Why does this demonstrate the advantage of Adjusted  $R^2$ ?

## 5 Dummy Variables and the Dummy Variable Trap

**Question 5.1:** Create a male dummy and try to include both female and male in the regression.

```

1 # Create male dummy
2 wage1$male <- 1 - wage1$female
3
4 # Try to include both
5 reg_trap <- lm(wage ~ educ + exper + female + male, data = wage1)
6 summary(reg_trap)

```

- What does R report for the coefficient on male? Why?
- Write down the mathematical reason: if  $\text{Female}_i + \text{Male}_i = 1$  for all  $i$ , show algebraically why this creates a problem.
- What is the general rule for including dummy variables for a categorical variable with  $k$  categories?

**Question 5.2:** The dataset has region indicators: northcen, south, and west (the omitted category is Northeast). Run the full model with region dummies.

```

1 reg4 <- lm(wage ~ educ + exper + tenure + female + nonwhite +
2           married + northcen + south + west, data = wage1)
3 summary(reg4)

```

- Which region is the **reference category**? How do you know?
- Interpret the coefficient on south. What comparison is being made?
- Are any of the region dummies individually statistically significant at the 5% level?

## 6 Multicollinearity

What happens when regressors are highly correlated with each other? Let's explore with two experiments.

**Question 6.1: Perfect multicollinearity.** In this dataset, a worker's age is approximately:  $\text{Age} = \text{Educ} + 6 + \text{Exper}$  (assuming school starts at age 6). Create this variable and try to include it in the regression.

```

1 # Construct age (exact linear function of educ and exper)
2 wage1$age <- wage1$educ + 6 + wage1$exper
3
4 # Try to include all three
5 reg_perfect_mc <- lm(wage ~ educ + exper + age, data = wage1)
6 summary(reg_perfect_mc)

```

What happens? Why can't R estimate all three coefficients?

**Question 6.2: Imperfect multicollinearity.** Now add random noise to age (simulating measurement error in age). This breaks the exact linear relationship but keeps the variables highly correlated.

```

1 set.seed(42) # For reproducibility
2 wage1$age_noisy <- wage1$age + rnorm(nrow(wage1), 0, 2)
3
4 cat("Correlation(age_noisy, age):",
5     round(cor(wage1$age_noisy, wage1$age), 4), "\n\n")
6
7 # Run the regression with the noisy age variable
8 reg_noisy <- lm(wage ~ educ + exper + age_noisy, data = wage1)
9 summary(reg_noisy)

```

- R can now estimate all coefficients. But look at the standard errors on `educ` and `exper`—how do they compare to Model 2 (which had no multicollinearity)?
- Are the coefficients on `educ` and `exper` still statistically significant?
- What happened to the magnitudes of the coefficients? Are they reasonable?

**Question 6.3:** Compute the Variance Inflation Factor (VIF) to formally diagnose multicollinearity. Recall:  $VIF(\hat{\beta}_j) = \frac{1}{1-R_j^2}$ , where  $R_j^2$  is the  $R^2$  from regressing  $X_j$  on all other regressors.

```
1 # Manual VIF function
2 compute_vif <- function(model) {
3   vars <- attr(terms(model), "term.labels")
4   dat <- model$model
5   vifs <- numeric(length(vars))
6   names(vifs) <- vars
7   for (i in seq_along(vars)) {
8     f <- as.formula(paste(vars[i], "~",
9       paste(vars[-i], collapse = " + "))
10    r2 <- summary(lm(f, data = dat))$r.squared
11    vifs[i] <- 1 / (1 - r2)
12  }
13  return(round(vifs, 2))
14 }
15
16 cat("VIF for Model 2 (no multicollinearity):\n")
17 print(compute_vif(reg2))
18
19 cat("\nVIF for noisy age model (severe multicollinearity):\n")
20 print(compute_vif(reg_noisy))
```

Which variables have  $VIF > 5$ ? What does this tell you? Would you keep `age_noisy` in the model?

## 7 Hypothesis Testing and F-Tests

**Question 7.1:** Using Model 4 (the full model), test whether the coefficient on education is statistically significant.

```

1 # Extract info for education coefficient
2 beta_educ <- coef(reg4)["educ"]
3 se_educ   <- summary(reg4)$coefficients["educ", "Std. Error"]
4 t_stat    <- beta_educ / se_educ
5 p_value   <- 2 * pnorm(-abs(t_stat))
6
7 cat("beta_hat (educ):", round(beta_educ, 4), "\n")
8 cat("SE:", round(se_educ, 4), "\n")
9 cat("t-statistic:", round(t_stat, 3), "\n")
10 cat("p-value:", p_value, "\n")
11 cat("95% CI: [", round(beta_educ - 1.96 * se_educ, 4), ",",
12     round(beta_educ + 1.96 * se_educ, 4), "]\n")

```

- State  $H_0$  and  $H_A$ .
- Is education significant at the 5% level? At the 1% level?
- Interpret the 95% confidence interval. Does it include zero?

**Question 7.2: The F-Test.** None of the three region dummies is individually significant at the 5% level. But are they *jointly* significant? That is, do we gain anything by including all three region dummies together?

$$H_0 : \beta_{\text{northcen}} = \beta_{\text{south}} = \beta_{\text{west}} = 0 \quad (\text{regions don't matter})$$

$$H_A : \text{At least one } \beta_j \neq 0 \quad (\text{at least one region matters})$$

```

1 # Restricted model (without region dummies)
2 reg_restricted <- lm(wage ~ educ + exper + tenure + female +
3                     nonwhite + married, data = wage1)
4
5 # Unrestricted model (with region dummies)
6 reg_unrestricted <- reg4 # already estimated above
7
8 # R-squared values
9 R2_unr <- summary(reg_unrestricted)$r.squared
10 R2_res <- summary(reg_restricted)$r.squared
11
12 cat("R2 (unrestricted):", round(R2_unr, 4), "\n")
13 cat("R2 (restricted):", round(R2_res, 4), "\n")
14

```

```
15 # F-test by hand
16 q <- 3      # number of restrictions
17 k <- 9      # regressors in unrestricted model
18 n <- nrow(wage1)
19
20 F_stat <- ((R2_unr - R2_res) / q) / ((1 - R2_unr) / (n - k - 1))
21 p_value_F <- 1 - pf(F_stat, q, n - k - 1)
22
23 cat("\nF-statistic (by hand):", round(F_stat, 4), "\n")
24 cat("Critical value (5%, q=3):", round(qf(0.95, q, n-k-1), 4),
25     "\n")
26 cat("p-value:", round(p_value_F, 4), "\n")
27
28 # Verify with R's built-in anova()
29 cat("\nanova() output:\n")
30 print(anova(reg_restricted, reg_unrestricted))
```

- (a) Compute the F-statistic by hand using the  $R^2$  formula. Show your work.
- (b) Compare it to the critical value. Do you reject  $H_0$  at the 5% level?
- (c) Explain the paradox: no individual region dummy is significant at 5%, yet the F-test rejects at 5%. Why can individual t-tests miss joint significance?

## 8 Putting It All Together

**Question 8.1:** Fill in the comparison table for our four main models.

	Model 1 Educ only	Model 2 + Exper, Tenure	Model 3 + Female, Nonwhite	Model 4 Full
$\hat{\beta}_{\text{educ}}$				
$SE(\hat{\beta}_{\text{educ}})$				
$R^2$				
Adjusted $R^2$				
SER				

```

1 cat("Side-by-side comparison:\n\n")
2 models_all <- list(reg1, reg2, reg3, reg4)
3 names_all <- c("(1) Educ", "(2)+Exp,Ten",
4               "(3)+Fem,NW", "(4) Full")
5
6 for (i in 1:4) {
7   s <- summary(models_all[[i]])
8   n <- nrow(models_all[[i]]$model)
9   k <- length(coef(models_all[[i]])) - 1
10  RSS <- sum(residuals(models_all[[i]])^2)
11  cat(names_all[i], "\n")
12  cat("  educ coef:", round(coef(models_all[[i]])["educ"], 4), "\n")
13  cat("  SE(educ):",
14      round(s$coefficients["educ", "Std. Error"], 4), "\n")
15  cat("  R2:", round(s$r.squared, 4), "\n")
16  cat("  Adj R2:", round(s$adj.r.squared, 4), "\n")
17  cat("  SER:", round(sqrt(RSS / (n - k - 1)), 4), "\n\n")
18 }

```

**Question 8.2:** Based on the table, discuss:

- How does the education coefficient change as we add controls? What does this tell you about omitted variable bias in the simple regression?
- Which model would you choose as your “preferred” specification? Why?
- Even with all these controls, can we claim that education *causes* higher wages? What concerns remain?

**Question 8.3:** Visualize the gender wage gap, controlling for education.

```
1 ggplot(wage1, aes(x = educ, y = wage, color = factor(female,
2     labels = c("Male", "Female")))) +
3   geom_jitter(alpha = 0.4, width = 0.3) +
4   geom_smooth(method = "lm", se = FALSE, linewidth = 1.2) +
5   labs(title = "Wage vs. Education by Gender",
6     x = "Years of Education", y = "Hourly Wage ($)",
7     color = "Gender") +
8   theme_minimal()
```

What does this plot reveal about the gender wage gap? Is it constant across education levels, or does it change?