

Week 2: Multiple Regression, Test Exercise 2

Coursera/Erasmus U., Econometric Methods and Applications

Anthony Nguyen

Notes:

- See website for how to submit your answers and how feedback is organized.
- For parts (e) and (f), you need regression results discussed in Lectures 2.1 and 2.5.

Goals and skills being used:

- Use matrix methods in the econometric analysis of multiple regression.
- Employ matrices and statistical methods in multiple regression analysis.
- Give numerical verification of mathematical results.

Questions:

This test exercise is of a theoretical nature. In our discussion of the F-test, the total set of explanatory factors was split in two parts. The factors in X_1 are always included in the model, whereas those in X_2 are possibly removed. In questions (a), (b), and (c) you derive relations between the two OLS estimates of the effects of X_1 on y , one in the large model and the other in the small model. In parts (d), (e), and (f), you check the relation of question (c) numerically for the wage data of our lectures.

We use the notation of Lecture 2.4.2 and assume that the standard regression assumptions A1-A6 are satisfied for the unrestricted model. The restricted model is obtained by deleting the set of g explanatory factors collected in the last g columns X_2 of X . We wrote the model with $X = (X_1 X_2)$ and corresponding partitioning of the OLS estimator b in b_1 and b_2 as $y = X_1\beta_1 + X_2\beta_2 + \varepsilon = X_1b_1 + X_2b_2 + e$. We denote by b_R the OLS estimator of β_1 obtained by regressing y on X_1 , so that $b_R = (X_1'X_1)^{-1}X_1'y$. Further, let $P = (X_1'X_1)^{-1}X_1'X_2$.

- (a) Prove that $E(b_R) = \beta_1 + P\beta_2$.
- (b) Prove that $\text{var}(b_R) = \sigma^2(X_1'X_1)^{-1}$.
- (c) Prove that $b_R = b_1 + Pb_2$.

Now consider the wage data of Lectures 2.1 and 2.5. Let y be log-wage (500×1 vector), and let X_1 be the (500×2) matrix for the constant term and the variable 'Female'. Further let X_2 be the (500×3) matrix with observations of the variables 'Age', 'Educ' and 'Parttime'. The values of b_R were given in Lecture 2.1, and those of b in Lecture 2.5.

- (d) Argue that the columns of the (2×3) matrix P are obtained by regressing each of the variables 'Age', 'Educ', and 'Parttime' on a constant term and the variable 'Female'.
- (e) Determine the values of P from the results in Lecture 2.1.

- (f) Check the numerical validity of the result in part (c). Note: This equation will not hold exactly because the coefficients have been rounded to two or three decimals; preciser results would have been obtained for higher precision coefficients.
-

Answers

To recap:

Linear regression assumptions:

A1: Data generating process is linear in parameters: $y_i = \alpha + \beta x_i + \varepsilon_i$.

A2: The n observations of x_i are fixed numbers.

A3: The n error terms, ε_i are random, with $E(\varepsilon_i) = 0$.

A4: The variance of n errors is fixed, $E(\varepsilon_i^2) = \sigma^2$.

A5: The errors are uncorrelated, $E(\varepsilon_i \varepsilon_j) = 0$ for all $i \neq j$.

A6: α and β are unknown, but fixed for all n observations.

A7: $\varepsilon_1, \dots, \varepsilon_n$ are jointly normally distributed; with A3, A4, A5: $\varepsilon \sim NID(0, \sigma^2)$.

(a) Prove that $E(b_R) = \beta_1 + P\beta_2$.

$$\begin{aligned} b_R &= (X_1' X_1)^{-1} X_1' y \\ &= (X_1' X_1)^{-1} X_1' (X_1 \beta_1 + X_2 \beta_2 + \varepsilon) \\ &= (X_1' X_1)^{-1} X_1' X_1 \beta_1 + (X_1' X_1)^{-1} X_1' X_2 \beta_2 + (X_1' X_1)^{-1} X_1' \varepsilon \\ &= I \beta_1 + (X_1' X_1)^{-1} X_1' X_2 \beta_2 + (X_1' X_1)^{-1} X_1' \varepsilon \\ &= \beta_1 + P \beta_2 + (X_1' X_1)^{-1} X_1' \varepsilon \end{aligned}$$

Therefore:

$$\begin{aligned} E(b_R) &= E(\beta_1 + P\beta_2 + (X_1' X_1)^{-1} X_1' \varepsilon) \\ &= \beta_1 + P\beta_2 + (X_1' X_1)^{-1} X_1' E(\varepsilon) \\ &= \beta_1 + P\beta_2 \end{aligned}$$

(b) Prove that $\text{var}(b_R) = \sigma^2 (X_1' X_1)^{-1}$.

The variance $\text{var}(X)$ of a random variable X is defined by:

$$\text{var}(X) = E[(X - E[X])^2]$$

Therefore we can find the $\text{var}(b_R)$ by setting:

$$\text{var}(b_R) = E[(b_R - E[b_R])^2]$$

Therefore:

$$\begin{aligned} b_R - E(b_R) &= [\beta_1 + P\beta_2 + (X_1' X_1)^{-1} X_1' \varepsilon] \\ &= (X_1' X_1)^{-1} X_1' \varepsilon \end{aligned}$$

So that:

$$\begin{aligned}
\text{var}(b_R) &= E([b_R - E(b_R)][b_R - E(b_R)]') \\
&= E([(X_1'X_1)^{-1}X_1'\varepsilon][(X_1'X_1)^{-1}X_1'\varepsilon]') \\
&= E([(X_1'X_1)^{-1}X_1'\varepsilon][\varepsilon'(X_1')'((X_1'X_1)^{-1})']) \\
&= E([(X_1'X_1)^{-1}X_1'\varepsilon][\varepsilon'X_1(X_1'X_1)^{-1}]) \\
&= E((X_1'X_1)^{-1}X_1'\varepsilon\varepsilon'X_1(X_1'X_1)^{-1}) \\
&= (X_1'X_1)^{-1}X_1'E(\varepsilon\varepsilon')X_1(X_1'X_1)^{-1} \\
&= (X_1'X_1)^{-1}X_1'\sigma^2IX_1(X_1'X_1)^{-1} \\
&= \sigma^2(X_1'X_1)^{-1}X_1'X_1(X_1'X_1)^{-1} \\
&= \sigma^2(X_1'X_1)^{-1}
\end{aligned}$$

(c) Prove that $b_R = b_1 + Pb_2$.

$$\begin{aligned}
b_R &= (X_1'X_1)^{-1}X_1'y \\
&= (X_1'X_1)^{-1}X_1'(X_1b_1 + X_2b_2 + e) \\
&= (X_1'X_1)^{-1}X_1'X_1b_1 + (X_1'X_1)^{-1}X_1'X_2b_2 + (X_1'X_1)^{-1}X_1'e \\
&= Ib_1 + (X_1'X_1)^{-1}X_1'X_2b_2 + (X_1'X_1)^{-1}X_1'e \\
&= b_1 + Pb_2 + (X_1'X_1)^{-1}X_1'e \\
&= b_1 + Pb_2
\end{aligned}$$

(d) Argue that the columns of the (2×3) matrix P are obtained by regressing each of the variables 'Age', 'Educ', and 'Parttime' on a constant term and the variable 'Female'.

If X_1 is an $(n \times 2)$ matrix, with columns for the constant term and the variable 'Female',
and X_1' is $(2 \times n)$,
and X_2 is an $(n \times 3)$ matrix with columns for 'Age', 'Educ', and 'Parttime',
then $(X_1'X_1)^{-1}$ is a (2×2) matrix,
and $(X_1'X_1)^{-1}X_1'$ is a $(2 \times n)$, therefore $P = (X_1'X_1)^{-1}X_1'X_2$ is a (2×3) matrix.

(e) Determine the values of P from the results in Lecture 2.1.

After loading the `TrainExer21` data set that we used in the previous exercises, we can split up the respective columns to create the X_1 and X_2 matrices needed to construct P .

```
X1 <- TrainExer21 %>% mutate("Constant" = rep(1, len = nrow(TrainExer21))) %>%
  select(Constant, Female) %>% as.matrix()
```

This results in an X_1 matrix with a dimension: (500, 2)

```
X2 <- TrainExer21 %>% select(Age, Educ, Parttime) %>% as.matrix()
```

This results in an X_2 matrix with a dimension: (500, 3)

We construct our P matrix using the formula:

$$P = (X_1'X_1)^{-1}X_1'X_2$$

```
P1 <- solve(t(X1)%*%X1) %*% (t(X1)%*%X2)
```

We can see that P is a (2, 3) dimension matrix, and takes on the following values:

Table 1:

	Age	Educ	Parttime
Constant	40.051	2.259	0.196
Female	-0.110	-0.493	0.249

(f) Check the numerical validity of the result in part (c). Note: This equation will not hold exactly because the coefficients have been rounded to two or three decimals; preciser results would have been obtained for higher precision coefficients.

From the lectures, our model is defined as:

$$\log(Wage)_i = \beta_1 + \beta_2 Female_i + \beta_3 Age_i + \beta_4 Educ_i + \beta_5 Parttime_i + \varepsilon_i$$

Thus, our OLS regression is determined to be:

$$\log(Wage)_i = 3.05 - 0.04 Female_i + 0.03 Age_i + 0.23 Educ_i - 0.37 Parttime_i + e_i$$

Using our findings from part(c), we can show that:

$$b_R = \begin{pmatrix} 3.05 \\ -0.04 \end{pmatrix} + \begin{pmatrix} 40.05 & 2.26 & 0.20 \\ -0.11 & -0.49 & 0.25 \end{pmatrix} \begin{pmatrix} 0.03 \\ 0.23 \\ 0.37 \end{pmatrix}$$

b_1 P b_2

```
b1 <- as.matrix(c(3.05, -0.04))
P2 <- cbind(c(40.05, -0.11), c(2.26, -0.49), c(0.20, 0.25))
b2 <- as.matrix(c(0.03, .23, -0.37))

bR <- b1 + P2%*%b2
print(bR)
```

```
##           [,1]
## [1,]  4.6973
## [2,] -0.2485
```

Carrying out our calculations, we find that b_R has a value: (4.70, -0.25)

This result approximates the more precise values of b_R given in lecture 2.1, which were: (4.73, -0.25)