

Chapter 12 Lecture Notes

Spring 2023

The research context is that two or more independent variables and one dependent variable have been observed for each of n participants. Here, I will discuss two independent variables x_{1i} and x_{2i} , $i = 1, \dots, n$ and one dependent variable y_i , $i = 1, \dots, n$.

The mathematics and analysis for more independent variables generalize routinely.

The research team then has a spreadsheet with n vectors of observations

(x_{1i}, x_{2i}, y_i) , $i = 1, \dots, n$. As in Chapter 11, one of the variables (here y) is the outcome variable or dependent variable. This is the variable hypothesized to be affected by the other variables in scientific research. The other variables (here x_{1i} and x_{2i} , $i = 1, \dots, n$) are the independent variables. They may be hypothesized to predict the outcome variable or to cause a change in the outcome variable. The research task is to document the association between independent and dependent variables.

As before, a recommended first step is to create the scatterplots of observations, with the vertical axis representing the dependent variable and the horizontal axis representing one of the independent variables. The “pencil test” can be used again. If the plot passes this test, then it is reasonable to assume that a **linear model** (such as $\beta_0 + \beta_1 x_1 + \beta_2 x_2$) describes the data. The linear model is reasonable for many data sets in observational studies. Specifically, the model for Chapter 12 is

$Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \sigma_{Y \bullet 12} Z_i$. The parameters $(\beta_0, \beta_1, \beta_2)$ are fixed but unknown. The parameter $\sigma_{Y \bullet 12}$ is the unknown conditional standard deviation of Y_i controlling for x_{1i} and x_{2i} , $i = 1, \dots, n$. The standard deviation of Y_i is assumed to be equal for each observation. The random errors Z_i are assumed to be independent. The independence of the random errors (and hence independence of Y_i) is important.

The assumption of a linear regression function (that is, $E(Y_i | x_{1i}, x_{2i}) = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}$) is also important. As in Chapter 11, this is equivalent to the joint distribution of the dependent variable values being $NID(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}, \sigma_{Y \bullet 12}^2)$.

Estimating the Linear Model Parameters

Again, OLS (ordinary least squares) is the most commonly used method to estimate the parameters of the linear model. A linear model with arbitrary arguments $b_0 + b_1x_1 + b_2x_2$ is used as a *fit* for the dependent variable values. The method uses the *residual* $y_i - b_0 - b_1x_{1i} - b_2x_{2i}$. As in Chapter 11, the fitting model is judged by how small the set of residuals is. Here OLS minimizes the sum of squares function $SS(b_0, b_1, b_2) = \sum_{i=1}^n (y_i - b_0 - b_1x_{1i} - b_2x_{2i})^2$. The OLS method is to find the arguments $(\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2)$ that make $SS(b_0, b_1, b_2)$ as small as possible. This minimization is a standard calculus problem. Step 1 is to calculate the partial derivatives of $SS(b_0, b_1, b_2)$ with respect to each argument. First, the partial with respect to b_0 :

$$\begin{aligned}\frac{\partial SS(b_0, b_1, b_2)}{\partial b_0} &= \frac{\partial}{\partial b_0} \sum_{i=1}^n (y_i - b_0 - b_1x_{1i} - b_2x_{2i})^2 = \sum_{i=1}^n \frac{\partial}{\partial b_0} (y_i - b_0 - b_1x_{1i} - b_2x_{2i})^2 \\ &= \sum_{i=1}^n 2(y_i - b_0 - b_1x_{1i} - b_2x_{2i}) \frac{\partial (y_i - b_0 - b_1x_{1i} - b_2x_{2i})}{\partial b_0}.\end{aligned}$$

$$\text{Simplifying, } \frac{\partial SS(b_0, b_1, b_2)}{\partial b_0} = \sum_{i=1}^n (-2)(y_i - b_0 - b_1x_{1i} - b_2x_{2i}).$$

Second, the partial with respect to b_1 :

$$\begin{aligned}\frac{\partial SS(b_0, b_1, b_2)}{\partial b_1} &= \frac{\partial}{\partial b_1} \sum_{i=1}^n (y_i - b_0 - b_1x_{1i} - b_2x_{2i})^2 = \sum_{i=1}^n \frac{\partial}{\partial b_1} (y_i - b_0 - b_1x_{1i} - b_2x_{2i})^2 \\ &= \sum_{i=1}^n 2(y_i - b_0 - b_1x_{1i} - b_2x_{2i}) \frac{\partial (y_i - b_0 - b_1x_{1i} - b_2x_{2i})}{\partial b_1}.\end{aligned}$$

$$\text{Then, } \frac{\partial SS(b_0, b_1, b_2)}{\partial b_1} = \sum_{i=1}^n (-2)(y_i - b_0 - b_1x_{1i} - b_2x_{2i})x_{1i}$$

Third, the partial with respect to b_2 is parallel, with the result that

$$\frac{\partial SS(b_0, b_1, b_2)}{\partial b_2} = \sum_{i=1}^n (-2)(y_i - b_0 - b_1x_{1i} - b_2x_{2i})x_{2i}$$

Step 2 is to find the arguments $(\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2)$ that make the three partial derivatives simultaneously zero. The resulting equations are still called the *normal equations*:

$$\sum_{i=1}^n (-2)(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i}) = 0,$$

$$\sum_{i=1}^n (-2)(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i})x_{1i} = 0, \text{ and}$$

$$\sum_{i=1}^n (-2)(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i})x_{2i} = 0, .$$

These equations still have a very important mathematical interpretation. Let

$r_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i}, i = 1, \dots, n$. The first normal equation is equivalent to $\sum_{i=1}^n r_i = 0$;

the second is $\sum_{i=1}^n r_i x_{1i} = 0$; and the third is $\sum_{i=1}^n r_i x_{2i} = 0$ That is, there are three

constraints on the n residuals. The OLS residuals must sum to zero, and the OLS residuals are orthogonal to the two independent variable values. The n residuals then have $n - 3$ degrees of freedom.

Step 3 is to solve this three linear equation system in three unknowns. There is a more general approach to solving systems like this. The first equation is

$$\sum_{i=1}^n (-2)(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i}) = 0, \text{ which can be written } \sum_{i=1}^n y_i = \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i}) \text{ and}$$

$$\sum_{i=1}^n (1 \times y_i) = [\sum_{i=1}^n (1 \times 1)]\hat{\beta}_0 + [\sum_{i=1}^n (1 \times x_{1i})]\hat{\beta}_1 + [\sum_{i=1}^n (1 \times x_{2i})]\hat{\beta}_2 . \text{ Similarly, the second normal}$$

equation can be written $\sum_{i=1}^n (x_{1i} \times y_i) = [\sum_{i=1}^n (1 \times x_{1i})]\hat{\beta}_0 + [\sum_{i=1}^n (x_{1i} \times x_{1i})]\hat{\beta}_1 + [\sum_{i=1}^n (x_{1i} \times x_{2i})]\hat{\beta}_2 ;$

and the third

$$\sum_{i=1}^n (x_{2i} \times y_i) = [\sum_{i=1}^n (1 \times x_{2i})]\hat{\beta}_0 + [\sum_{i=1}^n (x_{2i} \times x_{1i})]\hat{\beta}_1 + [\sum_{i=1}^n (x_{2i} \times x_{2i})]\hat{\beta}_2 .$$

While these look like complicated equations, matrix algebra leads to a

simpler expression. Let $y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$ be the $n \times 1$ column vector of dependent variable

values, and let $X = \begin{bmatrix} 1 & x_{11} & x_{21} \\ 1 & x_{12} & x_{22} \\ \vdots & \vdots & \vdots \\ 1 & x_{1n} & x_{2n} \end{bmatrix}$ be the $n \times 3$ matrix of coefficients of the parameters

$(\beta_0, \beta_1, \beta_2)$. From matrix algebra, $X^T X = \begin{bmatrix} \sum_{i=1}^n (1 \times 1) & \sum_{i=1}^n (1 \times x_{1i}) & \sum_{i=1}^n (1 \times x_{2i}) \\ \sum_{i=1}^n (1 \times x_{1i}) & \sum_{i=1}^n (x_{1i} \times x_{1i}) & \sum_{i=1}^n (x_{1i} \times x_{2i}) \\ \sum_{i=1}^n (1 \times x_{2i}) & \sum_{i=1}^n (x_{1i} \times x_{2i}) & \sum_{i=1}^n (x_{2i} \times x_{2i}) \end{bmatrix}$, and

$$X^T Y = \begin{bmatrix} \sum_{i=1}^n (1 \times y_i) \\ \sum_{i=1}^n (x_{1i} \times y_i) \\ \sum_{i=1}^n (x_{2i} \times y_i) \end{bmatrix}.$$

Recall the three normal equations above:

$$\sum_{i=1}^n (1 \times y_i) = \left[\sum_{i=1}^n (1 \times 1) \right] \hat{\beta}_0 + \left[\sum_{i=1}^n (1 \times x_{1i}) \right] \hat{\beta}_1 + \left[\sum_{i=1}^n (1 \times x_{2i}) \right] \hat{\beta}_2$$

$$\sum_{i=1}^n (x_{1i} \times y_i) = \left[\sum_{i=1}^n (1 \times x_{1i}) \right] \hat{\beta}_0 + \left[\sum_{i=1}^n (x_{1i} \times x_{1i}) \right] \hat{\beta}_1 + \left[\sum_{i=1}^n (x_{1i} \times x_{2i}) \right] \hat{\beta}_2$$

$$\sum_{i=1}^n (x_{2i} \times y_i) = \left[\sum_{i=1}^n (1 \times x_{2i}) \right] \hat{\beta}_0 + \left[\sum_{i=1}^n (x_{1i} \times x_{2i}) \right] \hat{\beta}_1 + \left[\sum_{i=1}^n (x_{2i} \times x_{2i}) \right] \hat{\beta}_2$$

The left-hand side terms are the same as the terms of $X^T Y$, and the coefficients of the OLS estimators match with the terms of $X^T X$. For this problem, then, the normal equations can be written in matrix form as $(X^T X) \hat{\beta} = X^T Y$. This result also holds for three or more independent variables. The proof is exactly the same as for the two independent variable case.

If $(X^T X)^{-1}$ exists, then $\hat{\beta} = (X^T X)^{-1} X^T Y$. The existence of $(X^T X)^{-1}$ is the usual case in observational studies using multiple regression. If $(X^T X)^{-1}$ does not exist, then the OLS estimators exist but are not unique.

Distribution of $\hat{\beta} = (X^T X)^{-1} X^T Y$

Let $Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$ be the vector of the random outcome variables. That is, the data will be

collected in the future as opposed to having the data in hand as we assumed in our OLS estimator derivation. The probabilistic model for the data can be written in matrix form $Y = X\beta + \sigma_{Y \cdot 12} Z$, where Z is the column vector of random errors Z_i that are assumed to be independent. From now on, we consider the general case with $p-1$ independent variables. The vector β of parameters in the regression function is then $p \times 1$, remembering that there is an intercept term in our model. The matrix X of coefficients of the parameters of the regression function is now $n \times p$, $p < n$, with rank p . Then $E(Y) = E(X\beta + \sigma_{Y \cdot 12} Z) = E(X\beta) + E(\sigma_{Y \cdot 12} Z) = X\beta + \sigma_{Y \cdot 12} E(Z) = X\beta$, and $\text{vcv}(Y) = \sigma_{Y \cdot 12}^2 I_{n \times n}$. An equivalent description is to say that Y is multivariate normal with dimension n ; that is, Y has the distribution $MVN_n(X\beta, \sigma_{Y \cdot 12}^2 I_{n \times n})$.

When the matrix X has rank p , $(X^T X)^{-1}$ exists. Then the vector of OLS estimators is

$\hat{\beta} = (X^T X)^{-1} X^T Y$. The expected value is given by

$$E(\hat{\beta}) = E[(X^T X)^{-1} X^T Y] = (X^T X)^{-1} X^T E(Y) = (X^T X)^{-1} X^T (X\beta) = [(X^T X)^{-1} (X^T X)]\beta = I_{p \times p} \beta = \beta.$$

The variance-covariance matrix of $\hat{\beta} = (X^T X)^{-1} X^T Y$ is calculated by

$\text{vcv}(\hat{\beta}) = \text{vcv}[(X^T X)^{-1} X^T Y] = (X^T X)^{-1} X^T \text{vcv}(Y) [(X^T X)^{-1} X^T]^T$. This can be simplified using the matrix algebra result that $(AB)^T = B^T A^T$ so that

$[(X^T X)^{-1} X^T]^T = (X^T)^T [(X^T X)^{-1}]^T = X(X^T X)^{-1}$. Recall that the transpose of the transpose of a matrix is just the matrix so that $(X^T)^T = X$. Further, a matrix is symmetric if its transpose is the matrix itself. That is, $(X^T X)^T = X^T (X^T)^T = X^T X$.

The inverse of a symmetric matrix is symmetric so that $[(X^T X)^{-1}]^T = (X^T X)^{-1}$. Using

these results in $\text{vcv}(\hat{\beta}) = (X^T X)^{-1} X^T \text{vcv}(Y) [(X^T X)^{-1} X^T]^T = (X^T X)^{-1} X^T \sigma_{Y \bullet X}^2 I_{n \times n} X (X^T X)^{-1}$,
 $\text{vcv}(\hat{\beta}) = (X^T X)^{-1} X^T \sigma_{Y \bullet X}^2 I_{n \times n} X (X^T X)^{-1} = \sigma_{Y \bullet X}^2 \{(X^T X)^{-1} [X^T X]\} (X^T X)^{-1} = \sigma_{Y \bullet X}^2 \{I_{p \times p}\} (X^T X)^{-1} = \sigma_{Y \bullet X}^2 (X^T X)^{-1}$.

The distribution of $\hat{\beta} = (X^T X)^{-1} X^T Y$ is $MVN_p(\beta, \sigma_{Y \bullet X}^2 (X^T X)^{-1})$.

Fisher's Decomposition of the (Uncorrected) Total Sum of Squares

The uncorrected total sum of squares of the dependent variable is defined to be $Y^T Y$ with n degrees of freedom. In Chapter 11, the (corrected) total sum of squares was used. This is $TSS = \sum_{i=1}^n (Y_i - \bar{Y}_n)^2 = Y^T Y - n\bar{Y}_n^2$ with $n-1$ degrees of freedom. First,

Fisher's decomposition of the uncorrected total sum of squares follows from

$$\begin{aligned} Y^T Y &= (Y - X\hat{\beta} + X\hat{\beta})^T (Y - X\hat{\beta} + X\hat{\beta}) \\ &= (Y - X\hat{\beta})^T (Y - X\hat{\beta}) + (X\hat{\beta})^T (Y - X\hat{\beta}) + (Y - X\hat{\beta})^T X\hat{\beta} + (X\hat{\beta})^T (X\hat{\beta}). \end{aligned}$$

This result can be simplified using

$$\begin{aligned} (Y - X\hat{\beta})^T X\hat{\beta} &= (Y - X(X^T X)^{-1} X^T Y)^T X(X^T X)^{-1} X^T Y = Y^T (I_{n \times n} - X(X^T X)^{-1} X^T)^T X(X^T X)^{-1} X^T Y \\ &= Y^T \{(I_{n \times n})^T - [X(X^T X)^{-1} X^T]^T\} X(X^T X)^{-1} X^T Y \\ &= Y^T [I_{n \times n} - X(X^T X)^{-1} X^T] X(X^T X)^{-1} X^T Y \\ &= Y^T [X(X^T X)^{-1} X^T - X(X^T X)^{-1} X^T X(X^T X)^{-1} X^T] Y \\ &= Y^T [X(X^T X)^{-1} X^T - X(X^T X)^{-1} \{(X^T X)(X^T X)^{-1}\} X^T] Y \\ &= Y^T [X(X^T X)^{-1} X^T - X(X^T X)^{-1} \{I_{p \times p}\} X^T] Y = 0. \end{aligned}$$

Of course, $(X\hat{\beta})^T (Y - X\hat{\beta}) = 0$.

Then

$$\begin{aligned} Y^T Y &= (Y - X\hat{\beta} + X\hat{\beta})^T (Y - X\hat{\beta} + X\hat{\beta}) \\ &= (Y - X\hat{\beta})^T (Y - X\hat{\beta}) + (X\hat{\beta})^T (X\hat{\beta}). \end{aligned}$$

The residuals R are defined to be the $n \times 1$ vector $R = Y - X\hat{\beta}$ on $n-p$ degrees of freedom, and the fitted values $\hat{Y} = X\hat{\beta}$ on p degrees of freedom. Then the uncorrected total sum of squares is $Y^T Y = R^T R + \hat{Y}^T \hat{Y}$. The error sum of squares is defined to be $R^T R$ with $n-p$ degrees of freedom. The uncorrected sum of squares due to regression is defined to be $\hat{Y}^T \hat{Y}$ with p degrees of freedom. Statistical computing programs subtract the correction $n\bar{Y}_n^2$ with 1 degree of freedom from both the uncorrected total sum of squares and uncorrected regression sum of squares. That is, the programs display the corrected total sum of squares

$$TSS = \sum_{i=1}^n (Y_i - \bar{Y}_n)^2 = Y^T Y - n\bar{Y}_n^2$$

and the corrected regression sum of squares in the Analysis of Variance Table as below:

Analysis of Variance Table

$p - 1$ Predictor Multiple Linear Regression

Source	DF	Sum of Squares	Mean Square	F
Regression	$p - 1$	$(X\hat{\beta})^T (X\hat{\beta}) - n\bar{Y}_n^2$	$\frac{(X\hat{\beta})^T (X\hat{\beta}) - n\bar{Y}_n^2}{p - 1}$	$\frac{MS_{REG}}{MSE}$
Error	$n - p$	$R^T R$	$\frac{R^T R}{(n - p)}$	
Total	$n - 1$	$TSS = (n - 1)s_{DV}^2$		

11.3 Inferences

The probabilistic model for the data is $Y = X\beta + \sigma_{Y \bullet 1 \dots (p-1)} Z$. The outcome or dependent (random) variables $Y_i, i = 1, \dots, n$ are each assumed to be the sum of the linear regression expected value $\beta_0 + \beta_1 x_{1i} + \dots + \beta_{p-1} x_{(p-1)i}$ and a random error term $\sigma_{Y \bullet 1 \dots (p-1)} Z_i$. The random variables $Z_i, i = 1, \dots, n$ are assumed to be independent standard normal random variables. The parameter β_0 is the intercept parameter and is fixed but unknown. The parameters $\beta_1, \dots, \beta_{p-1}$ are partial regression coefficient parameters and are also fixed but unknown. These parameters are the focus of the statistical analysis. The parameter $\sigma_{Y \bullet 1 \dots (p-1)}$ is also fixed but unknown. Another description of this model is that $Y_i, i = 1, \dots, n$ are independent normally distributed random variables with $Y_{n \times 1}$ having the distribution $MVN(X\beta, \sigma_{Y \bullet 1 \dots (p-1)}^2 I_{n \times n})$.

Again, there are four assumptions. The two important assumptions are that the outcome variables $Y_i, i = 1, \dots, n$ are independent and that the regression function is $\beta_0 + \beta_1 x_{1i} + \dots + \beta_{p-1} x_{(p-1)i}$ $i = 1, \dots, n$. Homoscedasticity is less important. The assumption that $Y_i, i = 1, \dots, n$ are normally distributed random variables is least important.

Testing null hypotheses about the partial regression coefficients

The mathematical analysis of the general problem is complicated. The analysis for two independent variables, however, is more manageable—particularly the problem of sequential tests. As before, the model for the data is

$$Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \sigma_{Y \cdot 12} Z_i.$$

The research problem is to consider a sequence of models. The first model is that $Y_i = \beta_0 + \beta_1 x_{1i} + \sigma_{Y \cdot 1} Z_i$, with null hypothesis $H_0 : \beta_1 = 0$. This is a Chapter 11 problem.

The second model is that $Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \sigma_{Y \cdot 12} Z_i$, with null hypothesis $H_0 : \beta_2 = 0$. This is an example of a sequential test. That is, the second hypothesis is tested after the first one. These tests require the definition of the partial correlation coefficient.

Partial correlation coefficient

Let the correlation matrix of (Y, x_1, x_2) be

$$\begin{pmatrix} 1 & \rho(y, x_1) = \rho_{y1} & \rho(y, x_2) = \rho_{y2} \\ \rho(y, x_1) = \rho_{y1} & 1 & \rho(x_1, x_2) = \rho_{12} \\ \rho(y, x_2) = \rho_{y2} & \rho(x_1, x_2) = \rho_{12} & 1 \end{pmatrix}$$

The *partial correlation* between Y and x_2 controlling for x_1 is defined to be

$$\rho_{y2.1} = \frac{\rho_{y2} - \rho_{y1}\rho_{12}}{\sqrt{(1 - \rho_{y1}^2)(1 - \rho_{12}^2)}}. \text{ Analogous definitions hold for the Pearson product}$$

moment correlations. That is, $r_{y2.1} = \frac{r_{y2} - r_{y1}r_{12}}{\sqrt{(1 - r_{y1}^2)(1 - r_{12}^2)}}$

Analysis of variance table for a sequential test

The (corrected) total sum of squares is always $TSS = (n-1)s_{DV}^2$. The first model is that $Y_i = \gamma_0 + \gamma_1 x_{1i} + \sigma_{Y \cdot 1} Z_i$, with null hypothesis $H_0 : \gamma_1 = 0$. The sum of squares due to the regression on x_1 is then $[r(x_1, y)]^2 TSS$ on 1 degree of freedom. The regression on

x_1 has error sum of squares $\{1 - [r(x_1, y)]^2\}TSS$ with $n - 2$ degrees of freedom. The new aspect of multiple regression is that there is a second independent (predictor) variable. The regression on x_2 after x_1 has been entered explains an additional $r_{y2 \cdot 1}^2$ of the $\{1 - [r(x_1, y)]^2\}TSS$ that was not explained by x_1 . That is, the sum of squares due to the regression on $x_2 | x_1$ is $[r_{y2 \cdot 1}^2 \{1 - [r(x_1, y)]^2\}TSS$ with one degree of freedom. The error sum of squares is obtained by subtracting both the sum of squares due to the regression on x_1 and the the sum of squares due to the regression on $x_2 | x_1$. The analysis of variance table below summarizes these results.

Analysis of variance table

Multiple regression of Y on x_1 and $x_2 | x_1$

Source	DF	Sum of Squares	Mean Square	
Reg on x_1	1	$[r(x_1, y)]^2 TSS$	$[r(x_1, y)]^2 TSS$	
Reg on $x_2 x_1$	1	$[r_{y2 \cdot 1}^2 \{1 - [r(x_1, y)]^2\} TSS$	$[r_{y2 \cdot 1}^2 \{1 - [r(x_1, y)]^2\} TSS$	
Error	$n - 3$	Subtraction	MSE	
Total (corrected)	$n - 1$	$TSS = (n - 1)s_{DV}^2$		

The test of $H_0 : \beta_2 = 0$ against the alternative hypothesis that $H_1 : \beta_2 \neq 0$ uses the test statistic $F_{2 \cdot 1} = \frac{[r_{y2 \cdot 1}^2 \{1 - [r(x_1, y)]^2\} TSS}{MSE}$, which has 1 numerator and $n - 3$ denominator degrees of freedom.

This presentation of the models has disguised the complexity of the coefficients. The first model was $Y_i = \gamma_0 + \gamma_1 x_{1i} + \sigma_{Y \cdot 1} Z_i$. The specification of the γ_1 parameter requires taking expectation to get $E(Y | x_1) = \gamma_0 + \gamma_1 x_1$. Then, $\gamma_1 = \frac{\partial}{\partial x_1} E(Y_i | x_1)$; that is, γ_1 is the expected increase in the value of the dependent variable associated with a unit increase in x_1 . The extended model was $Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \sigma_{Y \cdot 12} Z_i$, so that $E(Y | x_1, x_2) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$. Then

$\beta_2 = \frac{\partial}{\partial x_2} \{E(Y | x_1, x_2)\} |_{x_1 \text{ fixed}}$. The coefficient $\beta_2 = \frac{\partial}{\partial x_2} \{E(Y | x_1, x_2)\} |_{x_1 \text{ fixed}}$ is the expected increase in the value of the dependent variable associated with a unit increase in x_2 , controlling for x_1 being held constant (sometimes called *ceteris paribus*). These coefficients are partial derivatives of the conditional expectation of the dependent variable.

Example Examination Problem: A study collects the values of (Y, x_1, x_2) on 400 subjects. The total sum of squares for Y is 1000. The correlation between Y and x_1 is 0.67; the correlation between Y and x_2 is 0.50; and the correlation between x_1 and x_2 is 0.25.

- Compute the analysis of variance table for the multiple regression analysis of Y . Include the sum of squares due to the regression on x_1 and the sum of squares due to the regression on x_2 after including x_1 .
- Test the null hypothesis that both $\beta_2 = 0$ and $\beta_1 = 0$; that is, the null hypothesis is that there is no association between Y and these two independent variables.
- Test the null hypothesis that the variable x_2 does not improve the fit of the model once x_1 has been included against the alternative that the variable does improve the fit of the model. Report whether the test is significant at the 0.10, 0.05, 0.01 levels of significance.

Solution: a. The sum of squares due to the regression on x_1 has 1 degree of freedom and is equal to $SS(x_1) = r_{y1}^2 SS(Total) = (0.67)^2 \times 1000 = 448.9$. The partial correlation coefficient of x_2 with y after controlling for x_1 is

$$r_{y2.1} = \frac{r_{y2} - r_{y1}r_{12}}{\sqrt{(1-r_{y1}^2)(1-r_{12}^2)}} = \frac{0.50 - 0.67 \times 0.25}{\sqrt{(1-0.67^2)(1-0.25^2)}} = \frac{0.3325}{\sqrt{0.5511 \times 0.9375}} = \frac{0.3325}{0.7188} = 0.463$$

The sum

of squares due to the regression on x_2 after including x_1 also has 1 degree of freedom is $SS(x_2 | x_1) = r_{y2.1}^2 (1 - r_{y1}^2) SS(Total) = (0.463)^2 \times 0.5511 \times 1000 = 118.1$. The sum of squares for error has $400 - 3 = 397$ degrees of freedom and is $SS(Error) = SS(Total) - SS(x_1) - SS(x_2 | x_1) = 1000 - 448.9 - 118.1 = 1000 - 567.0 = 433$

b. The sum of squares using both x_1 and x_2 has 2 degrees of freedom and is equal to $SS(x_1, x_2) = SS(x_1) + SS(x_2 | x_1) = 448.9 + 118.1 = 567.0$. The F -test for this hypothesis is $F = \frac{SS(x_1, x_2)/2}{SS(Error)/397} = \frac{567.0/2}{433.0/397} = \frac{283.5}{1.091} = 259.9$. The critical value for the 0.01

level of significance is more than 4.61 (for two numerator and infinite denominator degrees of freedom) and less than 4.69 (for two numerator and 240 denominator degrees of freedom). Excel reports that the critical value is 4.659 for 2 numerator and 397 denominator degrees of freedom. I reject the null hypothesis that there is no association between Y and these two independent random variables.

c. The F -test for this hypothesis is $F = \frac{MS(x_2 | x_1)}{MS(Error)} = \frac{118.1/1}{433.0/397} = \frac{118.1}{1.091} = 108.3$ with 1

numerator and 397 denominator degrees of freedom. The critical value for the 0.01 level of significance is more than 6.63 (for one numerator and infinite denominator degrees of freedom) and less than 6.74 (for one numerator and 240 denominator degrees of freedom). Excel reports that the critical value is 6.700 for 1 numerator and 397 denominator degrees of freedom. I reject the null hypothesis that x_2 does not improve the fit of the model once x_1 has been included.

Complete Mediation and Complete Explanation Causal Models

In analyzing research data from engineering or physical sciences studies, the independent variables typically operate at the same time. Given this, the fact that a partial regression coefficient is an estimate of a partial derivative strongly indicates to the user that caution is warranted in the interpretation of a partial regression coefficient. In social science and epidemiological research, however, the independent variables may operate at different points of time. For example, x_1 may describe a variable measured when the participant was between ages 5 and 6, and x_2 may describe a variable measured when the participant was between the ages of 8 and 9. The time-ordering of the independent variables is a crucial consideration in the interpretation of partial regression coefficients.

For example, often one sees that ρ_{y2} appears significant (that is, x_2 has a significant F statistic in a multiple regression analysis or the r_{y2} , the Pearson product moment correlation, is significant) but that $\rho_{y2.1}$ does not appear

significant. That is, in multiple regression analysis, the variable x_2 does not have a significant F-to-enter once x_1 is in the regression equation. There is a fundamental paper (Simon, 1954, available on JSTOR and on the Blackboard site) that you should download and read it.

Simon points out that when one has a common cause model (or *explanation*), the independent variable x_1 precedes both x_2 and y with regard to operation impact. Then if x_1 “causes” x_2 and if x_1 “causes” y , then there will be a “spurious” correlation ρ_{y2} (this correlation will be non-zero even though x_2 has no causal relation to y) and $\rho_{y2.1}$ will be zero. For example, consider G. B. Shaw’s correlation between the number of suicides in England in a given year and the number of churches of England in the same year.

In a causal chain model, the independent variable x_2 operates before and causes x_1 , and x_1 operates before y and causes y . Simon also points out that, when the model is a causal chain (or *mediation*), one also observes that ρ_{y2} will be non-zero and $\rho_{y2.1}$ will be zero (even though x_2 causes y through the mediation of x_1). Both causal modeling situations have the same empirical fact that a partial correlation is near 0. Deciding which interpretation is valid requires clarifying the sequence of operation of the variables. In practice, the relevant partial correlation may not be essentially 0. In this event, researchers speak of partial explanation and partial mediation.

Example Past Examination Questions

Common Information for Questions 1, 2, and 3

A research team sought to estimate the model $E(Y) = \beta_0 + \beta_1 x + \beta_2 w$. The variable Y was a measure of depression of a participant observed at age 25; the variable x was a measure of anxiety shown by the participant at age 18; and the variable w was a measure of the extent of traumatic events experienced by the participant before age 15. They observed values of y , x , and w on $n = 800$ subjects. They found that the standard deviation of Y , where the variance estimator used division by $n - 1$, was 12.2. The correlation between Y and w was 0.31; the

correlation between Y and x was 0.14; and the correlation between x and w was 0.41.

1. Compute the partial correlation coefficients $r_{Yx \cdot w}$ and $r_{Yw \cdot x}$.

Answer: $r_{Yx \cdot w} = 0.0149$ and $r_{Yw \cdot x} = 0.2797$. For example,

$$r_{yx \cdot w} = \frac{r_{yx} - r_{yw}r_{xw}}{\sqrt{(1-r_{yw}^2)(1-r_{xw}^2)}} = \frac{0.14 - 0.31 \times 0.41}{\sqrt{(1-0.31^2)(1-0.41^2)}} = \frac{0.0129}{\sqrt{0.9039 \times 0.8319}} = \frac{0.0129}{0.8672} = 0.0149.$$

The partial correlation of y with x conditioning on the variable w is close to zero.

Since w operates before x , this suggests an explanation model.

2. Compute the analysis of variance table for the multiple regression analysis of Y . Include the sum of squares due to the regression on w and the sum of squares due to the regression on x after including w . Test the null hypothesis that $\beta_1 = 0$ against the alternative that the coefficient is not equal to zero.

That is, test whether x adds significant additional explanation after using w . Report whether the test is significant at the 0.10, 0.05, and 0.01 levels of significance.

Answer: The analysis of variance table is given by

Analysis of Variance Table

Source	DF	SS	MS	F Statistic
Regression on w	1	11428.52	11428.52	
Regression on $x w$	1	23.86	23.86	0.18
Error	797	107470.78	134.84	
Total	799	118923.16		

The value of the test statistic is $F_{x|w} = 0.18$. Since $F_{0.10,1,797} = 2.71^+ = 2.712$, $F_{0.05,1,797} = 3.84^+ = 3.853$ and $F_{0.01,1,797} = 6.63^+ = 6.667$, we accept the null hypothesis that x does not add significant explanation after including w at the 0.10 level.

3. What interpretations can you make of these results in terms of causal models?

Answer: It is an explanation model.

End of application of common information

Common Information for Questions 4, 5, and 6

A research team sought to estimate the model $E(Y) = \beta_0 + \beta_1 x + \beta_2 w$. The variable Y was a measure of the extent of criminal behavior of a participant observed at age 30; the variable x was a measure of the rebelliousness shown by the participant at age 12; and the variable w was a measure of delinquency shown at age 18. They observed values of y , x , and w on $n = 1500$ subjects. They found that the standard deviation of Y , where the variance estimator used division by $n - 1$, was 15.7. The correlation between Y and w is 0.62; the correlation between Y and x is 0.35; and the correlation between x and w is 0.58.

4. Compute the partial correlation coefficients $r_{Yx \cdot w}$ and $r_{Yw \cdot x}$.

Answer: $r_{Yx \cdot w} = -0.015$ and $r_{Yw \cdot x} = 0.5465$

The partial correlation of y with x conditioning on the variable w is close to zero. Since w operates after x and before y , this partial correlation suggests a mediation model.

5. Compute the analysis of variance table for the multiple regression analysis of Y . Include the sum of squares due to the regression on w and the sum of squares due to the regression on x after including w . Test the null hypothesis that $\beta_1 = 0$ against the alternative that the coefficient is not equal to zero. That is, test whether x adds significant additional explanation after using w . Report whether the test is significant at the 0.10, 0.05, and 0.01 levels of significance.

Answer: The analysis of variance table is given by:

Analysis of Variance Table

Source	DF	SS	MS	F Statistic
Regression on w	1	142031.38	142031.38	
Regression on $x w$	1	51.18	51.18	0.34
Error	1497	227405.95	151.91	
Total	1499	369488.51		

The value of the test statistic is $F_{x|w} = 0.34$. Since the critical value for (1,1497) degrees of freedom is 2.708, we accept the null hypothesis that x does not add significant explanation after including w at the 0.10 level.

6. What, if any, interpretations can you make of these results in terms of causal models?

Answer: It is a mediation model.

End of application of common information

12.8 Logistic Regression

The outcome variable in a regression analysis is a binary (0 or 1) variable. Then $E(Y|x) = \Pr\{Y = 1|x\} = p(x)$.

One can use OLS to analyze the data. The estimated regression coefficients are still unbiased. The variance-covariance matrix of Y will not be proportional to the identity matrix. The bigger problem is that the fitted values, however, will not be between 0 and 1. Reviewers of the OLS results will not accept these estimates and will insist upon the application of logistic regression.

The odds ratio may take values between 0 and infinity, and the logarithm of the odds ratio may take values between -infinity and +infinity. The logarithm of the odds ratio is a natural transformation to use. The simple (univariate) logistic regression model is

$$\ln\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + \beta_1 x,$$

or equivalently

$$p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

The parameters are estimated by a method discussed in AMS 412 called maximum likelihood.

The interpretation of the parameters is complex.

Interpretation of β_0

The parameter β_0 is related to $p(0)$.

That is,

$$p(0) = \frac{e^{\beta_0 + \beta_1 0}}{1 + e^{\beta_0 + \beta_1 0}} = \frac{e^{\beta_0}}{1 + e^{\beta_0}}.$$

Equivalently, $\beta_0 = \ln\left(\frac{p(0)}{1-p(0)}\right)$.

Interpretation of β_1

The parameter β_1 described the association between the probability of the event occurring and x .

When $\beta_1 = 0$, $p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}} = \frac{e^{\beta_0}}{1 + e^{\beta_0}}$ for all values of x . A positive value of β_1 implies that $\lim_{x \rightarrow \infty} p(x) = 1$. A negative value of β_1 implies that $\lim_{x \rightarrow \infty} p(x) = 0$.

Since

$$\ln\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + \beta_1 x,$$

$$\frac{p(x)}{1-p(x)} = \exp(\beta_0 + \beta_1 x) = \exp(\beta_0) [\exp(\beta_1)]^x.$$

Then a unit increase in x is associated with odds of the event being multiplied by $\exp(\beta_1)$.

$$\text{Equivalently } \beta_1 = \ln\left(\frac{p(x+1)}{1-p(x+1)}\right) - \ln\left(\frac{p(x)}{1-p(x)}\right) \text{ or } \exp(\beta_1) = \frac{\left[\frac{p(x+1)}{1-p(x+1)}\right]}{\left[\frac{p(x)}{1-p(x)}\right]}.$$

This approach can handle multiple independent variables. One can propose the model

$$\ln\left(\frac{p(x_1, x_2)}{1-p(x_1, x_2)}\right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2.$$

One can use even more independent variables.