

It is common for more than one factor to influence an outcome. Fitting regression models to data involving two or more predictors is one of the most widely used statistical procedures. In this chapter we consider multiple linear regression problems involving modeling the relationship between a dependent variable,  $Y$  and two or more predictor variables  $X_1, X_2, X_3$ , etc.

### Multiple linear regression model

A multiple linear regression (MLR) model that relates a response variable  $y$  to  $p-1$  predictor variables is written as:

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_{p-1} x_{i,p-1} + \varepsilon_i, \quad i = 1, 2, \dots, n$$

where

$$\varepsilon_i \sim \text{Normal}(0, \sigma^2), \quad \text{independent.}$$

- The subscript  $i$  refers to the  $i^{th}$  individual (or unit) in the population.
- In  $x_{ij}$ , the first subscript  $i$  denotes the individual and the second subscript  $j$  denotes which predictor variable it is.
- The word “linear” in multiple linear regression refers to linearity in the parameters  $\beta_j$ .
- Predictors  $x_{ij}$  can be raw variables or transformations (e.g.,  $x^2$ , interaction terms).
- The model includes  $p-1$  predictors but  $p$  parameters (including the intercept  $\beta_0$ ).

### Interpretation of the parameters

- Each  $\beta_j$  represents the change in the mean response  $E(Y)$  per unit increase in  $x_j$ , when all other predictors are held constant.
- $\beta_0$  represents the mean response when all predictors are equal to zero (may or may not be meaningful).
- For example: Consider the multiple regression function with two predictors:

$$E(Y) = \beta_0 + \beta_1 X_1 + \beta_2 X_2.$$

- a. The parameter  $\beta_0$  is the  $Y$ -intercept of the regression plane. If the scope of the model includes  $X_1 = 0$  and  $X_2 = 0$ , then  $\beta_0$  represents the mean response  $E(Y)$  at that point. Otherwise,  $\beta_0$  does not generally have any practical meaning as a separate term.
- b. The parameter  $\beta_1$  indicates the change in the mean response  $E(Y)$  per unit increase in  $X_1$  when  $X_2$  is held constant. Say,  $\beta_1 = 2$ , then the mean response  $E(Y)$  increases by 2 units for every unit increase in  $X_1$ , regardless of the level of  $X_2$ .
- c. The parameter  $\beta_2$  indicates the change in the mean response  $E(Y)$  per unit increase in  $X_2$  when  $X_1$  is held constant. Similarly, say,  $\beta_2 = 5$  indicates that the mean response increases by 5 units for every unit increase in  $X_2$ , when  $X_1$  is held constant.

### Matrix formulation

We now present the principal results for the multiple linear regression model in matrix terms. Define the following matrices:

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}_{n \times 1}, \quad \mathbf{X} = \begin{bmatrix} 1 & X_{11} & X_{12} & \cdots & X_{1,p-1} \\ 1 & X_{21} & X_{22} & \cdots & X_{2,p-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \cdots & X_{n,p-1} \end{bmatrix}_{n \times p},$$

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{bmatrix}_{p \times 1}, \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}_{n \times 1}.$$

– The multiple linear regression model can then be written concisely as:

$$\mathbf{Y} = \mathbf{X}\beta + \varepsilon,$$

where:

- $\mathbf{Y}$  is the  $n \times 1$  vector of responses,
  - $\mathbf{X}$  is the  $n \times p$  design matrix (including a column of ones for the intercept and one column for each of the  $p - 1$  predictors),
  - $\beta$  is the  $p \times 1$  vector of regression coefficients, and
  - $\varepsilon$  is the  $n \times 1$  vector of errors.
- We assume:
- $$E(\varepsilon) = \mathbf{0}, \quad \text{Var}(\varepsilon) = \sigma^2 \mathbf{I}_n, \quad \begin{bmatrix} 1 & & & 0 \\ & \ddots & & \\ & & \ddots & \\ 0 & & & 1 \end{bmatrix}_{n \times n}$$
- where  $\mathbf{I}_n$  is the  $n \times n$  identity matrix.
- Consequently, the random vector  $\mathbf{Y}$  has

$$E(\mathbf{Y}) = \mathbf{X}\beta, \text{ and } \text{Var}(\mathbf{Y}) = \sigma^2 \mathbf{I}_n.$$

$$\mathbf{A}^T \mathbf{A}$$

## Estimation

– The least squares method is again used to estimate the regression coefficients. The estimator  $\hat{\beta}$  minimizes the sum of squared errors:

$$SSE = (\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta).$$

$$\begin{aligned} \frac{\partial SSE}{\partial \beta} &= 0 \\ -2 \quad \mathbf{x}^T (\mathbf{Y} - \mathbf{x}\beta) &= 0 \\ \mathbf{x}^T \mathbf{Y} - \mathbf{x}^T \mathbf{x} \beta &= 0 \\ \mathbf{x}^T \mathbf{x} \beta &= \mathbf{x}^T \mathbf{Y} \end{aligned}$$

$$\hat{\beta} = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{Y}$$

– Setting derivatives to zero yields the **normal equations**:

$$\mathbf{x}^T \mathbf{x} \hat{\beta} = \mathbf{x}^T \mathbf{Y}.$$

– Provided  $\mathbf{x}^T \mathbf{x}$  is nonsingular, the least squares solution is:

$$\hat{\beta} = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{Y}.$$

## Properties of least squares estimates

- The expected value of  $\hat{\beta}$  is

$$\begin{aligned}
 E(\hat{\beta} | \mathbf{X}) &= E\left[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \mid \mathbf{X}\right] \\
 &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top E(\mathbf{Y} | \mathbf{X}) \\
 &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\beta) \\
 &= \beta.
 \end{aligned}$$

- The variance of  $\hat{\beta}$  is

$$\begin{aligned}
 \text{Var}(\hat{\beta} | \mathbf{X}) &= \text{Var}\left[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \mid \mathbf{X}\right] \\
 &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \text{Var}(\mathbf{Y} | \mathbf{X}) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\
 &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\sigma^2 \mathbf{I}_n) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\
 &= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}.
 \end{aligned}$$

- Since  $\mathbf{Y} \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n)$  and  $\hat{\beta}$  is a linear function of  $\mathbf{Y}$ , it follows that:

$$\hat{\beta} \sim N(\beta, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}).$$

This distributional result provides the basis for statistical inference on regression coefficients, including the construction of confidence intervals and hypothesis testing ( $t$ -tests and  $F$ -tests).

## Model evaluation

- Recall the following sum of squares from the simple linear regression

- SST: the total corrected sum of squares of the  $Y$ 's is given by  $SST = \sum_{i=1}^n (y_i - \bar{y})^2$ .
- SSE: the residual sum of squares (error sum of squares) is given by  $SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$ .
- SSR: the regression sum of squares (i.e., sum of squares explained by the regression model) is given by  $SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$ .

- Define:

- $\bar{Y} = \frac{1}{n} J Y$  where  $J$  is the  $n \times n$  matrix of ones.

$$J = \begin{bmatrix} 1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 1 \end{bmatrix}$$

- $\hat{Y} = HY$  where  $H = X(X^\top X)^{-1} X^\top$  is the hat matrix.

- Residuals:  $e = Y - \hat{Y} = (I - H)Y$ .

- We have:

- $(I - \frac{1}{n} J)^\top (I - \frac{1}{n} J) = (I - \frac{1}{n} J)$

- $H^\top H = H$   $(AB)^\top = B^\top A^\top$ ,  $(X^\top X)^\top = X^\top X$ ,  $((X^\top X)^\top)^\top = (X^\top X)^\top$   
 $H^\top H = (X(X^\top X)^\top X^\top)^\top (X(X^\top X)^\top X^\top) = X(X^\top X)^\top X^\top \cdot (X^\top X)^\top X^\top = X(X^\top X)^\top X^\top = X^\top X$

- $(I - H)^\top (I - H) = (I - H)$

$$\begin{aligned}
 (I - H)^\top (I - H) &= I^\top I - I^\top H - H^\top I + H^\top H \\
 &= I - H - H^\top + H \\
 &= I - H
 \end{aligned}$$

– Using the matrix formulation, the three sums of squares for the multiple regression model are

$$a. SST = (Y - \bar{Y})^T (Y - \bar{Y}) = (Y - \frac{1}{n}JY)^T (Y - \frac{1}{n}JY) = ((I - \frac{1}{n}J)Y)^T (I - \frac{1}{n}J)Y = Y^T (I - \frac{1}{n}J)^T (I - \frac{1}{n}J)Y = Y^T (I - \frac{1}{n}J)Y$$

$$b. SSR = (HY - \bar{Y})^T (HY - \bar{Y}) = (HY - \frac{1}{n}JY)^T (HY - \frac{1}{n}JY) = ((H - \frac{1}{n}J)Y)^T (H - \frac{1}{n}J)Y = Y^T (H - \frac{1}{n}J)^T (H - \frac{1}{n}J)Y = Y^T (H - \frac{1}{n}J)Y$$

$$c. SSE = (Y - HY)^T (Y - HY) = ((I - H)Y)^T (I - H)Y = Y^T (I - H)^T (I - H)Y = Y^T (I - H)Y$$

– The ANOVA table for the multiple regression model is

|            | SS                              | DF      | MS                        |
|------------|---------------------------------|---------|---------------------------|
| Regression | $SSR = Y^T (H - \frac{1}{n}J)Y$ | $p - 1$ | $MSR = \frac{SSR}{p - 1}$ |
| Error      | $SSE = Y^T (I - H)Y$            | $n - p$ | $MSE = \frac{SSE}{n - p}$ |
| Total      | $SST = Y^T (I - \frac{1}{n}J)Y$ | $n - 1$ |                           |

– For the simple linear regression model, there is only one slope parameter about which one can perform hypothesis tests. And the  $F$  test from the ANOVA table is testing the slope parameter. However, for the multiple linear regression model, there are three different hypothesis tests for slopes that one could conduct. They are:

a. a hypothesis test for testing that one slope parameter is 0  $H_0: \beta_1 = 0$  vs.  $H_a: \beta_1 \neq 0$  (+-test)

b. a hypothesis test for testing that all of the slope parameters are 0  $H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0$  vs.  $H_a: \text{at least one is not zero}$

c. a hypothesis test for testing that a subset (more than one, but not all) of the slope parameters are 0  $H_0: \beta_2 = \beta_3 = 0$  vs.  $H_a: \text{at least one of } \beta_2, \beta_3 \text{ is not zero}$  (F-test)

– The above hypotheses can be tested using the analysis of variance  $F$ -test. The **general linear  $F$ -test** involves three basic steps:

- Define a larger full model: By “larger,” we mean one with more parameters (i.e., containing additional predictors).
- Define a smaller reduced model: By “smaller,” we mean one with fewer parameters (a restricted version of the full model).
- Construct the  $F$ -statistic: to decide whether to reject the smaller (reduced) model in favor of the larger (full) model.

– In this framework, the null hypothesis  $H_0$  always pertains to the *reduced model*, while the alternative hypothesis  $H_a$  always pertains to the *full model*.

– Let  $SSE_R$  and  $SSE_F$  denote the error sums of squares for the reduced and full models, respectively. If the reduced model has  $df_R$  error degrees of freedom and the full model has  $df_F$  error degrees of freedom, then the test statistic is

$$F^* = \frac{(SSE_R - SSE_F) / (df_R - df_F)}{MSE_F}, \quad MSE_F = \frac{SSE_F}{df_F}.$$

$16198 - 13322 = 2875.7$   
 $\uparrow \quad \uparrow \quad \uparrow \quad \uparrow$   
 $36 - 34 = 2$   
 $34 \quad 34$   
 $13322$   
 $34$

– Under  $H_0$ , this statistic follows an  $F$  distribution with  $(df_R - df_F)$  and  $df_F$  degrees of freedom. The decision rule is:

Reject  $H_0$  if  $F^* > F_{1-\alpha; df_R - df_F, df_F}$ .

### Example 1

Are a person's brain size and body size predictive of his or her intelligence? To investigate this question, researchers (Willerman, et al., 1991) collected the following data (`IQsize.csv`) on a sample of  $n = 38$  college students: (1) performance IQ scores (PIQ) from the revised Wechsler Adult Intelligence Scale, which served as the investigator's measure of the individual's intelligence, (2) brain size based on the count obtained from MRI scans (measured in units of count/10,000), (3) height in inches, and (4) weight in pounds.

- a. Fit a multiple regression model above using the dataset. Report the estimated regression equation.

$$\hat{PIQ} = 111 + 2.06 \times \text{Brain} - 2.73 \times \text{Ht} + .0006 \times \text{wt}$$

- b. Interpret the estimated coefficient for brain size.

Holding height and weight constant, each additional unit of brain size (10,000 MRI counts) is associated with an estimated increase of about 2.06 points in PIQ.

- c. Perform an overall  $F$ -test to determine whether at least one of the predictors (brain size, height, weight) is significantly associated with PIQ. State the null and alternative hypotheses and report your conclusion.

$$H_0: \beta_1 = \beta_2 = \beta_3 = 0 \quad \text{v.s.} \quad H_a: \text{at least one of } \beta_1, \beta_2, \beta_3 \text{ is not zero.}$$

$$F\text{-test statistic} = 4.741 \sim F_{(3, 34)} \quad p\text{-value} = .007215 < \alpha = .05$$

$\downarrow$  3 predictors  $\rightarrow n-p = 38-4 = 34$

Yes, at least one of the is significant

- d. Determine whether brain size is a significant predictor of PIQ, controlling for height and weight. Clearly state the hypotheses, p-value, and conclusion at the  $\alpha = 0.05$  level.

$$H_0: \beta_1 = 0 \quad \text{v.s.} \quad H_a: \beta_1 \neq 0 \quad \text{Yes!}$$

$$t\text{-stat} = 3.657 \quad p\text{-value} = .000856 < \alpha = .05$$

- e. Determine whether height is a significant predictor of PIQ, controlling for brain size and weight. Clearly state the hypotheses, p-value, and conclusion at the  $\alpha = 0.05$  level.

$$H_0: \beta_2 = 0 \quad \text{v.s.} \quad H_a: \beta_2 \neq 0$$

$\downarrow$  res

$$t\text{-stat} = -2.22 \quad p\text{-value} = .033 < \alpha = .05$$

- f. Determine whether height and weight jointly contribute to explaining PIQ after accounting for brain size. Clearly state the hypotheses, p-value, and conclusion at the  $\alpha = 0.05$  level.

$$H_0: \beta_2 = \beta_3 = 0 \quad \text{v.s.} \quad H_a: \text{at least one } \beta_2, \beta_3 \text{ isn't zero.}$$

$$F\text{-stat} = 3.6696, \quad p\text{-value} = .03606$$

- g. Report the coefficient of determination and interpret it in the context of this study.

This means about 29.49% of the variation in PIQ scores is explained by brain size, height, and weight.

- Final notes:

1. Brain size is a strong and reliable predictor of performance IQ in this sample. Height adds a weaker, negative effect once brain size is accounted for. Weight provides no predictive value.
2. The model explains about 30% of the variation in PIQ scores, leaving considerable unexplained variation, so intelligence is clearly influenced by many factors beyond body and brain size.

## Example 2

A young chef is opening a new Italian restaurant in Manhattan, to be featured on an international reality TV show. The restaurant's goals are to provide the highest quality Italian food, incorporate state-of-the-art décor, and set a new standard for service. The location will be no further south than the Flatiron District and will be situated either east or west of Fifth Avenue. Your role is to determine a competitive dinner menu price by analyzing survey data from 168 Italian restaurants in the area. The dataset ([NYC.csv](#)) includes the following variables:

- Price: the cost of dinner in US dollars, including one drink and tip
- Food: customer rating of the food out of 30
- Décor: customer rating of the décor out of 30
- Service: customer rating of the service out of 30
- East: dummy variable equal to 1 if east of Fifth Avenue and 0 if west

- Fit a multiple regression model using all of the predictors. Report the estimated regression equation.
- Among Food, Décor, and Service, which has the largest estimated *effect* on Price? Is the largest effect also the most statistically significant?
- If the aim is to choose the location to maximize the achievable dinner price, should the restaurant be on the east or west side of Fifth Avenue? Justify using the sign, magnitude, and significance of the coefficient.
- Does it seem possible to achieve a price premium for “setting a new standard for high-quality service in Manhattan” for Italian restaurants?
- Test whether Food, Décor, and Service jointly contribute to Price after controlling for East. Clearly state the hypotheses, p-value, and conclusion at the  $\alpha = 0.05$  level.
- Test whether East and Service jointly add explanatory power beyond Food and Décor. Clearly state the hypotheses, p-value, and conclusion at the  $\alpha = 0.05$  level.