

	TV	radio	newspaper	sales
TV	1.0000	0.0548	0.0567	0.7822
radio		1.0000	0.3541	0.5762
newspaper			1.0000	0.2283
sales				1.0000

**TABLE 3.5.** Correlation matrix for TV, radio, newspaper, and sales for the Advertising data.

### 3.2.2 Some Important Questions

When we perform multiple linear regression, we usually are interested in answering a few important questions.

1. Is at least one of the predictors  $X_1, X_2, \dots, X_p$  useful in predicting the response?
2. Do all the predictors help to explain  $Y$ , or is only a subset of the predictors useful?
3. How well does the model fit the data?
4. Given a set of predictor values, what response value should we predict, and how accurate is our prediction?

We now address each of these questions in turn.

#### One: Is There a Relationship Between the Response and Predictors?

Recall that in the simple linear regression setting, in order to determine whether there is a relationship between the response and the predictor we can simply check whether  $\beta_1 = 0$ . In the multiple regression setting with  $p$  predictors, we need to ask whether all of the regression coefficients are zero, i.e. whether  $\beta_1 = \beta_2 = \dots = \beta_p = 0$ . As in the simple linear regression setting, we use a hypothesis test to answer this question. We test the null hypothesis,

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$$

versus the alternative

$$H_a : \text{at least one } \beta_j \text{ is non-zero.}$$

This hypothesis test is performed by computing the  $F$ -statistic,

$$F = \frac{(\text{TSS} - \text{RSS})/p}{\text{RSS}/(n - p - 1)}, \quad (3.23)$$

F-statistic

Quantity	Value
Residual standard error	1.69
$R^2$	0.897
F-statistic	570

**TABLE 3.6.** More information about the least squares model for the regression of number of units sold on TV, newspaper, and radio advertising budgets in the Advertising data. Other information about this model was displayed in Table 3.4.

where, as with simple linear regression,  $TSS = \sum (y_i - \bar{y})^2$  and  $RSS = \sum (y_i - \hat{y}_i)^2$ . If the linear model assumptions are correct, one can show that

$$E\{RSS/(n - p - 1)\} = \sigma^2$$

and that, provided  $H_0$  is true,

$$E\{(TSS - RSS)/p\} = \sigma^2.$$

Hence, when there is no relationship between the response and predictors, one would expect the F-statistic to take on a value close to 1. On the other hand, if  $H_a$  is true, then  $E\{(TSS - RSS)/p\} > \sigma^2$ , so we expect  $F$  to be greater than 1.

The F-statistic for the multiple linear regression model obtained by regressing sales onto radio, TV, and newspaper is shown in Table 3.6. In this example the F-statistic is 570. Since this is far larger than 1, it provides compelling evidence against the null hypothesis  $H_0$ . In other words, the large F-statistic suggests that at least one of the advertising media must be related to sales. However, what if the F-statistic had been closer to 1? How large does the F-statistic need to be before we can reject  $H_0$  and conclude that there is a relationship? It turns out that the answer depends on the values of  $n$  and  $p$ . When  $n$  is large, an F-statistic that is just a little larger than 1 might still provide evidence against  $H_0$ . In contrast, a larger F-statistic is needed to reject  $H_0$  if  $n$  is small. When  $H_0$  is true and the errors  $\epsilon_i$  have a normal distribution, the F-statistic follows an F-distribution.<sup>6</sup> For any given value of  $n$  and  $p$ , any statistical software package can be used to compute the p-value associated with the F-statistic using this distribution. Based on this p-value, we can determine whether or not to reject  $H_0$ . For the advertising data, the p-value associated with the F-statistic in Table 3.6 is essentially zero, so we have extremely strong evidence that at least one of the media is associated with increased sales.

In (3.23) we are testing  $H_0$  that all the coefficients are zero. Sometimes we want to test that a particular subset of  $q$  of the coefficients are zero. This corresponds to a null hypothesis

$$H_0 : \beta_{p-q+1} = \beta_{p-q+2} = \dots = \beta_p = 0,$$

<sup>6</sup>Even if the errors are not normally-distributed, the F-statistic approximately follows an F-distribution provided that the sample size  $n$  is large.

where for convenience we have put the variables chosen for omission at the end of the list. In this case we fit a second model that uses all the variables *except* those last  $q$ . Suppose that the residual sum of squares for that model is  $\text{RSS}_0$ . Then the appropriate F-statistic is

$$F = \frac{(\text{RSS}_0 - \text{RSS})/q}{\text{RSS}/(n - p - 1)}. \quad (3.24)$$

Notice that in Table 3.4, for each individual predictor a t-statistic and a p-value were reported. These provide information about whether each individual predictor is related to the response, after adjusting for the other predictors. It turns out that each of these are exactly equivalent<sup>7</sup> to the F-test that omits that single variable from the model, leaving all the others in—i.e.  $q=1$  in (3.24). So it reports the *partial effect* of adding that variable to the model. For instance, as we discussed earlier, these p-values indicate that **TV** and **radio** are related to **sales**, but that there is no evidence that **newspaper** is associated with **sales**, in the presence of these two.

Given these individual p-values for each variable, why do we need to look at the overall F-statistic? After all, it seems likely that *if any one of the p-values for the individual variables is very small, then at least one of the predictors is related to the response*. However, this logic is flawed, especially *when the number of predictors  $p$  is large*.

For instance, consider an example in which  $p = 100$  and  $H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$  is true, so no variable is truly associated with the response. In this situation, about 5% of the p-values associated with each variable (of the type shown in Table 3.4) will be below 0.05 by chance. *In other words, we expect to see approximately five small p-values even in the absence of any true association between the predictors and the response*. In fact, we are almost guaranteed that we will observe at least one p-value below 0.05 by chance! Hence, if we use the individual t-statistics and associated p-values in order to decide whether or not there is any association between the variables and the response, there is a very high chance that we will incorrectly conclude that there is a relationship. *However, the F-statistic does not suffer from this problem because it adjusts for the number of predictors*. Hence, if  $H_0$  is true, there is only a 5% chance that the F-statistic will result in a p-value below 0.05, regardless of the number of predictors or the number of observations.

The approach of using an **F-statistic to test for any association between the predictors and the response works when  $p$  is relatively small, and certainly small compared to  $n$** . However, sometimes we have a very large number of variables. If  $p > n$  then there are more coefficients  $\beta_j$  to estimate than observations from which to estimate them. In this case we cannot even fit the multiple linear regression model using least squares, so the

---

<sup>7</sup>The square of each t-statistic is the corresponding F-statistic.

F-statistic cannot be used, and neither can most of the other concepts that we have seen so far in this chapter. When  $p$  is large, some of the approaches discussed in the next section, such as *forward selection*, can be used. This *high-dimensional* setting is discussed in greater detail in Chapter 6.

high-  
dimensional

### Two: Deciding on Important Variables

As discussed in the previous section, the first step in a multiple regression analysis is to compute the F-statistic and to examine the associated p-value. If we conclude on the basis of that p-value that at least one of the predictors is related to the response, then it is natural to wonder which are the guilty ones! We could look at the individual p-values as in Table 3.4, but as discussed, if  $p$  is large we are likely to make some false discoveries.

It is possible that all of the predictors are associated with the response, but it is more often the case that the response is only related to a subset of the predictors. The task of determining which predictors are associated with the response, in order to fit a single model involving only those predictors, is referred to as *variable selection*. The variable selection problem is studied extensively in Chapter 6, and so here we will provide only a brief outline of some classical approaches.

variable  
selection

Ideally, we would like to perform variable selection by trying out a lot of different models, each containing a different subset of the predictors. For instance, if  $p = 2$ , then we can consider four models: (1) a model containing no variables, (2) a model containing  $X_1$  only, (3) a model containing  $X_2$  only, and (4) a model containing both  $X_1$  and  $X_2$ . We can then select the best model out of all of the models that we have considered. How do we determine which model is best? Various statistics can be used to judge the quality of a model. These include Mallow's  $C_p$ , Akaike information criterion (AIC), Bayesian information criterion (BIC), and adjusted  $R^2$ . These are discussed in more detail in Chapter 6. We can also determine which model is best by plotting various model outputs, such as the residuals, in order to search for patterns.

Mallow's  $C_p$   
Akaike  
information  
criterion  
Bayesian  
information  
criterion  
adjusted  $R^2$

Unfortunately, there are a total of  $2^p$  models that contain subsets of  $p$  variables. This means that even for moderate  $p$ , trying out every possible subset of the predictors is infeasible. For instance, we saw that if  $p = 2$ , then there are  $2^2 = 4$  models to consider. But if  $p = 30$ , then we must consider  $2^{30} = 1,073,741,824$  models! This is not practical. Therefore, unless  $p$  is very small, we cannot consider all  $2^p$  models, and instead we need an automated and efficient approach to choose a smaller set of models to consider. There are three classical approaches for this task:

- *Forward selection.* We begin with the null model—a model that contains an intercept but no predictors. We then fit  $p$  simple linear regressions and add to the null model the variable that results in the lowest RSS. We then add to that model the variable that results

forward  
selection  
null model

in the lowest RSS for the new two-variable model. This approach is continued until some stopping rule is satisfied.

- **Backward selection.** We start with all variables in the model, and remove the variable with the largest p-value—that is, the variable that is the least statistically significant. The new  $(p - 1)$ -variable model is fit, and the variable with the largest p-value is removed. This procedure continues until a stopping rule is reached. For instance, we may stop when all remaining variables have a p-value below some threshold. backward selection
- **Mixed selection.** This is a combination of forward and backward selection. We start with no variables in the model, and as with forward selection, we add the variable that provides the best fit. We continue to add variables one-by-one. Of course, as we noted with the Advertising example, the p-values for variables can become larger as new predictors are added to the model. Hence, if at any point the p-value for one of the variables in the model rises above a certain threshold, then we remove that variable from the model. We continue to perform these forward and backward steps until all variables in the model have a sufficiently low p-value, and all variables outside the model would have a large p-value if added to the model. mixed selection

Backward selection cannot be used if  $p > n$ , while forward selection can always be used. Forward selection is a greedy approach, and might include variables early that later become redundant. Mixed selection can remedy this.

### Three: Model Fit

Two of the most common numerical measures of model fit are the **RSE** and  **$R^2$** , the **fraction of variance explained**. These quantities are computed and interpreted in the same fashion as for simple linear regression.

Recall that in simple regression,  **$R^2$  is the square of the correlation of the response and the variable**. In multiple linear regression, it turns out that it equals  $\text{Cor}(Y, \hat{Y})^2$ , the square of the correlation between the response and the fitted linear model; in fact one property of the fitted linear model is that it maximizes this correlation among all possible linear models.

An  **$R^2$  value close to 1 indicates that the model explains a large portion of the variance in the response variable**. As an example, we saw in Table 3.6 that for the Advertising data, the model that uses all three advertising media to predict sales has an  $R^2$  of 0.8972. On the other hand, the model that uses only TV and radio to predict sales has an  $R^2$  value of 0.89719. In other words, there is a *small* increase in  $R^2$  if we include newspaper advertising in the model that already contains TV and radio advertising, even though we saw earlier that the p-value for newspaper advertising in Table 3.4 is not significant. It turns out that  **$R^2$  will always increase when more variables**

are added to the model, even if those variables are only weakly associated with the response. This is due to the fact that adding another variable to the least squares equations must allow us to fit the training data (though not necessarily the testing data) more accurately. Thus, the  $R^2$  statistic, which is also computed on the training data, must increase. The fact that adding newspaper advertising to the model containing only TV and radio advertising leads to just a tiny increase in  $R^2$  provides additional evidence that newspaper can be dropped from the model. Essentially, newspaper provides no real improvement in the model fit to the training samples, and its inclusion will likely lead to poor results on independent test samples due to overfitting.

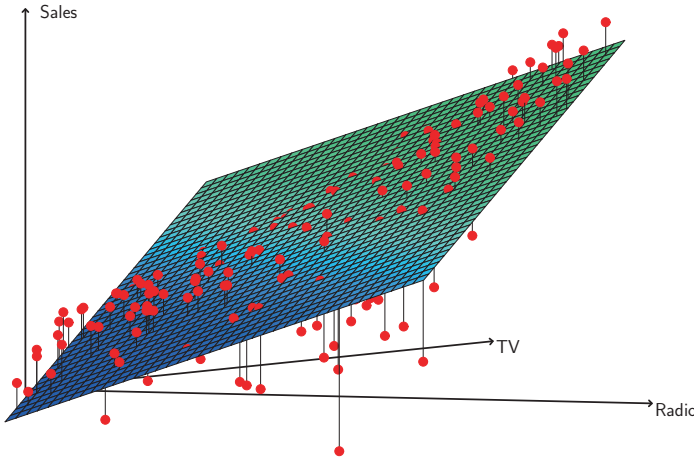
In contrast, the model containing only TV as a predictor had an  $R^2$  of 0.61 (Table 3.2). Adding radio to the model leads to a substantial improvement in  $R^2$ . This implies that a model that uses TV and radio expenditures to predict sales is substantially better than one that uses only TV advertising. We could further quantify this improvement by looking at the p-value for the radio coefficient in a model that contains only TV and radio as predictors.

The model that contains only TV and radio as predictors has an RSE of 1.681, and the model that also contains newspaper as a predictor has an RSE of 1.686 (Table 3.6). In contrast, the model that contains only TV has an RSE of 3.26 (Table 3.2). This corroborates our previous conclusion that a model that uses TV and radio expenditures to predict sales is much more accurate (on the training data) than one that only uses TV spending. Furthermore, given that TV and radio expenditures are used as predictors, there is no point in also using newspaper spending as a predictor in the model. The observant reader may wonder how RSE can increase when newspaper is added to the model given that RSS must decrease. In general RSE is defined as

$$\text{RSE} = \sqrt{\frac{1}{n - p - 1} \text{RSS}}, \quad (3.25)$$

which simplifies to (3.15) for a simple linear regression. Thus, models with more variables can have higher RSE if the decrease in RSS is small relative to the increase in  $p$ .

In addition to looking at the RSE and  $R^2$  statistics just discussed, it can be useful to plot the data. Graphical summaries can reveal problems with a model that are not visible from numerical statistics. For example, Figure 3.5 displays a three-dimensional plot of TV and radio versus sales. We see that some observations lie above and some observations lie below the least squares regression plane. In particular, the linear model seems to overestimate sales for instances in which most of the advertising money was spent exclusively on either TV or radio. It underestimates sales for instances where the budget was split between the two media. This pronounced non-linear pattern cannot be modeled accurately using linear re-



**FIGURE 3.5.** For the Advertising data, a linear regression fit to sales using TV and radio as predictors. From the pattern of the residuals, we can see that there is a pronounced non-linear relationship in the data. The positive residuals (those visible above the surface), tend to lie along the 45-degree line, where TV and Radio budgets are split evenly. The negative residuals (most not visible), tend to lie away from this line, where budgets are more lopsided.

gression. It suggests a *synergy* or *interaction* effect between the advertising media, whereby combining the media together results in a bigger boost to sales than using any single medium. In Section 3.3.2, we will discuss extending the linear model to accommodate such synergistic effects through the use of interaction terms.

#### Four: Predictions

Once we have fit the multiple regression model, it is straightforward to apply (3.21) in order to predict the response  $Y$  on the basis of a set of values for the predictors  $X_1, X_2, \dots, X_p$ . However, there are three sorts of uncertainty associated with this prediction.

1. The coefficient estimates  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$  are estimates for  $\beta_0, \beta_1, \dots, \beta_p$ . That is, the *least squares plane*

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_p X_p$$

is only an estimate for the *true population regression plane*

$$f(X) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p.$$

The inaccuracy in the coefficient estimates is related to the *reducible error* from Chapter 2. We can compute a *confidence interval* in order to determine how close  $\hat{Y}$  will be to  $f(X)$ .

2. Of course, in practice assuming a linear model for  $f(X)$  is almost always an approximation of reality, so there is an additional source of potentially reducible error which we call *model bias*. So when we use a linear model, we are in fact estimating the best linear approximation to the true surface. However, here we will ignore this discrepancy, and operate as if the linear model were correct.
3. Even if we knew  $f(X)$ —that is, even if we knew the true values for  $\beta_0, \beta_1, \dots, \beta_p$ —the response value cannot be predicted perfectly because of the random error  $\epsilon$  in the model (3.21). In Chapter 2, we referred to this as the *irreducible error*. How much will  $Y$  vary from  $\hat{Y}$ ? We use *prediction intervals* to answer this question. Prediction intervals are always wider than confidence intervals, because they incorporate both the error in the estimate for  $f(X)$  (the reducible error) and the uncertainty as to how much an individual point will differ from the population regression plane (the irreducible error).

We use a *confidence interval* to quantify the uncertainty surrounding the average sales over a large number of cities. For example, given that \$100,000 is spent on TV advertising and \$20,000 is spent on radio advertising in each city, the 95 % confidence interval is [10,985, 11,528]. We interpret this to mean that 95 % of intervals of this form will contain the true value of  $f(X)$ .<sup>8</sup> On the other hand, a *prediction interval* can be used to quantify the uncertainty surrounding sales for a particular city. Given that \$100,000 is spent on TV advertising and \$20,000 is spent on radio advertising in that city the 95 % prediction interval is [7,930, 14,580]. We interpret this to mean that 95 % of intervals of this form will contain the true value of  $Y$  for this city. Note that both intervals are centered at 11,256, but that the prediction interval is substantially wider than the confidence interval, reflecting the increased uncertainty about sales for a given city in comparison to the average sales over many locations.

confidence  
interval

prediction  
interval

## 3.3 Other Considerations in the Regression Model

### 3.3.1 Qualitative Predictors

In our discussion so far, we have assumed that all variables in our linear regression model are *quantitative*. But in practice, this is not necessarily the case; often some predictors are *qualitative*.

---

<sup>8</sup>In other words, if we collect a large number of data sets like the Advertising data set, and we construct a confidence interval for the average sales on the basis of each data set (given \$100,000 in TV and \$20,000 in radio advertising), then 95 % of these confidence intervals will contain the true value of average sales.