

	Coefficient	Std. error	<i>t</i> -statistic	<i>p</i> -value
<b>Intercept</b>	531.00	46.32	11.464	< 0.0001
<b>region[South]</b>	−12.50	56.68	−0.221	0.8260
<b>region[West]</b>	−18.69	65.02	−0.287	0.7740

**TABLE 3.8.** Least squares coefficient estimates associated with the regression of **balance** onto **region** in the **Credit** data set. The linear model is given in (3.30). That is, **region** is encoded via two dummy variables (3.28) and (3.29).

From Table 3.8, we see that the estimated **balance** for the baseline, East, is \$531.00. It is estimated that those in the South will have \$18.69 less debt than those in the East, and that those in the West will have \$12.50 less debt than those in the East. However, the *p*-values associated with the coefficient estimates for the two dummy variables are very large, suggesting no statistical evidence of a real difference in average credit card balance between South and East or between West and East.<sup>12</sup> Once again, the level selected as the baseline category is arbitrary, and the final predictions for each group will be the same regardless of this choice. However, the coefficients and their *p*-values do depend on the choice of dummy variable coding. Rather than rely on the individual coefficients, we can use an *F*-test to test  $H_0 : \beta_1 = \beta_2 = 0$ ; this does not depend on the coding. This *F*-test has a *p*-value of 0.96, indicating that we cannot reject the null hypothesis that there is no relationship between **balance** and **region**.

Using this dummy variable approach presents no difficulties when incorporating both quantitative and qualitative predictors. For example, to regress **balance** on both a quantitative variable such as **income** and a qualitative variable such as **student**, we must simply create a dummy variable for **student** and then fit a multiple regression model using **income** and the dummy variable as predictors for credit card balance.

There are many different ways of coding qualitative variables besides the dummy variable approach taken here. All of these approaches lead to equivalent model fits, but the coefficients are different and have different interpretations, and are designed to measure particular *contrasts*. This topic is beyond the scope of the book.

contrast

### 3.3.2 Extensions of the Linear Model

The standard linear regression model (3.19) provides interpretable results and works quite well on many real-world problems. However, it makes several highly restrictive assumptions that are often violated in practice. Two of the most important assumptions state that the relationship between the predictors and response are *additive* and *linear*. The additivity assumption means that the association between a predictor  $X_j$  and the response  $Y$  does not depend on the values of the other predictors. The linearity assumption states that the change in the response  $Y$  associated with a one-unit change in  $X_j$  is constant, regardless of the value of  $X_j$ . In later chapters of this book, we examine a number of sophisticated methods that relax these two

additive  
linear

<sup>12</sup>There could still in theory be a difference between South and West, although the data here does not suggest any difference.

assumptions. Here, we briefly examine some common classical approaches for extending the linear model.

### Removing the Additive Assumption

In our previous analysis of the **Advertising** data, we concluded that both **TV** and **radio** seem to be associated with **sales**. The linear models that formed the basis for this conclusion assumed that the effect on **sales** of increasing one advertising medium is independent of the amount spent on the other media. For example, the linear model (3.20) states that the average increase in **sales** associated with a one-unit increase in **TV** is always  $\beta_1$ , regardless of the amount spent on **radio**.

However, this simple model may be incorrect. Suppose that spending money on radio advertising actually increases the effectiveness of TV advertising, so that the slope term for **TV** should increase as **radio** increases. In this situation, given a fixed budget of \$100,000, spending half on **radio** and half on **TV** may increase **sales** more than allocating the entire amount to either **TV** or to **radio**. In marketing, this is known as a *synergy* effect, and in statistics it is referred to as an *interaction* effect. Figure 3.5 suggests that such an effect may be present in the advertising data. Notice that when levels of either **TV** or **radio** are low, then the true **sales** are lower than predicted by the linear model. But when advertising is split between the two media, then the model tends to underestimate **sales**.

Consider the standard linear regression model with two variables,

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon.$$

According to this model, a one-unit increase in  $X_1$  is associated with an average increase in  $Y$  of  $\beta_1$  units. Notice that the presence of  $X_2$  does not alter this statement—that is, regardless of the value of  $X_2$ , a one-unit increase in  $X_1$  is associated with a  $\beta_1$ -unit increase in  $Y$ . One way of extending this model is to include a third predictor, called an *interaction term*, which is constructed by computing the product of  $X_1$  and  $X_2$ . This results in the model

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \epsilon. \quad (3.31)$$

How does inclusion of this interaction term relax the additive assumption? Notice that (3.31) can be rewritten as

$$\begin{aligned} Y &= \beta_0 + (\beta_1 + \beta_3 X_2) X_1 + \beta_2 X_2 + \epsilon \\ &= \beta_0 + \tilde{\beta}_1 X_1 + \beta_2 X_2 + \epsilon \end{aligned} \quad (3.32)$$

where  $\tilde{\beta}_1 = \beta_1 + \beta_3 X_2$ . Since  $\tilde{\beta}_1$  is now a function of  $X_2$ , the association between  $X_1$  and  $Y$  is no longer constant: a change in the value of  $X_2$  will change the association between  $X_1$  and  $Y$ . A similar argument shows that a change in the value of  $X_1$  changes the association between  $X_2$  and  $Y$ .

For example, suppose that we are interested in studying the productivity of a factory. We wish to predict the number of **units** produced on the basis of the number of production **lines** and the total number of **workers**. It seems likely that the effect of increasing the number of production lines

	Coefficient	Std. error	<i>t</i> -statistic	<i>p</i> -value
Intercept	6.7502	0.248	27.23	< 0.0001
TV	0.0191	0.002	12.70	< 0.0001
radio	0.0289	0.009	3.24	0.0014
TV×radio	0.0011	0.000	20.73	< 0.0001

**TABLE 3.9.** For the **Advertising** data, least squares coefficient estimates associated with the regression of **sales** onto **TV** and **radio**, with an interaction term, as in (3.33).

will depend on the number of workers, since if no workers are available to operate the lines, then increasing the number of lines will not increase production. This suggests that it would be appropriate to include an interaction term between **lines** and **workers** in a linear model to predict **units**. Suppose that when we fit the model, we obtain

$$\begin{aligned}\text{units} &\approx 1.2 + 3.4 \times \text{lines} + 0.22 \times \text{workers} + 1.4 \times (\text{lines} \times \text{workers}) \\ &= 1.2 + (3.4 + 1.4 \times \text{workers}) \times \text{lines} + 0.22 \times \text{workers}.\end{aligned}$$

In other words, adding an additional line will increase the number of units produced by  $3.4 + 1.4 \times \text{workers}$ . Hence the more **workers** we have, the stronger will be the effect of **lines**.

We now return to the **Advertising** example. A linear model that uses **radio**, **TV**, and an interaction between the two to predict **sales** takes the form

$$\begin{aligned}\text{sales} &= \beta_0 + \beta_1 \times \text{TV} + \beta_2 \times \text{radio} + \beta_3 \times (\text{radio} \times \text{TV}) + \epsilon \\ &= \beta_0 + (\beta_1 + \beta_3 \times \text{radio}) \times \text{TV} + \beta_2 \times \text{radio} + \epsilon.\end{aligned}\quad (3.33)$$

We can interpret  $\beta_3$  as the increase in the effectiveness of TV advertising associated with a one-unit increase in radio advertising (or vice-versa). The coefficients that result from fitting the model (3.33) are given in Table 3.9.

The results in Table 3.9 strongly suggest that the model that includes the interaction term is superior to the model that contains only *main effects*. The *p*-value for the interaction term, **TV×radio**, is extremely low, indicating that there is strong evidence for  $H_a : \beta_3 \neq 0$ . In other words, it is clear that the true relationship is not additive. The  $R^2$  for the model (3.33) is 96.8%, compared to only 89.7% for the model that predicts **sales** using **TV** and **radio** without an interaction term. This means that  $(96.8 - 89.7)/(100 - 89.7) = 69\%$  of the variability in **sales** that remains after fitting the additive model has been explained by the interaction term. The coefficient estimates in Table 3.9 suggest that an increase in TV advertising of \$1,000 is associated with increased sales of  $(\hat{\beta}_1 + \hat{\beta}_3 \times \text{radio}) \times 1,000 = 19 + 1.1 \times \text{radio}$  units. And an increase in radio advertising of \$1,000 will be associated with an increase in sales of  $(\hat{\beta}_2 + \hat{\beta}_3 \times \text{TV}) \times 1,000 = 29 + 1.1 \times \text{TV}$  units.

In this example, the *p*-values associated with **TV**, **radio**, and the interaction term all are statistically significant (Table 3.9), and so it is obvious that all three variables should be included in the model. However, it is sometimes the case that an interaction term has a very small *p*-value, but the associated main effects (in this case, **TV** and **radio**) do not. The *hierarchical principle* states that if we include an interaction in a model, we

main effect

hierarchical principle

should also include the main effects, even if the  $p$ -values associated with their coefficients are not significant. In other words, if the interaction between  $X_1$  and  $X_2$  seems important, then we should include both  $X_1$  and  $X_2$  in the model even if their coefficient estimates have large  $p$ -values. The rationale for this principle is that if  $X_1 \times X_2$  is related to the response, then whether or not the coefficients of  $X_1$  or  $X_2$  are exactly zero is of little interest. Also  $X_1 \times X_2$  is typically correlated with  $X_1$  and  $X_2$ , and so leaving them out tends to alter the meaning of the interaction.

In the previous example, we considered an interaction between **TV** and **radio**, both of which are quantitative variables. However, the concept of interactions applies just as well to qualitative variables, or to a combination of quantitative and qualitative variables. In fact, an interaction between a qualitative variable and a quantitative variable has a particularly nice interpretation. Consider the **Credit** data set from Section 3.3.1, and suppose that we wish to predict **balance** using the **income** (quantitative) and **student** (qualitative) variables. In the absence of an interaction term, the model takes the form

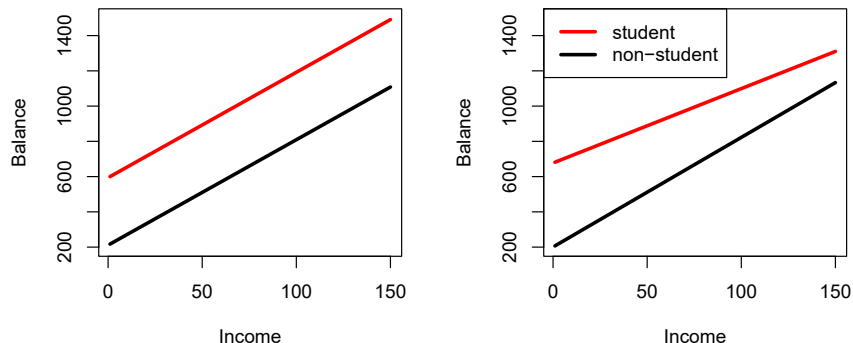
$$\begin{aligned} \text{balance}_i &\approx \beta_0 + \beta_1 \times \text{income}_i + \begin{cases} \beta_2 & \text{if } i\text{th person is a student} \\ 0 & \text{if } i\text{th person is not a student} \end{cases} \\ &= \beta_1 \times \text{income}_i + \begin{cases} \beta_0 + \beta_2 & \text{if } i\text{th person is a student} \\ \beta_0 & \text{if } i\text{th person is not a student.} \end{cases} \end{aligned} \quad (3.34)$$

Notice that this amounts to fitting two parallel lines to the data, one for students and one for non-students. The lines for students and non-students have different intercepts,  $\beta_0 + \beta_2$  versus  $\beta_0$ , but the same slope,  $\beta_1$ . This is illustrated in the left-hand panel of Figure 3.7. The fact that the lines are parallel means that the average effect on **balance** of a one-unit increase in **income** does not depend on whether or not the individual is a student. This represents a potentially serious limitation of the model, since in fact a change in **income** may have a very different effect on the credit card balance of a student versus a non-student.

This limitation can be addressed by adding an interaction variable, created by multiplying **income** with the dummy variable for **student**. Our model now becomes

$$\begin{aligned} \text{balance}_i &\approx \beta_0 + \beta_1 \times \text{income}_i + \begin{cases} \beta_2 + \beta_3 \times \text{income}_i & \text{if student} \\ 0 & \text{if not student} \end{cases} \\ &= \begin{cases} (\beta_0 + \beta_2) + (\beta_1 + \beta_3) \times \text{income}_i & \text{if student} \\ \beta_0 + \beta_1 \times \text{income}_i & \text{if not student.} \end{cases} \end{aligned} \quad (3.35)$$

Once again, we have two different regression lines for the students and the non-students. But now those regression lines have different intercepts,  $\beta_0 + \beta_2$  versus  $\beta_0$ , as well as different slopes,  $\beta_1 + \beta_3$  versus  $\beta_1$ . This allows for the possibility that changes in income may affect the credit card balances of students and non-students differently. The right-hand panel of Figure 3.7



**FIGURE 3.7.** For the **Credit** data, the least squares lines are shown for prediction of **balance** from **income** for students and non-students. Left: The model (3.34) was fit. There is no interaction between **income** and **student**. Right: The model (3.35) was fit. There is an interaction term between **income** and **student**.

shows the estimated relationships between **income** and **balance** for students and non-students in the model (3.35). We note that the slope for students is lower than the slope for non-students. This suggests that increases in income are associated with smaller increases in credit card balance among students as compared to non-students.

### Non-linear Relationships

As discussed previously, the linear regression model (3.19) assumes a linear relationship between the response and predictors. But in some cases, the true relationship between the response and the predictors may be non-linear. Here we present a very simple way to directly extend the linear model to accommodate non-linear relationships, using *polynomial regression*. In later chapters, we will present more complex approaches for performing non-linear fits in more general settings.

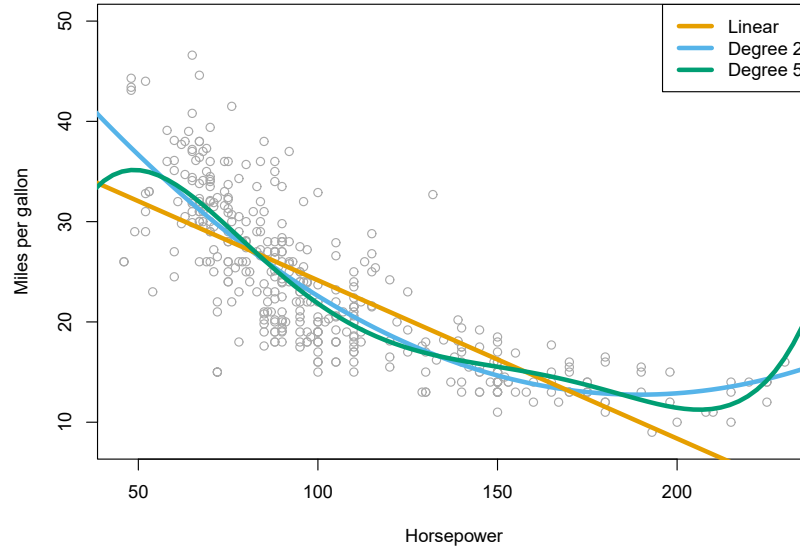
polynomial  
regression

Consider Figure 3.8, in which the **mpg** (gas mileage in miles per gallon) versus **horsepower** is shown for a number of cars in the **Auto** data set. The orange line represents the linear regression fit. There is a pronounced relationship between **mpg** and **horsepower**, but it seems clear that this relationship is in fact non-linear: the data suggest a curved relationship. A simple approach for incorporating non-linear associations in a linear model is to include transformed versions of the predictors. For example, the points in Figure 3.8 seem to have a *quadratic* shape, suggesting that a model of the form

quadratic

$$\text{mpg} = \beta_0 + \beta_1 \times \text{horsepower} + \beta_2 \times \text{horsepower}^2 + \epsilon \quad (3.36)$$

may provide a better fit. Equation 3.36 involves predicting **mpg** using a non-linear function of **horsepower**. *But it is still a linear model!* That is, (3.36) is simply a multiple linear regression model with  $X_1 = \text{horsepower}$  and  $X_2 = \text{horsepower}^2$ . So we can use standard linear regression software to estimate  $\beta_0, \beta_1$ , and  $\beta_2$  in order to produce a non-linear fit. The blue curve in Figure 3.8 shows the resulting quadratic fit to the data. The quadratic



**FIGURE 3.8.** The `Auto` data set. For a number of cars, `mpg` and `horsepower` are shown. The linear regression fit is shown in orange. The linear regression fit for a model that includes `horsepower`<sup>2</sup> is shown as a blue curve. The linear regression fit for a model that includes all polynomials of `horsepower` up to fifth-degree is shown in green.

	Coefficient	Std. error	<i>t</i> -statistic	<i>p</i> -value
Intercept	56.9001	1.8004	31.6	< 0.0001
<code>horsepower</code>	−0.4662	0.0311	−15.0	< 0.0001
<code>horsepower</code> <sup>2</sup>	0.0012	0.0001	10.1	< 0.0001

**TABLE 3.10.** For the `Auto` data set, least squares coefficient estimates associated with the regression of `mpg` onto `horsepower` and `horsepower`<sup>2</sup>.

fit appears to be substantially better than the fit obtained when just the linear term is included. The  $R^2$  of the quadratic fit is 0.688, compared to 0.606 for the linear fit, and the  $p$ -value in Table 3.10 for the quadratic term is highly significant.

If including `horsepower`<sup>2</sup> led to such a big improvement in the model, why not include `horsepower`<sup>3</sup>, `horsepower`<sup>4</sup>, or even `horsepower`<sup>5</sup>? The green curve in Figure 3.8 displays the fit that results from including all polynomials up to fifth degree in the model (3.36). The resulting fit seems unnecessarily wiggly—that is, it is unclear that including the additional terms really has led to a better fit to the data.

The approach that we have just described for extending the linear model to accommodate non-linear relationships is known as *polynomial regression*, since we have included polynomial functions of the predictors in the regression model. We further explore this approach and other non-linear extensions of the linear model in Chapter 7.

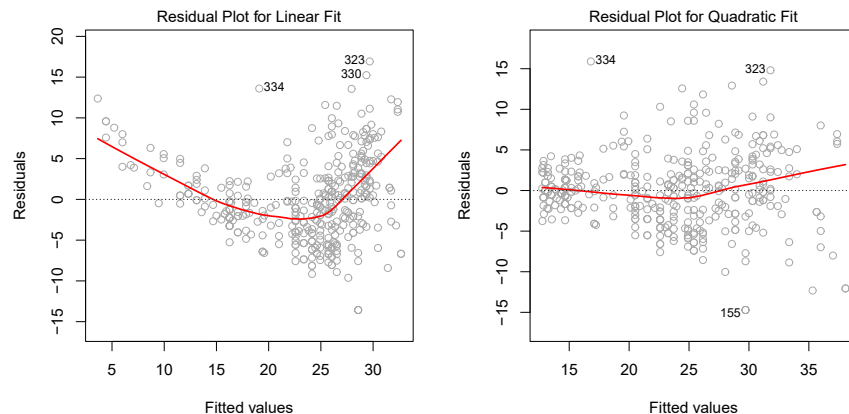
### 3.3.3 Potential Problems

When we fit a linear regression model to a particular data set, many problems may occur. Most common among these are the following:

1. *Non-linearity of the response-predictor relationships.*
2. *Correlation of error terms.*
3. *Non-constant variance of error terms.*
4. *Outliers.*
5. *High-leverage points.*
6. *Collinearity.*

In practice, identifying and overcoming these problems is as much an art as a science. Many pages in countless books have been written on this topic. Since the linear regression model is not our primary focus here, we will provide only a brief summary of some key points.

#### 1. Non-linearity of the Data



**FIGURE 3.9.** *Plots of residuals versus predicted (or fitted) values for the **Auto** data set. In each plot, the red line is a smooth fit to the residuals, intended to make it easier to identify a trend. Left: A linear regression of **mpg** on **horsepower**. A strong pattern in the residuals indicates non-linearity in the data. Right: A linear regression of **mpg** on **horsepower** and **horsepower**<sup>2</sup>. There is little pattern in the residuals.*

The linear regression model assumes that there is a straight-line relationship between the predictors and the response. If the true relationship is far from linear, then virtually all of the conclusions that we draw from the fit are suspect. In addition, the prediction accuracy of the model can be significantly reduced.

*Residual plots* are a useful graphical tool for identifying non-linearity. Given a simple linear regression model, we can plot the residuals,  $e_i =$

residual plot

$y_i - \hat{y}_i$ , versus the predictor  $x_i$ . In the case of a multiple regression model, since there are multiple predictors, we instead plot the residuals versus the predicted (or *fitted*) values  $\hat{y}_i$ . Ideally, the residual plot will show no discernible pattern. The presence of a pattern may indicate a problem with some aspect of the linear model.

The left panel of Figure 3.9 displays a residual plot from the linear regression of `mpg` onto `horsepower` on the `Auto` data set that was illustrated in Figure 3.8. The red line is a smooth fit to the residuals, which is displayed in order to make it easier to identify any trends. The residuals exhibit a clear U-shape, which provides a strong indication of non-linearity in the data. In contrast, the right-hand panel of Figure 3.9 displays the residual plot that results from the model (3.36), which contains a quadratic term. There appears to be little pattern in the residuals, suggesting that the quadratic term improves the fit to the data.

If the residual plot indicates that there are non-linear associations in the data, then a simple approach is to use non-linear transformations of the predictors, such as  $\log X$ ,  $\sqrt{X}$ , and  $X^2$ , in the regression model. In the later chapters of this book, we will discuss other more advanced non-linear approaches for addressing this issue.

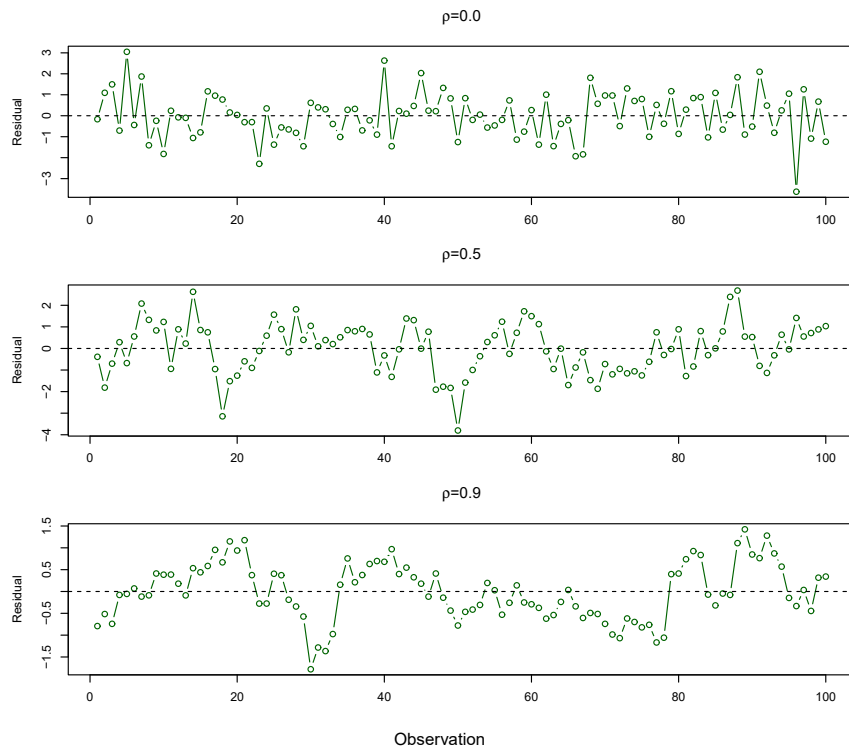
## 2. Correlation of Error Terms

An important assumption of the linear regression model is that the error terms,  $\epsilon_1, \epsilon_2, \dots, \epsilon_n$ , are uncorrelated. What does this mean? For instance, if the errors are uncorrelated, then the fact that  $\epsilon_i$  is positive provides little or no information about the sign of  $\epsilon_{i+1}$ . The standard errors that are computed for the estimated regression coefficients or the fitted values are based on the assumption of uncorrelated error terms. If in fact there is correlation among the error terms, then the estimated standard errors will tend to underestimate the true standard errors. As a result, confidence and prediction intervals will be narrower than they should be. For example, a 95% confidence interval may in reality have a much lower probability than 0.95 of containing the true value of the parameter. In addition,  $p$ -values associated with the model will be lower than they should be; this could cause us to erroneously conclude that a parameter is statistically significant. In short, if the error terms are correlated, we may have an unwarranted sense of confidence in our model.

As an extreme example, suppose we accidentally doubled our data, leading to observations and error terms identical in pairs. If we ignored this, our standard error calculations would be as if we had a sample of size  $2n$ , when in fact we have only  $n$  samples. Our estimated parameters would be the same for the  $2n$  samples as for the  $n$  samples, but the confidence intervals would be narrower by a factor of  $\sqrt{2}$ !

Why might correlations among the error terms occur? Such correlations frequently occur in the context of *time series* data, which consists of observations for which measurements are obtained at discrete points in time. In many cases, observations that are obtained at adjacent time points will have positively correlated errors. In order to determine if this is the case for a given data set, we can plot the residuals from our model as a function of time. If the errors are uncorrelated, then there should be no discernible pat-



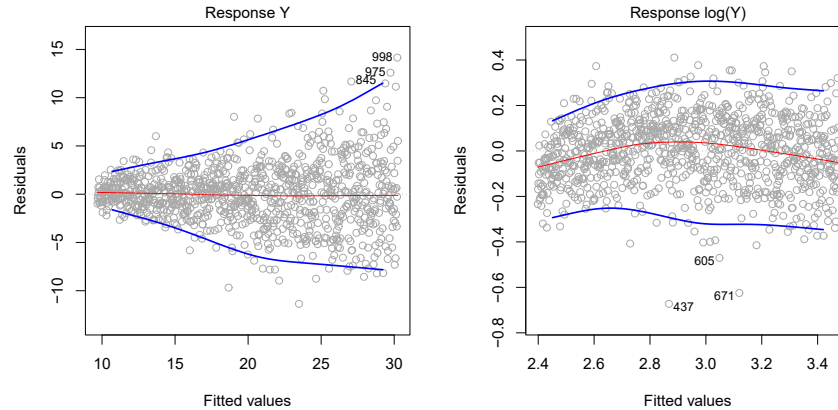


**FIGURE 3.10.** Plots of residuals from simulated time series data sets generated with differing levels of correlation  $\rho$  between error terms for adjacent time points.

tern. On the other hand, if the error terms are positively correlated, then we may see *tracking* in the residuals—that is, adjacent residuals may have similar values. Figure 3.10 provides an illustration. In the top panel, we see the residuals from a linear regression fit to data generated with uncorrelated errors. There is no evidence of a time-related trend in the residuals. In contrast, the residuals in the bottom panel are from a data set in which adjacent errors had a correlation of 0.9. Now there is a clear pattern in the residuals—adjacent residuals tend to take on similar values. Finally, the center panel illustrates a more moderate case in which the residuals had a correlation of 0.5. There is still evidence of tracking, but the pattern is less clear.

tracking

Many methods have been developed to properly take account of correlations in the error terms in time series data. Correlation among the error terms can also occur outside of time series data. For instance, consider a study in which individuals' heights are predicted from their weights. The assumption of uncorrelated errors could be violated if some of the individuals in the study are members of the same family, eat the same diet, or have been exposed to the same environmental factors. In general, the assumption of uncorrelated errors is extremely important for linear regression as well as for other statistical methods, and good experimental design is crucial in order to mitigate the risk of such correlations.



**FIGURE 3.11.** *Residual plots. In each plot, the red line is a smooth fit to the residuals, intended to make it easier to identify a trend. The blue lines track the outer quantiles of the residuals, and emphasize patterns. Left: The funnel shape indicates heteroscedasticity. Right: The response has been log transformed, and there is now no evidence of heteroscedasticity.*

### 3. Non-constant Variance of Error Terms

Another important assumption of the linear regression model is that the error terms have a constant variance,  $\text{Var}(\epsilon_i) = \sigma^2$ . The standard errors, confidence intervals, and hypothesis tests associated with the linear model rely upon this assumption.

Unfortunately, it is often the case that the variances of the error terms are non-constant. For instance, the variances of the error terms may increase with the value of the response. One can identify non-constant variances in the errors, or *heteroscedasticity*, from the presence of a *funnel shape* in the residual plot. An example is shown in the left-hand panel of Figure 3.11, in which the magnitude of the residuals tends to increase with the fitted values. When faced with this problem, one possible solution is to transform the response  $Y$  using a concave function such as  $\log Y$  or  $\sqrt{Y}$ . Such a transformation results in a greater amount of shrinkage of the larger responses, leading to a reduction in heteroscedasticity. The right-hand panel of Figure 3.11 displays the residual plot after transforming the response using  $\log Y$ . The residuals now appear to have constant variance, though there is some evidence of a slight non-linear relationship in the data.

hetero-  
scedasticity

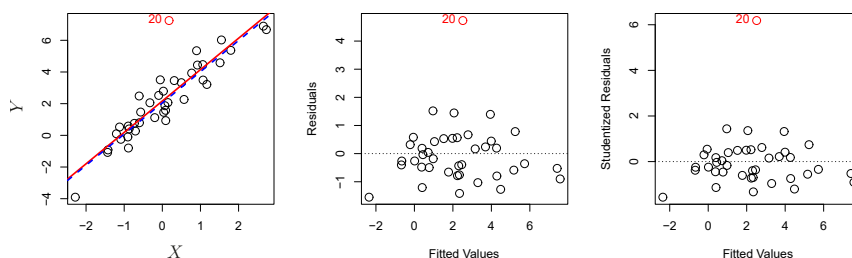
Sometimes we have a good idea of the variance of each response. For example, the  $i$ th response could be an average of  $n_i$  raw observations. If each of these raw observations is uncorrelated with variance  $\sigma^2$ , then their average has variance  $\sigma_i^2 = \sigma^2/n_i$ . In this case a simple remedy is to fit our model by *weighted least squares*, with weights proportional to the inverse variances—i.e.  $w_i = n_i$  in this case. Most linear regression software allows for observation weights.

weighted  
least squares

### 4. Outliers

An *outlier* is a point for which  $y_i$  is far from the value predicted by the

outlier



**FIGURE 3.12.** Left: The least squares regression line is shown in red, and the regression line after removing the outlier is shown in blue. Center: The residual plot clearly identifies the outlier. Right: The outlier has a studentized residual of 6; typically we expect values between  $-3$  and  $3$ .

model. Outliers can arise for a variety of reasons, such as incorrect recording of an observation during data collection.

The red point (observation 20) in the left-hand panel of Figure 3.12 illustrates a typical outlier. The red solid line is the least squares regression fit, while the blue dashed line is the least squares fit after removal of the outlier. In this case, removing the outlier has little effect on the least squares line: it leads to almost no change in the slope, and a miniscule reduction in the intercept. It is typical for an outlier that does not have an unusual predictor value to have little effect on the least squares fit. However, even if an outlier does not have much effect on the least squares fit, it can cause other problems. For instance, in this example, the RSE is 1.09 when the outlier is included in the regression, but it is only 0.77 when the outlier is removed. Since the RSE is used to compute all confidence intervals and  $p$ -values, such a dramatic increase caused by a single data point can have implications for the interpretation of the fit. Similarly, inclusion of the outlier causes the  $R^2$  to decline from 0.892 to 0.805.

Residual plots can be used to identify outliers. In this example, the outlier is clearly visible in the residual plot illustrated in the center panel of Figure 3.12. But in practice, it can be difficult to decide how large a residual needs to be before we consider the point to be an outlier. To address this problem, instead of plotting the residuals, we can plot the *studentized residuals*, computed by dividing each residual  $e_i$  by its estimated standard error. Observations whose studentized residuals are greater than 3 in absolute value are possible outliers. In the right-hand panel of Figure 3.12, the outlier's studentized residual exceeds 6, while all other observations have studentized residuals between  $-2$  and  $2$ .

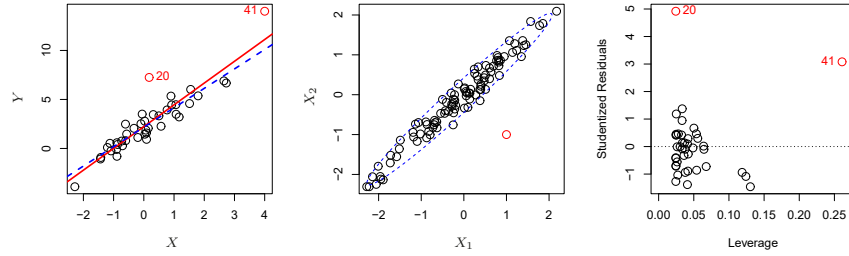
studentized  
residual

If we believe that an outlier has occurred due to an error in data collection or recording, then one solution is to simply remove the observation. However, care should be taken, since an outlier may instead indicate a deficiency with the model, such as a missing predictor.

### 5. High Leverage Points

We just saw that outliers are observations for which the response  $y_i$  is unusual given the predictor  $x_i$ . In contrast, observations with *high leverage* have an unusual value for  $x_i$ . For example, observation 41 in the left-hand

high  
leverage



**FIGURE 3.13.** Left: Observation 41 is a high leverage point, while 20 is not. The red line is the fit to all the data, and the blue line is the fit with observation 41 removed. Center: The red observation is not unusual in terms of its  $X_1$  value or its  $X_2$  value, but still falls outside the bulk of the data, and hence has high leverage. Right: Observation 41 has a high leverage and a high residual.

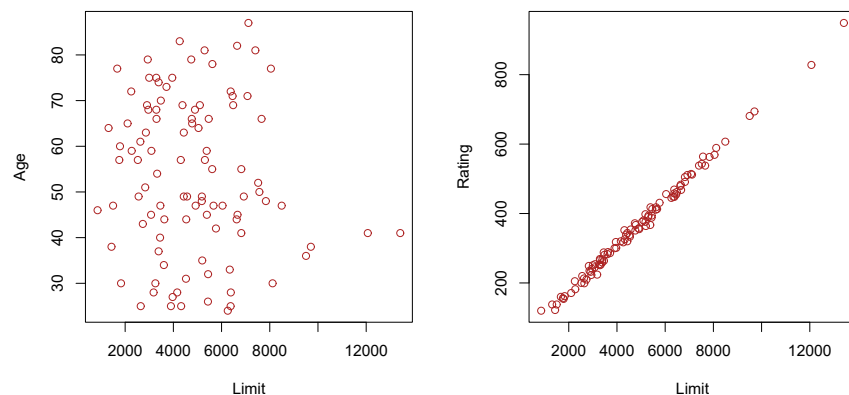
panel of Figure 3.13 has high leverage, in that the predictor value for this observation is large relative to the other observations. (Note that the data displayed in Figure 3.13 are the same as the data displayed in Figure 3.12, but with the addition of a single high leverage observation.) The red solid line is the least squares fit to the data, while the blue dashed line is the fit produced when observation 41 is removed. Comparing the left-hand panels of Figures 3.12 and 3.13, we observe that removing the high leverage observation has a much more substantial impact on the least squares line than removing the outlier. In fact, high leverage observations tend to have a sizable impact on the estimated regression line. It is cause for concern if the least squares line is heavily affected by just a couple of observations, because any problems with these points may invalidate the entire fit. For this reason, it is important to identify high leverage observations.

In a simple linear regression, high leverage observations are fairly easy to identify, since we can simply look for observations for which the predictor value is outside of the normal range of the observations. But in a multiple linear regression with many predictors, it is possible to have an observation that is well within the range of each individual predictor's values, but that is unusual in terms of the full set of predictors. An example is shown in the center panel of Figure 3.13, for a data set with two predictors,  $X_1$  and  $X_2$ . Most of the observations' predictor values fall within the blue dashed ellipse, but the red observation is well outside of this range. But neither its value for  $X_1$  nor its value for  $X_2$  is unusual. So if we examine just  $X_1$  or just  $X_2$ , we will fail to notice this high leverage point. This problem is more pronounced in multiple regression settings with more than two predictors, because then there is no simple way to plot all dimensions of the data simultaneously.

In order to quantify an observation's leverage, we compute the *leverage statistic*. A large value of this statistic indicates an observation with high leverage. For a simple linear regression,

leverage  
statistic

$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{i'=1}^n (x_{i'} - \bar{x})^2}. \quad (3.37)$$



**FIGURE 3.14.** Scatterplots of the observations from the **Credit** data set. Left: A plot of **age** versus **limit**. These two variables are not collinear. Right: A plot of **rating** versus **limit**. There is high collinearity.

It is clear from this equation that  $h_i$  increases with the distance of  $x_i$  from  $\bar{x}$ . There is a simple extension of  $h_i$  to the case of multiple predictors, though we do not provide the formula here. The leverage statistic  $h_i$  is always between  $1/n$  and 1, and the average leverage for all the observations is always equal to  $(p+1)/n$ . So if a given observation has a leverage statistic that greatly exceeds  $(p+1)/n$ , then we may suspect that the corresponding point has high leverage.

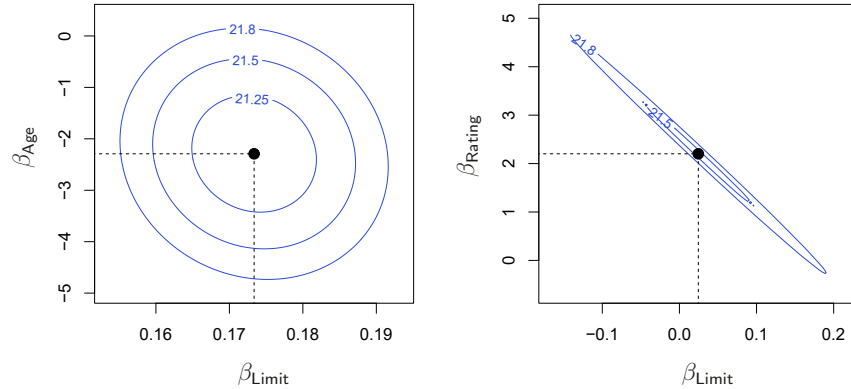
The right-hand panel of Figure 3.13 provides a plot of the studentized residuals versus  $h_i$  for the data in the left-hand panel of Figure 3.13. Observation 41 stands out as having a very high leverage statistic as well as a high studentized residual. In other words, it is an outlier as well as a high leverage observation. This is a particularly dangerous combination! This plot also reveals the reason that observation 20 had relatively little effect on the least squares fit in Figure 3.12: it has low leverage.

## 6. Collinearity

*Collinearity* refers to the situation in which two or more predictor variables are closely related to one another. The concept of collinearity is illustrated in Figure 3.14 using the **Credit** data set. In the left-hand panel of Figure 3.14, the two predictors **limit** and **age** appear to have no obvious relationship. In contrast, in the right-hand panel of Figure 3.14, the predictors **limit** and **rating** are very highly correlated with each other, and we say that they are *collinear*. The presence of collinearity can pose problems in the regression context, since it can be difficult to separate out the individual effects of collinear variables on the response. In other words, since **limit** and **rating** tend to increase or decrease together, it can be difficult to determine how each one separately is associated with the response, **balance**.

Figure 3.15 illustrates some of the difficulties that can result from collinearity. The left-hand panel of Figure 3.15 is a contour plot of the RSS (3.22) associated with different possible coefficient estimates for the regression of **balance** on **limit** and **age**. Each ellipse represents a set of coefficients

collinearity



**FIGURE 3.15.** Contour plots for the RSS values as a function of the parameters  $\beta$  for various regressions involving the **Credit** data set. In each plot, the black dots represent the coefficient values corresponding to the minimum RSS. Left: A contour plot of RSS for the regression of **balance** onto **age** and **limit**. The minimum value is well defined. Right: A contour plot of RSS for the regression of **balance** onto **rating** and **limit**. Because of the collinearity, there are many pairs  $(\beta_{\text{Limit}}, \beta_{\text{Rating}})$  with a similar value for RSS.

that correspond to the same RSS, with ellipses nearest to the center taking on the lowest values of RSS. The black dots and associated dashed lines represent the coefficient estimates that result in the smallest possible RSS—in other words, these are the least squares estimates. The axes for **limit** and **age** have been scaled so that the plot includes possible coefficient estimates that are up to four standard errors on either side of the least squares estimates. Thus the plot includes all plausible values for the coefficients. For example, we see that the true **limit** coefficient is almost certainly somewhere between 0.15 and 0.20.

In contrast, the right-hand panel of Figure 3.15 displays contour plots of the RSS associated with possible coefficient estimates for the regression of **balance** onto **limit** and **rating**, which we know to be highly collinear. Now the contours run along a narrow valley; there is a broad range of values for the coefficient estimates that result in equal values for RSS. Hence a small change in the data could cause the pair of coefficient values that yield the smallest RSS—that is, the least squares estimates—to move anywhere along this valley. This results in a great deal of uncertainty in the coefficient estimates. Notice that the scale for the **limit** coefficient now runs from roughly  $-0.2$  to  $0.2$ ; this is an eight-fold increase over the plausible range of the **limit** coefficient in the regression with **age**. Interestingly, even though the **limit** and **rating** coefficients now have much more individual uncertainty, they will almost certainly lie somewhere in this contour valley. For example, we would not expect the true value of the **limit** and **rating** coefficients to be  $-0.1$  and  $1$  respectively, even though such a value is plausible for each coefficient individually.

Since collinearity reduces the accuracy of the estimates of the regression coefficients, it causes the standard error for  $\hat{\beta}_j$  to grow. Recall that the  $t$ -statistic for each predictor is calculated by dividing  $\hat{\beta}_j$  by its standard

		Coefficient	Std. error	<i>t</i> -statistic	<i>p</i> -value
Model 1	Intercept	−173.411	43.828	−3.957	< 0.0001
	age	−2.292	0.672	−3.407	0.0007
	limit	0.173	0.005	34.496	< 0.0001
Model 2	Intercept	−377.537	45.254	−8.343	< 0.0001
	rating	2.202	0.952	2.312	0.0213
	limit	0.025	0.064	0.384	0.7012

**TABLE 3.11.** The results for two multiple regression models involving the **Credit** data set are shown. Model 1 is a regression of **balance** on **age** and **limit**, and Model 2 a regression of **balance** on **rating** and **limit**. The standard error of  $\hat{\beta}_{\text{limit}}$  increases 12-fold in the second regression, due to collinearity.

error. Consequently, collinearity results in a decline in the *t*-statistic. As a result, in the presence of collinearity, we may fail to reject  $H_0 : \beta_j = 0$ . This means that the *power* of the hypothesis test—the probability of correctly detecting a *non-zero* coefficient—is reduced by collinearity. power

Table 3.11 compares the coefficient estimates obtained from two separate multiple regression models. The first is a regression of **balance** on **age** and **limit**, and the second is a regression of **balance** on **rating** and **limit**. In the first regression, both **age** and **limit** are highly significant with very small *p*-values. In the second, the collinearity between **limit** and **rating** has caused the standard error for the **limit** coefficient estimate to increase by a factor of 12 and the *p*-value to increase to 0.701. In other words, the importance of the **limit** variable has been masked due to the presence of collinearity. To avoid such a situation, it is desirable to identify and address potential collinearity problems while fitting the model.

A simple way to detect collinearity is to look at the correlation matrix of the predictors. An element of this matrix that is large in absolute value indicates a pair of highly correlated variables, and therefore a collinearity problem in the data. Unfortunately, not all collinearity problems can be detected by inspection of the correlation matrix: it is possible for collinearity to exist between three or more variables even if no pair of variables has a particularly high correlation. We call this situation *multicollinearity*. Instead of inspecting the correlation matrix, a better way to assess multicollinearity is to compute the *variance inflation factor* (VIF). The VIF is the ratio of the variance of  $\hat{\beta}_j$  when fitting the full model divided by the variance of  $\hat{\beta}_j$  if fit on its own. The smallest possible value for VIF is 1, which indicates the complete absence of collinearity. Typically in practice there is a small amount of collinearity among the predictors. As a rule of thumb, a VIF value that exceeds 5 or 10 indicates a problematic amount of collinearity. The VIF for each variable can be computed using the formula multi-collinearity  
variance  
inflation  
factor

$$\text{VIF}(\hat{\beta}_j) = \frac{1}{1 - R_{X_j|X_{-j}}^2},$$

where  $R_{X_j|X_{-j}}^2$  is the  $R^2$  from a regression of  $X_j$  onto all of the other predictors. If  $R_{X_j|X_{-j}}^2$  is close to one, then collinearity is present, and so the VIF will be large.

In the `Credit` data, a regression of `balance` on `age`, `rating`, and `limit` indicates that the predictors have VIF values of 1.01, 160.67, and 160.59. As we suspected, there is considerable collinearity in the data!

When faced with the problem of collinearity, there are two simple solutions. The first is to drop one of the problematic variables from the regression. This can usually be done without much compromise to the regression fit, since the presence of collinearity implies that the information that this variable provides about the response is redundant in the presence of the other variables. For instance, if we regress `balance` onto `age` and `limit`, without the `rating` predictor, then the resulting VIF values are close to the minimum possible value of 1, and the  $R^2$  drops from 0.754 to 0.75. So dropping `rating` from the set of predictors has effectively solved the collinearity problem without compromising the fit. The second solution is to combine the collinear variables together into a single predictor. For instance, we might take the average of standardized versions of `limit` and `rating` in order to create a new variable that measures *credit worthiness*.

## 3.4 The Marketing Plan

We now briefly return to the seven questions about the `Advertising` data that we set out to answer at the beginning of this chapter.

1. *Is there a relationship between sales and advertising budget?*

This question can be answered by fitting a multiple regression model of `sales` onto `TV`, `radio`, and `newspaper`, as in (3.20), and testing the hypothesis  $H_0 : \beta_{\text{TV}} = \beta_{\text{radio}} = \beta_{\text{newspaper}} = 0$ . In Section 3.2.2, we showed that the  $F$ -statistic can be used to determine whether or not we should reject this null hypothesis. In this case the  $p$ -value corresponding to the  $F$ -statistic in Table 3.6 is very low, indicating clear evidence of a relationship between advertising and sales.

2. *How strong is the relationship?*

We discussed two measures of model accuracy in Section 3.1.3. First, the RSE estimates the standard deviation of the response from the population regression line. For the `Advertising` data, the RSE is 1.69 units while the mean value for the response is 14.022, indicating a percentage error of roughly 12%. Second, the  $R^2$  statistic records the percentage of variability in the response that is explained by the predictors. The predictors explain almost 90% of the variance in `sales`. The RSE and  $R^2$  statistics are displayed in Table 3.6.

3. *Which media are associated with sales?*

To answer this question, we can examine the  $p$ -values associated with each predictor's  $t$ -statistic (Section 3.1.2). In the multiple linear regression displayed in Table 3.4, the  $p$ -values for `TV` and `radio` are low, but the  $p$ -value for `newspaper` is not. This suggests that only `TV` and `radio` are related to `sales`. In Chapter 6 we explore this question in greater detail.



4. *How large is the association between each medium and sales?*

We saw in Section 3.1.2 that the standard error of  $\hat{\beta}_j$  can be used to construct confidence intervals for  $\beta_j$ . For the **Advertising** data, we can use the results in Table 3.4 to compute the 95 % confidence intervals for the coefficients in a multiple regression model using all three media budgets as predictors. The confidence intervals are as follows: (0.043, 0.049) for **TV**, (0.172, 0.206) for **radio**, and (−0.013, 0.011) for **newspaper**. The confidence intervals for **TV** and **radio** are narrow and far from zero, providing evidence that these media are related to **sales**. But the interval for **newspaper** includes zero, indicating that the variable is not statistically significant given the values of **TV** and **radio**.

We saw in Section 3.3.3 that collinearity can result in very wide standard errors. Could collinearity be the reason that the confidence interval associated with **newspaper** is so wide? The VIF scores are 1.005, 1.145, and 1.145 for **TV**, **radio**, and **newspaper**, suggesting no evidence of collinearity.

In order to assess the association of each medium individually on sales, we can perform three separate simple linear regressions. Results are shown in Tables 3.1 and 3.3. There is evidence of an extremely strong association between **TV** and **sales** and between **radio** and **sales**. There is evidence of a mild association between **newspaper** and **sales**, when the values of **TV** and **radio** are ignored.

5. *How accurately can we predict future sales?*

The response can be predicted using (3.21). The accuracy associated with this estimate depends on whether we wish to predict an individual response,  $Y = f(X) + \epsilon$ , or the average response,  $f(X)$  (Section 3.2.2). If the former, we use a prediction interval, and if the latter, we use a confidence interval. Prediction intervals will always be wider than confidence intervals because they account for the uncertainty associated with  $\epsilon$ , the irreducible error.

6. *Is the relationship linear?*

In Section 3.3.3, we saw that residual plots can be used in order to identify non-linearity. If the relationships are linear, then the residual plots should display no pattern. In the case of the **Advertising** data, we observe a non-linear effect in Figure 3.5, though this effect could also be observed in a residual plot. In Section 3.3.2, we discussed the inclusion of transformations of the predictors in the linear regression model in order to accommodate non-linear relationships.

7. *Is there synergy among the advertising media?*

The standard linear regression model assumes an additive relationship between the predictors and the response. An additive model is easy to interpret because the association between each predictor and the response is unrelated to the values of the other predictors. However, the additive assumption may be unrealistic for certain data sets. In Section 3.3.2, we showed how to include an interaction term

in the regression model in order to accommodate non-additive relationships. A small  $p$ -value associated with the interaction term indicates the presence of such relationships. Figure 3.5 suggested that the **Advertising** data may not be additive. Including an interaction term in the model results in a substantial increase in  $R^2$ , from around 90 % to almost 97 %.

### 3.5 Comparison of Linear Regression with $K$ -Nearest Neighbors

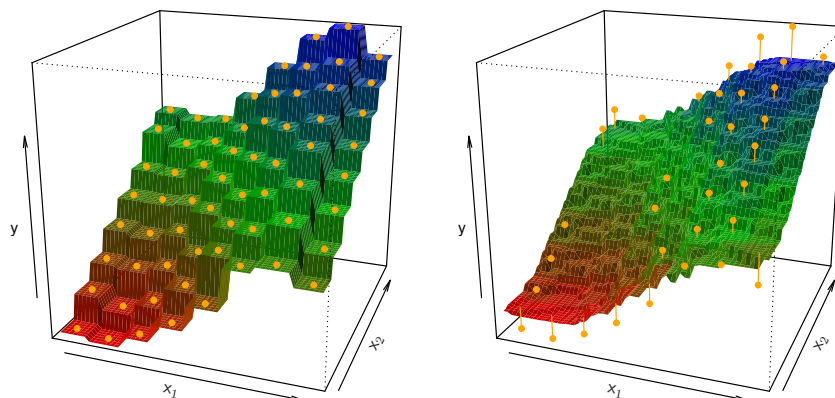
As discussed in Chapter 2, linear regression is an example of a *parametric* approach because it assumes a linear functional form for  $f(X)$ . Parametric methods have several advantages. They are often easy to fit, because one need estimate only a small number of coefficients. In the case of linear regression, the coefficients have simple interpretations, and tests of statistical significance can be easily performed. But parametric methods do have a disadvantage: by construction, they make strong assumptions about the form of  $f(X)$ . If the specified functional form is far from the truth, and prediction accuracy is our goal, then the parametric method will perform poorly. For instance, if we assume a linear relationship between  $X$  and  $Y$  but the true relationship is far from linear, then the resulting model will provide a poor fit to the data, and any conclusions drawn from it will be suspect.

In contrast, *non-parametric* methods do not explicitly assume a parametric form for  $f(X)$ , and thereby provide an alternative and more flexible approach for performing regression. We discuss various non-parametric methods in this book. Here we consider one of the simplest and best-known non-parametric methods,  *$K$ -nearest neighbors regression* (KNN regression). The KNN regression method is closely related to the KNN classifier discussed in Chapter 2. Given a value for  $K$  and a prediction point  $x_0$ , KNN regression first identifies the  $K$  training observations that are closest to  $x_0$ , represented by  $\mathcal{N}_0$ . It then estimates  $f(x_0)$  using the average of all the training responses in  $\mathcal{N}_0$ . In other words,

*$K$ -nearest  
neighbors  
regression*

$$\hat{f}(x_0) = \frac{1}{K} \sum_{x_i \in \mathcal{N}_0} y_i.$$

Figure 3.16 illustrates two KNN fits on a data set with  $p = 2$  predictors. The fit with  $K = 1$  is shown in the left-hand panel, while the right-hand panel corresponds to  $K = 9$ . We see that when  $K = 1$ , the KNN fit perfectly interpolates the training observations, and consequently takes the form of a step function. When  $K = 9$ , the KNN fit still is a step function, but averaging over nine observations results in much smaller regions of constant prediction, and consequently a smoother fit. In general, the optimal value for  $K$  will depend on the *bias-variance tradeoff*, which we introduced in Chapter 2. A small value for  $K$  provides the most flexible fit, which will have low bias but high variance. This variance is due to the fact that the prediction in a given region is entirely dependent on just one observation.

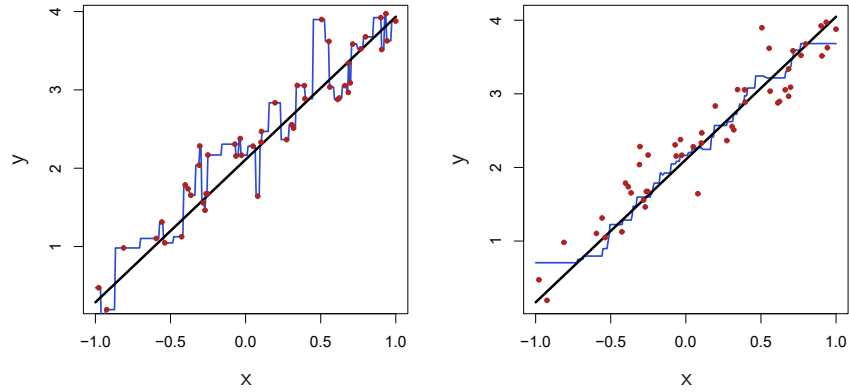


**FIGURE 3.16.** Plots of  $\hat{f}(X)$  using KNN regression on a two-dimensional data set with 64 observations (orange dots). Left:  $K = 1$  results in a rough step function fit. Right:  $K = 9$  produces a much smoother fit.

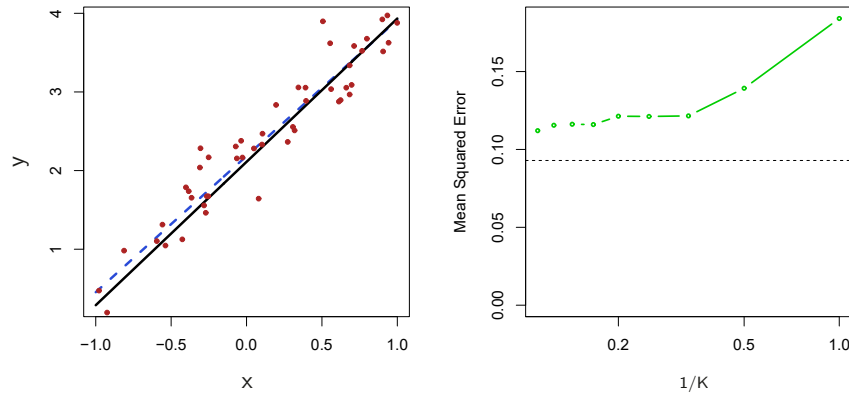
In contrast, larger values of  $K$  provide a smoother and less variable fit; the prediction in a region is an average of several points, and so changing one observation has a smaller effect. However, the smoothing may cause bias by masking some of the structure in  $f(X)$ . In Chapter 5, we introduce several approaches for estimating test error rates. These methods can be used to identify the optimal value of  $K$  in KNN regression.

In what setting will a parametric approach such as least squares linear regression outperform a non-parametric approach such as KNN regression? The answer is simple: *the parametric approach will outperform the non-parametric approach if the parametric form that has been selected is close to the true form of  $f$* . Figure 3.17 provides an example with data generated from a one-dimensional linear regression model. The black solid lines represent  $f(X)$ , while the blue curves correspond to the KNN fits using  $K = 1$  and  $K = 9$ . In this case, the  $K = 1$  predictions are far too variable, while the smoother  $K = 9$  fit is much closer to  $f(X)$ . However, since the true relationship is linear, it is hard for a non-parametric approach to compete with linear regression: a non-parametric approach incurs a cost in variance that is not offset by a reduction in bias. The blue dashed line in the left-hand panel of Figure 3.18 represents the linear regression fit to the same data. It is almost perfect. The right-hand panel of Figure 3.18 reveals that linear regression outperforms KNN for this data. The green solid line, plotted as a function of  $1/K$ , represents the test set mean squared error (MSE) for KNN. The KNN errors are well above the black dashed line, which is the test MSE for linear regression. When the value of  $K$  is large, then KNN performs only a little worse than least squares regression in terms of MSE. It performs far worse when  $K$  is small.

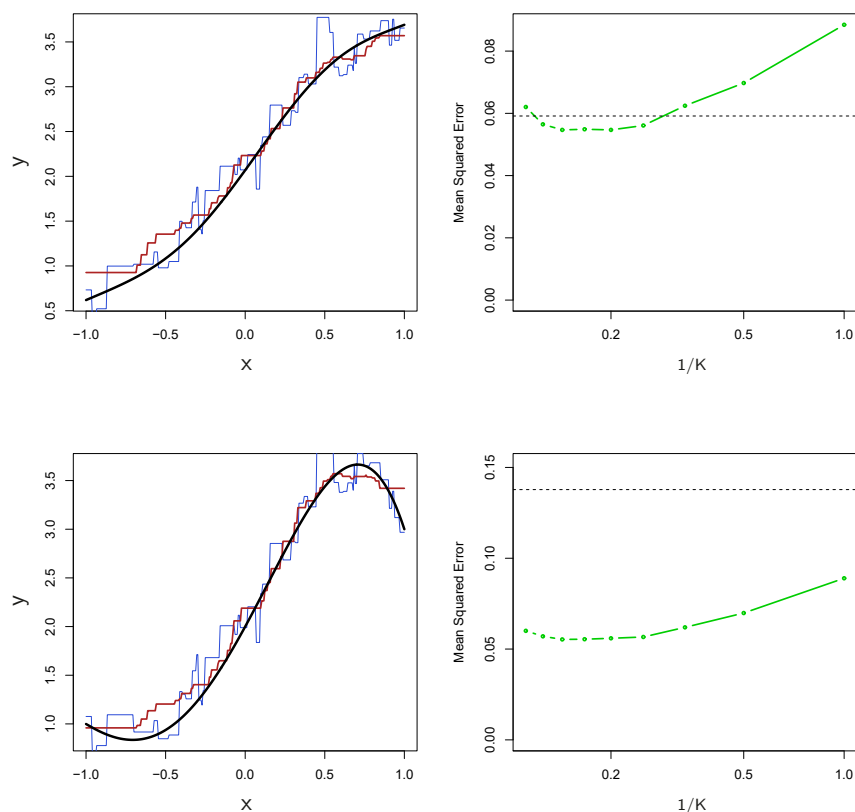
In practice, the true relationship between  $X$  and  $Y$  is rarely exactly linear. Figure 3.19 examines the relative performances of least squares regression and KNN under increasing levels of non-linearity in the relationship between  $X$  and  $Y$ . In the top row, the true relationship is nearly linear. In this case we see that the test MSE for linear regression is still superior



**FIGURE 3.17.** Plots of  $\hat{f}(X)$  using KNN regression on a one-dimensional data set with 50 observations. The true relationship is given by the black solid line. Left: The blue curve corresponds to  $K = 1$  and interpolates (i.e. passes directly through) the training data. Right: The blue curve corresponds to  $K = 9$ , and represents a smoother fit.



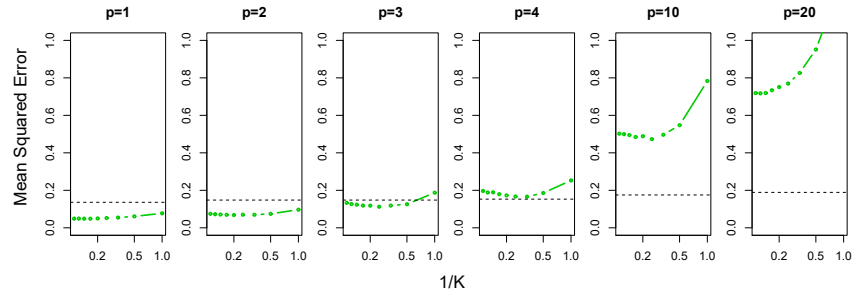
**FIGURE 3.18.** The same data set shown in Figure 3.17 is investigated further. Left: The blue dashed line is the least squares fit to the data. Since  $f(X)$  is in fact linear (displayed as the black line), the least squares regression line provides a very good estimate of  $f(X)$ . Right: The dashed horizontal line represents the least squares test set MSE, while the green solid line corresponds to the MSE for KNN as a function of  $1/K$  (on the log scale). Linear regression achieves a lower test MSE than does KNN regression, since  $f(X)$  is in fact linear. For KNN regression, the best results occur with a very large value of  $K$ , corresponding to a small value of  $1/K$ .



**FIGURE 3.19.** Top Left: In a setting with a slightly non-linear relationship between  $X$  and  $Y$  (solid black line), the KNN fits with  $K = 1$  (blue) and  $K = 9$  (red) are displayed. Top Right: For the slightly non-linear data, the test set MSE for least squares regression (horizontal black) and KNN with various values of  $1/K$  (green) are displayed. Bottom Left and Bottom Right: As in the top panel, but with a strongly non-linear relationship between  $X$  and  $Y$ .

to that of KNN for low values of  $K$ . However, for  $K \geq 4$ , KNN outperforms linear regression. The second row illustrates a more substantial deviation from linearity. In this situation, KNN substantially outperforms linear regression for all values of  $K$ . Note that as the extent of non-linearity increases, there is little change in the test set MSE for the non-parametric KNN method, but there is a large increase in the test set MSE of linear regression.

Figures 3.18 and 3.19 display situations in which KNN performs slightly worse than linear regression when the relationship is linear, but much better than linear regression for nonlinear situations. In a real life situation in which the true relationship is unknown, one might suspect that KNN should be favored over linear regression because it will at worst be slightly inferior to linear regression if the true relationship is linear, and may give substantially better results if the true relationship is non-linear. But in reality, even when the true relationship is highly non-linear, KNN may still provide inferior results to linear regression. In particular, both Figures 3.18



**FIGURE 3.20.** Test MSE for linear regression (black dashed lines) and KNN (green curves) as the number of variables  $p$  increases. The true function is non-linear in the first variable, as in the lower panel in Figure 3.19, and does not depend on the additional variables. The performance of linear regression deteriorates slowly in the presence of these additional noise variables, whereas KNN's performance degrades much more quickly as  $p$  increases.

and 3.19 illustrate settings with  $p = 1$  predictor. But in higher dimensions, KNN often performs worse than linear regression.

Figure 3.20 considers the same strongly non-linear situation as in the second row of Figure 3.19, except that we have added additional *noise* predictors that are not associated with the response. When  $p = 1$  or  $p = 2$ , KNN outperforms linear regression. But for  $p = 3$  the results are mixed, and for  $p \geq 4$  linear regression is superior to KNN. In fact, the increase in dimension has only caused a small deterioration in the linear regression test set MSE, but it has caused more than a ten-fold increase in the MSE for KNN. This decrease in performance as the dimension increases is a common problem for KNN, and results from the fact that in higher dimensions there is effectively a reduction in sample size. In this data set there are 50 training observations; when  $p = 1$ , this provides enough information to accurately estimate  $f(X)$ . However, spreading 50 observations over  $p = 20$  dimensions results in a phenomenon in which a given observation has no *nearby neighbors*—this is the so-called *curse of dimensionality*. That is, the  $K$  observations that are nearest to a given test observation  $x_0$  may be very far away from  $x_0$  in  $p$ -dimensional space when  $p$  is large, leading to a very poor prediction of  $f(x_0)$  and hence a poor KNN fit. As a general rule, parametric methods will tend to outperform non-parametric approaches when there is a small number of observations per predictor.

curse of dimensionality

Even when the dimension is small, we might prefer linear regression to KNN from an interpretability standpoint. If the test MSE of KNN is only slightly lower than that of linear regression, we might be willing to forego a little bit of prediction accuracy for the sake of a simple model that can be described in terms of just a few coefficients, and for which  $p$ -values are available.