

## Predictive Modeling for Regression

which the response variable does have numerical significance. Examples are the amount a retail store might earn from a given customer over a ten-year period, the rate of fuel consumption of a given type of car under normal conditions, the number of people who might access a particular Web site in a given month, and so on. The variables to be used as input for prediction will be called *predictor variables* and the variable to be predicted is the *response variable*. Other authors sometimes use the terms *dependent* or *target* for the response variable, and *independent*, *explanatory*, or *regressor* for the predictor variables. Other names used in the classification context were mentioned in [chapter 10](#). Note that the predictor variables can be numerical, but they need not be. Our aim, then, is to use a sample of objects, for which both the response variable and the predictor variables are known, to construct a model that will allow prediction of the numerical value of the response variable for a new case for which only the predictor variables are known. This is essentially the same problem as in [chapter 10](#), the only difference being the numerical instead of nominal nature of the response variable. In fact, as we will see later in this chapter, we can also treat prediction of nominal variables (that is, classification) within this general framework of regression.

Accuracy of prediction is one of the most important properties of such models, so various measures of accuracy have been devised. These measures may also be used for choosing between alternative models, and for choosing the values of parameters in models. In the terminology introduced earlier, these measures are *score functions*, by which different models may be compared.

Predictive accuracy is a critical aspect of models, but it is not the only aspect. For example, we might use the model to shed insight into which of the predictor variables are most important. We might even insist that some variables be included in the model, because we know they should be there on substantive grounds, even though they lead to only small predictive improvement. Contrariwise, we might omit variables that we feel would enhance our predictive performance. (An example of this situation arises in credit scoring, in which, in many countries, it is illegal to include sex or race as a predictor variable.) We might be interested in whether predictor variables interact, in the sense that the effect that one has on the response variable depends on the values taken by others. For obvious reasons, we might be interested in whether good prediction can be achieved by a simple model. Sometimes we might even be willing to sacrifice some predictive accuracy in exchange for substantially reduced model complexity. Though predictive accuracy is perhaps the most important component of the performance of a predictive model, this has to be tempered by the context in which the model is to be applied.

### 11.2 Linear Models and Least Squares Fitting

[Chapter 6](#) introduced the idea of linear models, so called because they are linear in the parameters. The simplest such model yields predicted values,  $\hat{y}$ , of the response variable  $y$ , that are also a linear combination of the predictor variables  $x_j$ :

$$(11.1) \quad \hat{y} = a_0 + \sum_{j=1}^p a_j x_j$$

In fact, of course, we will not normally be able to predict the response variable perfectly (life is seldom so simple) and a common aim is to predict the mean value that  $y$  takes at each vector of the predictor variables—so  $\hat{y}$  is our predicted estimate of the mean value at  $\mathbf{x} = (x_1, \dots, x_p)$ . Models of this form are known as *linear regression models*. In the simplest case of a single predictor variable (*simple regression*), we have a *regression line* in the space spanned by the response and predictor variables. More generally (*multiple regression*) we have a *regression plane*. Such models are the oldest, most important, and single most widely used form of predictive model. One reason for this is their evident simplicity; a simple weighted sum is very easy both to compute and to understand. Another compelling reason is that they often perform very well—even in circumstances in which we know enough to be confident that the true relationship between the predictor and response variables cannot be linear. This is not altogether surprising: when we expand continuous mathematical functions in a Taylor series we

often find that the lowest order terms—the linear terms—are the most important, so that the best simple approximation is obtained by using a linear model.

It is extremely rare that the chosen model is exactly right. This is especially true in data mining situations, where our model is generally empirical rather than being based on an underlying theory (see [chapter 9](#)). The model may not include all of the predictor variables that are needed for perfect prediction (many may not have been measured or even be measurable); it may not include certain functions of the predictor variables

(maybe  $x_1^2$  is needed as well as  $x_1$ , or maybe products of the predictor variables are needed because they interact in their effect on  $y$ ); and, in any case, no measurement is perfect; the  $y$  variable will have errors associated with it so that each vector  $(x_1, \dots, x_p)$  will be associated with a distribution of possible  $y$  values, as we have noted above.

All of this means that the actual  $y$  values in a sample will differ from the predicted values. The differences between observed and predicted values are called *residuals*, and we denote these by  $e$ :

$$(11.2) \quad y(i) = \hat{y}(i) + e(i) = a_0 + \sum_{j=1}^p a_j x_j(i) + e(i), \quad 1 \leq i \leq n.$$

In matrix terms, if we denote the observed  $y$  measurements on the  $n$  objects in the training sample by the vector  $\mathbf{y}$  and the  $p$  measurements of the predictor variables on the  $n$  objects by the  $n$  by  $p+1$  matrix  $\mathbf{X}$  (an additional column of 1s are added to incorporate the intercept term  $a_0$  in the model), we can express the relationship between the observed response and predictor measurements, in terms of our model, as

$$(11.3) \quad \mathbf{y} = \mathbf{X}\mathbf{a} + \mathbf{e},$$

where  $\mathbf{y}$  is an  $n \times 1$  matrix of response values,  $\mathbf{a} = (a_0, \dots, a_p)$  represents the  $(p+1) \times 1$  vector of parameter values, and the  $n \times 1$  vector  $\mathbf{e} = (e(1), \dots, e(n))$  contains the residuals. Clearly we want to choose the parameters in our model (the values in the  $p+1$  vector  $\mathbf{a}$ ) so as to yield predictions that are as accurate as possible. Put another way, we must find estimates for the  $a_j$  that minimize the  $e$  discrepancies in some way. To do this, we combine the elements of  $\mathbf{e}$  in such a way as to yield a single numerical measure that we can minimize. Various ways of combining the  $e(i)$  have been proposed, but by far the most popular method is to sum their squares—that is, the sum of squared errors score function. Thus we seek the values for the parameter vector  $\mathbf{a}$  that minimizes

$$(11.4) \quad \sum_{i=1}^n e(i)^2 = \sum_{i=1}^n \left( y(i) - \sum_{j=0}^p a_j x_j(i) \right)^2$$

In this expression,  $y(i)$  is the observed  $y$  value for the  $i$ th training sample point and

$$(x_0(i), x_1(i), \dots, x_p(i)) = (1, x_1(i), \dots, x_p(i))$$

is the vector of predictor variables for this point. For obvious reasons, this method is known as the *least squares method*. For simplicity, we will denote the parameter vector that minimizes this by  $(a_0, \dots, a_p)$ . (It would be more correct, of course, if we used some notation to indicate that it is an estimate, such as  $(\hat{a}_0, \dots, \hat{a}_p)$ , but our notation has the merit of simplicity.) In matrix terms, the values of the parameters that minimize [equation 11.4](#) can be shown to be

$$(11.5) \quad \mathbf{a} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

In linear regression in general, the  $a$  parameters are called *regression coefficients*. Once the parameters have been estimated, they are used in [equation 11.1](#) to yield predictions. The predicted value of  $y$ ,  $y_k$ , for a vector of predictor variables  $\mathbf{x}_k$ , is given by

$$\hat{y}_k = \mathbf{x}_k^T \mathbf{a} = \mathbf{a}^T \mathbf{x}_k.$$

### 11.2.1 Computational Issues in Fitting the Model

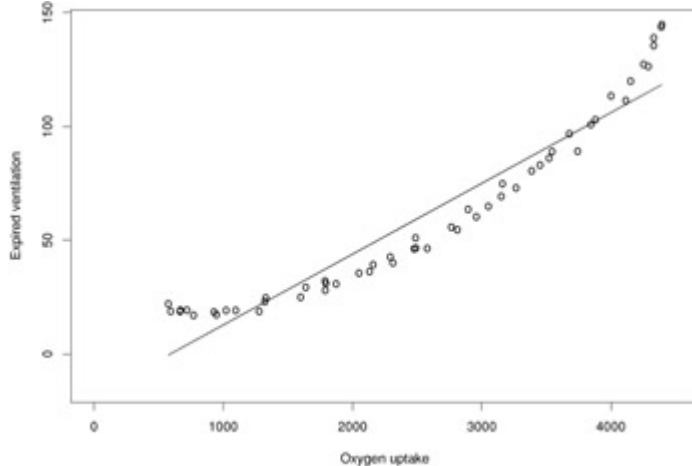
Solving [equation 11.5](#) directly requires that the matrix  $\mathbf{X}^T \mathbf{X}$  be invertible. Problems will arise if the sample size  $n$  is small (rare in data mining situations) or if there are linear dependencies between the measured values of the predictor variables (not so rare). In the latter case, modern software packages normally issue warnings, and appropriate action can be taken, such as dropping some of the predictor variables.

A rather more subtle problem arises when the measured values of the predictor variables are not exactly linearly dependent, but are almost so. Now the matrix can be inverted, but the solution will be unstable. This means that slight alterations to the observed  $\mathbf{X}$  values would lead to substantial differences in the estimated values of  $\mathbf{a}$ . Different

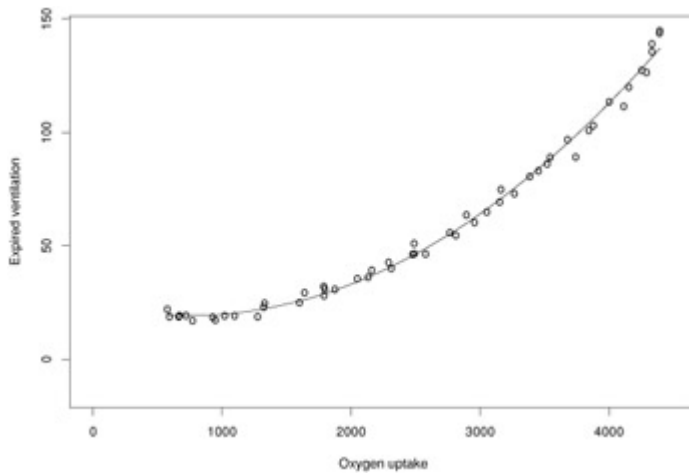
measurement errors or a slightly different training sample would have led to different parameter estimates. This problem is termed *multicollinearity*. The instability in the estimated parameters is a problem if these values are the focus of interest—for example, if we want to know which of the variables is most important in the model. However, it will not normally be a problem as far as predictive accuracy is concerned: although substantially different  $\mathbf{a}$  vectors may be produced by slight variations of the data, all of these vectors will lead to similar predictions for *most*  $\mathbf{x}_k$  vectors.

Solving [equation 11.5](#) is usually carried out by numerical linear algebra techniques for equation solving (such as the LU decomposition or the singular value decomposition (SVD)), which tend to have better numerical stability than that achieved by inverting the matrix  $\mathbf{X}^T \mathbf{X}$  directly. The underlying computational complexity is typically the same no matter which particular technique is used, namely,  $O(p^2 n + p^3)$ . The  $p^2 n$  term comes from the  $n$  multiplications required to calculate each element in the  $p \times p$  matrix  $\mathbf{C} = \mathbf{X}^T \mathbf{X}$ . The  $p^3$  term comes from then solving  $\mathbf{C}\mathbf{a} = \mathbf{X}^T \mathbf{y}$  for  $\mathbf{a}$ .

In [chapter 6](#) we remarked that the additive nature of the regression model could be retained while permitting more flexible model forms by including transformations of the raw  $x_i$  as well as the raw variables themselves. [Figure 11.1](#) shows a plot of data collected in an experiment in which a subject performed a physical task at a gradually increasing level of difficulty. The vertical axis shows a measure on the gases expired from the lungs while the horizontal axis shows the oxygen uptake. The nonlinearity of the relationship between these two variables is quite clear from the plot. A straight line  $y = a_0 + a_1 x$  provides a poor fit—as is shown in the figure. The predicted values from this model would be accurate only for  $x$  (oxygen uptake) values just above 1000 and just below 4000. (Despite this, the model is not grossly inaccurate—the point made earlier about models linear in  $x$  providing reasonable approximations is clearly true.) However, the model  $y = a_0 + a_1 x + a_2 x^2$  gives the fitted line shown in [figure 11.2](#). This model is still linear in the parameters, so that these can be easily estimated using the same standard matrix manipulation shown above in [equation 11.5](#). It is clear that the predictions obtained from this model are about as good as they can be. The remaining inaccuracy in the model is the irreducible measurement error associated with the variance of  $y$  about its mean at each value of  $x$ .



**Figure 11.1:** Expired Ventilation Plotted Against Oxygen Uptake in a Series of Trials, with Fitted Straight Line.



**Figure 11.2:** The Data From Figure 11.1 with a Model that Includes a Term in  $x^2$ .

### 11.2.2 A Probabilistic Interpretation of Linear Regression

This informal data analytic route allows us to fit a regression model to any data set involving a response variable and a set of predictor variables, and to obtain a vector of estimated regression coefficients. If our aim were merely to produce a convenient summary of the training data (as, very occasionally, it is) we could stop there. However, this chapter is concerned with predictive models. Our aim is to go beyond the training data to predict  $y$  values for other "out-of-sample" objects. Goodness of fit to the given data is all very well, but we are really interested in fit to future data that arise from the same process, so that our future predictions are as accurate as possible. In order to explore this, we need to embed the model-building process in a more formal inferential context. To do this, we suppose that each observed value  $y(i)$  is produced as a sum of weighted predictor variables  $\mathbf{a}^T \mathbf{x}(i)$  and a random term  $\epsilon(i)$  that follows a  $N(0, s^2)$  distribution independent of other values. (Note that implicit in this is the assumption that the variances of the random terms are all the same— $s^2$  is the same for all possible values of the vector of predictor variables. We will discuss this assumption further below.) The  $n \times 1$  random vector  $\mathbf{Y}$  thus takes the form  $\mathbf{Y} = \mathbf{X}\mathbf{a} + \boldsymbol{\epsilon}$ . The observed  $n \times 1$  vector in [equation 11.3](#) is a realization from this distribution. The components of the  $n \times 1$  vector  $\boldsymbol{\epsilon}$  are often called *errors*. Note that they are different from the residuals,  $\mathbf{e}$ . An "error" is a random realization from a given distribution, whereas a residual is a difference between a fitted model and an observed  $y$  value. Note also that  $\mathbf{a}$  is different from  $\hat{\mathbf{a}}$ .  $\mathbf{a}$  represents the underlying and unknown truth, whereas  $\hat{\mathbf{a}}$  gives the values used in a model of the truth.

It turns out that within this framework the least squares estimate  $\hat{\mathbf{a}}$  is also the maximum likelihood estimate of  $\mathbf{a}$ . Furthermore, the covariance matrix of the estimate  $\hat{\mathbf{a}}$  obtained above is  $(\mathbf{X}^T \mathbf{X})^{-1} s^2$ , where this covariance matrix expresses the uncertainty in our parameter estimates  $\hat{\mathbf{a}}$ . In the case of a single predictor variable, this gives

$$(11.6) \quad \left( 1 + \frac{n\bar{x}^2}{\sum_i (x(i) - \bar{x})^2} \right) \frac{\sigma^2}{n}$$

for the variance of the intercept term and

$$(11.7) \quad \frac{\sigma^2}{\sum_i (x(i) - \bar{x})^2}$$

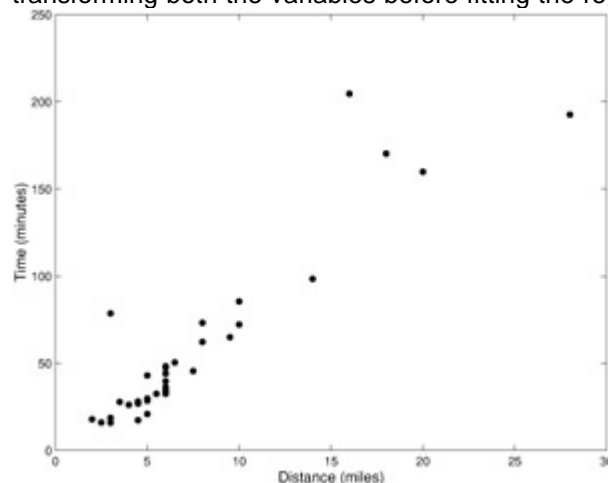
for the variance of the slope. Here  $\bar{x}$  is the sample mean of the single predictor variable.

The diagonal elements of the covariance matrix for  $\hat{\mathbf{a}}$  above give the variances of the regression coefficients—which can be used to test whether the individual regression coefficients are significantly different from zero. If  $v_j$  is the  $j$ th diagonal element of  $(\mathbf{X}^T \mathbf{X})^{-1} s^2$ , then the ratio  $\hat{a}_j / \sqrt{v_j}$  can be compared with a  $t(n - p - 1)$  distribution to see whether the regression coefficient is zero. However, as we discuss below, this test makes sense only in the context of the other variables included in the model, and alternative methods, also discussed below, are available for more elaborate model-building exercises. If  $\mathbf{x}$  is the vector of predictor variables for a new object, with predicted  $y$  value  $\hat{y}$ , then the variance

of  $y$  is  $\mathbf{x}^T(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{x}s^2$ . With one predictor variable, this reduces to  $\sigma^2\left(\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum(x(i) - \bar{x})^2}\right)$ . Note that this variance is greater the further  $\mathbf{x}$  is from the mean of the training sample. That is, the least accurate predictions, in terms of variance, are those in the tails of the distribution of predictor variables. Note also that confidence intervals (see [chapter 4](#)) based on this variance are confidence values for the *predicted value* of  $y$ . We may also be interested in (what are somewhat confusingly called) *prediction intervals*, telling us a range of plausible values for the observed  $y$  at a given value of  $x$ , not a range of plausible values for the predicted value. Prediction intervals must include the uncertainty arising from our prediction and also that arising from the variability of  $y$  about our predicted value. This means that the variance above is increased by an extra term  $s^2$ , yielding  $\sigma^2\left(1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{\sum(x(i) - \bar{x})^2}\right)$ .

### Example 11.1

The most important special case of linear regression arises when there is just one predictor variable. [Figure 11.3](#) shows a plot of the record time (in 1984, in minutes) against the distance (in miles) for 35 Scottish hill races. We can use regression to attempt to predict record time from distance. A simple linear regression of the data gives an estimated intercept value of -4.83 and an estimated regression coefficient of 8.33. Most modern data analytic packages will give the associated standard errors of the estimates, along with significance tests of the null hypotheses that the true parameters that led to the data are zero. In this case, the standard errors are 5.76 and 0.62, respectively, yielding significance probabilities of 0.41 and  $< 0.01$ . From this we would conclude that there is strong evidence that the positive linear relationship is real, but no evidence of a non-zero intercept. The plot in [figure 11.3](#) shows marked skewness in both variables (they become more sparsely spread towards the top and right of the figure). It is clear that the position of the regression line will be much more sensitive to the precise position of points to the right of the figure than it will be to the position of points to the left. Points that can have a big effect on the conclusion are called points of high *leverage*—they are points at the extreme values of estimated relative performance in [figure 11.3](#). Points that actually do have a big effect are called *influential points*. For example, if the rightmost point in [figure 11.3](#) had time of 100 (while still having distance around 28), it would clearly have a big effect on the regression line. The asymmetry of the leverage of the points in the figure might be regarded as undesirable. We might try to overcome this by reducing the skewness—for example, by log transforming both the variables before fitting the regression line.



**Figure 11.3:** A Plot of Record Time (in Minutes) Against Distance (in Miles) for 35 Scottish Hill Races From 1984.

### 11.2.3 Interpreting the Fitted Model

The coefficients in a multiple regression model can be interpreted as follows: if the  $j$ th predictor variable,  $x_j$ , is increased by one unit, while all the other predictor variables are kept fixed, then the response variable  $y$  will increase by  $a_j$ . The regression coefficients thus tell us the *conditional* effect of each predictor variable, conditional on keeping the other predictor variables constant. This is an important aspect of the interpretation. In particular, the size of the regression coefficient associated with the  $j$ th variable will depend on what other variables are in the model. This is clearly especially important if we are constructing models in a sequential manner: add another variable and the coefficients of those already in the model will change. (There is an exception to this. If the predictor variables are orthogonal, then the estimated regression coefficients are unaffected by the presence or absence of others in the model. However, this situation is most common in designed experiments, and is rare in the kinds of secondary data analyses encountered in data mining.) The sizes of the regression coefficients tell us the relative importance of the variables, in the sense that we can compare the effects of unit changes. Note also that the size of the effects depends on the chosen units of measurement for the predictor variables. If we measure  $x_1$  in kilometers instead of millimeters, then its associated regression coefficient will be multiplied by a million. This can make comparisons between variables difficult, so people often work with standardized variables—measuring each predictor variable relative to its standard deviation.

We used the sum of squared errors between the predictions and the observed  $y$  values as a criterion through which to choose the values of the parameters in the model. This is the *residual sum of squares* or the *sum of squared residuals*,  $\sum e(i)^2 = \sum (y(i) - \hat{y}(i))^2$ . In a sense, the worst model would be obtained if we simply predicted all of the  $y$  values by the value  $\bar{y}$ , the mean of the sample of  $y$  values that is constant relative to the  $x$  values (thus effectively ignoring the inputs to the model and always guessing the output to be the mean of  $y$ ). The *total sum of squares* is defined as the sum of squared errors for this worst model,  $\sum (y(i) - \bar{y})^2$ . The difference between the residual sum of squares from a model and the total sum of squares is the sum of squares that can be attributed to the regression for that model—it is the *regression sum of squares*. This is the sum of squared differences of the predicted values,  $\hat{y}(i)$ , from the overall mean,  $\bar{y}$ ,  $\sum (\hat{y}(i) - \bar{y})^2$ . The symbol  $R^2$  is often used for the "multiple correlation coefficient," the ratio of the regression sum of squares to total sum of squares:

$$(11.8) \quad R^2 = \frac{\sum (\hat{y}(i) - \bar{y})^2}{\sum (y(i) - \bar{y})^2}.$$

A value near 1 tells us that the model explains most of the  $y$  variation in the data. The number of independent components contributing to each sum of squares is called the number of *degrees of freedom* for that sum of squares. The degrees of freedom for the total sum of squares is  $n - 1$  (one less than the sample size, since the components are all calculated relative to the mean). The degrees of freedom for the residual sum of squares is  $n - 1 - p$  (although there are  $n$  terms in the summation,  $p + 1$  regression coefficients are calculated). The degrees of freedom for the regression sum of squares is  $p$ , the difference between the total and residual degrees of freedom. These sums of squares and their associated degrees of freedom are usefully put together in an analysis of variance table, as in [table 11.1](#), summarizing the decomposition of the totals into components. The meaning of the final column is described below.

**Table 11.1: The Analysis of Variance Decomposition Table for a Regression.**

| Source of variation | Sum of squares                  | Degrees of freedom | Mean square                                |
|---------------------|---------------------------------|--------------------|--|
| Regression          | $\sum (\hat{y}(i) - \bar{y})^2$ | $p$                | $\sum (y(i) - \bar{y})^2$                  |
| Residual            | $\sum (y(i) - \hat{y}(i))^2$    | $n - p - 1$        | $\sum (y(i) - \hat{y}(i))^2 / (n - p - 1)$ |
| Total               | $\sum (y(i) - \bar{y})^2$       | $n - 1$            |  |



### 11.2.4 Inference and Generalization

We have already noted that our real aim in building predictive models is one of inference: we want to make statements (predictions) about objects for which we do not know the  $y$  values. This means that goodness of fit to the training data is not our real objective. In particular, for example, merely because we have obtained nonzero estimated regression coefficients, this does not necessarily mean that the variables are related: it could be merely that our model has captured chance idiosyncrasies of the training sample. This is particularly relevant in the context of data mining where many models may be explored and fit to the data in a relatively automated fashion. As discussed earlier, we need some way to *test* the model, to see how easily the observed data could have arisen by chance, even if there was no structure in the population the data were collected from. In this case, we need to test whether the population regression coefficients are really zero. (Of course, this is not the only test we might be interested in, but it is the one most often required.) It can be shown that if the values of  $a_j$  are actually all zero (and still making the assumption that the  $\epsilon(i)$  are independently distributed as  $N(0, s^2)$ ),

$$(11.9) \quad \frac{\sum (\hat{y}(i) - \bar{y})^2 / p}{\sum (y(i) - \bar{y})^2 / (n - p - 1)}$$

has an  $F(p, n - p - 1)$  distribution. This is just the ratio of the two mean squares given in [table 11.1](#). The test is carried out by comparing the value of this ratio with the upper critical level of the  $F(p, n - p - 1)$  distribution. If the ratio exceeds this value the test is significant—and we would conclude that there is a linear relationship between the  $y$  and  $x_j$  variables (or that a very unlikely event has occurred). If the ratio is less than the critical value we have no evidence to reject the null hypothesis that the population regression coefficients are all zero.

### 11.2.5 Model Search and Model Building

We have described an overall test to see whether the regression coefficients in a given model are all zero. However, we are more often involved in a situation of searching over model space—or *model building*—in which we examine a sequence of models to find one that is "best" in some sense. In particular, we often need to examine the effect of adding a set of predictor variables to a set we have already included. Note that this includes the special case of adding just one extra variable, and that the idea is applied in reverse, it can also handle the situation of removing variables from a model.

In order to compare models we need a score function. Once again, the obvious one is the sum of squared errors between the predictions and the observed  $y$  values. Suppose we are comparing two models: a model with  $p$  predictor variables (model  $M$ ) and the largest model we are prepared to contemplate, with  $q$  variables (these will include all the untransformed predictor variables we think might be relevant, along with any transformations of them we think might be relevant), model  $M^*$ . Each of these models will have an associated residual sum of squares, and the difference between them will tell us how much better the larger model fits the data than the smaller model. (Equivalently, we could calculate the difference between the regression sums of squares. Since the residual and regression sum of squares sum to the total sum of squares, which is the same for both models, the two calculations will yield the same result.) The degrees of freedom associated with the difference between the residual sums of squares for the two models is  $q - p$ , the extra number of regression coefficients computed in fitting the larger model,  $M^*$ . The ratio between the difference of the residual sums of squares and the difference of degrees of freedom again gives us a mean square—now a mean square for the difference between the two models. Comparison of this with the residual mean square for model  $M^*$  gives us an  $F$ -test of whether the difference between the models is real or not. [Table 11.2](#) illustrates this extension. From this table, the ratio

$$\left[ \frac{(SS(M^*) - SS(M))}{(q - p)} \right] / \left[ \frac{(SS(T) - SS(M^*))}{(n - q - 1)} \right]$$

is compared with the critical value of an  $F(q - p, n - q - 1)$  distribution.

|   |
|---|
| <b>Table 11.2: The Analysis of Variance Decomposition Table for Model Building.</b> |
|---|

| Source of variation   | Sum of squares    | Degrees of freedom | Mean square                           |
|-----------------------|-------------------|--------------------|---------------------------------------|
| Regression Model 1    | $SS(M)$           | $p$                | $SS(M)/p$                             |
| Regression Full Model | $SS(M^*)$         | $q$                | $SS(M^*)/q$                           |
| Difference            | $SS(M^*) - SS(M)$ | $q - p$            | $\frac{SS(M^*) - SS(M)}{(q - p)}$     |
| Residual              | $SS(T) - SS(M^*)$ | $n - p - 1$        | $\frac{SS(T) - SS(M^*)}{(n - q - 1)}$ |
| Total                 | $SS(T)$           | $n - 1$            |                                       |

This is fine if we have just a few models we want to compare, but data mining problems are such that often we need to rely on automatic model building processes. Such automatic methods are available in most modern data mining computer packages. There are various strategies that may be adopted. A basic form is a *forward selection* method, mentioned in [chapter 8](#), in which variables are added one at a time to an existing model. At each step that variable is chosen from the set of potential variables that leads to the greatest increase in predictive power (measured in terms of reduction of sum of squared residuals), provided the increase exceeds some specified threshold. Ideally, the addition would be made as long as the increase in predictive power was statistically significant, but in practice this is complicated to ensure: the variable selection process necessarily involves carrying out many tests, not all independent, so that computing correct significance values is a nontrivial process. The simple significance level based on [table 11.2](#) does not apply when multiple dependent tests are made. (The implication of this is that if the significance level is being used to choose variables, then it is being used as a score function, and should not be given a probabilistic interpretation.)

We can, of course, in principle use any of the score functions discussed in [chapter 7](#) for model selection in regression, such as BIC, minimum description length, cross-validation, or more Bayesian methods. These provide an alternative to the hypothesis-testing framework that measures the statistical significance of adding and deleting terms on a model-by-model basis. Penalized score functions such as BIC, and variations on cross-validation tailored specifically to regression, are commonly used in practice as score functions for model selection in regression.

A strategy opposite to that of forward selection is *backward elimination*. We begin with the most complex model we might contemplate (the "largest model,"  $M^*$ , above) and progressively eliminate variables, selecting them on the basis that eliminating them leads to the least increase in sum of squared residuals (again, subject to some threshold). Other variants include combinations of forward selection and backward elimination. For example, we might add two variables, eliminate one, add two, remove one, and so on. For data sets where the number of variables  $p$  is very large, it may be much more practical computationally to build the model in the forward direction than in the backward direction. Stepwise methods are attempts to restrict the search of the space of all possible sets of predictor variables, so that the search is manageable. But if the search is restricted, it is possible that some highly effective combination of variables may be overlooked. Very occasionally (if the set of potential predictor variables is small), we can examine all possible sets of variables (although, with  $p$  variables, there are  $(2^p - 1)$  possible subsets). The size of problems for which all possible subsets can be examined has been expanded by the use of strategies such as branch and bound, which rely on the monotonicity of the residual sum of squares criterion (see [chapter 8](#)).

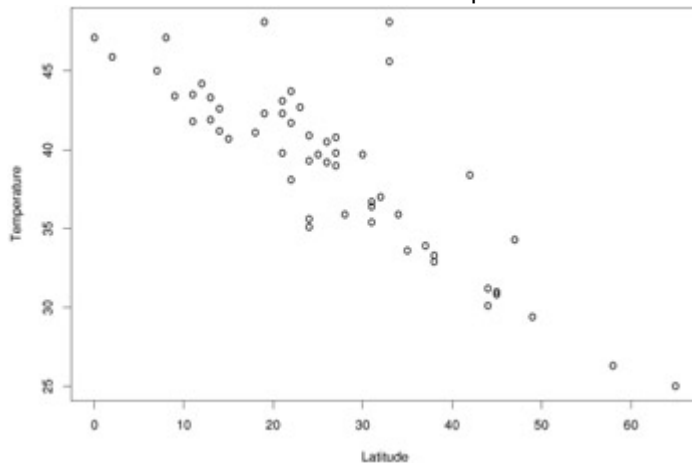
A couple of cautionary comments are worth making here. First, as we have noted, the coefficients of variables already in the model will generally change as new variables are added. A variable that is important for one model may become less so when the model is extended. Second, as we have discussed in earlier chapters, if too elaborate a search is carried out there is a high chance of overfitting the training set—that is, of obtaining a



model that provides a good fit to the training set (small residual sum of squares) but does not predict new data very well.

### 11.2.6 Diagnostics and Model Inspection

Although multiple regression is a very powerful and widely used technique, some of the assumptions might be regarded as restrictive. The assumption that the variance of the  $y$  distribution is the same at each vector  $\mathbf{x}$  is often inappropriate. (This assumption of equal variances is called *homoscedasticity*. The converse is *heteroscedasticity*.) For example, [figure 11.4](#) shows the normal average January minimum temperature (in deg F) plotted against the latitude (deg N) for 56 cities in the United States. There is evidence that, for smaller latitudes, at least, the variance of the temperature increases with increasing latitude (although the mean temperature seems to decrease). We can still apply the standard least squares algorithm above to estimate parameters in this new situation (and the resulting estimates would still be unbiased if the model form were correct), but we could do better in the sense that it is possible to find estimators with smaller variance.



**Figure 11.4:** Temperature (Degrees F) Against Latitude (Degrees N) for 56 Cities in the United States.

To do this we need to modify the basic method. Essentially, we need to arrange things so that those values of  $\mathbf{x}$  associated with  $y$  values with larger variance are weighted less heavily in the model fitting process. This makes perfect sense—it means that the estimator is more influenced by the more accurate values. Formally, this idea leads to a modification of the solution [equation 11.5](#). Suppose that the covariance matrix of the  $n \times 1$  random vector  $\epsilon$  is the  $n \times n$  matrix  $s^2 \mathbf{V}$  (previously we took  $\mathbf{V} = \mathbf{I}$ ). The case of unequal variances means that  $\mathbf{V}$  is diagonal with terms that are not all equal. Now it is possible (see any standard text on linear algebra) to find a unique nonsingular matrix  $\mathbf{P}$  such that  $\mathbf{P}^T \mathbf{P} = \mathbf{V}$ . We can use this to define a new random vector  $\mathbf{f} = \mathbf{P}^{-1} \epsilon$ , and it is easy to show that the covariance matrix of  $\mathbf{f}$  is  $s^2 \mathbf{I}$ . Using this idea, we form a new model by premultiplying the old one by  $\mathbf{P}^{-1}$ :

$$(11.10) \quad \mathbf{P}^{-1} \mathbf{Y} = \mathbf{P}^{-1} \mathbf{X} \alpha + \mathbf{P}^{-1} \epsilon$$

or

$$(11.11) \quad \mathbf{Z} = \mathbf{W} \beta + \mathbf{f},$$

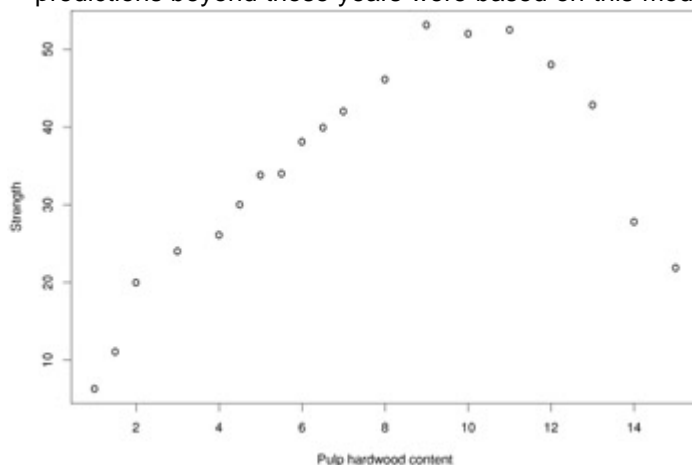
now of the form required to apply the standard least squares algorithm. If we do this, and then convert the solution back into the original variables  $\mathbf{Y}$ , we obtain:

$$(11.12) \quad \mathbf{a} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}) \mathbf{X} \mathbf{V}^{-1} \mathbf{y},$$

a *weighted* least squares solution. The variance of this estimated parameter vector  $\mathbf{a}$  is  $(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} s^2$ .

Unequal variances of the  $y$  distributions for different  $\mathbf{x}$  vectors is one way in which the assumptions of basic multiple regression can break down. There are others. What we really need are ways to explore the quality of the model and tools that will enable us to detect where and why the model deviates from the assumptions. That is, we require *diagnostic* tools. In simple regression, where there is only one predictor variable, we can see the quality of the model from a plot of  $y$  against  $x$  (see [figures 11.1](#), [11.2](#) and [11.4](#)). More generally, however, when there is more than one predictor variable, such a simple plot is not possible, and more sophisticated methods are needed. In general, the key

features for examining the quality of a regression model are the residuals, the components of the vector  $\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}}$ . If there is a pattern to these, it tells us that the model is failing to explain the distribution of the data. Various plots involving the residuals are used, including plotting the residuals against the fitted values, plotting standardized residuals (obtained by dividing the residuals by their standard errors) against the fitted values, and plotting the standardized residuals against standard normal quantiles. (The latter are "normal probability plots." If the residuals are approximately normally distributed, the points in this plot should lie roughly on a straight line.) Of course, interpreting some of the diagnostic plots requires practice and experience. One general cautionary comment, applies to all predictive models: such models are valid only within the bounds of the data. It can be very risky to extrapolate beyond the data. A very simple example is given in [figure 11.5](#). This shows a plot of the tensile strength of paper plotted against the percentage of hardwood in the pulp from which the paper was made. But suppose only those samples with pulp values between 1 and 9 had been measured. The figure shows that a straight line would provide quite a good fit to this subset of the data. For new samples of paper, with pulp values lying between 1 and 9, quite good prediction of the strength could legitimately be expected. But the figure also shows, strikingly clearly, that our model would produce predictions that were seriously amiss if we used it to predict the strength of paper with pulp values greater than 9. Only within the bounds of our data is the model trustworthy. We saw another example of this sort of thing in [chapter 3](#), where we showed the number of credit cards in circulation each year. A straight line fitted to years 1985 to 1990 provided a good fit—but if predictions beyond those years were based on this model, disaster would follow.



**Figure 11.5:** A Plot of Tensile Strength of Paper against the Percentage of Hardwood in the Pulp.

These examples are particularly clear—but they involve just a few data points and a single predictor variable. In data mining applications, with large data sets and many variables, things may not be so clear. Caution needs to be exercised when we make predictions.

### 11.3 Generalized Linear Models

[Section 11.2](#) described the *linear model*, in which the response variable was decomposed into two parts: a weighted sum of the predictor variables and a random component:  $Y(i) = \beta_0 + \beta_1 x_1(i) + \dots + \beta_p x_p(i) + \epsilon(i)$ . For inferential purposes we also assumed that the  $\epsilon(i)$  were independently distributed as  $N(0, s^2)$ . We can write this another way, which permits convenient generalization, splitting the description of the model into three parts:

- The  $Y(i)$  are independent random variables, with distribution  $N(\mu(i), s^2)$ .
- The parameters enter the model in a linear way via the sum  $\mu(i) = \beta_0 + \beta_1 x_1(i) + \dots + \beta_p x_p(i)$ .
- The  $\mu(i)$  and  $\beta_j$  are linked by  $\mu(i) = \beta_0 + \beta_1 x_1(i) + \dots + \beta_p x_p(i)$ .

This permits two immediate generalizations, while retaining the advantages of the linear combination of the parameters. First, in (i) we can relax the requirement that the random variables follow a normal distribution. Second, we can generalize the link expressed in

(iii), so that some other *link* function  $g(\mu(i)) = \eta(i)$  relates the parameter of the distribution to the linear term  $\eta(i) = \sum_j a_j x_j(i)$ . These extensions result in what are called *generalized linear models*. They are one of the most important advances in data analysis of the last two decades. As we shall see, such models can also be regarded as fundamental components of feed forward neural networks.

One of the most important kinds of generalized linear model for data mining is *logistic regression*. We have already encountered this in [chapter 10](#) in the form of logistic discrimination, but we describe it in rather more detail here, and use it as an illustration of the ideas underlying generalized linear models. In many situations the response variable is not continuous, as we have assumed above, but is a proportion: the number of flies from a given sample that die when exposed to an insecticide, the proportion of questions people get correct in a test, the proportion of oranges in a carton that are rotten. The extreme of this arises when the proportion is out of 1, that is, the observed response is binary: whether or not an individual insect dies, whether or not a person gets a particular one of the questions right, whether or not an individual orange is rotten. This is exactly the situation we discussed in [chapter 10](#), though here we embed it in a more general context. We are now dealing with a binary response variable, with the random variable taking values 0 or 1 corresponding to the two possible outcomes. We will assume that the probability that the  $i$ th individual yields the value 1 is  $p(i)$ , and that the responses of different individuals are independent. This means that the response for the  $i$ th individual follows a Bernoulli distribution:

$$(11.13) \quad p(Y(i) = y(i)) = p(i)^{y(i)} (1 - p(i))^{1-y(i)},$$

where here  $y(i) \in \{0, 1\}$ . For logistic regression, this is the generalization of (i) above: the Bernoulli distribution is replacing the normal distribution.

Our aim is to formulate a model for the probability that an object with predictor vector  $\mathbf{x}$  will take value 1. That is, we want a model for the mean value of the response, the probability  $p(y = 1|\mathbf{x})$ . We could use a linear model—a weighted sum of the predictor variables. However, this would not be ideal. Most obviously, a linear model can take values less than 0 and greater than 1 (if the  $\mathbf{x}$  values are extreme enough). This suggests that we need to modify the model to include a nonlinear aspect. We achieve this by transforming the probability, nonlinearly, so that it can be modeled by a linear combination. That is, we use a nonlinear link function in (iii). A suitable function (not the only possible one) is a *logistic* (or *logit*) link function, in which

$$(11.14) \quad g(p(y = 1|\mathbf{x})) = \log \frac{p(y = 1|\mathbf{x})}{1 - p(y = 1|\mathbf{x})},$$

where  $g(p(y = 1|\mathbf{x}))$  is modeled as  $\sum_j a_j x_j$ . As  $p$  varies from 0 to 1,  $\log(p/1 - p)$  clearly varies from  $-\infty$  to  $\infty$ , matching the potential range of  $g(p) = \sum_j a_j x_j(i)$ . One of the advantages of the logistic link function over alternatives is that it permits convenient interpretation. For example:

- The ratio  $\frac{p(y=1|\mathbf{x})}{1-p(y=1|\mathbf{x})}$  in the transformation is the familiar *odds* that a 1 will be observed and  $\log \frac{p(y=1|\mathbf{x})}{1-p(y=1|\mathbf{x})}$  is the *log odds*.
- Given a new vector of predictor variables  $\mathbf{x} = (x_1, \dots, x_n)$ , the predicted probability of observing a 1 is derived from  $\log \frac{p(y=1|\mathbf{x})}{1-p(y=1|\mathbf{x})}$ . The effect on this of changing the  $j$ th predictor variable by one unit is simply  $a_j$ . Thus the coefficients tell us the difference in log odds—or, equivalently, the log odds ratio resulting from the two values. From this it is easy to see that  $e^{a_j}$  is the factor by which the odds changes when the  $j$ th predictor variable changes by one unit (see the discussion of the effect of a unit change of one variable in the multiple regression case discussed in [section 11.2](#)).

### Example 11.2

Two minutes into its flight on January 29, 1986, the space shuttle Challenger exploded, killing everyone on board. The two booster rockets for the shuttle are made of several pieces, with each of three joints sealed with a rubber "O-ring," making six rings in total. It was known that these O-rings were sensitive to temperature. Records of the proportion of O-rings damaged in previous flights were available, along with the temperatures on those

days. The lowest previous temperature was 53degF. On the day of the flight the temperature was 31degF, so there was much discussion about whether the flight should go ahead. One argument was based on an analysis of the seven previous flights that had resulted in damage to at least one O-ring. A logistic regression to predict the probability of failure from temperature led to a slope estimate of 0.0014 with a standard error of 0.0498. From this, the predicted logit of the probability of failure at 31degF is 1.3466, yielding a predicted probability 0.206. The slope in this model is positive, suggesting that, if anything, the probability of failure is lower at low temperatures. However, this slope is not significantly different from zero, so that there is little evidence for a relationship between failure probability and temperature.

This analysis is far from ideal. First, 31degF is far below 53degF, so one is extrapolating beyond the data—a practice we warned against above. Secondly, there is valuable information in the 16 flights that had not resulted in O-ring damage. This is immediately obvious from a comparison of [figure 11.6\(a\)](#), which shows the numbers damaged for the seven flights above (vertical axis) against temperature (horizontal axis), and [figure 11.6\(b\)](#), which shows the number for all 23 flights. These 16 flights all took place at relatively high temperatures. The second figure suggests that the relationship might, in fact, have a negative slope. A logistic model fitted to the data in [figure 11.6\(b\)](#) gave a slope estimate of -0.1156, with a standard error of -2.46 (and an intercept estimate of 5.08 with standard error of 3.05). From this the predicted probability at 31degF is 0.817. This gives a rather different picture, one that could have been deduced before the flight if all the data had been studied.



**Figure 11.6:** Number of O-Rings Damaged (Vertical Axis) against Temperature on Day of Flight, (a) Data Examined before the Flight, and (b) The Complete Data.

Generalized linear models thus have three main features:

- i. The  $Y(i)$ ,  $i = 1, \dots, n$ , are independent random variables, with the same exponential family distribution (see below).
- ii. The predictor variables are combined in a form  $\eta(i) = \sum_j a_j x_j(i)$ , called the *linear predictor*, where the  $a_j$ s are *estimates* of the  $\beta_j$ s.
- iii. The mean  $\mu(i)$  of the distribution for a given predictor vector is related to the linear combination in (ii) through the link function  $\eta(i) = \mu(i) = \sum_j a_j x_j(i)$ .

The *exponential family* of distributions is an important family that includes the normal, the Poisson, the Bernoulli, and the binomial distributions. Members of this family can be expressed in the general form

$$(11.15) \quad f(y; \theta, \phi) = c \frac{\exp(\eta(\theta)T(y) - \phi(\theta))}{a(\phi)},$$

If  $\eta$  is known, then  $\theta$  is called the *natural* or *canonical* parameter. When, as is often the case,  $a(\eta) = \eta$ ,  $\phi$  is called the *dispersion* or *scale* parameter. A little algebra reveals that the mean of this distribution is given by  $\eta'(\theta)$  and variance by  $a(\phi) \eta''(\theta)$ . Note that the variance is related to the mean via  $\eta''(\theta)$ , and this, expressed in the form  $V(\eta)$ , is sometimes called the *variance function*. In the model as described in (i) to (iii) above, there are no restrictions on the link function. However (and this is where the exponential family comes in), things simplify if the link function is chosen to be the function expressing the canonical parameter for the distribution being used as a linear sum. For multiple regression this is simply the identity distribution and for logistic regression it is the logistic transformation presented above. For *Poisson regression*, in which the distribution in (i) is the Poisson distribution, the canonical link is the log link  $g(u) = \log(u)$ . Prediction from a generalized linear model requires the inversion of the relationship  $g(\mu(i)) = \sum_j a_j x_j(i)$ . The algorithms in least squares estimation were very straightforward, essentially involving only matrix inversion. For generalized linear models, however, things are more complicated: the non-linearity means that an iterative scheme has to be adopted. We will not go into details of the mathematics here, but it is not difficult to show that the maximum likelihood solution is given by solving the equations

$$(11.16) \quad \sum_{i=1}^n \frac{x_j(i) (y(i) - \mu(i))}{a_i(\phi) V(\mu(i)) g'(\mu(i))} = 0, \quad j = 1, \dots, p,$$

where the  $i$  indices for  $a_i(\phi)$  and  $\mu(i)$  are in recognition of the fact that these vary from data point to data point. Standard application of the Newton-Raphson method ([chapter 8](#)) leads to iteration of the equations

$$(11.17) \quad \mathbf{a}^{(s)} = \mathbf{a}^{(s-1)} - \mathbf{M}_{s-1}^{-1} \mathbf{u}_{s-1},$$

where  $\mathbf{a}^{(s)}$  represents the vector of values of  $(a_1, \dots, a_p)$  at the  $s$ th iteration,  $\mathbf{u}_{s-1}$  is the vector of first derivatives of the log likelihood, evaluated at  $\mathbf{a}^{(s-1)}$ , and  $\mathbf{M}_{s-1}$  is the matrix of second derivatives of the log likelihood, again evaluated at  $\mathbf{a}^{(s-1)}$ .

An alternative method, the method of "scoring" (this is a traditional name, and is not to be confused with our use of the word score in "score function," though the meaning is similar), replaces  $\mathbf{M}_{s-1}$  by the matrix of expected second derivatives. The iterative steps of this method can be expressed in a form similar to the weighted version, [equation 11.12](#), of the standard least squares matrix solution, [equation 11.5](#):

$$(11.18) \quad (\mathbf{X}' \mathbf{W}_{(s-1)} \mathbf{X}) \mathbf{a}^{(s)} = \mathbf{X}' \mathbf{W}_{(s-1)} \mathbf{z}_{(s-1)},$$

where  $\mathbf{W}_{(s-1)}$  is a diagonal matrix with  $i$ th element  $(\mu(i) - \eta(i))^2 / \text{var}(Y(i))$  evaluated at  $\mathbf{a}^{(s-1)}$  and  $\mathbf{z}_{(s-1)}$  is a vector with  $i$ th element  $x_j(i) a_j + (\eta(i) - \mu(i)) \eta'(i) / \mu(i)$  again evaluated at  $\mathbf{a}^{(s-1)}$ . Given the similarity of this to [equation 11.12](#) it will hardly be surprising to learn that this method is called *iteratively weighted least squares*. We need a measure of the goodness of fit of a generalized linear model, analogous to the sum of squares used for linear regression. Such a measure is the *deviance* of a model. In fact, the sum of squares is the special case of deviance when it is applied to linear models. Deviance is defined as  $D(M) = -2 (\log L(M; Y) - \log L(M^*; Y))$ , essentially the difference between the log likelihood of model  $M$  and the log likelihood of the largest model we are prepared to contemplate,  $M^*$ . Deviance can be decomposed like the sum of squares to permit exploration of classes of models.

### Example 11.3

In a study of ear infections in swimmers, 287 swimmers were asked if they were frequent ocean swimmers, whether they preferred beach or nonbeach, their age, their sex, and also the number of self-diagnosed ear infections they had had in a given period. The last variable here is the response variable, and a predictive model is sought, in which the number of ear infections can be predicted from the other variables. Clearly standard linear regression would be inappropriate: the response variable is discrete and, being a count, is unlikely to look remotely like a normal distribution. Likewise, it is not a proportion, it is not bounded between 0 and 1, so it would be inappropriate to model it using logistic regression. Instead, it is reasonable to assume that the response variable follows a Poisson distribution, with a parameter depending on the value of the predictor variables. Fitting a generalized linear model to predict the number of infections from the other variables, with the response following a Poisson distribution and using a log function for the link, led to the *analysis of deviance* [table 11.3](#).

**Table 11.3: Analysis of Deviance Table.**

|            | d.f. | deviance | mean deviance | deviance ratio |
|------------|------|----------|---------------|----------------|
| Regression | 4    | 1.67     | 0.4166        | 0.42           |
| Residual   | 282  | 47.11    | 0.1671        |                |
| Total      | 286  | 48.78    | 0.1706        |                |
| Change     | -4   | -1.67    | 0.4166        | 0.42           |

To test the null hypothesis of no predictive relationship between the response variable and the predictors, we compare the value of the regression deviance (1.67, from the top of the second column of numbers) with the chi-squared distribution with 4 degrees of freedom (given at the top of the first column of numbers). This gives a  $p$ -value of 0.7962. This is far from small, suggesting that there is little evidence that the response variable is related to the predictor variables. Not all data necessarily lead to a model that gives accurate predictions!

Before leaving this section, it is worth noting a property of equations 11.16. Although these were derived assuming that the random variables follow an exponential family distribution, examination reveals that these estimating equations make use of only the means  $\mu(i)$ ; the variances  $a_i(f) V(\mu(i))$ , as well as the link function and the data values. There is nothing about any other aspect of the distributions. This means that even if we are not prepared to make tighter distributional assumptions, we can still estimate the parameters in the linear predictor  $\eta(i) = \sum a_i x_i(i)$ . Because no full likelihood has to be formulated in this approach, it is termed *quasilikelihood estimation*. Once again, of course, iterative algorithms are needed.

## 11.4 Artificial Neural Networks

Artificial neural networks (ANNs) are one of a class of highly parameterized statistical models that have attracted considerable attention in recent years (other such models are outlined in later sections). In the present context, we will be concerned only with *feed-forward neural networks* or *multilayer perceptrons*, as originally discussed in [chapter 5](#). In this section, we can barely scratch the surface of this topic, and suitable further reading is suggested below. The fact that ANNs are highly parameterized makes them very flexible, so that they can accurately model relatively small irregularities in functions. On the other hand, as we have noted before, such flexibility means that there is a serious danger of overfitting. Indeed, early (by which is meant during the 1980s) work



was characterized by inflated claims when such networks were overfitted to training sets, with predictions of future performance being based on the training set performance. In recent years strategies have been developed for overcoming this problem, resulting in a very powerful class of predictive models.

To set ANNs in context, recall that the generalized linear models of the [previous section](#) formed a linear combination of the predictor variables, and transformed this via a nonlinear transformation. Feedforward ANNs adopt this as the basic element. However, instead of using just one such element, they use multiple layers of many such elements. The outputs from one layer—the transformed linear combinations from each basic element—serve as inputs to the next layer. In this next layer the inputs are combined in exactly the same way—each element forms a weighted sum that is then non-linearly transformed. Mathematically, for a network with just one layer of transformations between the input variables  $x$  and the output  $y$  (one *hidden* layer), we have

$$(11.19) \quad y = \sum_k w_k^{(2)} f_k \left( \sum_j w_j^{(1)} x_j \right).$$

Here the  $w$  are the weights in the linear combinations and the  $f_k$ s are the non-linear transformations. The nonlinearity of these transformations is essential, since otherwise the model reduces to a nested series of linear combinations of linear combinations—which is simply a linear combination. The term *network* derives from a graphical representation of this structure in which the predictor variables and each weighted sum are nodes, with edges connecting the terms in the summation to the node.

There is no limit to the number of layers that can be used, though it can be proven that a single hidden layer (with enough nodes in that layer) is sufficient to model any continuous functions. Of course, the practicality of this will depend on the available data, and it might be convenient for other reasons (such as interpretability) to use more than one hidden layer. There are also generalizations, in which layers are skipped, with inputs to a node coming not only from the layer immediately preceding it but also from other preceding layers.

The earliest forms of ANN used threshold logic units as the nonlinear transformations: the output was 0 if the weighted sum of inputs was below some threshold and 1 otherwise. However, there are mathematical advantages to be gained by adopting differentiable forms for these functions. In applications, the two most common forms are logistic  $f(x) = e^x / (1 + e^x)$  and hyperbolic tangent  $f(x) = \tanh(x)$  transformations of the weighted sums.

We saw, when we moved from simple linear models to generalized linear models, that estimating the parameters became more complicated. A further extra level of complication occurs when we move from generalized linear models to ANNs. This will probably not come as a surprise, given the number of parameters (these now being the weights in the linear combinations) in the model and the fundamental nonlinearity of the transformations. As a consequence of this, neural network models can be slow to train. This can limit their applicability in data mining problems involving large data sets. (But slow estimation and convergence is not all bad. There are stories within the ANN folklore relating how severe overfitting by a flexible model has been avoided by accident, simply because the estimation procedure was stopped early.) Various estimation algorithms have been proposed. A popular approach is to minimize the score function consisting of the sum of squared deviations (again!) between the output and predicted values by steepest descent on the weight parameters. This can be expressed as a sequence of steps in which the weights are updated, working from the output node(s) back to the input nodes. For this reason, the method is called *back-propagation*. Other criteria have also been used. When  $Y$  takes only two values (so that the problem is really one of supervised classification, as discussed in [chapter 10](#)) the sum of squared deviations is rather unnatural (since, as we have seen, the sum of squared deviations arises as a score function naturally from the log-likelihood for normal distributions). A more natural score function, based on log-likelihood for Bernoulli data, is

$$(11.20) \quad \sum_i \left[ y(i) \log \frac{\hat{y}(i)}{y(i)} - (1 - y(i)) \log \frac{(1 - \hat{y}(i))}{(1 - y(i))} \right].$$

As it happens, in practical applications with reasonably sized data sets, the precise choice of criterion seems to make little difference. The vast amount of work on neural networks in recent years, which has been carried out by a diverse range of intellectual communities, has led to the rediscovery of many concepts and phenomena already well known and understood in other areas. It has also led to the introduction of unnecessary new terminology.

Nonetheless, research in this area has also led to several novel general forms of models that we have not discussed here. For example, *radial basis function* networks replace the typical logistic nonlinearity of feedforward net-works with a "bump" function (a radial basis function). An example would be a set of  $p$ -dimensional Gaussian bumps in  $\mathbf{x}$  space, with specified widths. The output is approximated as a linear weighted combination of these bump functions. Model training consists of estimating the locations, widths, and weights of the bumps, in a manner reminiscent of mixture models described in [chapter 9](#).

## 11.5 Other Highly Parameterized Models

The characterizing feature of neural networks is that they provide a very flexible model with which to approximate functions. Partly because of this power and flexibility, but probably also partly because of the appeal of their name with its implied promise, they have attracted a great deal of media attention. However, they are not the only class of flexible models. Others, in some cases with an approximating power equivalent to that of neural net-works, have also been developed. Some of these have advantages as far as interpretation and estimation goes. In this section we briefly outline two of the more important classes of flexible model. Others are mentioned in [section 11.5.2](#).

### 11.5.1 Generalized Additive Models

We have seen how the generalized linear model extends the ideas of linear models. Yet further extension arises in the form of *generalized additive models*. These replace the simple weighted sums of the predictor variables by weighted sums of *transformed* versions of the predictor variables. To achieve greater flexibility, the relationships between the response variable and the predictor variables are estimated nonparametrically—for example, by kernel or spline smoothing (see [chapter 6](#)), so that the generalized linear model form  $g(\mu(i)) = \sum_j \alpha_j x_j(i)$  becomes  $g(\mu(i)) = \sum_j \alpha_j f_j(x_j(i))$ . The right-hand side here is sometimes termed the *additive predictor*. Such models take to the nonparametric limit the idea of extending the scope of linear models by transforming the predictor variables. Generalized additive models of this form retain the merits of linear and generalized linear models. In particular, how  $g$  changes with any particular predictor variable does not depend on how other predictor variables change; interpretation is eased. Of course, this is at the cost of assuming that such an additive form does provide a good approximation to the "true" surface. The model can be readily generalized by including multiple predictor variables within individual  $f$  components of the sum, but this is at the cost of relaxing the simple additive interpretation. The additive form also means that we can examine each smoothed predictor variable separately, to see how well it fits the data.

In the special case in which  $g$  is the identity function, appropriate smoothing functions can be found by a *backfitting* algorithm. If the additive model  $y(i) = \sum_j \alpha_j f_j(x_j(i)) + \epsilon(i)$  is correct, then

$$f_k(X_k) = E \left( Y - \sum_{j \neq k} \alpha_j f_j(X_j(i)) \mid X_k \right).$$

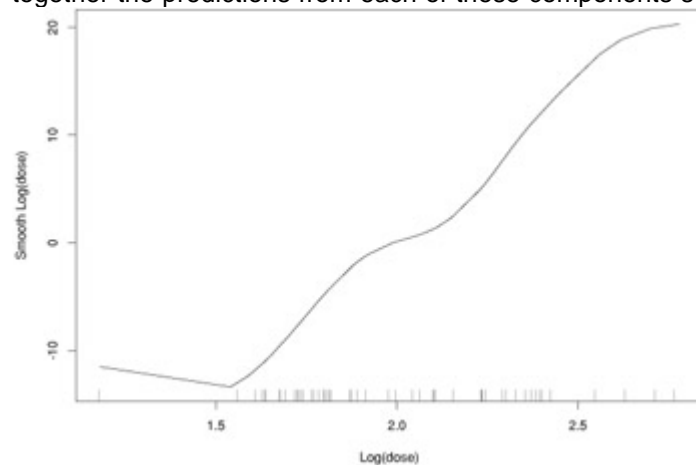
This leads to an iterative algorithm in which, at each step the "partial residuals"  $y - \sum_{j \neq k} \alpha_j f_j(x_j(i))$  for the  $k$ th predictor variable are smoothed, cycling through the predictor variables until the smoothed functions do not change. The precise details will, of course, depend on the choice of smoothing method: kernel, spline, or whatever.

To extend this from additive to generalized additive models, we make the same extension as above, where we extended the ideas from linear to generalized linear models. We have already outlined the iteratively weighted least squares algorithm for fitting generalized linear models. We showed that this was essentially an iteration of a

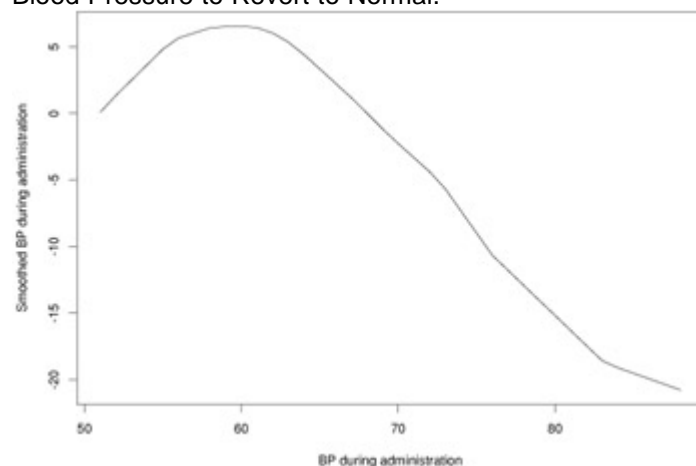
weighted least squares solution applied to an "adjusted" response variable, defined by  $\sum_j x_j(i)a_j + (y(i) - \mu(i)) \frac{\partial \eta(i)}{\partial \mu(i)}$ . For generalized additive models, instead of the weighted linear regression we adopt an algorithm for fitting a weighted additive model.

#### Example 11.4

Sometimes blood pressure is deliberately lowered during surgery, using drugs. Once the operation is completed, and the administration of the drug discontinued, it is desirable that the blood pressure return to normal as soon as possible. The data in this example relate to how soon (in minutes) systolic blood pressure returned to 100 mm of mercury after the medication was discontinued. There are two predictor variables: the log of the dose of the particular drug used and the average systolic blood pressure of the patient during administration of the drug. A generalized additive model was fitted, using splines (in fact, cubic B-splines) to effect the smoothing. [Figures 11.7](#) and [11.8](#) show, respectively, a plot of the transformed Log(dose) against observed Log(dose) values and a plot of the transformed blood pressure during administration against the observed values. (There is some nonlinearity evident in both these plots—although that in the Log(dose) plot seems to be attributable to a single point.) Predictions to new data points are made by adding to together the predictions from each of these components separately.



**Figure 11.7:** The Transformation Function of Log(dose) in the Model for Predicting Time for Blood Pressure to Revert to Normal.



**Figure 11.8:** The Transformation Function of Blood Pressure During Administration in the Model for Predicting Time for Blood Pressure to Revert to Normal.

### 11.5.2 Projection Pursuit Regression

Projection pursuit regression models can be proven to have the same ability to estimate arbitrary functions as neural networks, but they are not as widely used. This is perhaps unfortunate, since estimating their parameters can have advantages over the neural network situation. The additive models of the last section essentially focus on individual variables (albeit transformed versions of these). Such models can be extended so that each additive component involves several variables, but it is not clear how best to select such subsets. If the total number of available variables is large, then we may also be faced with a combinatorial explosion of possibilities. The basic projection pursuit regression model takes the form

$$(11.21) \quad Y = \alpha_0 + \sum f_k(\alpha_k^T \mathbf{X}) + \varepsilon.$$

This has obvious close similarities to the neural network model—it is a linear combination of (potentially nonlinear) transformations of linear combinations of the raw variables. Here, however, the  $f$  functions are not constrained (as in neural networks) to take a particular form, but are usually found by smoothing, as in generalized additive models. This makes them a generalization of neural networks. Various forms of smoothing have been used, including spline methods, Friedman's "supersmoother" (which makes a local linear fit about the point where the smooth is required), and various polynomial functions. The term *projection pursuit* arises from the viewpoint that one is projecting  $\mathbf{X}$  in direction  $\mathbf{a}_k$ , and then seeking directions of projection that are optimal for some purpose. (In this case, optimal as components in a predictive model.) Various algorithms have been developed to estimate the parameters. In one, components of the sum are added sequentially up to some maximum value, and then sequentially dropped, each time selecting on the basis of least squares fit of the model to the data. For a given number of terms, the model is fitted using standard iterative procedures to estimate the parameters in the  $\mathbf{a}_k$  vector. This fitting process is rather complex from a computational viewpoint, so that projection pursuit regression tends may not be practical for data sets that are massive (large  $n$ ) and high-dimensional (large  $p$ ).

## 11.6 Further Reading

Traditional linear regression is covered in depth in the classic book of [Draper and Smith \(1981\)](#), as well as in innumerable other texts. [Furnival and Wilson \(1974\)](#) describe the classic "leaps and bounds" algorithm, which efficiently searches for the best subset of predictors to include in a regression model. The seminal text on generalized linear models is that of [McCullagh and Nelder \(1989\)](#), and a comprehensive outline of generalized additive models is given in the book by [Hastie and Tibshirani \(1990\)](#). Projection pursuit regression (PPR) was introduced by [Friedman and Stuetzle \(1981\)](#), and theoretical approximation results are given in (for example) [Diaconis and Shashahani \(1984\)](#). A very flexible data-driven model for multivariate regression called MARS (Multivariate Adaptive Regression Splines) was introduced by [Friedman \(1991\)](#). [Breiman et al. \(1984\)](#) describe the application of tree-structure models to regression, and [Weiss and Indurkha \(1995\)](#) describe related techniques for rule-based regression models. The technique of boosting, mentioned in chapter 10 in the context of classification, can also be usefully applied to regression. Regression can of course be cast in a Bayesian context, e.g., [Gelman, Carlin, Stern, and Rubin \(1995\)](#). Techniques for *local regression*, analogous to kernel models for density estimation (chapter 9) and nearest neighbor methods for classification (chapter 10), rely on adaptive local fits to achieve a nonparametric regression function (for example, [Cleveland and Devlin \(1988\)](#) and [Atkeson, Schall and Moore \(1997\)](#)). Such techniques, however, can be quite computationally intensive and also are susceptible to the same estimation problems that plague local kernel methods in general in high dimensions. Good introductions to neural networks are given by [Bishop \(1995\)](#), [Ripley \(1996\)](#), [Golden \(1996\)](#), [Ballard \(1997\)](#), and [Fine \(1999\)](#). Ripley's text is particularly noteworthy in that it includes an integrated and extensive discussion of many techniques from the fields of neural networks, statistics, machine learning, and pattern recognition (unlike most texts which tend to focus on one or two of these areas). Bayesian approaches to neural network training are described in [MacKay \(1992\)](#) and [Neal \(1996\)](#).