

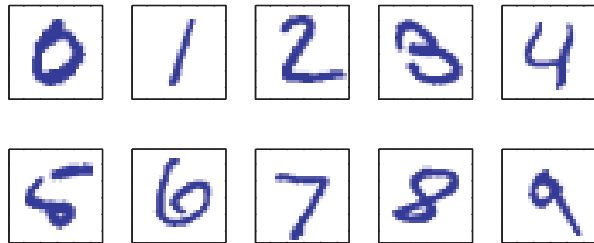
# 1

# Introduction

The problem of searching for patterns in data is a fundamental one and has a long and successful history. For instance, the extensive astronomical observations of Tycho Brahe in the 16<sup>th</sup> century allowed Johannes Kepler to discover the empirical laws of planetary motion, which in turn provided a springboard for the development of classical mechanics. Similarly, the discovery of regularities in atomic spectra played a key role in the development and verification of quantum physics in the early twentieth century. The field of pattern recognition is concerned with the automatic discovery of regularities in data through the use of computer algorithms and with the use of these regularities to take actions such as classifying the data into different categories.

Consider the example of recognizing handwritten digits, illustrated in Figure 1.1. Each digit corresponds to a  $28 \times 28$  pixel image and so can be represented by a vector  $\mathbf{x}$  comprising 784 real numbers. The goal is to build a machine that will take such a vector  $\mathbf{x}$  as input and that will produce the identity of the digit  $0, \dots, 9$  as the output. This is a nontrivial problem due to the wide variability of handwriting. It could be

**Figure 1.1** Examples of hand-written digits taken from US zip codes.



tackled using handcrafted rules or heuristics for distinguishing the digits based on the shapes of the strokes, but in practice such an approach leads to a proliferation of rules and of exceptions to the rules and so on, and invariably gives poor results.

Far better results can be obtained by adopting a machine learning approach in which a large set of  $N$  digits  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  called a *training set* is used to tune the parameters of an adaptive model. The categories of the digits in the training set are known in advance, typically by inspecting them individually and hand-labelling them. We can express the category of a digit using *target vector*  $\mathbf{t}$ , which represents the identity of the corresponding digit. Suitable techniques for representing categories in terms of vectors will be discussed later. Note that there is one such target vector  $\mathbf{t}$  for each digit image  $\mathbf{x}$ .

The result of running the machine learning algorithm can be expressed as a function  $\mathbf{y}(\mathbf{x})$  which takes a new digit image  $\mathbf{x}$  as input and that generates an output vector  $\mathbf{y}$ , encoded in the same way as the target vectors. The precise form of the function  $\mathbf{y}(\mathbf{x})$  is determined during the *training* phase, also known as the *learning* phase, on the basis of the training data. Once the model is trained it can then determine the identity of new digit images, which are said to comprise a *test set*. The ability to categorize correctly new examples that differ from those used for training is known as *generalization*. In practical applications, the variability of the input vectors will be such that the training data can comprise only a tiny fraction of all possible input vectors, and so generalization is a central goal in pattern recognition.

For most practical applications, the original input variables are typically *preprocessed* to transform them into some new space of variables where, it is hoped, the pattern recognition problem will be easier to solve. For instance, in the digit recognition problem, the images of the digits are typically translated and scaled so that each digit is contained within a box of a fixed size. This greatly reduces the variability within each digit class, because the location and scale of all the digits are now the same, which makes it much easier for a subsequent pattern recognition algorithm to distinguish between the different classes. This pre-processing stage is sometimes also called *feature extraction*. Note that new test data must be pre-processed using the same steps as the training data.

Pre-processing might also be performed in order to speed up computation. For example, if the goal is real-time face detection in a high-resolution video stream, the computer must handle huge numbers of pixels per second, and presenting these directly to a complex pattern recognition algorithm may be computationally infeasible. Instead, the aim is to find useful features that are fast to compute, and yet that

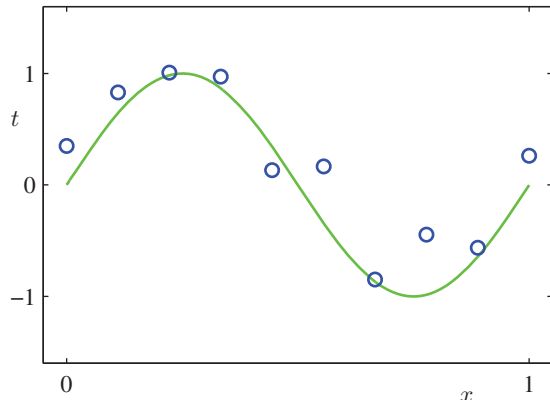
also preserve useful discriminatory information enabling faces to be distinguished from non-faces. These features are then used as the inputs to the pattern recognition algorithm. For instance, the average value of the image intensity over a rectangular subregion can be evaluated extremely efficiently (Viola and Jones, 2004), and a set of such features can prove very effective in fast face detection. Because the number of such features is smaller than the number of pixels, this kind of pre-processing represents a form of dimensionality reduction. Care must be taken during pre-processing because often information is discarded, and if this information is important to the solution of the problem then the overall accuracy of the system can suffer.

Applications in which the training data comprises examples of the input vectors along with their corresponding target vectors are known as *supervised learning* problems. Cases such as the digit recognition example, in which the aim is to assign each input vector to one of a finite number of discrete categories, are called *classification* problems. If the desired output consists of one or more continuous variables, then the task is called *regression*. An example of a regression problem would be the prediction of the yield in a chemical manufacturing process in which the inputs consist of the concentrations of reactants, the temperature, and the pressure.

In other pattern recognition problems, the training data consists of a set of input vectors  $\mathbf{x}$  without any corresponding target values. The goal in such *unsupervised learning* problems may be to discover groups of similar examples within the data, where it is called *clustering*, or to determine the distribution of data within the input space, known as *density estimation*, or to project the data from a high-dimensional space down to two or three dimensions for the purpose of *visualization*.

Finally, the technique of *reinforcement learning* (Sutton and Barto, 1998) is concerned with the problem of finding suitable actions to take in a given situation in order to maximize a reward. Here the learning algorithm is not given examples of optimal outputs, in contrast to supervised learning, but must instead discover them by a process of trial and error. Typically there is a sequence of states and actions in which the learning algorithm is interacting with its environment. In many cases, the current action not only affects the immediate reward but also has an impact on the reward at all subsequent time steps. For example, by using appropriate reinforcement learning techniques a neural network can learn to play the game of backgammon to a high standard (Tesauro, 1994). Here the network must learn to take a board position as input, along with the result of a dice throw, and produce a strong move as the output. This is done by having the network play against a copy of itself for perhaps a million games. A major challenge is that a game of backgammon can involve dozens of moves, and yet it is only at the end of the game that the reward, in the form of victory, is achieved. The reward must then be attributed appropriately to all of the moves that led to it, even though some moves will have been good ones and others less so. This is an example of a *credit assignment* problem. A general feature of reinforcement learning is the trade-off between *exploration*, in which the system tries out new kinds of actions to see how effective they are, and *exploitation*, in which the system makes use of actions that are known to yield a high reward. Too strong a focus on either exploration or exploitation will yield poor results. Reinforcement learning continues to be an active area of machine learning research. However, a

**Figure 1.2** Plot of a training data set of  $N = 10$  points, shown as blue circles, each comprising an observation of the input variable  $x$  along with the corresponding target variable  $t$ . The green curve shows the function  $\sin(2\pi x)$  used to generate the data. Our goal is to predict the value of  $t$  for some new value of  $x$ , without knowledge of the green curve.



detailed treatment lies beyond the scope of this book.

Although each of these tasks needs its own tools and techniques, many of the key ideas that underpin them are common to all such problems. One of the main goals of this chapter is to introduce, in a relatively informal way, several of the most important of these concepts and to illustrate them using simple examples. Later in the book we shall see these same ideas re-emerge in the context of more sophisticated models that are applicable to real-world pattern recognition applications. This chapter also provides a self-contained introduction to three important tools that will be used throughout the book, namely probability theory, decision theory, and information theory. Although these might sound like daunting topics, they are in fact straightforward, and a clear understanding of them is essential if machine learning techniques are to be used to best effect in practical applications.

## 1.1. Example: Polynomial Curve Fitting

We begin by introducing a simple regression problem, which we shall use as a running example throughout this chapter to motivate a number of key concepts. Suppose we observe a real-valued input variable  $x$  and we wish to use this observation to predict the value of a real-valued target variable  $t$ . For the present purposes, it is instructive to consider an artificial example using synthetically generated data because we then know the precise process that generated the data for comparison against any learned model. The data for this example is generated from the function  $\sin(2\pi x)$  with random noise included in the target values, as described in detail in Appendix A.

Now suppose that we are given a training set comprising  $N$  observations of  $x$ , written  $\mathbf{x} \equiv (x_1, \dots, x_N)^T$ , together with corresponding observations of the values of  $t$ , denoted  $\mathbf{t} \equiv (t_1, \dots, t_N)^T$ . Figure 1.2 shows a plot of a training set comprising  $N = 10$  data points. The input data set  $\mathbf{x}$  in Figure 1.2 was generated by choosing values of  $x_n$ , for  $n = 1, \dots, N$ , spaced uniformly in range  $[0, 1]$ , and the target data set  $\mathbf{t}$  was obtained by first computing the corresponding values of the function

$\sin(2\pi x)$  and then adding a small level of random noise having a Gaussian distribution (the Gaussian distribution is discussed in Section 1.2.4) to each such point in order to obtain the corresponding value  $t_n$ . By generating data in this way, we are capturing a property of many real data sets, namely that they possess an underlying regularity, which we wish to learn, but that individual observations are corrupted by random noise. This noise might arise from intrinsically stochastic (i.e. random) processes such as radioactive decay but more typically is due to there being sources of variability that are themselves unobserved.

Our goal is to exploit this training set in order to make predictions of the value  $\hat{t}$  of the target variable for some new value  $\hat{x}$  of the input variable. As we shall see later, this involves implicitly trying to discover the underlying function  $\sin(2\pi x)$ . This is intrinsically a difficult problem as we have to generalize from a finite data set. Furthermore the observed data are corrupted with noise, and so for a given  $\hat{x}$  there is uncertainty as to the appropriate value for  $\hat{t}$ . Probability theory, discussed in Section 1.2, provides a framework for expressing such uncertainty in a precise and quantitative manner, and decision theory, discussed in Section 1.5, allows us to exploit this probabilistic representation in order to make predictions that are optimal according to appropriate criteria.

For the moment, however, we shall proceed rather informally and consider a simple approach based on curve fitting. In particular, we shall fit the data using a polynomial function of the form

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j \quad (1.1)$$

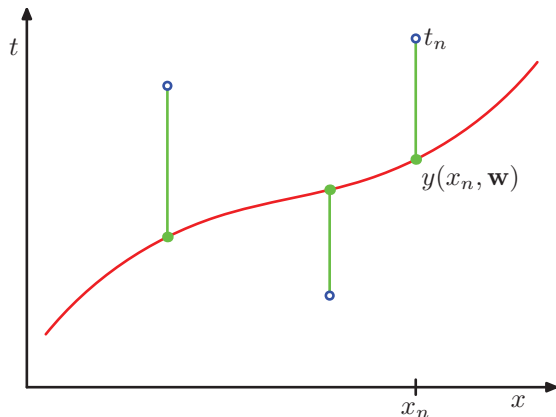
where  $M$  is the *order* of the polynomial, and  $x^j$  denotes  $x$  raised to the power of  $j$ . The polynomial coefficients  $w_0, \dots, w_M$  are collectively denoted by the vector  $\mathbf{w}$ . Note that, although the polynomial function  $y(x, \mathbf{w})$  is a nonlinear function of  $x$ , it is a linear function of the coefficients  $\mathbf{w}$ . Functions, such as the polynomial, which are linear in the unknown parameters have important properties and are called *linear models* and will be discussed extensively in Chapters 3 and 4.

The values of the coefficients will be determined by fitting the polynomial to the training data. This can be done by minimizing an *error function* that measures the misfit between the function  $y(x, \mathbf{w})$ , for any given value of  $\mathbf{w}$ , and the training set data points. One simple choice of error function, which is widely used, is given by the sum of the squares of the errors between the predictions  $y(x_n, \mathbf{w})$  for each data point  $x_n$  and the corresponding target values  $t_n$ , so that we minimize

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 \quad (1.2)$$

where the factor of  $1/2$  is included for later convenience. We shall discuss the motivation for this choice of error function later in this chapter. For the moment we simply note that it is a nonnegative quantity that would be zero if, and only if, the

**Figure 1.3** The error function (1.2) corresponds to (one half of) the sum of the squares of the displacements (shown by the vertical green bars) of each data point from the function  $y(x, \mathbf{w})$ .



function  $y(x, \mathbf{w})$  were to pass exactly through each training data point. The geometrical interpretation of the sum-of-squares error function is illustrated in Figure 1.3.

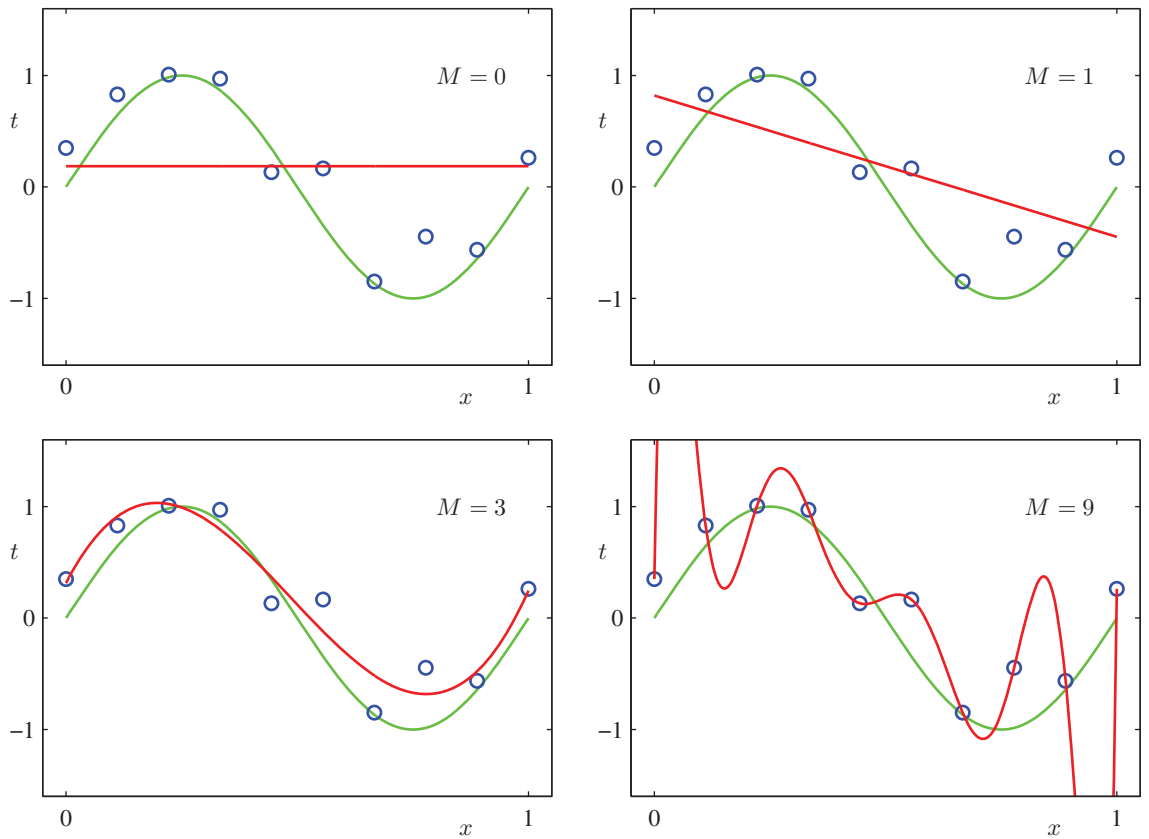
We can solve the curve fitting problem by choosing the value of  $\mathbf{w}$  for which  $E(\mathbf{w})$  is as small as possible. Because the error function is a quadratic function of the coefficients  $\mathbf{w}$ , its derivatives with respect to the coefficients will be linear in the elements of  $\mathbf{w}$ , and so the minimization of the error function has a unique solution, denoted by  $\mathbf{w}^*$ , which can be found in closed form. The resulting polynomial is given by the function  $y(x, \mathbf{w}^*)$ .

### Exercise 1.1

There remains the problem of choosing the order  $M$  of the polynomial, and as we shall see this will turn out to be an example of an important concept called *model comparison* or *model selection*. In Figure 1.4, we show four examples of the results of fitting polynomials having orders  $M = 0, 1, 3$ , and 9 to the data set shown in Figure 1.2.

We notice that the constant ( $M = 0$ ) and first order ( $M = 1$ ) polynomials give rather poor fits to the data and consequently rather poor representations of the function  $\sin(2\pi x)$ . The third order ( $M = 3$ ) polynomial seems to give the best fit to the function  $\sin(2\pi x)$  of the examples shown in Figure 1.4. When we go to a much higher order polynomial ( $M = 9$ ), we obtain an excellent fit to the training data. In fact, the polynomial passes exactly through each data point and  $E(\mathbf{w}^*) = 0$ . However, the fitted curve oscillates wildly and gives a very poor representation of the function  $\sin(2\pi x)$ . This latter behaviour is known as *over-fitting*.

As we have noted earlier, the goal is to achieve good generalization by making accurate predictions for new data. We can obtain some quantitative insight into the dependence of the generalization performance on  $M$  by considering a separate test set comprising 100 data points generated using exactly the same procedure used to generate the training set points but with new choices for the random noise values included in the target values. For each choice of  $M$ , we can then evaluate the residual value of  $E(\mathbf{w}^*)$  given by (1.2) for the training data, and we can also evaluate  $E(\mathbf{w}^*)$  for the test data set. It is sometimes more convenient to use the root-mean-square



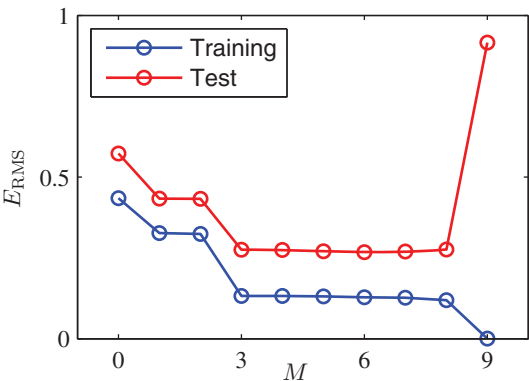
**Figure 1.4** Plots of polynomials having various orders  $M$ , shown as red curves, fitted to the data set shown in Figure 1.2.

(RMS) error defined by

$$E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N} \tag{1.3}$$

in which the division by  $N$  allows us to compare different sizes of data sets on an equal footing, and the square root ensures that  $E_{\text{RMS}}$  is measured on the same scale (and in the same units) as the target variable  $t$ . Graphs of the training and test set RMS errors are shown, for various values of  $M$ , in Figure 1.5. The test set error is a measure of how well we are doing in predicting the values of  $t$  for new data observations of  $x$ . We note from Figure 1.5 that small values of  $M$  give relatively large values of the test set error, and this can be attributed to the fact that the corresponding polynomials are rather inflexible and are incapable of capturing the oscillations in the function  $\sin(2\pi x)$ . Values of  $M$  in the range  $3 \leq M \leq 8$  give small values for the test set error, and these also give reasonable representations of the generating function  $\sin(2\pi x)$ , as can be seen, for the case of  $M = 3$ , from Figure 1.4.

**Figure 1.5** Graphs of the root-mean-square error, defined by (1.3), evaluated on the training set and on an independent test set for various values of  $M$ .



For  $M = 9$ , the training set error goes to zero, as we might expect because this polynomial contains 10 degrees of freedom corresponding to the 10 coefficients  $w_0, \dots, w_9$ , and so can be tuned exactly to the 10 data points in the training set. However, the test set error has become very large and, as we saw in Figure 1.4, the corresponding function  $y(x, \mathbf{w}^*)$  exhibits wild oscillations.

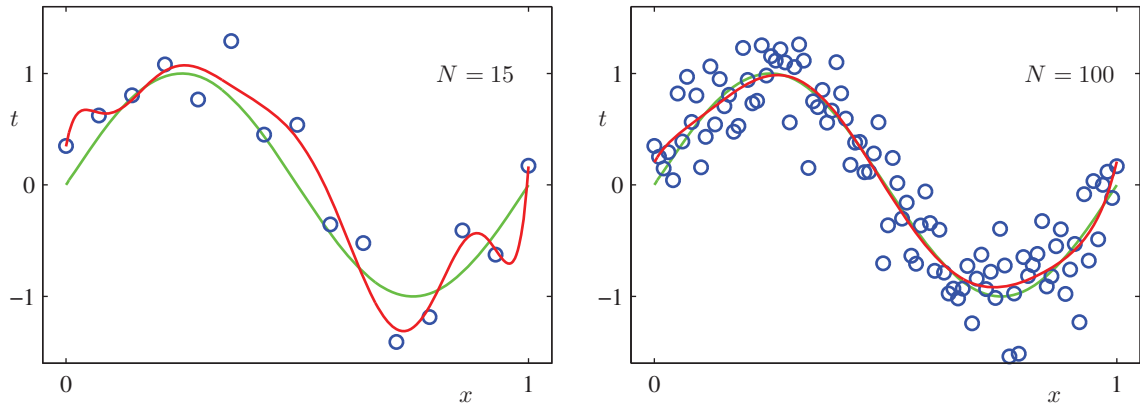
This may seem paradoxical because a polynomial of given order contains all lower order polynomials as special cases. The  $M = 9$  polynomial is therefore capable of generating results at least as good as the  $M = 3$  polynomial. Furthermore, we might suppose that the best predictor of new data would be the function  $\sin(2\pi x)$  from which the data was generated (and we shall see later that this is indeed the case). We know that a power series expansion of the function  $\sin(2\pi x)$  contains terms of all orders, so we might expect that results should improve monotonically as we increase  $M$ .

We can gain some insight into the problem by examining the values of the coefficients  $\mathbf{w}^*$  obtained from polynomials of various order, as shown in Table 1.1. We see that, as  $M$  increases, the magnitude of the coefficients typically gets larger. In particular for the  $M = 9$  polynomial, the coefficients have become finely tuned to the data by developing large positive and negative values so that the correspond-

**Table 1.1** Table of the coefficients  $\mathbf{w}^*$  for polynomials of various order. Observe how the typical magnitude of the coefficients increases dramatically as the order of the polynomial increases.

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
$w_0^*$	0.19	0.82	0.31	0.35
$w_1^*$		-1.27	7.99	232.37
$w_2^*$			-25.43	-5321.83
$w_3^*$			17.37	48568.31
$w_4^*$				-231639.30
$w_5^*$				640042.26
$w_6^*$				-1061800.52
$w_7^*$				1042400.18
$w_8^*$				-557682.99
$w_9^*$				125201.43





**Figure 1.6** Plots of the solutions obtained by minimizing the sum-of-squares error function using the  $M = 9$  polynomial for  $N = 15$  data points (left plot) and  $N = 100$  data points (right plot). We see that increasing the size of the data set reduces the over-fitting problem.

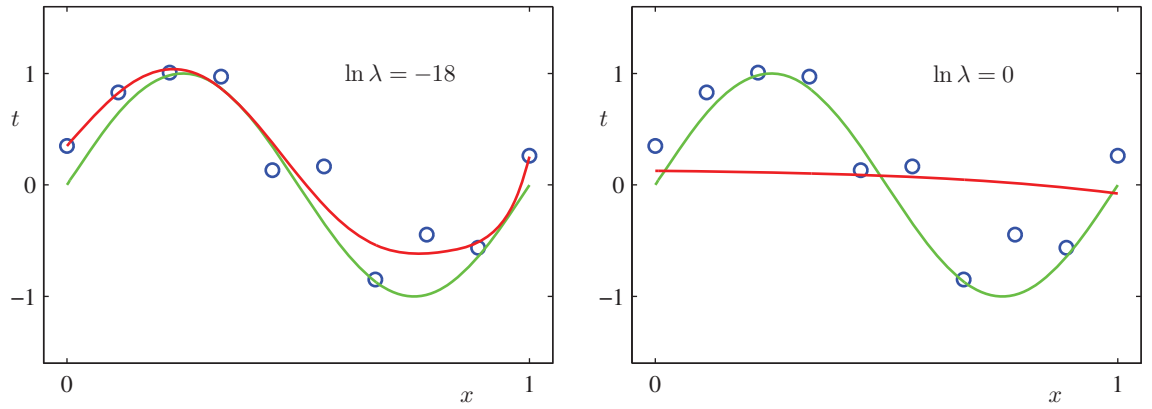
ing polynomial function matches each of the data points exactly, but between data points (particularly near the ends of the range) the function exhibits the large oscillations observed in Figure 1.4. Intuitively, what is happening is that the more flexible polynomials with larger values of  $M$  are becoming increasingly tuned to the random noise on the target values.

It is also interesting to examine the behaviour of a given model as the size of the data set is varied, as shown in Figure 1.6. We see that, for a given model complexity, the over-fitting problem become less severe as the size of the data set increases. Another way to say this is that the larger the data set, the more complex (in other words more flexible) the model that we can afford to fit to the data. One rough heuristic that is sometimes advocated is that the number of data points should be no less than some multiple (say 5 or 10) of the number of adaptive parameters in the model. However, as we shall see in Chapter 3, the number of parameters is not necessarily the most appropriate measure of model complexity.

Also, there is something rather unsatisfying about having to limit the number of parameters in a model according to the size of the available training set. It would seem more reasonable to choose the complexity of the model according to the complexity of the problem being solved. We shall see that the least squares approach to finding the model parameters represents a specific case of *maximum likelihood* (discussed in Section 1.2.5), and that the over-fitting problem can be understood as a general property of maximum likelihood. By adopting a *Bayesian* approach, the over-fitting problem can be avoided. We shall see that there is no difficulty from a Bayesian perspective in employing models for which the number of parameters greatly exceeds the number of data points. Indeed, in a Bayesian model the *effective* number of parameters adapts automatically to the size of the data set.

For the moment, however, it is instructive to continue with the current approach and to consider how in practice we can apply it to data sets of limited size where we

Section 3.4



**Figure 1.7** Plots of  $M = 9$  polynomials fitted to the data set shown in Figure 1.2 using the regularized error function (1.4) for two values of the regularization parameter  $\lambda$  corresponding to  $\ln \lambda = -18$  and  $\ln \lambda = 0$ . The case of no regularizer, i.e.,  $\lambda = 0$ , corresponding to  $\ln \lambda = -\infty$ , is shown at the bottom right of Figure 1.4.

may wish to use relatively complex and flexible models. One technique that is often used to control the over-fitting phenomenon in such cases is that of *regularization*, which involves adding a penalty term to the error function (1.2) in order to discourage the coefficients from reaching large values. The simplest such penalty term takes the form of a sum of squares of all of the coefficients, leading to a modified error function of the form

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad (1.4)$$

where  $\|\mathbf{w}\|^2 \equiv \mathbf{w}^T \mathbf{w} = w_0^2 + w_1^2 + \dots + w_M^2$ , and the coefficient  $\lambda$  governs the relative importance of the regularization term compared with the sum-of-squares error term. Note that often the coefficient  $w_0$  is omitted from the regularizer because its inclusion causes the results to depend on the choice of origin for the target variable (Hastie *et al.*, 2001), or it may be included but with its own regularization coefficient (we shall discuss this topic in more detail in Section 5.5.1). Again, the error function in (1.4) can be minimized exactly in closed form. Techniques such as this are known in the statistics literature as *shrinkage* methods because they reduce the value of the coefficients. The particular case of a quadratic regularizer is called *ridge regression* (Hoerl and Kennard, 1970). In the context of neural networks, this approach is known as *weight decay*.

Figure 1.7 shows the results of fitting the polynomial of order  $M = 9$  to the same data set as before but now using the regularized error function given by (1.4). We see that, for a value of  $\ln \lambda = -18$ , the over-fitting has been suppressed and we now obtain a much closer representation of the underlying function  $\sin(2\pi x)$ . If, however, we use too large a value for  $\lambda$  then we again obtain a poor fit, as shown in Figure 1.7 for  $\ln \lambda = 0$ . The corresponding coefficients from the fitted polynomials are given in Table 1.2, showing that regularization has the desired effect of reducing

**Table 1.2** Table of the coefficients  $w^*$  for  $M = 9$  polynomials with various values for the regularization parameter  $\lambda$ . Note that  $\ln \lambda = -\infty$  corresponds to a model with no regularization, i.e., to the graph at the bottom right in Figure 1.4. We see that, as the value of  $\lambda$  increases, the typical magnitude of the coefficients gets smaller.

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
$w_0^*$	0.35	0.35	0.13
$w_1^*$	232.37	4.74	-0.05
$w_2^*$	-5321.83	-0.77	-0.06
$w_3^*$	48568.31	-31.97	-0.05
$w_4^*$	-231639.30	-3.89	-0.03
$w_5^*$	640042.26	55.28	-0.02
$w_6^*$	-1061800.52	41.32	-0.01
$w_7^*$	1042400.18	-45.95	-0.00
$w_8^*$	-557682.99	-91.53	0.00
$w_9^*$	125201.43	72.68	0.01

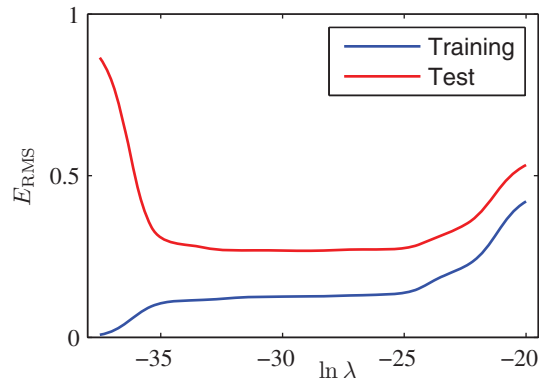
the magnitude of the coefficients.

The impact of the regularization term on the generalization error can be seen by plotting the value of the RMS error (1.3) for both training and test sets against  $\ln \lambda$ , as shown in Figure 1.8. We see that in effect  $\lambda$  now controls the effective complexity of the model and hence determines the degree of over-fitting.

The issue of model complexity is an important one and will be discussed at length in Section 1.3. Here we simply note that, if we were trying to solve a practical application using this approach of minimizing an error function, we would have to find a way to determine a suitable value for the model complexity. The results above suggest a simple way of achieving this, namely by taking the available data and partitioning it into a training set, used to determine the coefficients  $w$ , and a separate *validation* set, also called a *hold-out* set, used to optimize the model complexity (either  $M$  or  $\lambda$ ). In many cases, however, this will prove to be too wasteful of valuable training data, and we have to seek more sophisticated approaches.

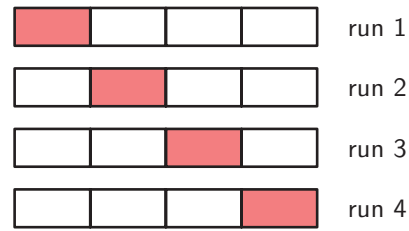
So far our discussion of polynomial curve fitting has appealed largely to intuition. We now seek a more principled approach to solving problems in pattern recognition by turning to a discussion of probability theory. As well as providing the foundation for nearly all of the subsequent developments in this book, it will also

**Figure 1.8** Graph of the root-mean-square error (1.3) versus  $\ln \lambda$  for the  $M = 9$  polynomial.



**Figure 1.18**

The technique of  $S$ -fold cross-validation, illustrated here for the case of  $S = 4$ , involves taking the available data and partitioning it into  $S$  groups (in the simplest case these are of equal size). Then  $S - 1$  of the groups are used to train a set of models that are then evaluated on the remaining group. This procedure is then repeated for all  $S$  possible choices for the held-out group, indicated here by the red blocks, and the performance scores from the  $S$  runs are then averaged.



data to assess performance. When data is particularly scarce, it may be appropriate to consider the case  $S = N$ , where  $N$  is the total number of data points, which gives the *leave-one-out* technique.

One major drawback of cross-validation is that the number of training runs that must be performed is increased by a factor of  $S$ , and this can prove problematic for models in which the training is itself computationally expensive. A further problem with techniques such as cross-validation that use separate data to assess performance is that we might have multiple complexity parameters for a single model (for instance, there might be several regularization parameters). Exploring combinations of settings for such parameters could, in the worst case, require a number of training runs that is exponential in the number of parameters. Clearly, we need a better approach. Ideally, this should rely only on the training data and should allow multiple hyperparameters and model types to be compared in a single training run. We therefore need to find a measure of performance which depends only on the training data and which does not suffer from bias due to over-fitting.

Historically various ‘information criteria’ have been proposed that attempt to correct for the bias of maximum likelihood by the addition of a penalty term to compensate for the over-fitting of more complex models. For example, the *Akaike information criterion*, or AIC (Akaike, 1974), chooses the model for which the quantity

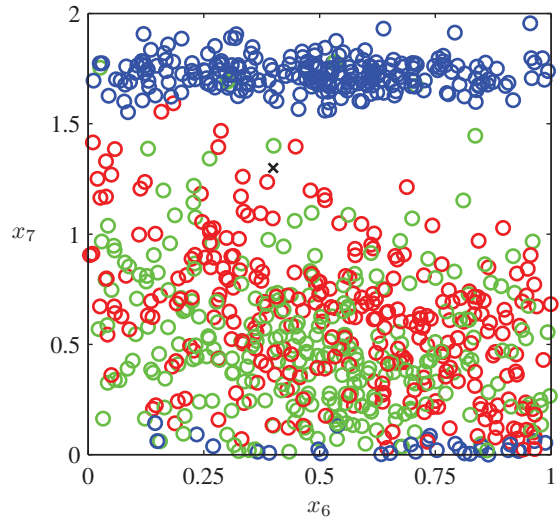
$$\ln p(\mathcal{D}|\mathbf{w}_{\text{ML}}) - M \quad (1.73)$$

is largest. Here  $p(\mathcal{D}|\mathbf{w}_{\text{ML}})$  is the best-fit log likelihood, and  $M$  is the number of adjustable parameters in the model. A variant of this quantity, called the *Bayesian information criterion*, or BIC, will be discussed in Section 4.4.1. Such criteria do not take account of the uncertainty in the model parameters, however, and in practice they tend to favour overly simple models. We therefore turn in Section 3.4 to a fully Bayesian approach where we shall see how complexity penalties arise in a natural and principled way.

## 1.4. The Curse of Dimensionality

In the polynomial curve fitting example we had just one input variable  $x$ . For practical applications of pattern recognition, however, we will have to deal with spaces

**Figure 1.19** Scatter plot of the oil flow data for input variables  $x_6$  and  $x_7$ , in which red denotes the ‘homogenous’ class, green denotes the ‘annular’ class, and blue denotes the ‘laminar’ class. Our goal is to classify the new test point denoted by ‘x’.

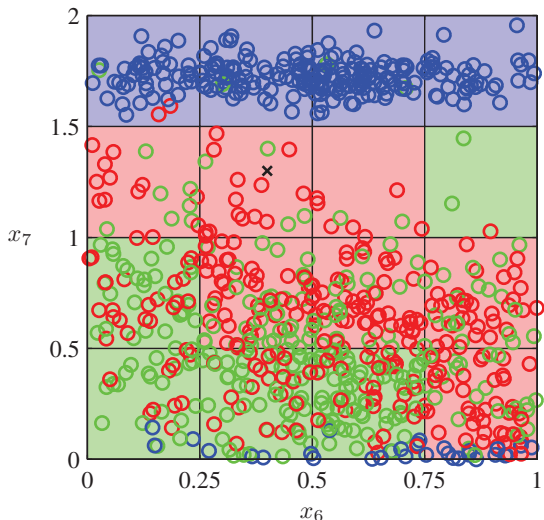


of high dimensionality comprising many input variables. As we now discuss, this poses some serious challenges and is an important factor influencing the design of pattern recognition techniques.

In order to illustrate the problem we consider a synthetically generated data set representing measurements taken from a pipeline containing a mixture of oil, water, and gas (Bishop and James, 1993). These three materials can be present in one of three different geometrical configurations known as ‘homogenous’, ‘annular’, and ‘laminar’, and the fractions of the three materials can also vary. Each data point comprises a 12-dimensional input vector consisting of measurements taken with gamma ray densitometers that measure the attenuation of gamma rays passing along narrow beams through the pipe. This data set is described in detail in Appendix A. Figure 1.19 shows 100 points from this data set on a plot showing two of the measurements  $x_6$  and  $x_7$  (the remaining ten input values are ignored for the purposes of this illustration). Each data point is labelled according to which of the three geometrical classes it belongs to, and our goal is to use this data as a training set in order to be able to classify a new observation  $(x_6, x_7)$ , such as the one denoted by the cross in Figure 1.19. We observe that the cross is surrounded by numerous red points, and so we might suppose that it belongs to the red class. However, there are also plenty of green points nearby, so we might think that it could instead belong to the green class. It seems unlikely that it belongs to the blue class. The intuition here is that the identity of the cross should be determined more strongly by nearby points from the training set and less strongly by more distant points. In fact, this intuition turns out to be reasonable and will be discussed more fully in later chapters.

How can we turn this intuition into a learning algorithm? One very simple approach would be to divide the input space into regular cells, as indicated in Figure 1.20. When we are given a test point and we wish to predict its class, we first decide which cell it belongs to, and we then find all of the training data points that

**Figure 1.20** Illustration of a simple approach to the solution of a classification problem in which the input space is divided into cells and any new test point is assigned to the class that has a majority number of representatives in the same cell as the test point. As we shall see shortly, this simplistic approach has some severe shortcomings.



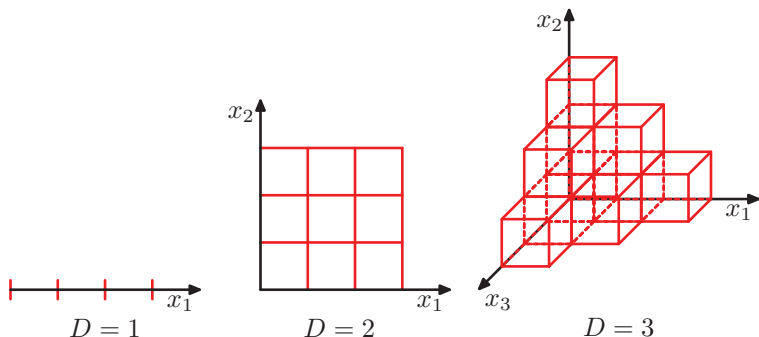
fall in the same cell. The identity of the test point is predicted as being the same as the class having the largest number of training points in the same cell as the test point (with ties being broken at random).

There are numerous problems with this naive approach, but one of the most severe becomes apparent when we consider its extension to problems having larger numbers of input variables, corresponding to input spaces of higher dimensionality. The origin of the problem is illustrated in Figure 1.21, which shows that, if we divide a region of a space into regular cells, then the number of such cells grows exponentially with the dimensionality of the space. The problem with an exponentially large number of cells is that we would need an exponentially large quantity of training data in order to ensure that the cells are not empty. Clearly, we have no hope of applying such a technique in a space of more than a few variables, and so we need to find a more sophisticated approach.

We can gain further insight into the problems of high-dimensional spaces by returning to the example of polynomial curve fitting and considering how we would

## Section 1.1

**Figure 1.21** Illustration of the curse of dimensionality, showing how the number of regions of a regular grid grows exponentially with the dimensionality  $D$  of the space. For clarity, only a subset of the cubical regions are shown for  $D = 3$ .



extend this approach to deal with input spaces having several variables. If we have  $D$  input variables, then a general polynomial with coefficients up to order 3 would take the form

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^D w_i x_i + \sum_{i=1}^D \sum_{j=1}^D w_{ij} x_i x_j + \sum_{i=1}^D \sum_{j=1}^D \sum_{k=1}^D w_{ijk} x_i x_j x_k. \quad (1.74)$$

As  $D$  increases, so the number of independent coefficients (not all of the coefficients are independent due to interchange symmetries amongst the  $x$  variables) grows proportionally to  $D^3$ . In practice, to capture complex dependencies in the data, we may need to use a higher-order polynomial. For a polynomial of order  $M$ , the growth in the number of coefficients is like  $D^M$ . Although this is now a power law growth, rather than an exponential growth, it still points to the method becoming rapidly unwieldy and of limited practical utility.

*Exercise 1.16*

Our geometrical intuitions, formed through a life spent in a space of three dimensions, can fail badly when we consider spaces of higher dimensionality. As a simple example, consider a sphere of radius  $r = 1$  in a space of  $D$  dimensions, and ask what is the fraction of the volume of the sphere that lies between radius  $r = 1 - \epsilon$  and  $r = 1$ . We can evaluate this fraction by noting that the volume of a sphere of radius  $r$  in  $D$  dimensions must scale as  $r^D$ , and so we write

$$V_D(r) = K_D r^D \quad (1.75)$$

*Exercise 1.18*

where the constant  $K_D$  depends only on  $D$ . Thus the required fraction is given by

$$\frac{V_D(1) - V_D(1 - \epsilon)}{V_D(1)} = 1 - (1 - \epsilon)^D \quad (1.76)$$

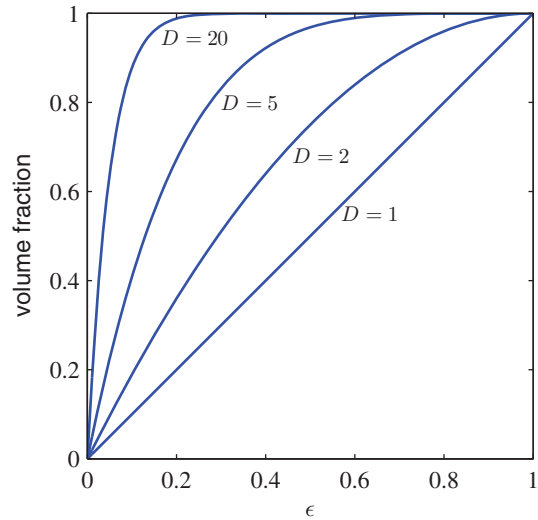
which is plotted as a function of  $\epsilon$  for various values of  $D$  in Figure 1.22. We see that, for large  $D$ , this fraction tends to 1 even for small values of  $\epsilon$ . Thus, in spaces of high dimensionality, most of the volume of a sphere is concentrated in a thin shell near the surface!

As a further example, of direct relevance to pattern recognition, consider the behaviour of a Gaussian distribution in a high-dimensional space. If we transform from Cartesian to polar coordinates, and then integrate out the directional variables, we obtain an expression for the density  $p(r)$  as a function of radius  $r$  from the origin. Thus  $p(r)\delta r$  is the probability mass inside a thin shell of thickness  $\delta r$  located at radius  $r$ . This distribution is plotted, for various values of  $D$ , in Figure 1.23, and we see that for large  $D$  the probability mass of the Gaussian is concentrated in a thin shell.

*Exercise 1.20*

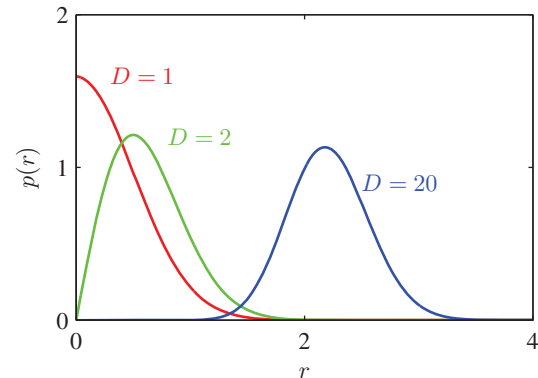
The severe difficulty that can arise in spaces of many dimensions is sometimes called the *curse of dimensionality* (Bellman, 1961). In this book, we shall make extensive use of illustrative examples involving input spaces of one or two dimensions, because this makes it particularly easy to illustrate the techniques graphically. The reader should be warned, however, that not all intuitions developed in spaces of low dimensionality will generalize to spaces of many dimensions.

**Figure 1.22** Plot of the fraction of the volume of a sphere lying in the range  $r = 1 - \epsilon$  to  $r = 1$  for various values of the dimensionality  $D$ .



Although the curse of dimensionality certainly raises important issues for pattern recognition applications, it does not prevent us from finding effective techniques applicable to high-dimensional spaces. The reasons for this are twofold. First, real data will often be confined to a region of the space having lower effective dimensionality, and in particular the directions over which important variations in the target variables occur may be so confined. Second, real data will typically exhibit some smoothness properties (at least locally) so that for the most part small changes in the input variables will produce small changes in the target variables, and so we can exploit local interpolation-like techniques to allow us to make predictions of the target variables for new values of the input variables. Successful pattern recognition techniques exploit one or both of these properties. Consider, for example, an application in manufacturing in which images are captured of identical planar objects on a conveyor belt, in which the goal is to determine their orientation. Each image is a point

**Figure 1.23** Plot of the probability density with respect to radius  $r$  of a Gaussian distribution for various values of the dimensionality  $D$ . In a high-dimensional space, most of the probability mass of a Gaussian is located within a thin shell at a specific radius.





in a high-dimensional space whose dimensionality is determined by the number of pixels. Because the objects can occur at different positions within the image and in different orientations, there are three degrees of freedom of variability between images, and a set of images will live on a three dimensional *manifold* embedded within the high-dimensional space. Due to the complex relationships between the object position or orientation and the pixel intensities, this manifold will be highly nonlinear. If the goal is to learn a model that can take an input image and output the orientation of the object irrespective of its position, then there is only one degree of freedom of variability within the manifold that is significant.

## 1.5. Decision Theory

---

We have seen in Section 1.2 how probability theory provides us with a consistent mathematical framework for quantifying and manipulating uncertainty. Here we turn to a discussion of decision theory that, when combined with probability theory, allows us to make optimal decisions in situations involving uncertainty such as those encountered in pattern recognition.

Suppose we have an input vector  $\mathbf{x}$  together with a corresponding vector  $\mathbf{t}$  of target variables, and our goal is to predict  $\mathbf{t}$  given a new value for  $\mathbf{x}$ . For regression problems,  $\mathbf{t}$  will comprise continuous variables, whereas for classification problems  $\mathbf{t}$  will represent class labels. The joint probability distribution  $p(\mathbf{x}, \mathbf{t})$  provides a complete summary of the uncertainty associated with these variables. Determination of  $p(\mathbf{x}, \mathbf{t})$  from a set of training data is an example of *inference* and is typically a very difficult problem whose solution forms the subject of much of this book. In a practical application, however, we must often make a specific prediction for the value of  $\mathbf{t}$ , or more generally take a specific action based on our understanding of the values  $\mathbf{t}$  is likely to take, and this aspect is the subject of decision theory.

Consider, for example, a medical diagnosis problem in which we have taken an X-ray image of a patient, and we wish to determine whether the patient has cancer or not. In this case, the input vector  $\mathbf{x}$  is the set of pixel intensities in the image, and output variable  $t$  will represent the presence of cancer, which we denote by the class  $C_1$ , or the absence of cancer, which we denote by the class  $C_2$ . We might, for instance, choose  $t$  to be a binary variable such that  $t = 0$  corresponds to class  $C_1$  and  $t = 1$  corresponds to class  $C_2$ . We shall see later that this choice of label values is particularly convenient for probabilistic models. The general inference problem then involves determining the joint distribution  $p(\mathbf{x}, C_k)$ , or equivalently  $p(\mathbf{x}, t)$ , which gives us the most complete probabilistic description of the situation. Although this can be a very useful and informative quantity, in the end we must decide either to give treatment to the patient or not, and we would like this choice to be optimal in some appropriate sense (Duda and Hart, 1973). This is the *decision* step, and it is the subject of decision theory to tell us how to make optimal decisions given the appropriate probabilities. We shall see that the decision stage is generally very simple, even trivial, once we have solved the inference problem.

Here we give an introduction to the key ideas of decision theory as required for

the rest of the book. Further background, as well as more detailed accounts, can be found in Berger (1985) and Bather (2000).

Before giving a more detailed analysis, let us first consider informally how we might expect probabilities to play a role in making decisions. When we obtain the X-ray image  $\mathbf{x}$  for a new patient, our goal is to decide which of the two classes to assign to the image. We are interested in the probabilities of the two classes given the image, which are given by  $p(C_k|\mathbf{x})$ . Using Bayes' theorem, these probabilities can be expressed in the form

$$p(C_k|\mathbf{x}) = \frac{p(\mathbf{x}|C_k)p(C_k)}{p(\mathbf{x})}. \quad (1.77)$$

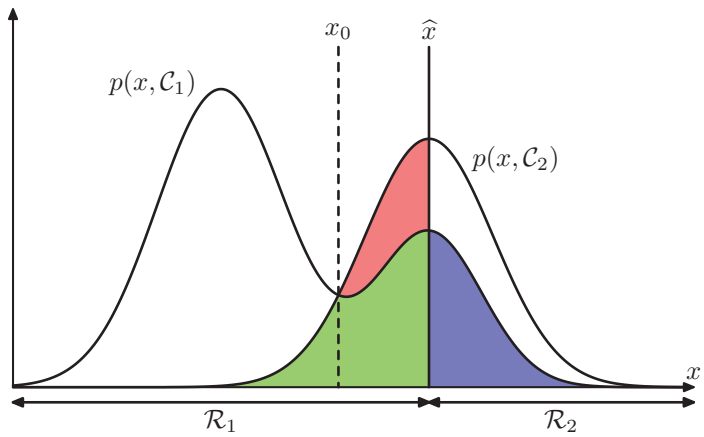
Note that any of the quantities appearing in Bayes' theorem can be obtained from the joint distribution  $p(\mathbf{x}, C_k)$  by either marginalizing or conditioning with respect to the appropriate variables. We can now interpret  $p(C_k)$  as the prior probability for the class  $C_k$ , and  $p(C_k|\mathbf{x})$  as the corresponding posterior probability. Thus  $p(C_1)$  represents the probability that a person has cancer, before we take the X-ray measurement. Similarly,  $p(C_1|\mathbf{x})$  is the corresponding probability, revised using Bayes' theorem in light of the information contained in the X-ray. If our aim is to minimize the chance of assigning  $\mathbf{x}$  to the wrong class, then intuitively we would choose the class having the higher posterior probability. We now show that this intuition is correct, and we also discuss more general criteria for making decisions.

### 1.5.1 Minimizing the misclassification rate

Suppose that our goal is simply to make as few misclassifications as possible. We need a rule that assigns each value of  $\mathbf{x}$  to one of the available classes. Such a rule will divide the input space into regions  $\mathcal{R}_k$  called *decision regions*, one for each class, such that all points in  $\mathcal{R}_k$  are assigned to class  $C_k$ . The boundaries between decision regions are called *decision boundaries* or *decision surfaces*. Note that each decision region need not be contiguous but could comprise some number of disjoint regions. We shall encounter examples of decision boundaries and decision regions in later chapters. In order to find the optimal decision rule, consider first of all the case of two classes, as in the cancer problem for instance. A mistake occurs when an input vector belonging to class  $C_1$  is assigned to class  $C_2$  or vice versa. The probability of this occurring is given by

$$\begin{aligned} p(\text{mistake}) &= p(\mathbf{x} \in \mathcal{R}_1, C_2) + p(\mathbf{x} \in \mathcal{R}_2, C_1) \\ &= \int_{\mathcal{R}_1} p(\mathbf{x}, C_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, C_1) d\mathbf{x}. \end{aligned} \quad (1.78)$$

We are free to choose the decision rule that assigns each point  $\mathbf{x}$  to one of the two classes. Clearly to minimize  $p(\text{mistake})$  we should arrange that each  $\mathbf{x}$  is assigned to whichever class has the smaller value of the integrand in (1.78). Thus, if  $p(\mathbf{x}, C_1) > p(\mathbf{x}, C_2)$  for a given value of  $\mathbf{x}$ , then we should assign that  $\mathbf{x}$  to class  $C_1$ . From the product rule of probability we have  $p(\mathbf{x}, C_k) = p(C_k|\mathbf{x})p(\mathbf{x})$ . Because the factor  $p(\mathbf{x})$  is common to both terms, we can restate this result as saying that the minimum



**Figure 1.24** Schematic illustration of the joint probabilities  $p(x, \mathcal{C}_k)$  for each of two classes plotted against  $x$ , together with the decision boundary  $x = \hat{x}$ . Values of  $x \geq \hat{x}$  are classified as class  $\mathcal{C}_2$  and hence belong to decision region  $\mathcal{R}_2$ , whereas points  $x < \hat{x}$  are classified as  $\mathcal{C}_1$  and belong to  $\mathcal{R}_1$ . Errors arise from the blue, green, and red regions, so that for  $x < \hat{x}$  the errors are due to points from class  $\mathcal{C}_2$  being misclassified as  $\mathcal{C}_1$  (represented by the sum of the red and green regions), and conversely for points in the region  $x \geq \hat{x}$  the errors are due to points from class  $\mathcal{C}_1$  being misclassified as  $\mathcal{C}_2$  (represented by the blue region). As we vary the location  $\hat{x}$  of the decision boundary, the combined areas of the blue and green regions remains constant, whereas the size of the red region varies. The optimal choice for  $\hat{x}$  is where the curves for  $p(x, \mathcal{C}_1)$  and  $p(x, \mathcal{C}_2)$  cross, corresponding to  $\hat{x} = x_0$ , because in this case the red region disappears. This is equivalent to the minimum misclassification rate decision rule, which assigns each value of  $x$  to the class having the higher posterior probability  $p(\mathcal{C}_k|x)$ .

probability of making a mistake is obtained if each value of  $\mathbf{x}$  is assigned to the class for which the posterior probability  $p(\mathcal{C}_k|\mathbf{x})$  is largest. This result is illustrated for two classes, and a single input variable  $x$ , in Figure 1.24.

For the more general case of  $K$  classes, it is slightly easier to maximize the probability of being correct, which is given by

$$\begin{aligned}
 p(\text{correct}) &= \sum_{k=1}^K p(\mathbf{x} \in \mathcal{R}_k, \mathcal{C}_k) \\
 &= \sum_{k=1}^K \int_{\mathcal{R}_k} p(\mathbf{x}, \mathcal{C}_k) \, d\mathbf{x}
 \end{aligned} \tag{1.79}$$

which is maximized when the regions  $\mathcal{R}_k$  are chosen such that each  $\mathbf{x}$  is assigned to the class for which  $p(\mathbf{x}, \mathcal{C}_k)$  is largest. Again, using the product rule  $p(\mathbf{x}, \mathcal{C}_k) = p(\mathcal{C}_k|\mathbf{x})p(\mathbf{x})$ , and noting that the factor of  $p(\mathbf{x})$  is common to all terms, we see that each  $\mathbf{x}$  should be assigned to the class having the largest posterior probability  $p(\mathcal{C}_k|\mathbf{x})$ .

**Figure 1.25** An example of a loss matrix with elements  $L_{kj}$  for the cancer treatment problem. The rows correspond to the true class, whereas the columns correspond to the assignment of class made by our decision criterion.

	cancer	normal
cancer	0	1000
normal	1	0

### 1.5.2 Minimizing the expected loss

For many applications, our objective will be more complex than simply minimizing the number of misclassifications. Let us consider again the medical diagnosis problem. We note that, if a patient who does not have cancer is incorrectly diagnosed as having cancer, the consequences may be some patient distress plus the need for further investigations. Conversely, if a patient with cancer is diagnosed as healthy, the result may be premature death due to lack of treatment. Thus the consequences of these two types of mistake can be dramatically different. It would clearly be better to make fewer mistakes of the second kind, even if this was at the expense of making more mistakes of the first kind.

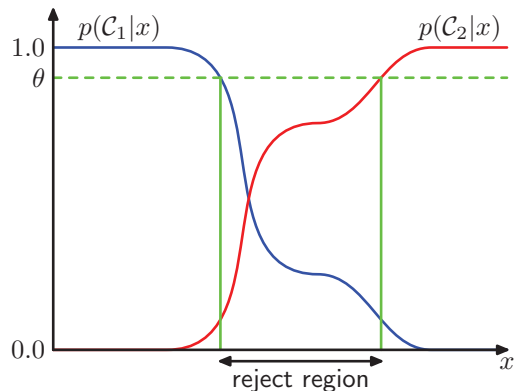
We can formalize such issues through the introduction of a *loss function*, also called a *cost function*, which is a single, overall measure of loss incurred in taking any of the available decisions or actions. Our goal is then to minimize the total loss incurred. Note that some authors consider instead a *utility function*, whose value they aim to maximize. These are equivalent concepts if we take the utility to be simply the negative of the loss, and throughout this text we shall use the loss function convention. Suppose that, for a new value of  $\mathbf{x}$ , the true class is  $\mathcal{C}_k$  and that we assign  $\mathbf{x}$  to class  $\mathcal{C}_j$  (where  $j$  may or may not be equal to  $k$ ). In so doing, we incur some level of loss that we denote by  $L_{kj}$ , which we can view as the  $k, j$  element of a *loss matrix*. For instance, in our cancer example, we might have a loss matrix of the form shown in Figure 1.25. This particular loss matrix says that there is no loss incurred if the correct decision is made, there is a loss of 1 if a healthy patient is diagnosed as having cancer, whereas there is a loss of 1000 if a patient having cancer is diagnosed as healthy.

The optimal solution is the one which minimizes the loss function. However, the loss function depends on the true class, which is unknown. For a given input vector  $\mathbf{x}$ , our uncertainty in the true class is expressed through the joint probability distribution  $p(\mathbf{x}, \mathcal{C}_k)$  and so we seek instead to minimize the average loss, where the average is computed with respect to this distribution, which is given by

$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x}. \tag{1.80}$$

Each  $\mathbf{x}$  can be assigned independently to one of the decision regions  $\mathcal{R}_j$ . Our goal is to choose the regions  $\mathcal{R}_j$  in order to minimize the expected loss (1.80), which implies that for each  $\mathbf{x}$  we should minimize  $\sum_k L_{kj} p(\mathbf{x}, \mathcal{C}_k)$ . As before, we can use the product rule  $p(\mathbf{x}, \mathcal{C}_k) = p(\mathcal{C}_k|\mathbf{x})p(\mathbf{x})$  to eliminate the common factor of  $p(\mathbf{x})$ . Thus the decision rule that minimizes the expected loss is the one that assigns each

**Figure 1.26** Illustration of the reject option. Inputs  $x$  such that the larger of the two posterior probabilities is less than or equal to some threshold  $\theta$  will be rejected.



new  $\mathbf{x}$  to the class  $j$  for which the quantity

$$\sum_k L_{kj} p(C_k | \mathbf{x}) \quad (1.81)$$

is a minimum. This is clearly trivial to do, once we know the posterior class probabilities  $p(C_k | \mathbf{x})$ .

### 1.5.3 The reject option

We have seen that classification errors arise from the regions of input space where the largest of the posterior probabilities  $p(C_k | \mathbf{x})$  is significantly less than unity, or equivalently where the joint distributions  $p(\mathbf{x}, C_k)$  have comparable values. These are the regions where we are relatively uncertain about class membership. In some applications, it will be appropriate to avoid making decisions on the difficult cases in anticipation of a lower error rate on those examples for which a classification decision is made. This is known as the *reject option*. For example, in our hypothetical medical illustration, it may be appropriate to use an automatic system to classify those X-ray images for which there is little doubt as to the correct class, while leaving a human expert to classify the more ambiguous cases. We can achieve this by introducing a threshold  $\theta$  and rejecting those inputs  $\mathbf{x}$  for which the largest of the posterior probabilities  $p(C_k | \mathbf{x})$  is less than or equal to  $\theta$ . This is illustrated for the case of two classes, and a single continuous input variable  $x$ , in Figure 1.26. Note that setting  $\theta = 1$  will ensure that all examples are rejected, whereas if there are  $K$  classes then setting  $\theta < 1/K$  will ensure that no examples are rejected. Thus the fraction of examples that get rejected is controlled by the value of  $\theta$ .

We can easily extend the reject criterion to minimize the expected loss, when a loss matrix is given, taking account of the loss incurred when a reject decision is made.

*Exercise 1.24*

### 1.5.4 Inference and decision

We have broken the classification problem down into two separate stages, the *inference stage* in which we use training data to learn a model for  $p(C_k | \mathbf{x})$ , and the

subsequent *decision* stage in which we use these posterior probabilities to make optimal class assignments. An alternative possibility would be to solve both problems together and simply learn a function that maps inputs  $\mathbf{x}$  directly into decisions. Such a function is called a *discriminant function*.

In fact, we can identify three distinct approaches to solving decision problems, all of which have been used in practical applications. These are given, in decreasing order of complexity, by:

- (a) First solve the inference problem of determining the class-conditional densities  $p(\mathbf{x}|\mathcal{C}_k)$  for each class  $\mathcal{C}_k$  individually. Also separately infer the prior class probabilities  $p(\mathcal{C}_k)$ . Then use Bayes' theorem in the form

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})} \quad (1.82)$$

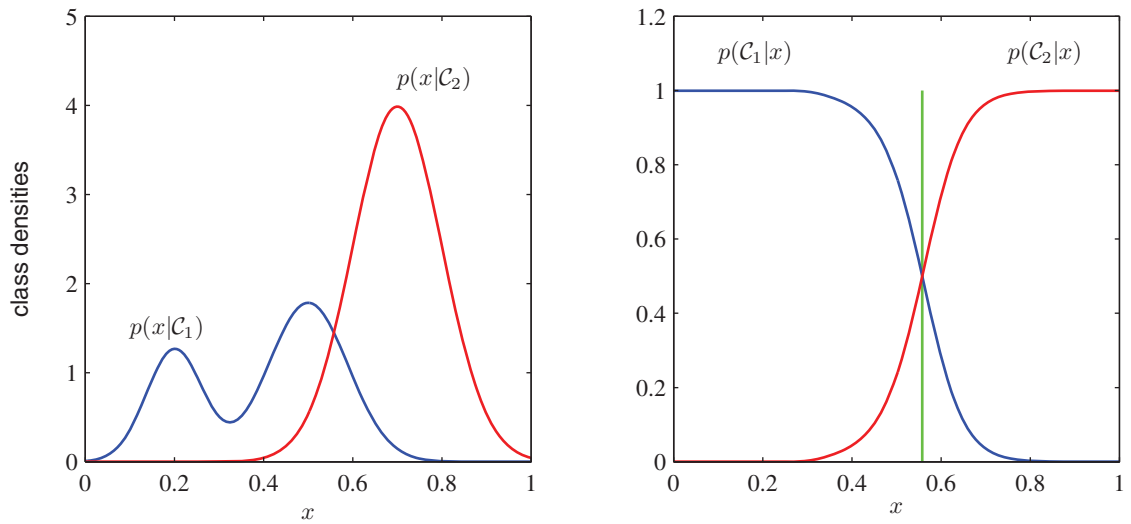
to find the posterior class probabilities  $p(\mathcal{C}_k|\mathbf{x})$ . As usual, the denominator in Bayes' theorem can be found in terms of the quantities appearing in the numerator, because

$$p(\mathbf{x}) = \sum_k p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k). \quad (1.83)$$

Equivalently, we can model the joint distribution  $p(\mathbf{x}, \mathcal{C}_k)$  directly and then normalize to obtain the posterior probabilities. Having found the posterior probabilities, we use decision theory to determine class membership for each new input  $\mathbf{x}$ . Approaches that explicitly or implicitly model the distribution of inputs as well as outputs are known as *generative models*, because by sampling from them it is possible to generate synthetic data points in the input space.

- (b) First solve the inference problem of determining the posterior class probabilities  $p(\mathcal{C}_k|\mathbf{x})$ , and then subsequently use decision theory to assign each new  $\mathbf{x}$  to one of the classes. Approaches that model the posterior probabilities directly are called *discriminative models*.
- (c) Find a function  $f(\mathbf{x})$ , called a discriminant function, which maps each input  $\mathbf{x}$  directly onto a class label. For instance, in the case of two-class problems,  $f(\cdot)$  might be binary valued and such that  $f = 0$  represents class  $\mathcal{C}_1$  and  $f = 1$  represents class  $\mathcal{C}_2$ . In this case, probabilities play no role.

Let us consider the relative merits of these three alternatives. Approach (a) is the most demanding because it involves finding the joint distribution over both  $\mathbf{x}$  and  $\mathcal{C}_k$ . For many applications,  $\mathbf{x}$  will have high dimensionality, and consequently we may need a large training set in order to be able to determine the class-conditional densities to reasonable accuracy. Note that the class priors  $p(\mathcal{C}_k)$  can often be estimated simply from the fractions of the training set data points in each of the classes. One advantage of approach (a), however, is that it also allows the marginal density of data  $p(\mathbf{x})$  to be determined from (1.83). This can be useful for detecting new data points that have low probability under the model and for which the predictions may



**Figure 1.27** Example of the class-conditional densities for two classes having a single input variable  $x$  (left plot) together with the corresponding posterior probabilities (right plot). Note that the left-hand mode of the class-conditional density  $p(x|\mathcal{C}_1)$ , shown in blue on the left plot, has no effect on the posterior probabilities. The vertical green line in the right plot shows the decision boundary in  $x$  that gives the minimum misclassification rate.

be of low accuracy, which is known as *outlier detection* or *novelty detection* (Bishop, 1994; Tarassenko, 1995).

However, if we only wish to make classification decisions, then it can be wasteful of computational resources, and excessively demanding of data, to find the joint distribution  $p(\mathbf{x}, \mathcal{C}_k)$  when in fact we only really need the posterior probabilities  $p(\mathcal{C}_k|\mathbf{x})$ , which can be obtained directly through approach (b). Indeed, the class-conditional densities may contain a lot of structure that has little effect on the posterior probabilities, as illustrated in Figure 1.27. There has been much interest in exploring the relative merits of generative and discriminative approaches to machine learning, and in finding ways to combine them (Jebara, 2004; Lasserre *et al.*, 2006).

An even simpler approach is (c) in which we use the training data to find a discriminant function  $f(\mathbf{x})$  that maps each  $\mathbf{x}$  directly onto a class label, thereby combining the inference and decision stages into a single learning problem. In the example of Figure 1.27, this would correspond to finding the value of  $x$  shown by the vertical green line, because this is the decision boundary giving the minimum probability of misclassification.

With option (c), however, we no longer have access to the posterior probabilities  $p(\mathcal{C}_k|\mathbf{x})$ . There are many powerful reasons for wanting to compute the posterior probabilities, even if we subsequently use them to make decisions. These include:

**Minimizing risk.** Consider a problem in which the elements of the loss matrix are subjected to revision from time to time (such as might occur in a financial

application). If we know the posterior probabilities, we can trivially revise the minimum risk decision criterion by modifying (1.81) appropriately. If we have only a discriminant function, then any change to the loss matrix would require that we return to the training data and solve the classification problem afresh.

**Reject option.** Posterior probabilities allow us to determine a rejection criterion that will minimize the misclassification rate, or more generally the expected loss, for a given fraction of rejected data points.

**Compensating for class priors.** Consider our medical X-ray problem again, and suppose that we have collected a large number of X-ray images from the general population for use as training data in order to build an automated screening system. Because cancer is rare amongst the general population, we might find that, say, only 1 in every 1,000 examples corresponds to the presence of cancer. If we used such a data set to train an adaptive model, we could run into severe difficulties due to the small proportion of the cancer class. For instance, a classifier that assigned every point to the normal class would already achieve 99.9% accuracy and it would be difficult to avoid this trivial solution. Also, even a large data set will contain very few examples of X-ray images corresponding to cancer, and so the learning algorithm will not be exposed to a broad range of examples of such images and hence is not likely to generalize well. A balanced data set in which we have selected equal numbers of examples from each of the classes would allow us to find a more accurate model. However, we then have to compensate for the effects of our modifications to the training data. Suppose we have used such a modified data set and found models for the posterior probabilities. From Bayes' theorem (1.82), we see that the posterior probabilities are proportional to the prior probabilities, which we can interpret as the fractions of points in each class. We can therefore simply take the posterior probabilities obtained from our artificially balanced data set and first divide by the class fractions in that data set and then multiply by the class fractions in the population to which we wish to apply the model. Finally, we need to normalize to ensure that the new posterior probabilities sum to one. Note that this procedure cannot be applied if we have learned a discriminant function directly instead of determining posterior probabilities.

**Combining models.** For complex applications, we may wish to break the problem into a number of smaller subproblems each of which can be tackled by a separate module. For example, in our hypothetical medical diagnosis problem, we may have information available from, say, blood tests as well as X-ray images. Rather than combine all of this heterogeneous information into one huge input space, it may be more effective to build one system to interpret the X-ray images and a different one to interpret the blood data. As long as each of the two models gives posterior probabilities for the classes, we can combine the outputs systematically using the rules of probability. One simple way to do this is to assume that, for each class separately, the distributions of inputs for the X-ray images, denoted by  $\mathbf{x}_I$ , and the blood data, denoted by  $\mathbf{x}_B$ , are



independent, so that

$$p(\mathbf{x}_I, \mathbf{x}_B | \mathcal{C}_k) = p(\mathbf{x}_I | \mathcal{C}_k) p(\mathbf{x}_B | \mathcal{C}_k). \quad (1.84)$$

This is an example of *conditional independence* property, because the independence holds when the distribution is conditioned on the class  $\mathcal{C}_k$ . The posterior probability, given both the X-ray and blood data, is then given by

$$\begin{aligned} p(\mathcal{C}_k | \mathbf{x}_I, \mathbf{x}_B) &\propto p(\mathbf{x}_I, \mathbf{x}_B | \mathcal{C}_k) p(\mathcal{C}_k) \\ &\propto p(\mathbf{x}_I | \mathcal{C}_k) p(\mathbf{x}_B | \mathcal{C}_k) p(\mathcal{C}_k) \\ &\propto \frac{p(\mathcal{C}_k | \mathbf{x}_I) p(\mathcal{C}_k | \mathbf{x}_B)}{p(\mathcal{C}_k)} \end{aligned} \quad (1.85)$$

Thus we need the class prior probabilities  $p(\mathcal{C}_k)$ , which we can easily estimate from the fractions of data points in each class, and then we need to normalize the resulting posterior probabilities so they sum to one. The particular conditional independence assumption (1.84) is an example of the *naive Bayes model*. Note that the joint marginal distribution  $p(\mathbf{x}_I, \mathbf{x}_B)$  will typically not factorize under this model. We shall see in later chapters how to construct models for combining data that do not require the conditional independence assumption (1.84).

### 1.5.5 Loss functions for regression

So far, we have discussed decision theory in the context of classification problems. We now turn to the case of regression problems, such as the curve fitting example discussed earlier. The decision stage consists of choosing a specific estimate  $y(\mathbf{x})$  of the value of  $t$  for each input  $\mathbf{x}$ . Suppose that in doing so, we incur a loss  $L(t, y(\mathbf{x}))$ . The average, or expected, loss is then given by

$$\mathbb{E}[L] = \iint L(t, y(\mathbf{x})) p(\mathbf{x}, t) d\mathbf{x} dt. \quad (1.86)$$

A common choice of loss function in regression problems is the squared loss given by  $L(t, y(\mathbf{x})) = \{y(\mathbf{x}) - t\}^2$ . In this case, the expected loss can be written

$$\mathbb{E}[L] = \iint \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt. \quad (1.87)$$

Our goal is to choose  $y(\mathbf{x})$  so as to minimize  $\mathbb{E}[L]$ . If we assume a completely flexible function  $y(\mathbf{x})$ , we can do this formally using the calculus of variations to give

$$\frac{\delta \mathbb{E}[L]}{\delta y(\mathbf{x})} = 2 \int \{y(\mathbf{x}) - t\} p(\mathbf{x}, t) dt = 0. \quad (1.88)$$

Solving for  $y(\mathbf{x})$ , and using the sum and product rules of probability, we obtain

$$y(\mathbf{x}) = \frac{\int t p(\mathbf{x}, t) dt}{p(\mathbf{x})} = \int t p(t | \mathbf{x}) dt = \mathbb{E}_t[t | \mathbf{x}] \quad (1.89)$$

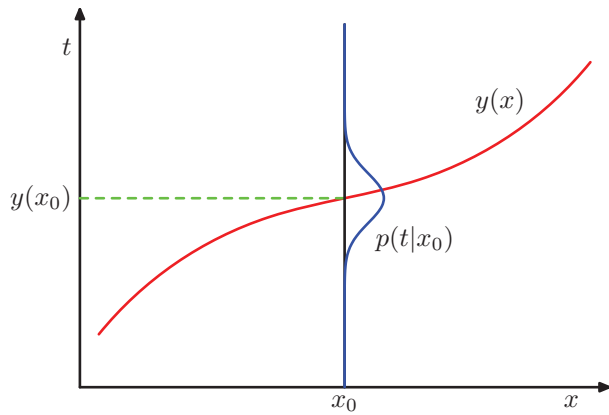
Section 8.2

Section 8.2.2

Section 1.1

Appendix D

**Figure 1.28** The regression function  $y(x)$ , which minimizes the expected squared loss, is given by the mean of the conditional distribution  $p(t|x)$ .



which is the conditional average of  $t$  conditioned on  $\mathbf{x}$  and is known as the *regression function*. This result is illustrated in Figure 1.28. It can readily be extended to multiple target variables represented by the vector  $\mathbf{t}$ , in which case the optimal solution is the conditional average  $\mathbf{y}(\mathbf{x}) = \mathbb{E}_{\mathbf{t}}[\mathbf{t}|\mathbf{x}]$ .

### Exercise 1.25

We can also derive this result in a slightly different way, which will also shed light on the nature of the regression problem. Armed with the knowledge that the optimal solution is the conditional expectation, we can expand the square term as follows

$$\begin{aligned} \{y(\mathbf{x}) - t\}^2 &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}] + \mathbb{E}[t|\mathbf{x}] - t\}^2 \\ &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 + 2\{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}\{\mathbb{E}[t|\mathbf{x}] - t\} + \{\mathbb{E}[t|\mathbf{x}] - t\}^2 \end{aligned}$$

where, to keep the notation uncluttered, we use  $\mathbb{E}[t|\mathbf{x}]$  to denote  $\mathbb{E}_t[t|\mathbf{x}]$ . Substituting into the loss function and performing the integral over  $t$ , we see that the cross-term vanishes and we obtain an expression for the loss function in the form

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 p(\mathbf{x}) d\mathbf{x} + \int \{\mathbb{E}[t|\mathbf{x}] - t\}^2 p(\mathbf{x}) d\mathbf{x}. \quad (1.90)$$

The function  $y(\mathbf{x})$  we seek to determine enters only in the first term, which will be minimized when  $y(\mathbf{x})$  is equal to  $\mathbb{E}[t|\mathbf{x}]$ , in which case this term will vanish. This is simply the result that we derived previously and that shows that the optimal least squares predictor is given by the conditional mean. The second term is the variance of the distribution of  $t$ , averaged over  $\mathbf{x}$ . It represents the intrinsic variability of the target data and can be regarded as noise. Because it is independent of  $y(\mathbf{x})$ , it represents the irreducible minimum value of the loss function.

As with the classification problem, we can either determine the appropriate probabilities and then use these to make optimal decisions, or we can build models that make decisions directly. Indeed, we can identify three distinct approaches to solving regression problems given, in order of decreasing complexity, by:

- (a) First solve the inference problem of determining the joint density  $p(\mathbf{x}, t)$ . Then normalize to find the conditional density  $p(t|\mathbf{x})$ , and finally marginalize to find the conditional mean given by (1.89).

(b) First solve the inference problem of determining the conditional density  $p(t|\mathbf{x})$ , and then subsequently marginalize to find the conditional mean given by (1.89).

(c) Find a regression function  $y(\mathbf{x})$  directly from the training data.

The relative merits of these three approaches follow the same lines as for classification problems above.

The squared loss is not the only possible choice of loss function for regression. Indeed, there are situations in which squared loss can lead to very poor results and where we need to develop more sophisticated approaches. An important example concerns situations in which the conditional distribution  $p(t|\mathbf{x})$  is multimodal, as often arises in the solution of inverse problems. Here we consider briefly one simple generalization of the squared loss, called the *Minkowski* loss, whose expectation is given by

$$\mathbb{E}[L_q] = \iint |y(\mathbf{x}) - t|^q p(\mathbf{x}, t) d\mathbf{x} dt \quad (1.91)$$

which reduces to the expected squared loss for  $q = 2$ . The function  $|y - t|^q$  is plotted against  $y - t$  for various values of  $q$  in Figure 1.29. The minimum of  $\mathbb{E}[L_q]$  is given by the conditional mean for  $q = 2$ , the conditional median for  $q = 1$ , and the conditional mode for  $q \rightarrow 0$ .

## 1.6. Information Theory

In this chapter, we have discussed a variety of concepts from probability theory and decision theory that will form the foundations for much of the subsequent discussion in this book. We close this chapter by introducing some additional concepts from the field of information theory, which will also prove useful in our development of pattern recognition and machine learning techniques. Again, we shall focus only on the key concepts, and we refer the reader elsewhere for more detailed discussions (Viterbi and Omura, 1979; Cover and Thomas, 1991; MacKay, 2003).

We begin by considering a discrete random variable  $x$  and we ask how much information is received when we observe a specific value for this variable. The amount of information can be viewed as the ‘degree of surprise’ on learning the value of  $x$ . If we are told that a highly improbable event has just occurred, we will have received more information than if we were told that some very likely event has just occurred, and if we knew that the event was certain to happen we would receive no information. Our measure of information content will therefore depend on the probability distribution  $p(x)$ , and we therefore look for a quantity  $h(x)$  that is a monotonic function of the probability  $p(x)$  and that expresses the information content. The form of  $h(\cdot)$  can be found by noting that if we have two events  $x$  and  $y$  that are unrelated, then the information gain from observing both of them should be the sum of the information gained from each of them separately, so that  $h(x, y) = h(x) + h(y)$ . Two unrelated events will be statistically independent and so  $p(x, y) = p(x)p(y)$ . From these two relationships, it is easily shown that  $h(x)$  must be given by the logarithm of  $p(x)$  and so we have



# 3

## Linear Models for Regression

The focus so far in this book has been on unsupervised learning, including topics such as density estimation and data clustering. We turn now to a discussion of supervised learning, starting with regression. The goal of regression is to predict the value of one or more continuous *target* variables  $t$  given the value of a  $D$ -dimensional vector  $\mathbf{x}$  of *input* variables. We have already encountered an example of a regression problem when we considered polynomial curve fitting in Chapter 1. The polynomial is a specific example of a broad class of functions called linear regression models, which share the property of being linear functions of the adjustable parameters, and which will form the focus of this chapter. The simplest form of linear regression models are also linear functions of the input variables. However, we can obtain a much more useful class of functions by taking linear combinations of a fixed set of nonlinear functions of the input variables, known as *basis functions*. Such models are linear functions of the parameters, which gives them simple analytical properties, and yet can be nonlinear with respect to the input variables.

Given a training data set comprising  $N$  observations  $\{\mathbf{x}_n\}$ , where  $n = 1, \dots, N$ , together with corresponding target values  $\{t_n\}$ , the goal is to predict the value of  $t$  for a new value of  $\mathbf{x}$ . In the simplest approach, this can be done by directly constructing an appropriate function  $y(\mathbf{x})$  whose values for new inputs  $\mathbf{x}$  constitute the predictions for the corresponding values of  $t$ . More generally, from a probabilistic perspective, we aim to model the predictive distribution  $p(t|\mathbf{x})$  because this expresses our uncertainty about the value of  $t$  for each value of  $\mathbf{x}$ . From this conditional distribution we can make predictions of  $t$ , for any new value of  $\mathbf{x}$ , in such a way as to minimize the expected value of a suitably chosen loss function. As discussed in Section 1.5.5, a common choice of loss function for real-valued variables is the squared loss, for which the optimal solution is given by the conditional expectation of  $t$ .

Although linear models have significant limitations as practical techniques for pattern recognition, particularly for problems involving input spaces of high dimensionality, they have nice analytical properties and form the foundation for more sophisticated models to be discussed in later chapters.

### 3.1. Linear Basis Function Models

The simplest linear model for regression is one that involves a linear combination of the input variables

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1x_1 + \dots + w_Dx_D \quad (3.1)$$

where  $\mathbf{x} = (x_1, \dots, x_D)^T$ . This is often simply known as *linear regression*. The key property of this model is that it is a linear function of the parameters  $w_0, \dots, w_D$ . It is also, however, a linear function of the input variables  $x_i$ , and this imposes significant limitations on the model. We therefore extend the class of models by considering linear combinations of fixed nonlinear functions of the input variables, of the form

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}) \quad (3.2)$$

where  $\phi_j(\mathbf{x})$  are known as *basis functions*. By denoting the maximum value of the index  $j$  by  $M - 1$ , the total number of parameters in this model will be  $M$ .

The parameter  $w_0$  allows for any fixed offset in the data and is sometimes called a *bias* parameter (not to be confused with ‘bias’ in a statistical sense). It is often convenient to define an additional dummy ‘basis function’  $\phi_0(\mathbf{x}) = 1$  so that

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) \quad (3.3)$$

where  $\mathbf{w} = (w_0, \dots, w_{M-1})^T$  and  $\boldsymbol{\phi} = (\phi_0, \dots, \phi_{M-1})^T$ . In many practical applications of pattern recognition, we will apply some form of fixed pre-processing,

or feature extraction, to the original data variables. If the original variables comprise the vector  $\mathbf{x}$ , then the features can be expressed in terms of the basis functions  $\{\phi_j(\mathbf{x})\}$ .

By using nonlinear basis functions, we allow the function  $y(\mathbf{x}, \mathbf{w})$  to be a nonlinear function of the input vector  $\mathbf{x}$ . Functions of the form (3.2) are called linear models, however, because this function is linear in  $\mathbf{w}$ . It is this linearity in the parameters that will greatly simplify the analysis of this class of models. However, it also leads to some significant limitations, as we discuss in Section 3.6.

The example of polynomial regression considered in Chapter 1 is a particular example of this model in which there is a single input variable  $x$ , and the basis functions take the form of powers of  $x$  so that  $\phi_j(x) = x^j$ . One limitation of polynomial basis functions is that they are global functions of the input variable, so that changes in one region of input space affect all other regions. This can be resolved by dividing the input space up into regions and fit a different polynomial in each region, leading to *spline functions* (Hastie *et al.*, 2001).

There are many other possible choices for the basis functions, for example

$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\} \quad (3.4)$$

where the  $\mu_j$  govern the locations of the basis functions in input space, and the parameter  $s$  governs their spatial scale. These are usually referred to as ‘Gaussian’ basis functions, although it should be noted that they are not required to have a probabilistic interpretation, and in particular the normalization coefficient is unimportant because these basis functions will be multiplied by adaptive parameters  $w_j$ .

Another possibility is the sigmoidal basis function of the form

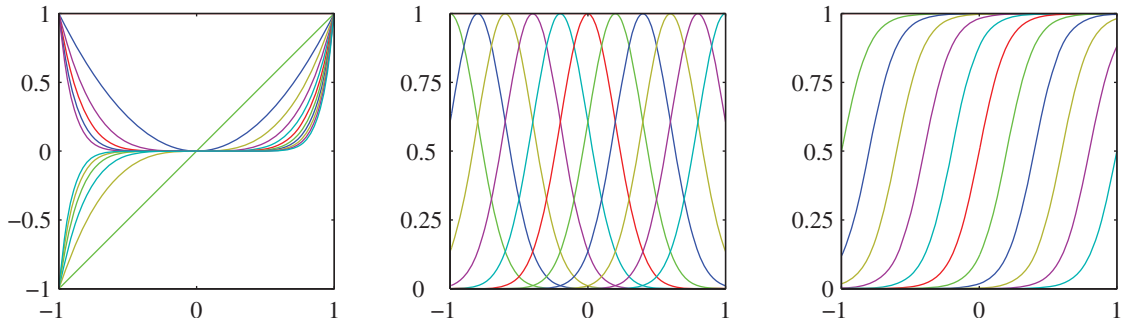
$$\phi_j(x) = \sigma \left( \frac{x - \mu_j}{s} \right) \quad (3.5)$$

where  $\sigma(a)$  is the logistic sigmoid function defined by

$$\sigma(a) = \frac{1}{1 + \exp(-a)}. \quad (3.6)$$

Equivalently, we can use the ‘tanh’ function because this is related to the logistic sigmoid by  $\tanh(a) = 2\sigma(a) - 1$ , and so a general linear combination of logistic sigmoid functions is equivalent to a general linear combination of ‘tanh’ functions. These various choices of basis function are illustrated in Figure 3.1.

Yet another possible choice of basis function is the Fourier basis, which leads to an expansion in sinusoidal functions. Each basis function represents a specific frequency and has infinite spatial extent. By contrast, basis functions that are localized to finite regions of input space necessarily comprise a spectrum of different spatial frequencies. In many signal processing applications, it is of interest to consider basis functions that are localized in both space and frequency, leading to a class of functions known as *wavelets*. These are also defined to be mutually orthogonal, to simplify their application. Wavelets are most applicable when the input values live



**Figure 3.1** Examples of basis functions, showing polynomials on the left, Gaussians of the form (3.4) in the centre, and sigmoidal of the form (3.5) on the right.

on a regular lattice, such as the successive time points in a temporal sequence, or the pixels in an image. Useful texts on wavelets include Ogden (1997), Mallat (1999), and Vidakovic (1999).

Most of the discussion in this chapter, however, is independent of the particular choice of basis function set, and so for most of our discussion we shall not specify the particular form of the basis functions, except for the purposes of numerical illustration. Indeed, much of our discussion will be equally applicable to the situation in which the vector  $\phi(\mathbf{x})$  of basis functions is simply the identity  $\phi(\mathbf{x}) = \mathbf{x}$ . Furthermore, in order to keep the notation simple, we shall focus on the case of a single target variable  $t$ . However, in Section 3.1.5, we consider briefly the modifications needed to deal with multiple target variables.

### 3.1.1 Maximum likelihood and least squares

In Chapter 1, we fitted polynomial functions to data sets by minimizing a sum-of-squares error function. We also showed that this error function could be motivated as the maximum likelihood solution under an assumed Gaussian noise model. Let us return to this discussion and consider the least squares approach, and its relation to maximum likelihood, in more detail.

As before, we assume that the target variable  $t$  is given by a deterministic function  $y(\mathbf{x}, \mathbf{w})$  with additive Gaussian noise so that

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad (3.7)$$

where  $\epsilon$  is a zero mean Gaussian random variable with precision (inverse variance)  $\beta$ . Thus we can write

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1}). \quad (3.8)$$

Recall that, if we assume a squared loss function, then the optimal prediction, for a new value of  $\mathbf{x}$ , will be given by the conditional mean of the target variable. In the case of a Gaussian conditional distribution of the form (3.8), the conditional mean

will be simply

$$\mathbb{E}[t|\mathbf{x}] = \int t p(t|\mathbf{x}) dt = y(\mathbf{x}, \mathbf{w}). \quad (3.9)$$

Note that the Gaussian noise assumption implies that the conditional distribution of  $t$  given  $\mathbf{x}$  is unimodal, which may be inappropriate for some applications. An extension to mixtures of conditional Gaussian distributions, which permit multimodal conditional distributions, will be discussed in Section 14.5.1.

Now consider a data set of inputs  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  with corresponding target values  $t_1, \dots, t_N$ . We group the target variables  $\{t_n\}$  into a column vector that we denote by  $\mathbf{t}$  where the typeface is chosen to distinguish it from a single observation of a multivariate target, which would be denoted  $\mathbf{t}$ . Making the assumption that these data points are drawn independently from the distribution (3.8), we obtain the following expression for the likelihood function, which is a function of the adjustable parameters  $\mathbf{w}$  and  $\beta$ , in the form

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \quad (3.10)$$

where we have used (3.3). Note that in supervised learning problems such as regression (and classification), we are not seeking to model the distribution of the input variables. Thus  $\mathbf{x}$  will always appear in the set of conditioning variables, and so from now on we will drop the explicit  $\mathbf{x}$  from expressions such as  $p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)$  in order to keep the notation uncluttered. Taking the logarithm of the likelihood function, and making use of the standard form (1.46) for the univariate Gaussian, we have

$$\begin{aligned} \ln p(\mathbf{t}|\mathbf{w}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w}) \end{aligned} \quad (3.11)$$

where the sum-of-squares error function is defined by

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2. \quad (3.12)$$

Having written down the likelihood function, we can use maximum likelihood to determine  $\mathbf{w}$  and  $\beta$ . Consider first the maximization with respect to  $\mathbf{w}$ . As observed already in Section 1.2.5, we see that maximization of the likelihood function under a conditional Gaussian noise distribution for a linear model is equivalent to minimizing a sum-of-squares error function given by  $E_D(\mathbf{w})$ . The gradient of the log likelihood function (3.11) takes the form

$$\nabla \ln p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\} \phi(\mathbf{x}_n)^T. \quad (3.13)$$



Setting this gradient to zero gives

$$0 = \sum_{n=1}^N t_n \phi(\mathbf{x}_n)^T - \mathbf{w}^T \left( \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \right). \quad (3.14)$$

Solving for  $\mathbf{w}$  we obtain

$$\mathbf{w}_{\text{ML}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad (3.15)$$

which are known as the *normal equations* for the least squares problem. Here  $\Phi$  is an  $N \times M$  matrix, called the *design matrix*, whose elements are given by  $\Phi_{nj} = \phi_j(\mathbf{x}_n)$ , so that

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}. \quad (3.16)$$

The quantity

$$\Phi^\dagger \equiv (\Phi^T \Phi)^{-1} \Phi^T \quad (3.17)$$

is known as the *Moore-Penrose pseudo-inverse* of the matrix  $\Phi$  (Rao and Mitra, 1971; Golub and Van Loan, 1996). It can be regarded as a generalization of the notion of matrix inverse to nonsquare matrices. Indeed, if  $\Phi$  is square and invertible, then using the property  $(\mathbf{A}\mathbf{B})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$  we see that  $\Phi^\dagger \equiv \Phi^{-1}$ .

At this point, we can gain some insight into the role of the bias parameter  $w_0$ . If we make the bias parameter explicit, then the error function (3.12) becomes

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left\{ t_n - w_0 - \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}_n) \right\}^2. \quad (3.18)$$

Setting the derivative with respect to  $w_0$  equal to zero, and solving for  $w_0$ , we obtain

$$w_0 = \bar{t} - \sum_{j=1}^{M-1} w_j \bar{\phi}_j \quad (3.19)$$

where we have defined

$$\bar{t} = \frac{1}{N} \sum_{n=1}^N t_n, \quad \bar{\phi}_j = \frac{1}{N} \sum_{n=1}^N \phi_j(\mathbf{x}_n). \quad (3.20)$$

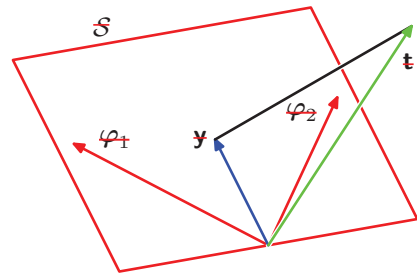
Thus the bias  $w_0$  compensates for the difference between the averages (over the training set) of the target values and the weighted sum of the averages of the basis function values.

We can also maximize the log likelihood function (3.11) with respect to the noise precision parameter  $\beta$ , giving

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{n=1}^N \{ t_n - \mathbf{w}_{\text{ML}}^T \phi(\mathbf{x}_n) \}^2 \quad (3.21)$$

**Figure 3.2**

Geometrical interpretation of the least-squares solution, in an  $N$ -dimensional space whose axes are the values of  $t_1, \dots, t_N$ . The least-squares regression function is obtained by finding the orthogonal projection of the data vector  $\mathbf{t}$  onto the subspace spanned by the basis functions  $\phi_j(\mathbf{x})$  in which each basis function is viewed as a vector  $\varphi_j$  of length  $N$  with elements  $\phi_j(\mathbf{x}_n)$ .



and so we see that the inverse of the noise precision is given by the residual variance of the target values around the regression function.

### 3.1.2 Geometry of least squares

At this point, it is instructive to consider the geometrical interpretation of the least-squares solution. To do this we consider an  $N$ -dimensional space whose axes are given by the  $t_n$ , so that  $\mathbf{t} = (t_1, \dots, t_N)^T$  is a vector in this space. Each basis function  $\phi_j(\mathbf{x}_n)$ , evaluated at the  $N$  data points, can also be represented as a vector in the same space, denoted by  $\varphi_j$ , as illustrated in Figure 3.2. Note that  $\varphi_j$  corresponds to the  $j^{\text{th}}$  column of  $\Phi$ , whereas  $\phi(\mathbf{x}_n)$  corresponds to the  $n^{\text{th}}$  row of  $\Phi$ . If the number  $M$  of basis functions is smaller than the number  $N$  of data points, then the  $M$  vectors  $\phi_j(\mathbf{x}_n)$  will span a linear subspace  $\mathcal{S}$  of dimensionality  $M$ . We define  $\mathbf{y}$  to be an  $N$ -dimensional vector whose  $n^{\text{th}}$  element is given by  $y(\mathbf{x}_n, \mathbf{w})$ , where  $n = 1, \dots, N$ . Because  $\mathbf{y}$  is an arbitrary linear combination of the vectors  $\varphi_j$ , it can live anywhere in the  $M$ -dimensional subspace. The sum-of-squares error (3.12) is then equal (up to a factor of  $1/2$ ) to the squared Euclidean distance between  $\mathbf{y}$  and  $\mathbf{t}$ . Thus the least-squares solution for  $\mathbf{w}$  corresponds to that choice of  $\mathbf{y}$  that lies in subspace  $\mathcal{S}$  and that is closest to  $\mathbf{t}$ . Intuitively, from Figure 3.2, we anticipate that this solution corresponds to the orthogonal projection of  $\mathbf{t}$  onto the subspace  $\mathcal{S}$ . This is indeed the case, as can easily be verified by noting that the solution for  $\mathbf{y}$  is given by  $\Phi \mathbf{w}_{\text{ML}}$ , and then confirming that this takes the form of an orthogonal projection.

In practice, a direct solution of the normal equations can lead to numerical difficulties when  $\Phi^T \Phi$  is close to singular. In particular, when two or more of the basis vectors  $\varphi_j$  are co-linear, or nearly so, the resulting parameter values can have large magnitudes. Such near degeneracies will not be uncommon when dealing with real data sets. The resulting numerical difficulties can be addressed using the technique of *singular value decomposition*, or *SVD* (Press *et al.*, 1992; Bishop and Nabney, 2008). Note that the addition of a regularization term ensures that the matrix is non-singular, even in the presence of degeneracies.

### 3.1.3 Sequential learning

Batch techniques, such as the maximum likelihood solution (3.15), which involve processing the entire training set in one go, can be computationally costly for large data sets. As we have discussed in Chapter 1, if the data set is sufficiently large, it may be worthwhile to use *sequential* algorithms, also known as *on-line* algorithms,

*Exercise 3.2*

in which the data points are considered one at a time, and the model parameters updated after each such presentation. Sequential learning is also appropriate for real-time applications in which the data observations are arriving in a continuous stream, and predictions must be made before all of the data points are seen.

We can obtain a sequential learning algorithm by applying the technique of *stochastic gradient descent*, also known as *sequential gradient descent*, as follows. If the error function comprises a sum over data points  $E = \sum_n E_n$ , then after presentation of pattern  $n$ , the stochastic gradient descent algorithm updates the parameter vector  $\mathbf{w}$  using

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n \quad (3.22)$$

where  $\tau$  denotes the iteration number, and  $\eta$  is a learning rate parameter. We shall discuss the choice of value for  $\eta$  shortly. The value of  $\mathbf{w}$  is initialized to some starting vector  $\mathbf{w}^{(0)}$ . For the case of the sum-of-squares error function (3.12), this gives

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \eta(t_n - \mathbf{w}^{(\tau)\top} \phi_n) \phi_n \quad (3.23)$$

where  $\phi_n = \phi(\mathbf{x}_n)$ . This is known as *least-mean-squares* or the *LMS algorithm*. The value of  $\eta$  needs to be chosen with care to ensure that the algorithm converges (Bishop and Nabney, 2008).

### 3.1.4 Regularized least squares

In Section 1.1, we introduced the idea of adding a regularization term to an error function in order to control over-fitting, so that the total error function to be minimized takes the form

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w}) \quad (3.24)$$

where  $\lambda$  is the regularization coefficient that controls the relative importance of the data-dependent error  $E_D(\mathbf{w})$  and the regularization term  $E_W(\mathbf{w})$ . One of the simplest forms of regularizer is given by the sum-of-squares of the weight vector elements

$$E_W(\mathbf{w}) = \frac{1}{2} \mathbf{w}^\top \mathbf{w}. \quad (3.25)$$

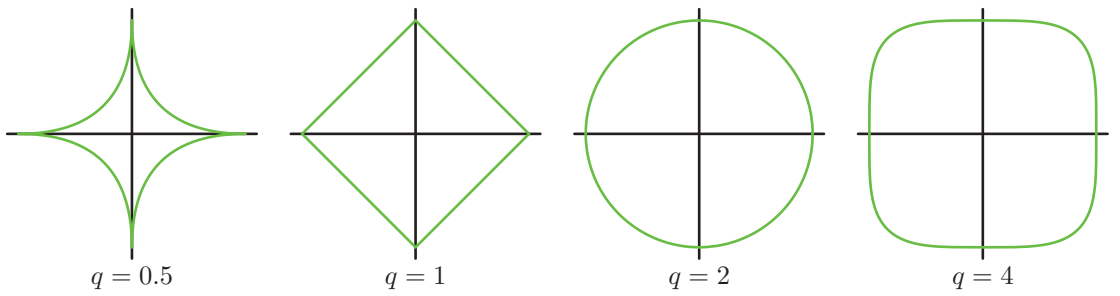
If we also consider the sum-of-squares error function given by

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^\top \phi(\mathbf{x}_n)\}^2 \quad (3.26)$$

then the total error function becomes

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^\top \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}. \quad (3.27)$$

This particular choice of regularizer is known in the machine learning literature as *weight decay* because in sequential learning algorithms, it encourages weight values to decay towards zero, unless supported by the data. In statistics, it provides an example of a *parameter shrinkage* method because it shrinks parameter values towards



**Figure 3.3** Contours of the regularization term in (3.29) for various values of the parameter  $q$ .

zero. It has the advantage that the error function remains a quadratic function of  $\mathbf{w}$ , and so its exact minimizer can be found in closed form. Specifically, setting the gradient of (3.27) with respect to  $\mathbf{w}$  to zero, and solving for  $\mathbf{w}$  as before, we obtain

$$\mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}. \quad (3.28)$$

This represents a simple extension of the least-squares solution (3.15).

A more general regularizer is sometimes used, for which the regularized error takes the form

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q \quad (3.29)$$

where  $q = 2$  corresponds to the quadratic regularizer (3.27). Figure 3.3 shows contours of the regularization function for different values of  $q$ .

The case of  $q = 1$  is known as the *lasso* in the statistics literature (Tibshirani, 1996). It has the property that if  $\lambda$  is sufficiently large, some of the coefficients  $w_j$  are driven to zero, leading to a *sparse* model in which the corresponding basis functions play no role. To see this, we first note that minimizing (3.29) is equivalent to minimizing the unregularized sum-of-squares error (3.12) subject to the constraint

$$\sum_{j=1}^M |w_j|^q \leq \eta \quad (3.30)$$

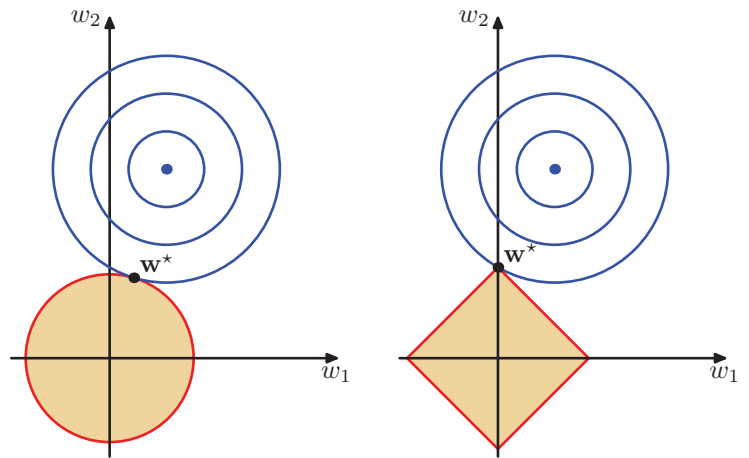
for an appropriate value of the parameter  $\eta$ , where the two approaches can be related using Lagrange multipliers. The origin of the sparsity can be seen from Figure 3.4, which shows that the minimum of the error function, subject to the constraint (3.30). As  $\lambda$  is increased, so an increasing number of parameters are driven to zero.

Regularization allows complex models to be trained on data sets of limited size without severe over-fitting, essentially by limiting the effective model complexity. However, the problem of determining the optimal model complexity is then shifted from one of finding the appropriate number of basis functions to one of determining a suitable value of the regularization coefficient  $\lambda$ . We shall return to the issue of model complexity later in this chapter.

*Exercise 3.5*

*Appendix E*

**Figure 3.4** Plot of the contours of the unregularized error function (blue) along with the constraint region (3.30) for the quadratic regularizer  $q = 2$  on the left and the lasso regularizer  $q = 1$  on the right, in which the optimum value for the parameter vector  $\mathbf{w}$  is denoted by  $\mathbf{w}^*$ . The lasso gives a sparse solution in which  $w_1^* = 0$ .



For the remainder of this chapter we shall focus on the quadratic regularizer (3.27) both for its practical importance and its analytical tractability.

### 3.1.5 Multiple outputs

So far, we have considered the case of a single target variable  $t$ . In some applications, we may wish to predict  $K > 1$  target variables, which we denote collectively by the target vector  $\mathbf{t}$ . This could be done by introducing a different set of basis functions for each component of  $\mathbf{t}$ , leading to multiple, independent regression problems. However, a more interesting, and more common, approach is to use the same set of basis functions to model all of the components of the target vector so that

$$\mathbf{y}(\mathbf{x}, \mathbf{w}) = \mathbf{W}^T \phi(\mathbf{x}) \quad (3.31)$$

where  $\mathbf{y}$  is a  $K$ -dimensional column vector,  $\mathbf{W}$  is an  $M \times K$  matrix of parameters, and  $\phi(\mathbf{x})$  is an  $M$ -dimensional column vector with elements  $\phi_j(\mathbf{x})$ , with  $\phi_0(\mathbf{x}) = 1$  as before. Suppose we take the conditional distribution of the target vector to be an isotropic Gaussian of the form

$$p(\mathbf{t}|\mathbf{x}, \mathbf{W}, \beta) = \mathcal{N}(\mathbf{t}|\mathbf{W}^T \phi(\mathbf{x}), \beta^{-1}\mathbf{I}). \quad (3.32)$$

If we have a set of observations  $\mathbf{t}_1, \dots, \mathbf{t}_N$ , we can combine these into a matrix  $\mathbf{T}$  of size  $N \times K$  such that the  $n^{\text{th}}$  row is given by  $\mathbf{t}_n^T$ . Similarly, we can combine the input vectors  $\mathbf{x}_1, \dots, \mathbf{x}_N$  into a matrix  $\mathbf{X}$ . The log-likelihood function is then given by

$$\begin{aligned} \ln p(\mathbf{T}|\mathbf{X}, \mathbf{W}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(\mathbf{t}_n|\mathbf{W}^T \phi(\mathbf{x}_n), \beta^{-1}\mathbf{I}) \\ &= \frac{NK}{2} \ln \left( \frac{\beta}{2\pi} \right) - \frac{\beta}{2} \sum_{n=1}^N \left\| \mathbf{t}_n - \mathbf{W}^T \phi(\mathbf{x}_n) \right\|^2. \end{aligned} \quad (3.33)$$

As before, we can maximize this function with respect to  $\mathbf{W}$ , giving

$$\mathbf{W}_{\text{ML}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{T}. \quad (3.34)$$

If we examine this result for each target variable  $t_k$ , we have

$$\mathbf{w}_k = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}_k = \Phi^\dagger \mathbf{t}_k \quad (3.35)$$

where  $\mathbf{t}_k$  is an  $N$ -dimensional column vector with components  $t_{nk}$  for  $n = 1, \dots, N$ . Thus the solution to the regression problem decouples between the different target variables, and we need only compute a single pseudo-inverse matrix  $\Phi^\dagger$ , which is shared by all of the vectors  $\mathbf{w}_k$ .

The extension to general Gaussian noise distributions having arbitrary covariance matrices is straightforward. Again, this leads to a decoupling into  $K$  independent regression problems. This result is unsurprising because the parameters  $\mathbf{W}$  define only the mean of the Gaussian noise distribution, and we know from Section 2.3.4 that the maximum likelihood solution for the mean of a multivariate Gaussian is independent of the covariance. From now on, we shall therefore consider a single target variable  $t$  for simplicity.

## 3.2. The Bias-Variance Decomposition

So far in our discussion of linear models for regression, we have assumed that the form and number of basis functions are both fixed. As we have seen in Chapter 1, the use of maximum likelihood, or equivalently least squares, can lead to severe over-fitting if complex models are trained using data sets of limited size. However, limiting the number of basis functions in order to avoid over-fitting has the side effect of limiting the flexibility of the model to capture interesting and important trends in the data. Although the introduction of regularization terms can control over-fitting for models with many parameters, this raises the question of how to determine a suitable value for the regularization coefficient  $\lambda$ . Seeking the solution that minimizes the regularized error function with respect to both the weight vector  $\mathbf{w}$  and the regularization coefficient  $\lambda$  is clearly not the right approach since this leads to the unregularized solution with  $\lambda = 0$ .

As we have seen in earlier chapters, the phenomenon of over-fitting is really an unfortunate property of maximum likelihood and does not arise when we marginalize over parameters in a Bayesian setting. In this chapter, we shall consider the Bayesian view of model complexity in some depth. Before doing so, however, it is instructive to consider a frequentist viewpoint of the model complexity issue, known as the *bias-variance* trade-off. Although we shall introduce this concept in the context of linear basis function models, where it is easy to illustrate the ideas using simple examples, the discussion has more general applicability.

In Section 1.5.5, when we discussed decision theory for regression problems, we considered various loss functions each of which leads to a corresponding optimal prediction once we are given the conditional distribution  $p(t|\mathbf{x})$ . A popular choice is

the squared loss function, for which the optimal prediction is given by the conditional expectation, which we denote by  $h(\mathbf{x})$  and which is given by

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int tp(t|\mathbf{x}) dt. \quad (3.36)$$

At this point, it is worth distinguishing between the squared loss function arising from decision theory and the sum-of-squares error function that arose in the maximum likelihood estimation of model parameters. We might use more sophisticated techniques than least squares, for example regularization or a fully Bayesian approach, to determine the conditional distribution  $p(t|\mathbf{x})$ . These can all be combined with the squared loss function for the purpose of making predictions.

We showed in Section 1.5.5 that the expected squared loss can be written in the form

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt. \quad (3.37)$$

Recall that the second term, which is independent of  $y(\mathbf{x})$ , arises from the intrinsic noise on the data and represents the minimum achievable value of the expected loss. The first term depends on our choice for the function  $y(\mathbf{x})$ , and we will seek a solution for  $y(\mathbf{x})$  which makes this term a minimum. Because it is nonnegative, the smallest that we can hope to make this term is zero. If we had an unlimited supply of data (and unlimited computational resources), we could in principle find the regression function  $h(\mathbf{x})$  to any desired degree of accuracy, and this would represent the optimal choice for  $y(\mathbf{x})$ . However, in practice we have a data set  $\mathcal{D}$  containing only a finite number  $N$  of data points, and consequently we do not know the regression function  $h(\mathbf{x})$  exactly.

If we model the  $h(\mathbf{x})$  using a parametric function  $y(\mathbf{x}, \mathbf{w})$  governed by a parameter vector  $\mathbf{w}$ , then from a Bayesian perspective the uncertainty in our model is expressed through a posterior distribution over  $\mathbf{w}$ . A frequentist treatment, however, involves making a point estimate of  $\mathbf{w}$  based on the data set  $\mathcal{D}$ , and tries instead to interpret the uncertainty of this estimate through the following thought experiment. Suppose we had a large number of data sets each of size  $N$  and each drawn independently from the distribution  $p(t, \mathbf{x})$ . For any given data set  $\mathcal{D}$ , we can run our learning algorithm and obtain a prediction function  $y(\mathbf{x}; \mathcal{D})$ . Different data sets from the ensemble will give different functions and consequently different values of the squared loss. The performance of a particular learning algorithm is then assessed by taking the average over this ensemble of data sets.

Consider the integrand of the first term in (3.37), which for a particular data set  $\mathcal{D}$  takes the form

$$\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2. \quad (3.38)$$

Because this quantity will be dependent on the particular data set  $\mathcal{D}$ , we take its average over the ensemble of data sets. If we add and subtract the quantity  $\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]$

inside the braces, and then expand, we obtain

$$\begin{aligned} & \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ &= \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 + \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ & \quad + 2\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}. \end{aligned} \quad (3.39)$$

We now take the expectation of this expression with respect to  $\mathcal{D}$  and note that the final term will vanish, giving

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2] \\ &= \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2}_{(\text{bias})^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2]}_{\text{variance}}. \end{aligned} \quad (3.40)$$

We see that the expected squared difference between  $y(\mathbf{x}; \mathcal{D})$  and the regression function  $h(\mathbf{x})$  can be expressed as the sum of two terms. The first term, called the squared *bias*, represents the extent to which the average prediction over all data sets differs from the desired regression function. The second term, called the *variance*, measures the extent to which the solutions for individual data sets vary around their average, and hence this measures the extent to which the function  $y(\mathbf{x}; \mathcal{D})$  is sensitive to the particular choice of data set. We shall provide some intuition to support these definitions shortly when we consider a simple example.

So far, we have considered a single input value  $\mathbf{x}$ . If we substitute this expansion back into (3.37), we obtain the following decomposition of the expected squared loss

$$\text{expected loss} = (\text{bias})^2 + \text{variance} + \text{noise} \quad (3.41)$$

where

$$(\text{bias})^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} \quad (3.42)$$

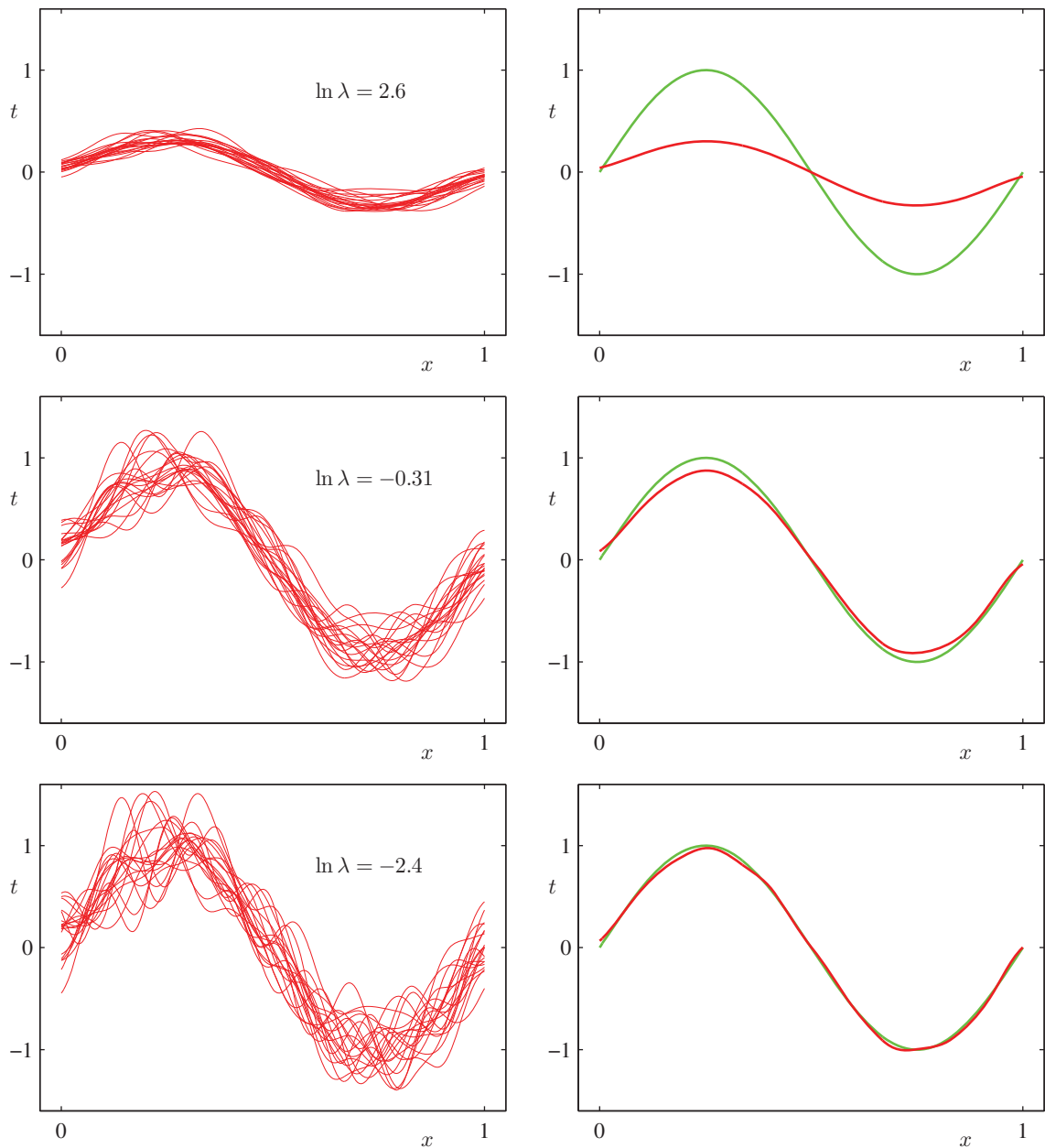
$$\text{variance} = \int \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) d\mathbf{x} \quad (3.43)$$

$$\text{noise} = \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt \quad (3.44)$$

and the bias and variance terms now refer to integrated quantities.

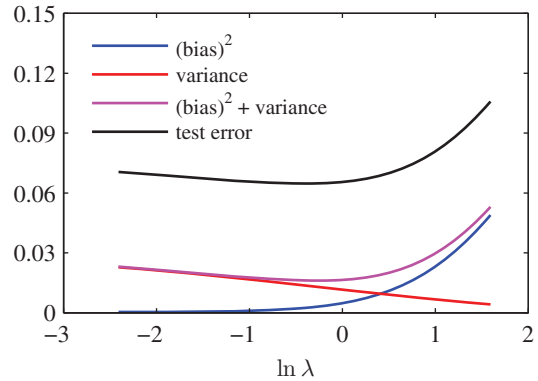
Our goal is to minimize the expected loss, which we have decomposed into the sum of a (squared) bias, a variance, and a constant noise term. As we shall see, there is a trade-off between bias and variance, with very flexible models having low bias and high variance, and relatively rigid models having high bias and low variance. The model with the optimal predictive capability is the one that leads to the best balance between bias and variance. This is illustrated by considering the sinusoidal data set from Chapter 1. Here we generate 100 data sets, each containing  $N = 25$  data points, independently from the sinusoidal curve  $h(x) = \sin(2\pi x)$ . The data sets are indexed by  $l = 1, \dots, L$ , where  $L = 100$ , and for each data set  $\mathcal{D}^{(l)}$  we





**Figure 3.5** Illustration of the dependence of bias and variance on model complexity, governed by a regularization parameter  $\lambda$ , using the sinusoidal data set from Chapter 1. There are  $L = 100$  data sets, each having  $N = 25$  data points, and there are 24 Gaussian basis functions in the model so that the total number of parameters is  $M = 25$  including the bias parameter. The left column shows the result of fitting the model to the data sets for various values of  $\ln \lambda$  (for clarity, only 20 of the 100 fits are shown). The right column shows the corresponding average of the 100 fits (red) along with the sinusoidal function from which the data sets were generated (green).

**Figure 3.6** Plot of squared bias and variance, together with their sum, corresponding to the results shown in Figure 3.5. Also shown is the average test set error for a test data set size of 1000 points. The minimum value of  $(\text{bias})^2 + \text{variance}$  occurs around  $\ln \lambda = -0.31$ , which is close to the value that gives the minimum error on the test data.



fit a model with 24 Gaussian basis functions by minimizing the regularized error function (3.27) to give a prediction function  $y^{(l)}(x)$  as shown in Figure 3.5. The top row corresponds to a large value of the regularization coefficient  $\lambda$  that gives low variance (because the red curves in the left plot look similar) but high bias (because the two curves in the right plot are very different). Conversely on the bottom row, for which  $\lambda$  is small, there is large variance (shown by the high variability between the red curves in the left plot) but low bias (shown by the good fit between the average model fit and the original sinusoidal function). Note that the result of averaging many solutions for the complex model with  $M = 25$  is a very good fit to the regression function, which suggests that averaging may be a beneficial procedure. Indeed, a weighted averaging of multiple solutions lies at the heart of a Bayesian approach, although the averaging is with respect to the posterior distribution of parameters, not with respect to multiple data sets.

We can also examine the bias-variance trade-off quantitatively for this example. The average prediction is estimated from

$$\bar{y}(x) = \frac{1}{L} \sum_{l=1}^L y^{(l)}(x) \quad (3.45)$$

and the integrated squared bias and integrated variance are then given by

$$(\text{bias})^2 = \frac{1}{N} \sum_{n=1}^N \{\bar{y}(x_n) - h(x_n)\}^2 \quad (3.46)$$

$$\text{variance} = \frac{1}{N} \sum_{n=1}^N \frac{1}{L} \sum_{l=1}^L \{y^{(l)}(x_n) - \bar{y}(x_n)\}^2 \quad (3.47)$$

where the integral over  $x$  weighted by the distribution  $p(x)$  is approximated by a finite sum over data points drawn from that distribution. These quantities, along with their sum, are plotted as a function of  $\ln \lambda$  in Figure 3.6. We see that small values of  $\lambda$  allow the model to become finely tuned to the noise on each individual

data set leading to large variance. Conversely, a large value of  $\lambda$  pulls the weight parameters towards zero leading to large bias.

Although the bias-variance decomposition may provide some interesting insights into the model complexity issue from a frequentist perspective, it is of limited practical value, because the bias-variance decomposition is based on averages with respect to ensembles of data sets, whereas in practice we have only the single observed data set. If we had a large number of independent training sets of a given size, we would be better off combining them into a single large training set, which of course would reduce the level of over-fitting for a given model complexity.

Given these limitations, we turn in the next section to a Bayesian treatment of linear basis function models, which not only provides powerful insights into the issues of over-fitting but which also leads to practical techniques for addressing the question model complexity.

### 3.3. Bayesian Linear Regression

---

In our discussion of maximum likelihood for setting the parameters of a linear regression model, we have seen that the effective model complexity, governed by the number of basis functions, needs to be controlled according to the size of the data set. Adding a regularization term to the log likelihood function means the effective model complexity can then be controlled by the value of the regularization coefficient, although the choice of the number and form of the basis functions is of course still important in determining the overall behaviour of the model.

This leaves the issue of deciding the appropriate model complexity for the particular problem, which cannot be decided simply by maximizing the likelihood function, because this always leads to excessively complex models and over-fitting. Independent hold-out data can be used to determine model complexity, as discussed in Section 1.3, but this can be both computationally expensive and wasteful of valuable data. We therefore turn to a Bayesian treatment of linear regression, which will avoid the over-fitting problem of maximum likelihood, and which will also lead to automatic methods of determining model complexity using the training data alone. Again, for simplicity we will focus on the case of a single target variable  $t$ . Extension to multiple target variables is straightforward and follows the discussion of Section 3.1.5.

#### 3.3.1 Parameter distribution

We begin our discussion of the Bayesian treatment of linear regression by introducing a prior probability distribution over the model parameters  $\mathbf{w}$ . For the moment, we shall treat the noise precision parameter  $\beta$  as a known constant. First note that the likelihood function  $p(\mathbf{t}|\mathbf{w})$  defined by (3.10) is the exponential of a quadratic function of  $\mathbf{w}$ . The corresponding conjugate prior is therefore given by a Gaussian distribution of the form

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0) \quad (3.48)$$

having mean  $\mathbf{m}_0$  and covariance  $\mathbf{S}_0$ .

Next we compute the posterior distribution, which is proportional to the product of the likelihood function and the prior. Due to the choice of a conjugate Gaussian prior distribution, the posterior will also be Gaussian. We can evaluate this distribution by the usual procedure of completing the square in the exponential, and then finding the normalization coefficient using the standard result for a normalized Gaussian. However, we have already done the necessary work in deriving the general result (2.116), which allows us to write down the posterior distribution directly in the form

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N) \quad (3.49)$$

where

$$\mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t}) \quad (3.50)$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi. \quad (3.51)$$

Note that because the posterior distribution is Gaussian, its mode coincides with its mean. Thus the maximum posterior weight vector is simply given by  $\mathbf{w}_{\text{MAP}} = \mathbf{m}_N$ . If we consider an infinitely broad prior  $\mathbf{S}_0 = \alpha^{-1} \mathbf{I}$  with  $\alpha \rightarrow 0$ , the mean  $\mathbf{m}_N$  of the posterior distribution reduces to the maximum likelihood value  $\mathbf{w}_{\text{ML}}$  given by (3.15). Similarly, if  $N = 0$ , then the posterior distribution reverts to the prior. Furthermore, if data points arrive sequentially, then the posterior distribution at any stage acts as the prior distribution for the subsequent data point, such that the new posterior distribution is again given by (3.49).

For the remainder of this chapter, we shall consider a particular form of Gaussian prior in order to simplify the treatment. Specifically, we consider a zero-mean isotropic Gaussian governed by a single precision parameter  $\alpha$  so that

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1} \mathbf{I}) \quad (3.52)$$

and the corresponding posterior distribution over  $\mathbf{w}$  is then given by (3.49) with

$$\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t} \quad (3.53)$$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi. \quad (3.54)$$

The log of the posterior distribution is given by the sum of the log likelihood and the log of the prior and, as a function of  $\mathbf{w}$ , takes the form

$$\ln p(\mathbf{w}|\mathbf{t}) = -\frac{\beta}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const.} \quad (3.55)$$

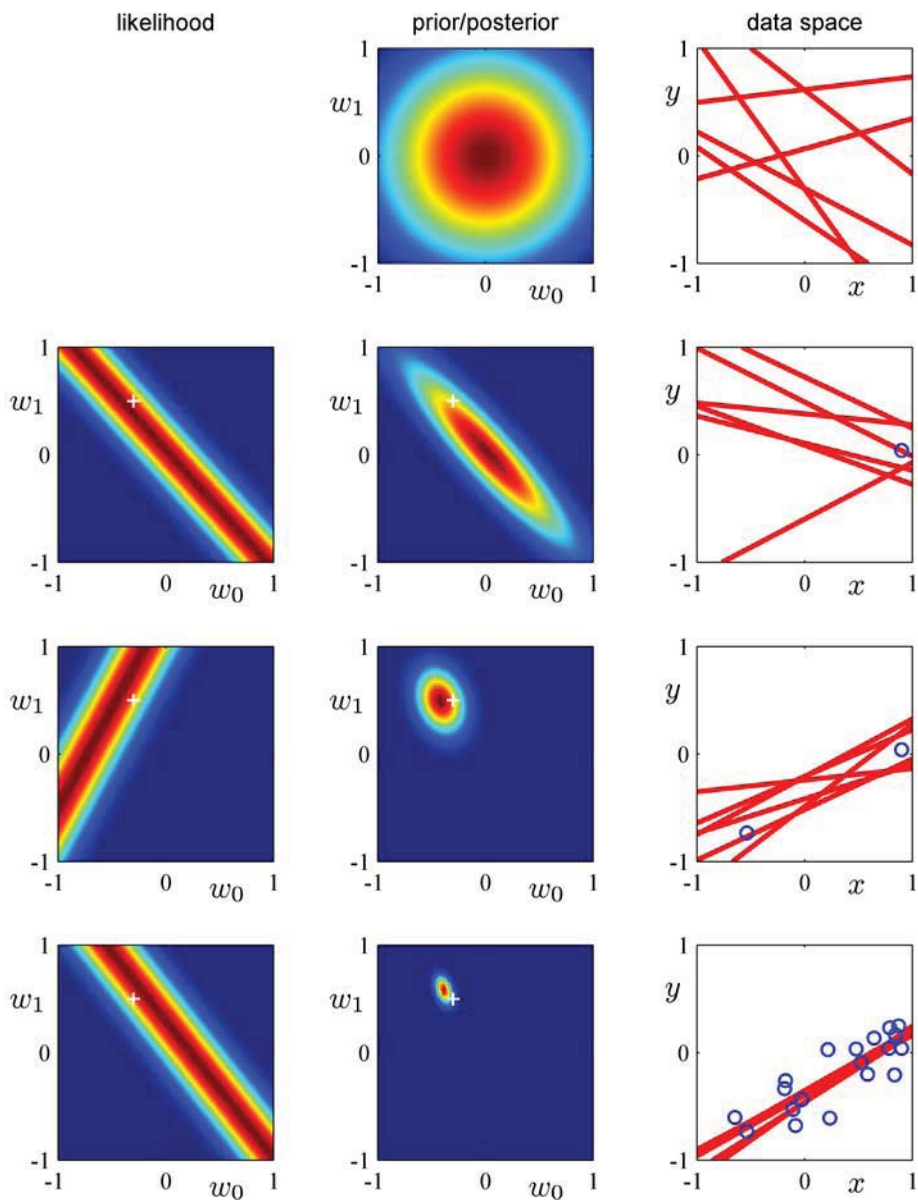
Maximization of this posterior distribution with respect to  $\mathbf{w}$  is therefore equivalent to the minimization of the sum-of-squares error function with the addition of a quadratic regularization term, corresponding to (3.27) with  $\lambda = \alpha/\beta$ .

We can illustrate Bayesian learning in a linear basis function model, as well as the sequential update of a posterior distribution, using a simple example involving straight-line fitting. Consider a single input variable  $x$ , a single target variable  $t$  and

*Exercise 3.7*

*Exercise 3.8*

a linear model of the form  $y(x, \mathbf{w}) = w_0 + w_1x$ . Because this has just two adaptive parameters, we can plot the prior and posterior distributions directly in parameter space. We generate synthetic data from the function  $f(x, \mathbf{a}) = a_0 + a_1x$  with parameter values  $a_0 = -0.3$  and  $a_1 = 0.5$  by first choosing values of  $x_n$  from the uniform distribution  $\mathcal{U}(x|-1, 1)$ , then evaluating  $f(x_n, \mathbf{a})$ , and finally adding Gaussian noise with standard deviation of 0.2 to obtain the target values  $t_n$ . Our goal is to recover the values of  $a_0$  and  $a_1$  from such data, and we will explore the dependence on the size of the data set. We assume here that the noise variance is known and hence we set the precision parameter to its true value  $\beta = (1/0.2)^2 = 25$ . Similarly, we fix the parameter  $\alpha$  to 2.0. We shall shortly discuss strategies for determining  $\alpha$  and  $\beta$  from the training data. Figure 3.7 shows the results of Bayesian learning in this model as the size of the data set is increased and demonstrates the sequential nature of Bayesian learning in which the current posterior distribution forms the prior when a new data point is observed. It is worth taking time to study this figure in detail as it illustrates several important aspects of Bayesian inference. The first row of this figure corresponds to the situation before any data points are observed and shows a plot of the prior distribution in  $\mathbf{w}$  space together with six samples of the function  $y(x, \mathbf{w})$  in which the values of  $\mathbf{w}$  are drawn from the prior. In the second row, we see the situation after observing a single data point. The location  $(x, t)$  of the data point is shown by a blue circle in the right-hand column. In the left-hand column is a plot of the likelihood function  $p(t|x, \mathbf{w})$  for this data point as a function of  $\mathbf{w}$ . Note that the likelihood function provides a soft constraint that the line must pass close to the data point, where close is determined by the noise precision  $\beta$ . For comparison, the true parameter values  $a_0 = -0.3$  and  $a_1 = 0.5$  used to generate the data set are shown by a white cross in the plots in the left column of Figure 3.7. When we multiply this likelihood function by the prior from the top row, and normalize, we obtain the posterior distribution shown in the middle plot on the second row. Samples of the regression function  $y(x, \mathbf{w})$  obtained by drawing samples of  $\mathbf{w}$  from this posterior distribution are shown in the right-hand plot. Note that these sample lines all pass close to the data point. The third row of this figure shows the effect of observing a second data point, again shown by a blue circle in the plot in the right-hand column. The corresponding likelihood function for this second data point alone is shown in the left plot. When we multiply this likelihood function by the posterior distribution from the second row, we obtain the posterior distribution shown in the middle plot of the third row. Note that this is exactly the same posterior distribution as would be obtained by combining the original prior with the likelihood function for the two data points. This posterior has now been influenced by two data points, and because two points are sufficient to define a line this already gives a relatively compact posterior distribution. Samples from this posterior distribution give rise to the functions shown in red in the third column, and we see that these functions pass close to both of the data points. The fourth row shows the effect of observing a total of 20 data points. The left-hand plot shows the likelihood function for the 20<sup>th</sup> data point alone, and the middle plot shows the resulting posterior distribution that has now absorbed information from all 20 observations. Note how the posterior is much sharper than in the third row. In the limit of an infinite number of data points, the



**Figure 3.7** Illustration of sequential Bayesian learning for a simple linear model of the form  $y(x, \mathbf{w}) = w_0 + w_1 x$ . A detailed description of this figure is given in the text.

posterior distribution would become a delta function centred on the true parameter values, shown by the white cross.

Other forms of prior over the parameters can be considered. For instance, we can generalize the Gaussian prior to give

$$p(\mathbf{w}|\alpha) = \left[ \frac{q}{2} \left( \frac{\alpha}{2} \right)^{1/q} \frac{1}{\Gamma(1/q)} \right]^M \exp \left( -\frac{\alpha}{2} \sum_{j=1}^M |w_j|^q \right) \quad (3.56)$$

in which  $q = 2$  corresponds to the Gaussian distribution, and only in this case is the prior conjugate to the likelihood function (3.10). Finding the maximum of the posterior distribution over  $\mathbf{w}$  corresponds to minimization of the regularized error function (3.29). In the case of the Gaussian prior, the mode of the posterior distribution was equal to the mean, although this will no longer hold if  $q \neq 2$ .

### 3.3.2 Predictive distribution

In practice, we are not usually interested in the value of  $\mathbf{w}$  itself but rather in making predictions of  $t$  for new values of  $\mathbf{x}$ . This requires that we evaluate the *predictive distribution* defined by

$$p(t|\mathbf{t}, \alpha, \beta) = \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w} \quad (3.57)$$

in which  $\mathbf{t}$  is the vector of target values from the training set, and we have omitted the corresponding input vectors from the right-hand side of the conditioning statements to simplify the notation. The conditional distribution  $p(t|\mathbf{x}, \mathbf{w}, \beta)$  of the target variable is given by (3.8), and the posterior weight distribution is given by (3.49). We see that (3.57) involves the convolution of two Gaussian distributions, and so making use of the result (2.115) from Section 8.1.4, we see that the predictive distribution takes the form

$$p(t|\mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t|\mathbf{m}_N^T \phi(\mathbf{x}), \sigma_N^2(\mathbf{x})) \quad (3.58)$$

where the variance  $\sigma_N^2(\mathbf{x})$  of the predictive distribution is given by

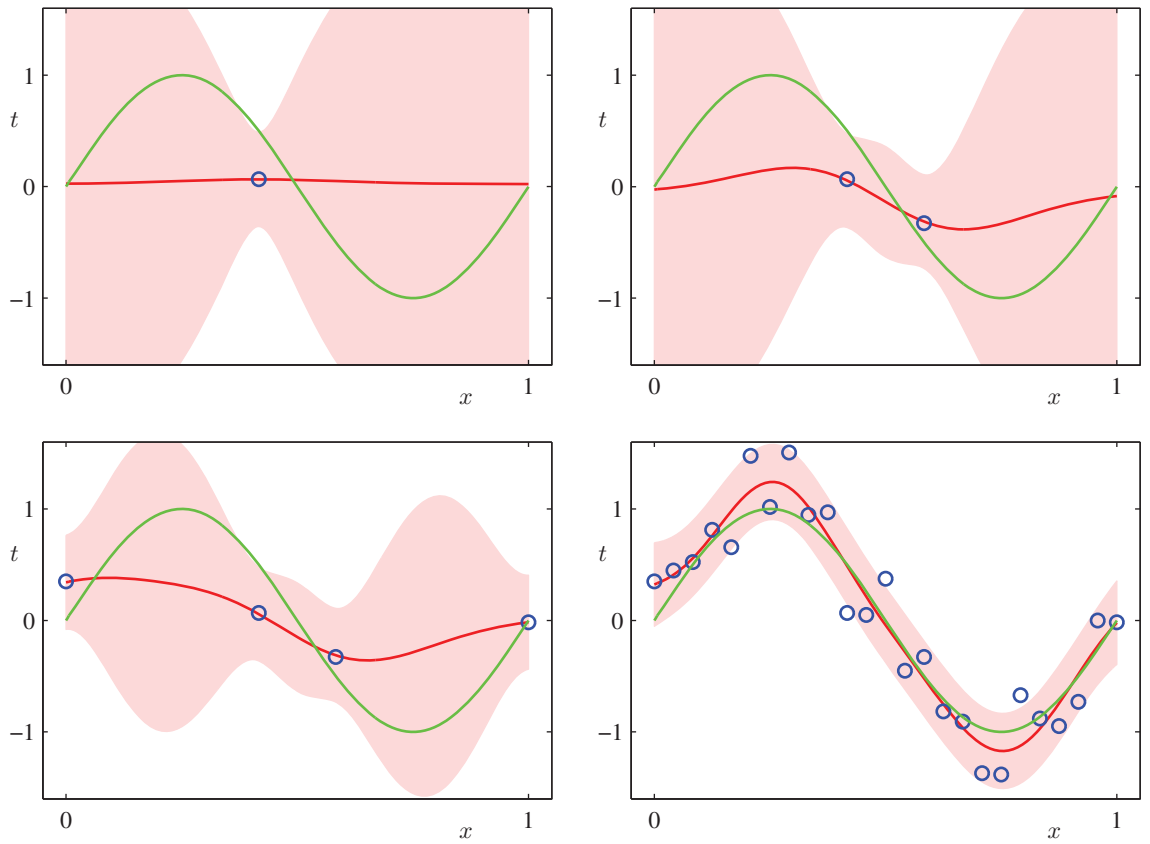
$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \phi(\mathbf{x})^T \mathbf{S}_N \phi(\mathbf{x}). \quad (3.59)$$

The first term in (3.59) represents the noise on the data whereas the second term reflects the uncertainty associated with the parameters  $\mathbf{w}$ . Because the noise process and the distribution of  $\mathbf{w}$  are independent Gaussians, their variances are additive. Note that, as additional data points are observed, the posterior distribution becomes narrower. As a consequence it can be shown (Qazaz *et al.*, 1997) that  $\sigma_{N+1}^2(\mathbf{x}) \leq \sigma_N^2(\mathbf{x})$ . In the limit  $N \rightarrow \infty$ , the second term in (3.59) goes to zero, and the variance of the predictive distribution arises solely from the additive noise governed by the parameter  $\beta$ .

As an illustration of the predictive distribution for Bayesian linear regression models, let us return to the synthetic sinusoidal data set of Section 1.1. In Figure 3.8,

*Exercise 3.10*

*Exercise 3.11*

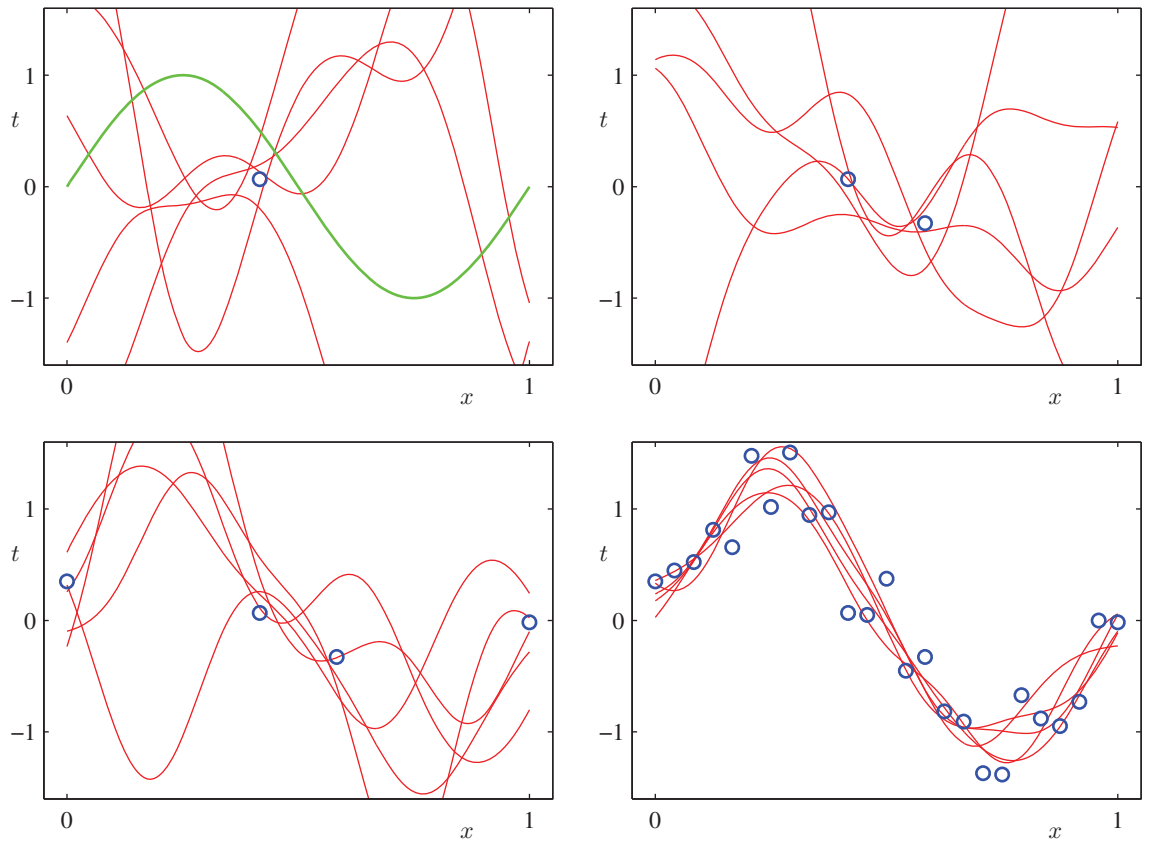


**Figure 3.8** Examples of the predictive distribution (3.58) for a model consisting of 9 Gaussian basis functions of the form (3.4) using the synthetic sinusoidal data set of Section 1.1. See the text for a detailed discussion.

we fit a model comprising a linear combination of Gaussian basis functions to data sets of various sizes and then look at the corresponding posterior distributions. Here the green curves correspond to the function  $\sin(2\pi x)$  from which the data points were generated (with the addition of Gaussian noise). Data sets of size  $N = 1$ ,  $N = 2$ ,  $N = 4$ , and  $N = 25$  are shown in the four plots by the blue circles. For each plot, the red curve shows the mean of the corresponding Gaussian predictive distribution, and the red shaded region spans one standard deviation either side of the mean. Note that the predictive uncertainty depends on  $x$  and is smallest in the neighbourhood of the data points. Also note that the level of uncertainty decreases as more data points are observed.

The plots in Figure 3.8 only show the point-wise predictive variance as a function of  $x$ . In order to gain insight into the covariance between the predictions at different values of  $x$ , we can draw samples from the posterior distribution over  $\mathbf{w}$ , and then plot the corresponding functions  $y(x, \mathbf{w})$ , as shown in Figure 3.9.





**Figure 3.9** Plots of the function  $y(x, \mathbf{w})$  using samples from the posterior distributions over  $\mathbf{w}$  corresponding to the plots in Figure 3.8.

If we used localized basis functions such as Gaussians, then in regions away from the basis function centres, the contribution from the second term in the predictive variance (3.59) will go to zero, leaving only the noise contribution  $\beta^{-1}$ . Thus, the model becomes very confident in its predictions when extrapolating outside the region occupied by the basis functions, which is generally an undesirable behaviour. This problem can be avoided by adopting an alternative Bayesian approach to regression known as a Gaussian process.

#### Section 6.4

*Exercise 3.12*

*Exercise 3.13*

Note that, if both  $\mathbf{w}$  and  $\beta$  are treated as unknown, then we can introduce a conjugate prior distribution  $p(\mathbf{w}, \beta)$  that, from the discussion in Section 2.3.6, will be given by a Gaussian-gamma distribution (Denison *et al.*, 2002). In this case, the predictive distribution is a Student's t-distribution.