

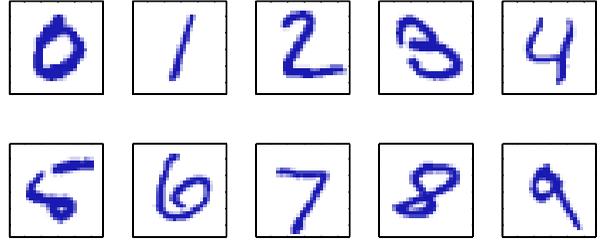
# 1

## 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

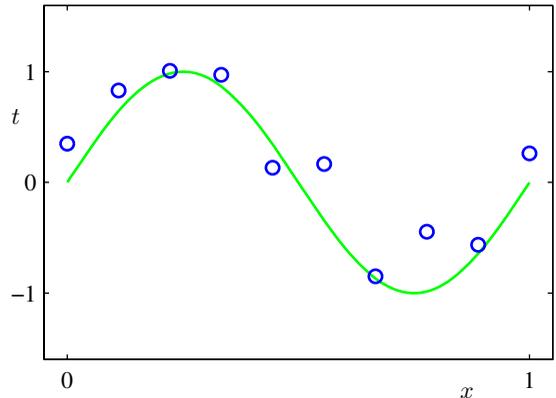
also preserve useful discriminatory information enabling faces to be distinguished from non-faces. These features are then used as the inputs to the pattern recognition algorithm. For instance, the average value of the image intensity over a rectangular subregion can be evaluated extremely efficiently (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

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

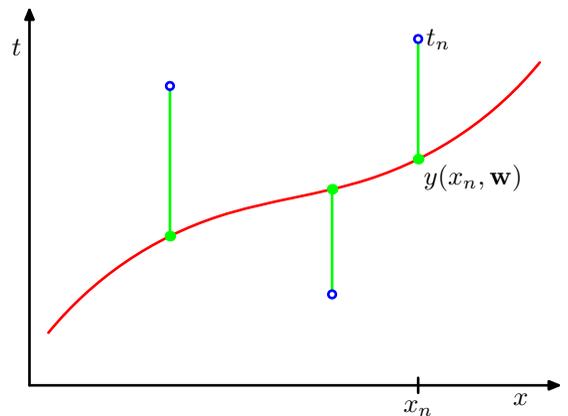
where  $M$  is the *order* of the polynomial, and  $x^j$  denotes  $x$  raised to the power of  $j$ . The polynomial coefficients  $w_0, \dots, w_M$  are collectively denoted by the vector  $\mathbf{w}$ . Note that, although the polynomial function  $y(x, \mathbf{w})$  is a nonlinear function of  $x$ , 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.

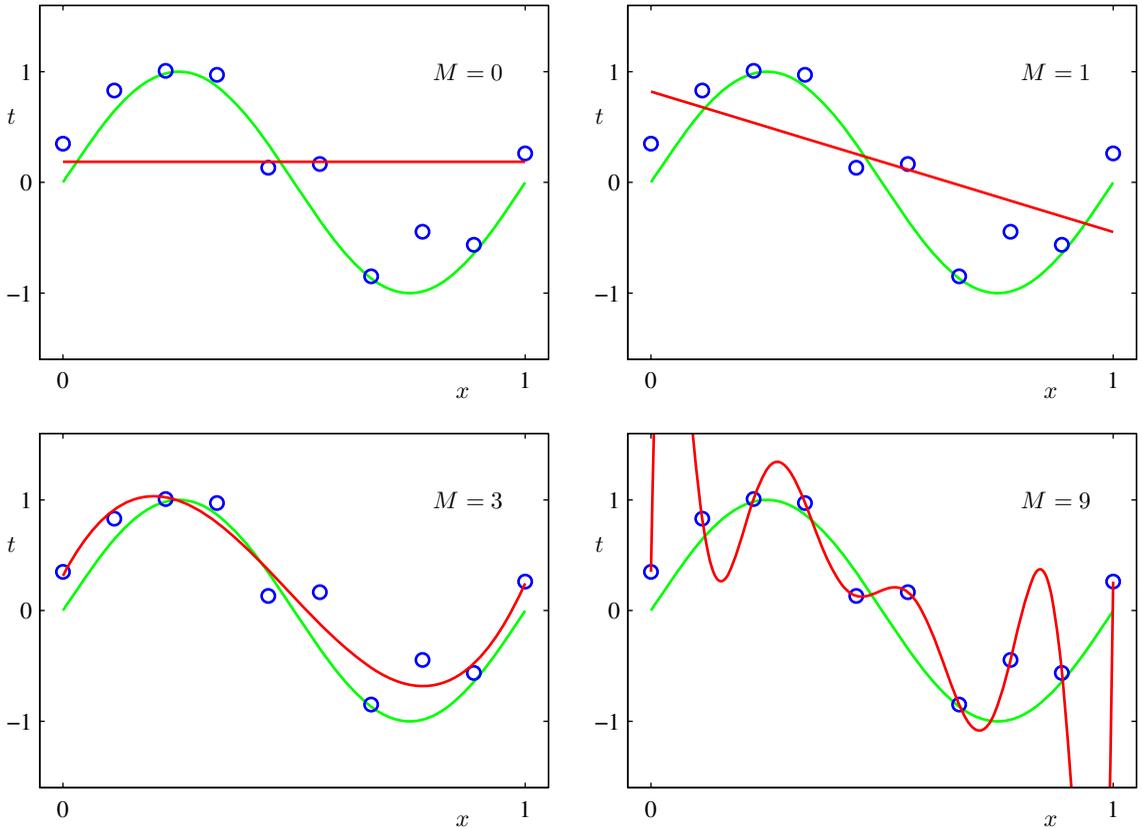
We can solve the curve fitting problem by choosing the value of  $\mathbf{w}$  for which  $E(\mathbf{w})$  is as small as possible. Because the error function is a quadratic function of the coefficients  $\mathbf{w}$ , its derivatives with respect to the coefficients will be linear in the elements of  $\mathbf{w}$ , 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



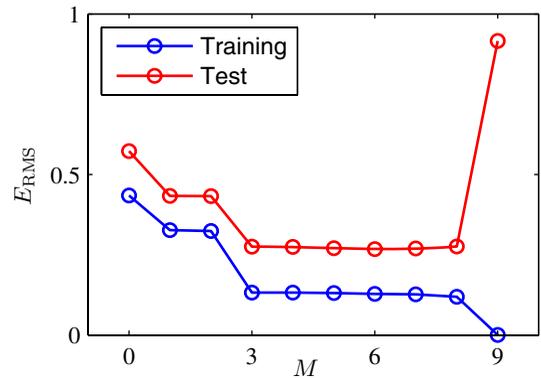
**Figure 1.4** Plots of polynomials having various orders  $M$ , shown as red curves, fitted to the data set shown in Figure 1.2.

(RMS) error defined by

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

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  $M$ .



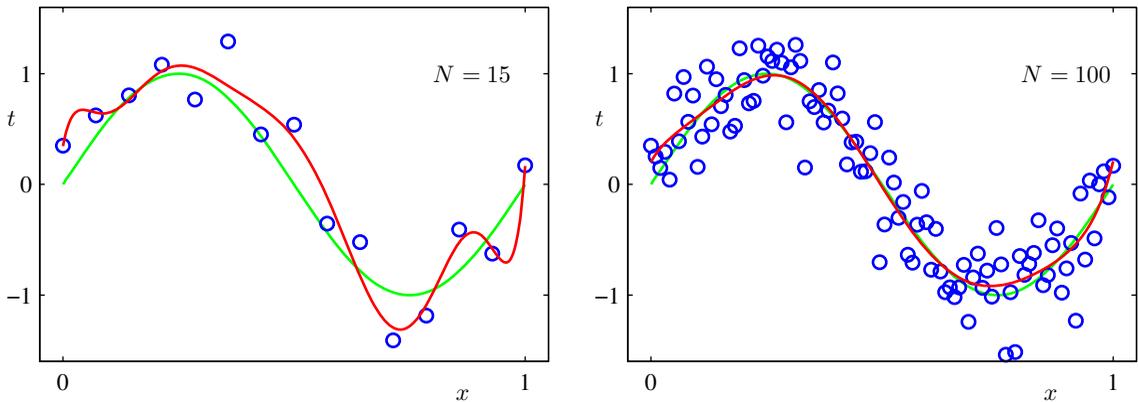
For  $M = 9$ , the training set error goes to zero, as we might expect because this polynomial contains 10 degrees of freedom corresponding to the 10 coefficients  $w_0, \dots, w_9$ , and so can be tuned exactly to the 10 data points in the training set. However, the test set error has become very large 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		0.31	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.

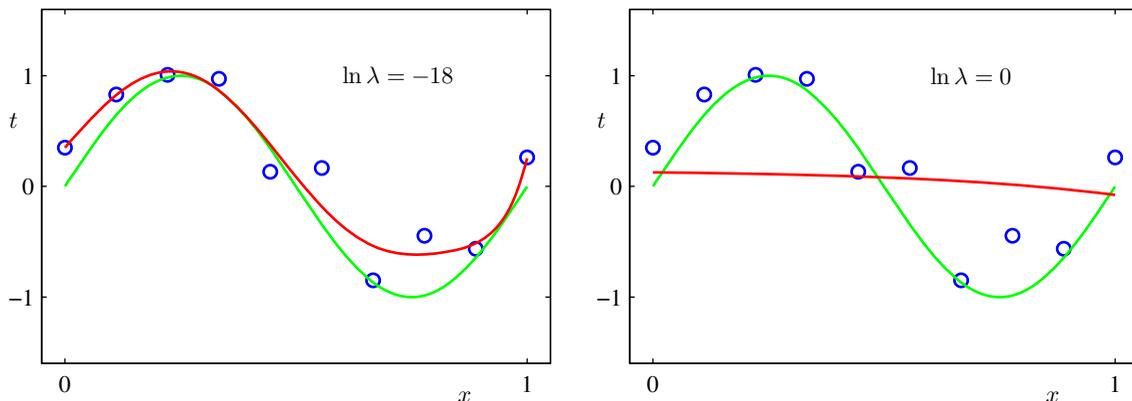
ing polynomial function matches each of the data points exactly, but between data points (particularly near the ends of the range) the function exhibits the large oscillations observed in Figure 1.4. Intuitively, what is happening is that the more flexible polynomials with larger values of  $M$  are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

### Exercise 1.2

**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.

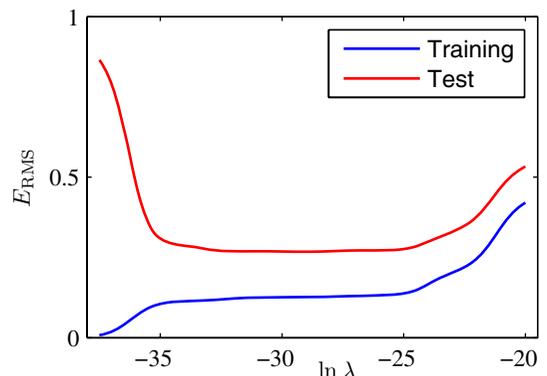
The impact of the regularization term on the generalization error can be seen by plotting the value of the RMS error (1.3) for both training and test sets against  $\ln \lambda$ , as shown in Figure 1.8. We see that in effect  $\lambda$  now controls the effective complexity of the model and hence determines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

### Section 1.3

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



give us some important insights into the concepts we have introduced in the context of polynomial curve fitting and will allow us to extend these to more complex situations.

## 1.2. Probability Theory

---

A key concept in the field of pattern recognition is that of uncertainty. It arises both through noise on measurements, as well as through the finite size of data sets. Probability theory provides a consistent framework for the quantification and manipulation of uncertainty and forms one of the central foundations for pattern recognition. When combined with decision theory, discussed in Section 1.5, it allows us to make optimal predictions given all the information available to us, even though that information may be incomplete or ambiguous.

We will introduce the basic concepts of probability theory by considering a simple example. Imagine we have two boxes, one red and one blue, and in the red box we have 2 apples and 6 oranges, and in the blue box we have 3 apples and 1 orange. This is illustrated in Figure 1.9. Now suppose we randomly pick one of the boxes and from that box we randomly select an item of fruit, and having observed which sort of fruit it is we replace it in the box from which it came. We could imagine repeating this process many times. Let us suppose that in so doing we pick the red box 40% of the time and we pick the blue box 60% of the time, and that when we remove an item of fruit from a box we are equally likely to select any of the pieces of fruit in the box.

In this example, the identity of the box that will be chosen is a random variable, which we shall denote by  $B$ . This random variable can take one of two possible values, namely  $r$  (corresponding to the red box) or  $b$  (corresponding to the blue box). Similarly, the identity of the fruit is also a random variable and will be denoted by  $F$ . It can take either of the values  $a$  (for apple) or  $o$  (for orange).

To begin with, we shall define the probability of an event to be the fraction of times that event occurs out of the total number of trials, in the limit that the total number of trials goes to infinity. Thus the probability of selecting the red box is  $4/10$

**Figure 1.9** We use a simple example of two coloured boxes each containing fruit (apples shown in green and oranges shown in orange) to introduce the basic ideas of probability.

