1 Introduction

1.1 What are functional data?

Figure 1.1 provides a prototype for the type of data that we shall consider. It shows the heights of 10 girls measured at a set of 31 ages in the Berkeley Growth Study (Tuddenham and Snyder, 1954). The ages are not equally spaced; there are four measurements while the child is one year old, annual measurements from two to eight years, followed by heights measured biannually. Although great care was taken in the measurement process, there is an uncertainty or noise in height values that has a standard deviation of about three millimeters. Even though each record involves only discrete values, these values reflect a smooth variation in height that could be assessed, in principle, as often as desired, and is therefore a height function. Thus, the data consist of a sample of 10 functional observations $\text{Height}_i(t)$.

There are features in this data too subtle to see in this type of plot. Figure 1.2 displays the acceleration curves $D^2 \text{Height}_i$ estimated from these data by Ramsay, Bock and Gasser (1995) using a technique discussed in Chapter 5. We use the notation D for differentiation, as in

$$D^2$$
Height = $\frac{d^2$ Height}{dt^2}.

In Figure 1.2 the pubertal growth spurt shows up as a pulse of strong positive acceleration followed by sharp negative deceleration. But most records also show a bump at around six years that is termed the mid-spurt. We therefore conclude that some of the variation from curve to curve can be explained at the level of certain derivatives. The fact that derivatives



Figure 1.1. The heights of 10 girls measured at 31 ages. The circles indicate the unequally spaced ages of measurement.



Figure 1.2. The estimated accelerations of height for 10 girls, measured in centimeters per year. The heavy dashed line is the cross-sectional mean, and is a rather poor summary of the curves.

are of interest is further reason to think of the records as functions, rather than vectors of observations in discrete time.

The ages themselves must also play an explicit role in our analysis, because they are not equally spaced. Although it might be mildly interesting to correlate heights at ages 9, 10 and 10.5, this would not take account of the fact that we expect the correlation for two ages separated by only half a year to be higher than that for a separation of one year. Indeed, although in this particular example the ages at which the observations are taken are nominally the same for each girl, there is no real need for this to be so; in general, the points at which the functions are observed may well vary from one record to another.

The replication of these height curves invites an exploration of the ways in which the curves vary. This is potentially complex. For example, the rapid growth during puberty is visible in all curves, but both the timing and the intensity of pubertal growth differ from girl to girl. Some type of principal components analysis would undoubtedly be helpful, but we must adapt the procedure to take account of the unequal age spacing and the smoothness of the underlying height functions. One objective might be to separate variation in timing of significant growth events, such as the pubertal growth spurt, from variation in the intensity of growth.

Not all functional data involves independent replications; we often have to work with a single long record. Figure 1.3 shows an important economic indicator, the nondurable goods manufacturing index for the United States. Data like these often show variation as multiple levels. There is a tendency for the index to show geometric or exponential increase over the whole century. But at a finer scale, we see departures from this trend due to the depression, World War II, the end of the Vietnam War and other more localized events. Moreover, at an even finer scale, there is a marked annual variation, and we can wonder whether this *seasonal trend* itself shows some longer term changes. Although there are no independent replications here, there is still a lot of repetition of information that we can exploit to obtain stable estimates of interesting curve features.

Functional data also arise as input/output pairs, such as in the data in Figure 1.4 collected at an oil refinery in Texas. The amount of a petroleum product at a certain level in a distillation column or cracking tower, shown in the top panel, reacts to the change in the flow of a vapor into the tray, shown in the bottom panel, at that level. How can we characterize this dependency?



Figure 1.3. The nondurable goods manufacturing index for the United States.



Figure 1.4. The top panel shows 193 measurements of tray level in a distillation column in an oil refinery. The bottom panel shows the flow of a vapor into the tray during an experiment.

1.2 Functional models for nonfunctional data

The data examples above seem to deserve the label "functional" since they so clearly reflect the smooth curves that we assume generated them. But not all data subject to a functional data analysis are themselves functional.

Consider the problem of estimating a probability density function p to describe the distribution of a sample of observations x_1, \ldots, x_n . The classic approach to this problem is to propose, after considering basic principles and closely studying the data, a *parametric model* with values $p(x|\theta)$ defined by a fixed and usually small number of parameters in the vector θ . For example, we might consider the normal distribution as appropriate for the data, so that $\theta = (\mu, \sigma^2)'$. The parameters themselves are usually chosen to be descriptors of the shape of the density, as in location and spread for the normal density, and are therefore the focus of the analysis.

But suppose that we do not want to assume in advance one of the many textbook density functions because, perhaps, none of them seem to capture features of the behavior of the data that we can see in histograms and other graphical displays. *Nonparametric density* estimation methods assume only smoothness, and permit as much flexibility in the estimated p(x) as the data require. To be sure, parameters are often involved, as in the density estimation method of Chapter 6, but the number of parameters is not fixed in advance of the data analysis, and our attention is focussed on the function p itself rather than on the estimated parameter values. Much of the technology for estimation of smooth *functional parameters* was originally developed and honed in the density estimation context, and Silverman (1986) can be consulted for further details.

Psychometrics or mental test theory also relies heavily on functional models for seemingly nonfunctional data. The data are usually zeros and ones indicating unsuccessful and correct answers to test items, but the model consists of a set of *item response functions*, one per test item, displaying the smooth relationship between the probability of success on an item and a presumed latent ability continuum. Figure 1.5 shows three such functional parameters for a test of mathematics estimated by the functional data analytic methods reported in Rossi, Wang and Ramsay (2002).

1.3 Some functional data analyses

Data in many fields come to us through a process naturally described as functional. To turn to a completely different context, consider Figure 1.6, where the mean monthly temperatures for four Canadian weather stations are plotted. It also shows estimates of the corresponding smooth temperature functions presumed to generate the observations. Montreal, with the warmest summer temperature, has a temperature pattern that appears to



Figure 1.5. Each panel shows an item response function relating an examinee's position θ on a latent ability continuum to the probability of getting a test item in a mathematics test correct.



Figure 1.6. Mean monthly temperatures for the Canadian weather stations. In descending order of the temperatures at the start of the year, the stations are Prince Rupert, Montreal, Edmonton, and Resolute.

be nicely sinusoidal. Edmonton, with the next warmest summer temperature, seems to have some distinctive departures from sinusoidal variation that might call for explanation. The marine climate of Prince Rupert is evident in the small amount of annual variation in temperature, and Resolute has bitterly cold but strongly sinusoidal temperature.

One expects temperature to be primarily sinusoidal in character, and certainly periodic over the annual cycle. There is some variation in phase,



Figure 1.7. The result of applying the differential operator $L = (\pi/6)^2 D + D^3$ to the estimated temperature functions in Figure 1.6. If the variation in temperature were purely sinusoidal, these curves would be exactly zero.

because the coldest day of the year seems to be later in Montreal and Resolute than in Edmonton and Prince Rupert. Consequently, a model of the form

$$\operatorname{Temp}_{i}(t) \approx c_{i1} + c_{i2}\sin(\pi t/6) + c_{i3}\cos(\pi t/6)$$
 (1.1)

should do rather nicely for these data, where Temp_i is the temperature function for the *i*th weather station, and (c_{i1}, c_{i2}, c_{i3}) is a vector of three parameters associated with that station.

In fact, there are clear departures from sinusoidal or simple harmonic behavior. One way to see this is to compute the function

$$L\text{Temp} = (\pi/6)^2 D\text{Temp} + D^3\text{Temp}.$$
 (1.2)

As we have already noted in Section 1.1, the notation D^m Temp means "take the *m*th derivative of function Temp," and the notation *L*Temp stands for the function which results from applying the linear differential operator $L = (\pi/6)^2 D + D^3$ to the function Temp. The resulting function, *L*Temp, is often called a *forcing function*. Now, if a temperature function is truly sinusoidal, then *L*Temp should be exactly zero, as it would be for any function of the form (1.1). That is, it would conform to the *differential equation*

$$D^3$$
Temp = $-(\pi/6)^2 D$ Temp.

But Figure 1.7 indicates that the functions $L\text{Temp}_i$ display systematic features that are especially strong in the spring and autumn months. Put



Figure 1.8. The angles in the sagittal plane formed by the hip and by the knee as 39 children go through a gait cycle. The interval [0, 1] is a single cycle, and the dotted curves show the periodic extension of the data beyond either end of the cycle.

another way, temperature at a particular weather station can be described as the solution of the *nonhomogeneous* differential equation corresponding to LTemp = u, where the forcing function u can be viewed as input from outside of the system, or an exogenous influence. Meteorologists suggest, for example, that these spring and autumn effects are partly due to the change in the reflectance of land when snow or ice melts, and this would be consistent with the fact that the least sinusoidal records are associated with continental stations well separated from large bodies of water.

Here, the point is that we may often find it interesting to remove effects of a simple character by applying a differential operator, rather than by simply subtracting them. This exploits the intrinsic smoothness in the process, and long experience in the natural and engineering sciences suggests that this may get closer to the underlying driving forces at work than just adding and subtracting effects, as one routinely does in multivariate data analysis. We will consider this idea in depth beginning with Chapter 18.

Functional data are often multivariate in a different sense. Our third example is in Figure 1.8. The Motion Analysis Laboratory at Children's Hospital, San Diego, collected these data, which consist of the angles formed by the hip and knee of each of 39 children over each child's gait cycle. See Olshen et al. (1989) for full details. Time is measured in terms of the individual gait cycle, so that every curve is given for values of t in [0, 1]. The cycle begins and ends at the point where the heel of the limb under observation strikes the ground. Both sets of functions are periodic, and are plotted as dotted curves somewhat beyond the interval for clarity. We see that the knee shows a two-phase process, while the hip motion is single-phase. What is harder to see is how the two joints interact; of course the figure does not indicate which hip curve is paired with which knee curve, and among many other things this example demonstrates the need for graphical ingenuity in functional data analysis.

Figure 1.9 shows the gait cycle for a single child by plotting knee angle against hip angle as time progresses round the cycle. The periodic nature of the process implies that this forms a closed curve. Also shown for reference purposes is the same relationship for the average across the 39 children. Now we see an interesting feature: a cusp occurring at the heel strike. The angular velocity is clearly visible in terms of the spacing between numbers, and it varies considerably as the cycle proceeds. The child whose gait is represented by the solid curve differs from the average in two principal ways. First, the portion of the gait pattern in the C–D part of the cycle shows an exaggeration of movement relative to the average, and second, in the part of the cycle where the hip is most bent, the amount by which the hip is bent is markedly less than average; interestingly, this is not accompanied by any strong effect on the knee angle. The overall shape of the cycle for the particular child is rather different from the average. The exploration of variability in these functional data must focus on features such as these.

Finally, in this introduction to types of functional data, we must not forget that they may come to our attention as full-blown functions, so that each record may consist of functions observed, for all practical purposes, everywhere. Sophisticated on-line sensing and monitoring equipment is now routinely used in research in medicine, seismology, meteorology, physiology, and many other fields.

1.4 The goals of functional data analysis

The goals of functional data analysis are essentially the same as those of any other branch of statistics. They include the following aims:

- to represent the data in ways that aid further analysis
- to display the data so as to highlight various characteristics
- to study important sources of pattern and variation among the data
- to explain variation in an outcome or dependent variable by using input or independent variable information



Figure 1.9. Solid line: The angles in the sagittal plane formed by the hip and by the knee for a single child plotted against each other. Dotted line: The corresponding plot for the average across children. The points indicate 20 equally spaced time points in the gait cycle, and the letters are plotted at intervals of one-fifth of the cycle, with A marking the heel strike.

• to compare two or more sets of data with respect to certain types of variation, where two sets of data can contain different sets of replicates of the same functions, or different functions for a common set of replicates.

Subsequent chapters explore each of these themes, and they are introduced only briefly here.

Each of these activities can be conducted with techniques appropriate to certain goals. Another way to characterize the strategy in a data analysis is as *exploratory*, *confirmatory*, or *predictive*. In exploratory mode, the questions put to the data tend to be rather open-ended in the sense that one expects the right technique to reveal new and interesting aspects of the data, as well as to shed light on known and obvious features. Exploratory investigations tend to consider only the data at hand, with less concern for statements about larger issues such as characteristics of populations or events not observed in the data. Confirmatory analyses, on the other hand, tend to be inferential and to be determined by specific questions about the data. Some type of structure is assumed to be present in the data, and one wants to know whether certain specific statements or hypotheses can be considered confirmed by the data. The dividing line between exploratory and confirmatory analyses tends to be the extent to which probability theory is used, in the sense that most confirmatory analyses are summarized by one or more probability statements. Predictive studies are somewhat less common, and focus on using the data at hand to make a statement about unobserved states, such as the future.

Functional principal components and canonical correlation analyses are mainly exploratory methods, and are covered in Chapters 8 to 11. Functional linear models, on the other hand, are often used in a confirmatory way, and in 12 to 17 we introduce confidence interval estimation. In general, prediction is beyond our scope, and is only considered here and there.

1.5 The first steps in a functional data analysis

1.5.1 Data representation: smoothing and interpolation

Assuming that a functional datum for replication i arrives as a set of discrete measured values, y_{i1}, \ldots, y_{in} , the first task is to convert these values to a function x_i with values $x_i(t)$ computable for any desired argument value t. If the discrete values are assumed to be errorless, then the process is *interpolation*, but if they have some observational error that needs removing, then the conversion from discrete data to functions may involve smoothing.

Chapters 3 to 6 offer a survey of these procedures. The *roughness penalty* smoothing method discussed in Chapter 5 will be used much more broadly in many contexts throughout the book, not merely for the purpose of estimating a function from a set of observed values. The daily precipitation data for Prince Rupert, one of the wettest places on the continent, is shown in Figure 1.10. The curve in the figure, which seems to capture the smooth variation in precipitation, was estimated using a penalty on the harmonic acceleration as measured by the differential operator (1.2).

The gait data in Figure 1.8 were converted to functions by the simplest of interpolation schemes: joining each pair of adjacent observations by a straight line segment. This approach would be inadequate if we require derivative information. However, one might perform a certain amount of smoothing while still respecting the periodicity of the data by fitting a Fourier series to each record: A constant plus three pairs of sine and cosine terms does a reasonable job for these data. The growth data in Figure 1.1 and the temperature data in Figure 1.6 were smoothed using smoothing splines, and this more sophisticated technique also provides high quality derivative information.

There are often conceptual constraints on the functions that we estimate. For example, a smooth of precipitation such as that in Figure 1.10 should logically never be negative. There is no danger of this happening for a station as moist as this, but a smooth of the data in Resolute, the driest



Figure 1.10. The points indicate average daily rainfall at Prince Rupert on the northern coast of British Columbia. The curve was fit to these data using a roughness penalty method.

place that we have data for, can easily violate this constraint. The growth curve fits should be strictly increasing, and we shall see that imposing this constraint results in a rather better estimate of the acceleration curves that we saw in Figure 1.2. Chapter 6 shows how to fit a variety of constrained functions to data.

1.5.2 Data registration or feature alignment

Figure 1.11 shows some biomechanical data. The curves in the figure are twenty records of the force exerted on a meter during a brief pinch by the thumb and forefinger. The subject was required to maintain a certain background force on a force meter and then to squeeze the meter aiming at a specified maximum value, returning afterwards to the background level. The purpose of the experiment was to study the neurophysiology of the thumb–forefinger muscle group. The data were collected at the MRC Applied Psychology Unit, Cambridge, by R. Flanagan; see Ramsay, Wang and Flanagan (1995).

These data illustrate a common problem in functional data analysis. The start of the pinch is located arbitrarily in time, and a first step is to align the records by some shift of the time axis. In Chapter 7 we take up the question of how to estimate this shift, and how to go further if necessary to estimate record-specific linear transformations of the argument, or even nonlinear transformations.



Figure 1.11. Twenty recordings of the force exerted by the thumb and forefinger where a constant background force of two newtons was maintained prior to a brief impulse targeted to reach 10 newtons. Force was sampled 500 times per second.

1.5.3 Data display

Displaying the results of a functional data analysis can be a challenge. With the gait data in Figures 1.8 and 1.9, we have already seen that different displays of data can bring out different features of interest, and that the standard plot of x(t) against t is not necessarily the most informative. It is impossible to be prescriptive about the best type of plot for a given set of data or procedure, but we shall give illustrations of various ways of plotting the results. These are intended to stimulate the reader's imagination rather than to lay down rigid rules.

1.5.4 Plotting pairs of derivatives

Helpful clues to the processes giving rise to functional data can often be found in the *relationships* between derivatives. For example, two functions exhibiting simple derivative relationships are frequently found as strong influences in functional data: the exponential function, $f(t) = C_1 + C_2 e^{\alpha t}$, satisfies the differential equation

$$Df = -\alpha(f - C_1)$$

and the sinusoid $f(t) = C_1 + C_2 \sin[\omega(t-\tau)]$ with phase constant τ satisfies

$$D^2 f = -\omega^2 (f - C_1).$$



Figure 1.12. The left panel gives the annual variation in mean temperature at Montreal. The times of the mid-months are indicated by the first letters of the months. The right panel displays the relationship between the second derivative of temperature and temperature less its annual mean. Strictly sinusoidal or harmonic variation in temperature would imply a linear relationship.

Plotting the first or second derivative against the function value explores the possibility of demonstrating a linear relationship corresponding to one of these differential equations. Of course, it is usually not difficult to spot these types of functional variation by plotting the data themselves. However, plotting the higher derivative against the lower is often more informative, partly because it is easier to detect departures from linearity than from other functional forms, and partly because the differentiation may expose effects not easily seen in the original functions.

Consider, for example, the variation in mean temperature Temp at Montreal displayed in the left panel of Figure 1.12. Casual inspection does indeed suggest a strongly sinusoidal relationship between temperature and month, but the right panel shows that things are not so simple. Although there is a broadly linear relationship between $-D^2$ Temp and Temp after subtracting the mean annual temperature, there is obviously an additional systematic trend, which is more evident in the summer through winter months than in the spring. This plot greatly enhances the small departures from sinusoidal behavior, and invites further attention.

Figure 1.13 plots the estimated derivatives for the logarithm of the U. S. nondurable goods index shown in Figure 1.3 for the year 1964. The second derivative or acceleration on the vertical axis is plotted against the first derivative or velocity on the horizontal axis in what is called a *phase plane plot*. The plot focuses attention on the interplay between Dx and D^2x by eliminating the explicit role of argument t, and reveals a fascinating cyclic structure that we will learn how to interpret in Chapter 2. Plotting derivatives as well as curve values is an essential part of functional data analysis.



Figure 1.13. A phase plane plot of the first two derivatives of the logarithm of the U. S. nondurable goods manufacturing index in Figure 1.3 over 1964.

1.6 Exploring variability in functional data

The examples considered so far offer a glimpse of ways in which the variability of a set of functional data can be interesting, but there is a need for more detailed and sophisticated ways of investigating variability, and these are a major theme of this book.

1.6.1 Functional descriptive statistics

Any data analysis begins with the basics: Estimating means and standard deviations. Functional versions of these elementary statistics are given in Chapter 2. But what is elementary for univariate and multivariate data turns out to be not always so simple for functional data. Chapter 7 returns to the functional data summary problem, and shows that *curve registration* or feature alignment may have to be applied in order to separate *amplitude variation* from *phase variation* before these statistics are used.

1.6.2 Functional principal components analysis

Most sets of data display a small number of dominant or substantial modes of variation, even after subtracting the mean function from each observation. An approach to identifying and exploring these, set out in Chapter 8, is to adapt the classical multivariate procedure of principal components

16 1. Introduction

analysis to functional data. In Chapter 9, techniques of smoothing or regularization are incorporated into the functional principal components analysis itself, thereby demonstrating that smoothing methods have a far wider rôle in functional data analysis than merely in the initial step of converting discrete observations to functional form. In Chapter 10, we show that functional principal components analysis can be made more selective and informative by considering specific types of variation in a special way. For example, we shall see that estimating a small shift of time for each temperature record and studying its variation will give a clearer understanding of record-to-record temperature variability.

1.6.3 Functional canonical correlation

How do two or more sets of records covary or depend on one another? As we saw in the cross-correlation plots, this is a question to pose for gait data, because relationships between record-to-record variation in hip angle and knee angle seem likely.

The functional linear modelling framework approaches this question by considering one of the sets of functional observations as a covariate and the other as a response variable, but in many cases, such as the gait data, it does not seem reasonable to impose this kind of asymmetry, and we shall develop two rather different methods that treat both sets of variables in an even-handed way. One method, described in Section 8.5, essentially treats the pair $(\text{Hip}_i, \text{Knee}_i)$ as a single vector-valued function, and then extends the functional principal components approach to perform an analysis. Chapter 11 takes another approach, a functional version of canonical correlation analysis, identifying components of variability in each of the two sets of observations which are highly correlated with one another.

For many of the methods we discuss, a naïve approach extending the classical multivariate method will usually give reasonable results, though regularization will often improve these. However, when a linear predictor is based on a functional observation, and also in functional canonical correlation analysis, regularization is not an optional extra but is an intrinsic and necessary part of the analysis; the reasons are discussed in Chapters 11, 15 and 16.

1.7 Functional linear models

The classical techniques of linear regression, analysis of variance, and linear modelling all investigate the way in which variability in observed data can be accounted for by other known or observed variables. They can all be placed within the framework of the general linear model

$$y = \mathbf{Z}\beta + \epsilon \tag{1.3}$$

where, in the simplest case, y is typically a vector of observations, β is a parameter vector, \mathbf{Z} is a matrix that defines a linear transformation from parameter space to observation space, and ϵ is an error vector with mean zero. The design matrix \mathbf{Z} incorporates observed covariates or independent variables.

To extend these ideas to the functional context, we retain the basic structure (1.3) but allow more general interpretations of the symbols within it. For example, we might ask of the Canadian weather data:

- If each weather station is broadly categorized as being Atlantic, Pacific, Continental or Arctic, in what way does the geographical category characterize the detailed temperature profile Temp and account for the different profiles observed? In Chapter 12 we introduce a functional analysis of variance methodology, where both the parameters and the observations become functions, but the matrix **Z** remains the same as in the classical multivariate case.
- Could a temperature record **Temp** be used to predict the logarithm of total annual precipitation? In Chapter 15 we extend the idea of linear regression to the case where the independent variable, or covariate, is a function, but the response variable (log total annual precipitation in this case) is not.
- Can the temperature record **Temp** be used as a predictor of the entire precipitation profile, not merely the total precipitation? This requires a fully functional linear model, where all the terms in the model have more general form than in the classical case. This topic is considered in Chapters 14 and 16.
- We considered earlier the many roles that derivatives play in functional data analysis. In the functional linear model, we may use derivatives as dependent and independent variables. Chapter 17 is a first look at this idea, and sets the stage for the following chapters on differential equations.

1.8 Using derivatives in functional data analysis

In Section 1.3 we have already had a taste of the ways in which derivatives and linear differential operators are useful in functional data analysis. The use of derivatives is important both in extending the range of simple graphical exploratory methods, and in the development of more detailed methodology. This is a theme that will be explored in much more detail in Chapters 18, 19 and 21, but some preliminary discussion is appropriate here.

Chapter 19 takes up the question, novel in functional data analysis, of how to use derivative information in studying components of variation. An approach called *principal differential analysis* identifies important variance components by estimating a linear differential operator that will annihilate them. Linear differential operators, whether estimated from data or constructed from external modelling considerations, also play an important part in developing regularization methods more general than those in common use. Some of their aspects and advantages will be discussed in Chapter 21.

1.9 Concluding remarks

The last chapter of the book, Chapter 22, includes a discussion of some historical perspectives and bibliographic references not included in the main part of our development.

In the course of the book, we shall describe a considerable number of techniques and algorithms, to explain how the methodology we develop can actually be used in practice. We shall also illustrate our methodology on a variety of data sets drawn from various fields, including where appropriate the examples we have already introduced in this chapter. However, it is not our intention to provide a cook-book for functional data analysis.

In broad terms, we have a grander aim: to encourage readers to think about and understand functional data in a new way. The methods we set out are hardly the last word in approaching the particular problems, and we believe that readers will gain more benefit by using the principles we have laid down than by following our own suggestions to the letter.

For those who would like access to the software we have used ourselves, a selection is available on the website:

http://www.functionaldata.org

This website will also be used to publicize related and future work by the authors and others, and to make available the data sets referred to in the book that we are permitted to release publicly.