

# INTRODUCTION

# 1

---

Installation	3
Object-oriented Programming	3
<b>Introductory Tutorial (Height Data)</b>	<b>4</b>
Selecting the Basis Functions	5
Smoothing	7
<b>A Linear Model for the Height Data</b>	<b>12</b>
Discussion	14
<b>Cluster Analysis of the Height Data</b>	<b>15</b>
Computing a Distance Matrix	15
Displaying the Cluster Mean Functions	17
Between Cluster Distances	18
Summary	19
Multidimensional Scaling	19
<b>FDA Flow Chart</b>	<b>21</b>

Functional data arise in many fields of research. Measurements are often best thought of as functions, even in cases where the data is gathered at a relatively small number of points. Examples include weather changes, stock prices, bone shapes, growth rates, health status indicators, and tumor size.

For time-dependent data, observations may be viewed as realizations of a smooth function  $y(t)$  of time that have been measured (with error) at specific time points  $t_j$ , but which could have been measured at any time. Spatial functional data is also common, e.g., the length of a bone along an axis, the concentration of a drug in a tissue as a function of depth, yearly mean temperature as a function of location.

Historically, functional data have been analyzed using multivariate or time-series methods at discrete measurement points. Analyzing functional data instead as functions has several advantages:

- Functions, unlike raw data, can be evaluated at any “time” point. This is important because it allows the use of statistical methods requiring evenly-spaced measurements and allows extrapolation for use in predictions or treatment decisions.
- Functional methods (e.g., functional principal components, functional canonical correlation) apply even when the data have been gathered at irregular intervals, or at different times on different subjects, when multivariate analogues of these methods are either inappropriate or unavailable.
- Derivatives and integrals of functions may provide important information about the underlying process. For example, knowledge of the direction and rate of change of a patient’s temperature may be more important than knowledge of the patient’s current temperature.

Functional methods can also be used when the parameters to be estimated are functions. Ramsay and Silverman (1997) use smoothing spline methods for density estimation, and to estimate the link function in generalized linear models. Another example is regression splines for fitting time-dependent hazard regression models (Koopman and Clarkson, 1997).

S+FDA integrates functional data analysis methods into S-PLUS. It includes a complete commercial implementation of the exploratory methods of Ramsay and Silverman (1997, 2002), featuring:

- methods for transforming observed data to a smoothed functional form,
- predicting a functional or nonfunctional variable  $y(t)$  as a function of one or more functional or nonfunctional variables,
- finding and rotating the functional “principal components” of a functional variable,
- finding the canonical correlations between two functional variables, and
- performing a “principal differential analysis”.

S+FDA also incorporates more recent innovations and extensions, such as allowing the use of functions with arbitrary bases, and providing methods for functional generalized linear models and functional cluster analysis.

## **Installation**

To install the software:

- **Go to the website:** <http://www.insightful.com/downloads/libraries/default.asp>
- Follow the on-screen Setup instructions; default settings are recommended.

## **Object-oriented Programming**

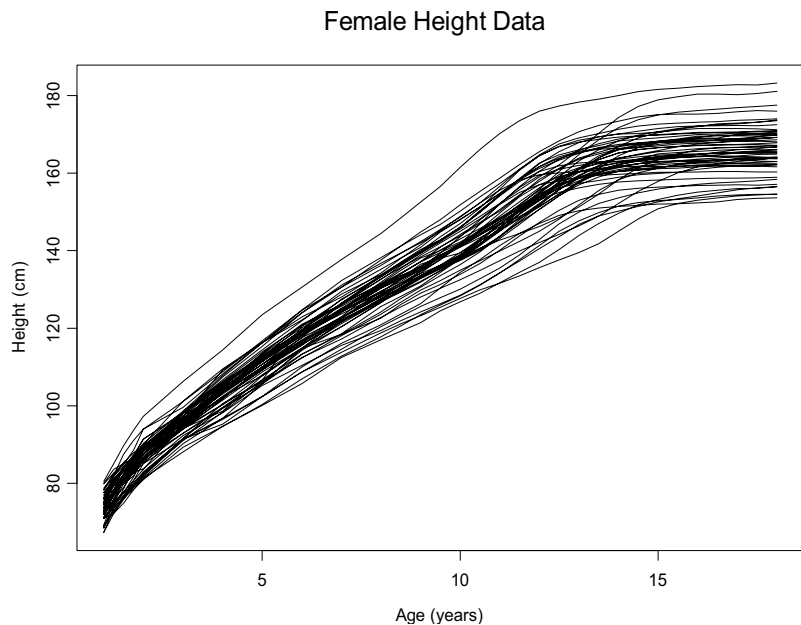
S+FDA makes use of the object-oriented capabilities of the S-PLUS language. In object-oriented programming, constructor functions create structured data “objects” that are assigned a class (which typically has the same name as the constructor). The object-oriented paradigm allows users to apply generic functions (such as `plot`) to these classed objects, the details of which are handled transparently through class-specific functions or “methods”. This simplifies programming by avoiding the need to explicitly invoke different functions or to have additional function arguments when generic operations are applied to objects of different structures.

## INTRODUCTORY TUTORIAL (HEIGHT DATA)

We illustrate some exploratory functional data analysis methods using the Berkeley height data (Tuddenham and Snyder, 1954). The corresponding data frame, `heightData`, is included in the S+FDA library. This data contains the heights of 54 female (columns 2 to 55) and 39 male (columns 56 to 94) children observed at 31 times from age 1 to age 18. The times of measurement are included as the variable `age` (column 1). We first inspect the data graphically by plotting the height curves as follows:

```
#Set up the plot and label
> plot(heightData$age, heightData[,2], type="n",
       ylim=range(unlist(heightData[,2:55])),
       xlab="Age (years)", ylab="Height (cm)",
       main="Female Height Data")
#draw the height curves
> matlines(heightData$age,as.matrix(heightData[,2:55]))
```

The result is shown in Figure 1.1.



**Figure 1.1:** *Female height data.*

Although the data appear as smooth curves, only 31 discrete values of height were measured. The curves are produced by connecting these discrete points with straight lines.

As a functional data analysis application, we fit a function to each height curve using linear least squares. The function is represented as a linear combination of basis functions  $b_j(t)$  and coefficients  $\beta_j$  that vary from one height function to the next:

$$f(x) = \sum_{j=1}^{n_b} \beta_j b_j(t)$$

There are a variety of choices for the basis functions, e.g., B-splines, Fourier series, and exponential series. Once the basis is chosen, the coefficients are estimated based on the observed data. In Figure 1.1, a *polygonal basis* of connected line segments is used to draw the curves.

Although the functional representation almost always differs from the data at the points of observation, these differences are assumed to be small in the sense that the coefficients  $\beta_j$  capture the information contained in the discretized curve. In most analyses, the raw data is ignored once the  $\beta_j$  have been estimated because it is simpler to work with the functional form. The assumption is that the within-subject variance in the  $\beta_j$  estimates is small compared to the between-subject variance.

**Warning**

When the number of observations for estimating the  $\beta_j$  is small to moderate or when the within-subject variance of the  $\beta_j$  estimates is large, a mixed-effects model may be preferred so that information may be combined across subjects.

**Selecting the Basis Functions**

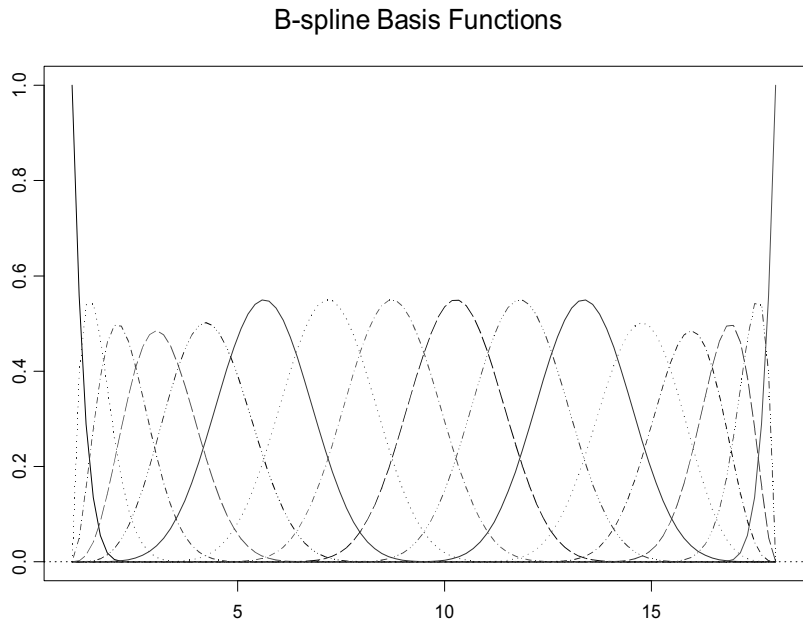
To perform a functional data analysis, we must first choose an appropriate set of basis functions. In the example above, 16 B-spline basis functions of order 6 were used. Since the order of a polynomial basis is the degree plus one, this basis consists of 16 piecewise polynomial splines of degree 5. By default, the *interior* knots for the 16

basis functions are equally spaced over the range of the independent variable (the two *exterior* knots are placed at the endpoints of the function domain). Since height is being viewed as a function of age, the appropriate domain for the basis functions is the age span of the data. The following forms an object of class “bsplineBasis” for the height data:

```
> heightBasis <- bsplineBasis(fDomain  
                             =range(heightData$age), nbasis=16,norder=6)
```

The basis functions, displayed in Figure 1.2, are equally spaced over the domain:

```
> plot(heightBasis, main="B-spline Basis Functions")
```



**Figure 1.2:** A set of 16 B-spline basis functions.

Now that we have defined a basis, we need to calculate the coefficients for each height curve. Since there are 93 subjects in this dataset, there should be 93 sets of coefficients (one set for each function). The S+FDA function `fVector` takes the basis, the data matrix, and the independent variable, and returns an object of class “fVector” containing the linear least-squares estimates of the coefficients. An “fVector” object has two additional attributes:

“basis”, which stores the basis used in the fit, and “fNames”, which stores labeling information for the data. In the code below, we also specify names for the independent variable (`age`), the subjects (`child`), and the units of the response (`height`). These names are used in the plotting and printing functions.

```
> fHgt <- fVector(object=heightBasis, y=heightData[,2:94],
  fArgs=heightData$age,
  fNames=list(age=heightData$age,
  child=names(heightData)[2:94], height='cm'))
```

Extract the estimated coefficients, basis functions, and function names from `fHgt` using the commands `getCoef(fHgt)`, `getBasis(fHgt)`, and `getNames(fHgt)`, respectively.

## Smoothing

Although the basis functions smooth the curves, additional smoothing may be beneficial. The S-PLUS functions for creating functional data objects allow specification of a smoothing penalty in the least-squares objective. The penalty also requires a smoothing parameter, `lambda`. You may estimate an optimal `lambda` by minimizing a generalized cross validation statistic. See section Generalized Cross Validation on page 82 for more details.

Smoothing techniques are largely exploratory in nature, and are discussed in more detail in Chapter 4 of this manual, as well as in Chapter 4 of Ramsay and Silverman (1997). We will have occasion to use smoothing techniques for most of the functional data analysis methods provided in S+FDA.

As an example, penalize the squared second derivative with a penalty parameter `lambda=0.001`:

```
> fHgt2 <- fVector(object=heightBasis,
  y=heightData[, 2:94], fArgs=heightData$age,
  penalty=list(lambda=0.001, linDop=fDop(2)),
  fNames=list(age=heightData$age,
  child=names(heightData)[2:94], height='cm'))
```

Compare with the original data of Figure 1.1 to see how closely the smoothed functions fit the data. The S-PLUS function `fEval` evaluates an object of class “fVector” at any point in the domain of the basis. Here, we evaluate the 54 spline curves for the females at the original

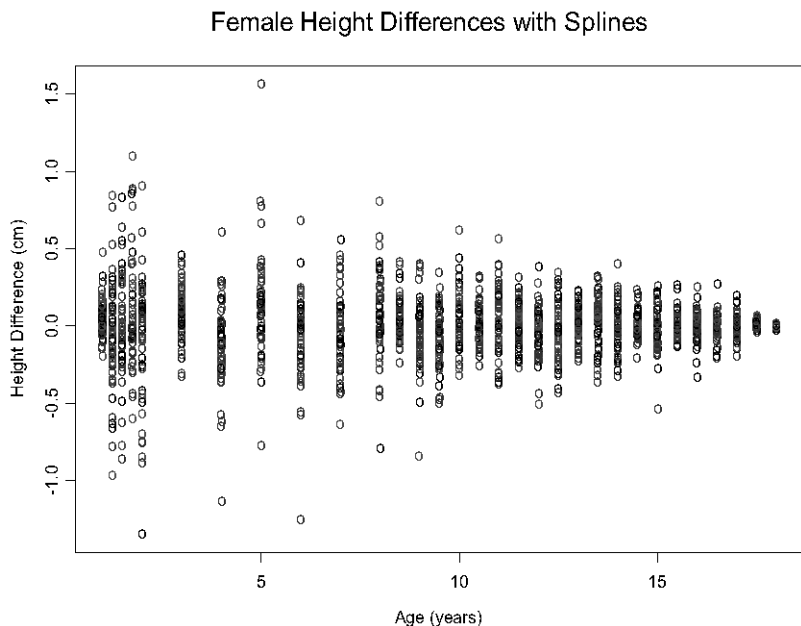
age values (`heightData$age`), calculate the difference between predicted and observed heights, and then plot the curve differences at the given ages:

```
> hgtFemale<-fEval(fHgt2[1:54], heightData$age)

> plot(heightData$age, hgtFemale[,1], type="n",
       ylim=range(hgtFemale-as.matrix(heightData[,2:55])),
       xlab="Age (years)", ylab="Height Difference (cm)",
       main="Female Height Differences with Splines")

> matpoints(heightData$age,
           hgtFemale-as.matrix(heightData[, 2:55]), pch="o")
```

The resulting plot is given in Figure 1.3.



**Figure 1.3:** *Difference between predicted and actual female height data when using cubic B-splines for function representation.*

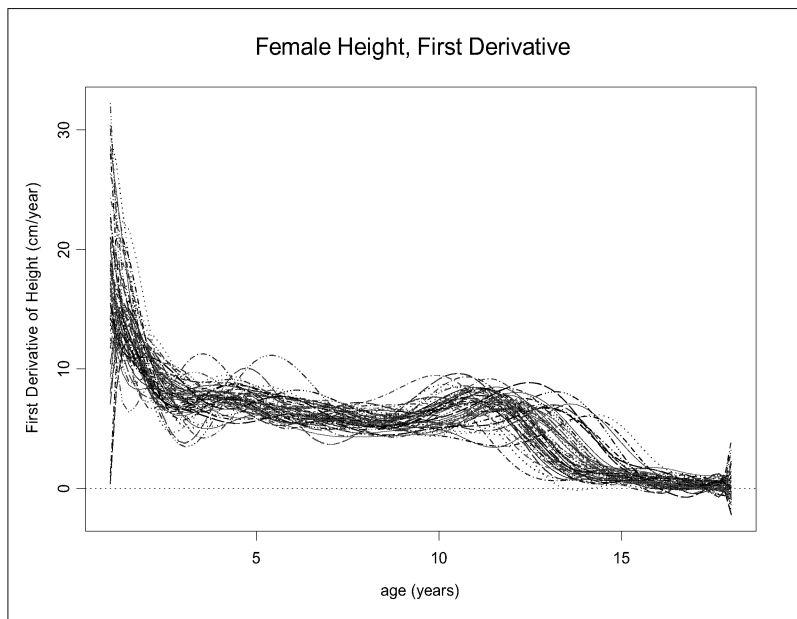
The maximum deviation between the spline approximation and the true heights is about 1.5 cm compared with height values of 80 cm or more (see Figure 1.1). These differences are small enough that we consider the smoothed functions to be acceptable for subsequent analysis.



Given a representation of the data as an `fVector` object, it is easy to conduct several kinds of exploratory analyses with `S+FDA`. Here, we compute the first two derivatives of height with respect to time. We begin with the first derivative:

```
> plot(fVector(fHgt2[1:54], linDop = fDop(1)),  
      xlab="age (years)",  
      ylab="First Derivative of Height (cm/year)",  
      main="Female Height, First Derivative")
```

The result is displayed in Figure 1.4.



**Figure 1.4:** *First derivatives of the functional representation of the female height data. The second derivative was penalized for smoothing, with penalty parameter 0.001.*

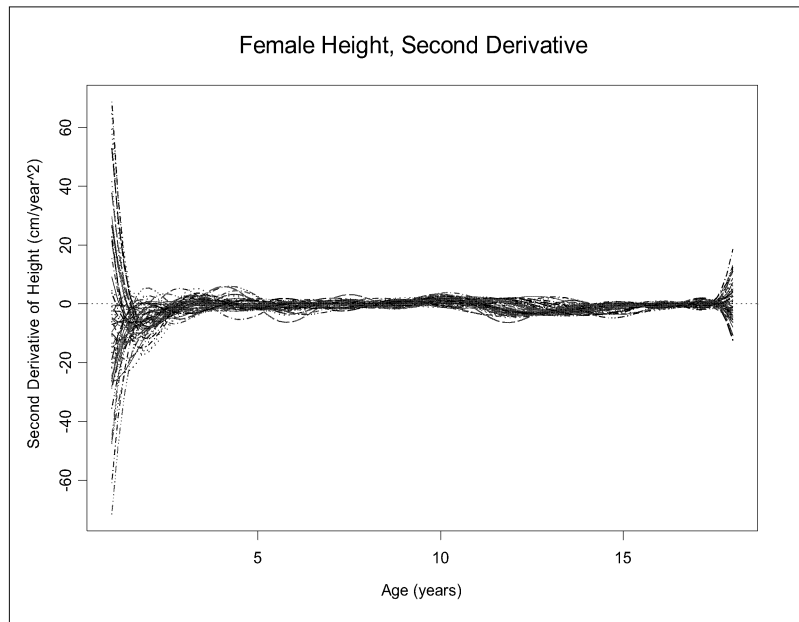
Despite the large number of curves in Figure 1.4, some general trends are apparent: there appears to be an acceleration in growth around age 4, with a second acceleration after age 10. Further exploratory analysis, such as plotting the mean of the 54 derivative functions, may help reveal more structure.

The plot of the second derivatives is produced in a similar fashion:

## Chapter 1 Introduction

```
> plot(fVector(fHgt2[1:54], linDop=fDop(2)),  
      xlab="Age (years)",  
      ylab="Second Derivative of Height (cm/year^2)",  
      main="Female Height, Second Derivative")
```

The result is displayed in Figure 1.5.



**Figure 1.5:** *Second derivatives of the functional representation of the female height data. The second derivative was penalized for smoothing, with penalty parameter 0.001.*

The large function values near the endpoints in both derivative plots are due to lack of information concerning values outside the interval. Smoothing by penalizing a higher derivative would reduce the variation at the endpoints, although possibly at the risk of oversmoothing the function. Such considerations are discussed in more detail in the chapter on smoothing.

Because we use splines of degree five (order 6) when fitting the functions, the second derivatives are smooth, cubic splines. Had we fit the raw data with cubic splines (order 4), the second derivative curves would have been piecewise linear. In general, if an analysis requires a

smooth  $k$ th derivative, and smoothness in higher derivatives is unimportant, splines of degree  $k+3$  (order  $k+4$ ) should be used to fit the functions so that the  $k$ th derivative will be a cubic spline.

The ease with which you can examine the derivatives is a direct consequence of the functional approach, and one of its main advantages. By regarding the height measurements for each person as a smooth curve, you are no longer constrained by discrete observation times.

## A LINEAR MODEL FOR THE HEIGHT DATA

Now consider a functional linear model for predicting sex in terms of the growth rate, the first derivative of the height curve. Since the dependent variable is binary, this model can also be considered a discriminant function for predicting sex in terms of the growth rate.

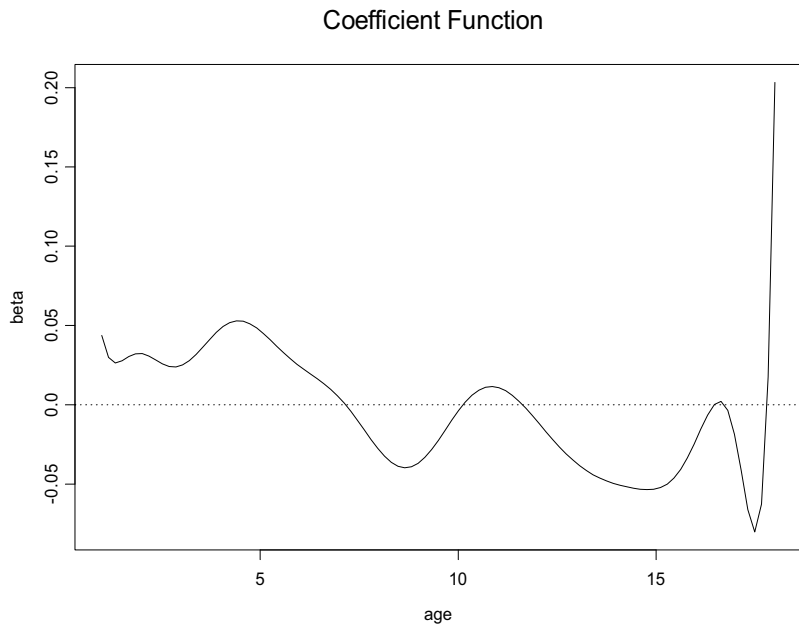
For the height data, fit a functional linear model as follows:

```
> predLm <- fLM(sex~-1+fVector(fHgt, linDop=fDop(1)),
  data.frame(fHgt=fHgt,
  sex=c(rep(1,54), rep(0,39))))
```

Here the -1 in the model formula eliminates the intercept, which is already contained in the B-splines. The coefficients in the resulting model are functional. The first coefficient estimate may be plotted as follows:

```
> plot(predLm$coef[[1]], xlab="age", ylab="beta",
  main="Coefficient Function")
```

The resulting plot is shown in Figure 1.6.



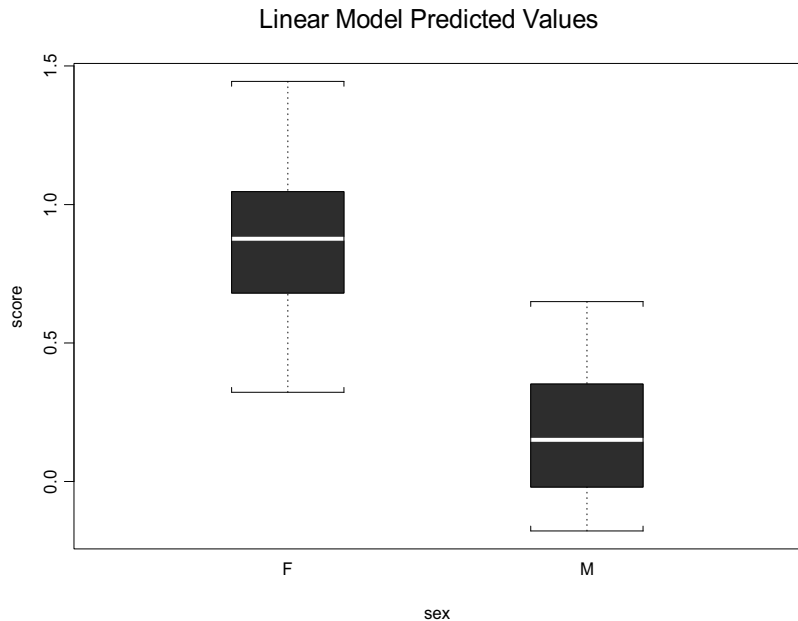
**Figure 1.6:** *The function of coefficients predicting sex in terms of the height function.*

The effect of the growth rate on the linear model prediction has a maximum around age 5, is positive again at around age 11, and is negative during the puberty growth spurt after age 11. The negative lobe after age 11 predicts maleness, when the boys have their growth spurts, but the girls are finished theirs.

To see how well the resulting model can discriminate between males and females, plot the fitted values:

```
> score <- getCoef(predLm$fitted.values)
> plot(as.factor(c(rep("F",54),rep("M",39))),
       score, main="Linear Model Predicted Values")
```

The results are displayed in Figure 1.7.



**Figure 1.7:** *Predicted height scores for each sex*

Most females have a score above 0.5, and most males have a score below 0.5, so that growth rates are an effective means of classifying the observations.

**Discussion**

In this simple example of a functional linear model, we have again used derivative information, this time to predict the sex of the individual. Although the results for this example are good, generally predictions based on functional linear models should be viewed with caution. When the independent variable is functional, so are the coefficient estimates, and outliers may significantly influence the outcome (overfitting). Methods to avoid overfitting, particularly smoothing methods, are discussed in more detail in the chapter on functional linear models.

## CLUSTER ANALYSIS OF THE HEIGHT DATA

One approach to cluster analysis is to search for natural groups of observations by examining “distances” between observations. For the height data discussed in the previous section, clustering can be based on a Euclidean or other distance measure between the observed heights at the observation times (the ages). Specifying these distances requires that all individuals be measured at the same times. This requirement can be met by first converting the observed data to functional form. Once this is accomplished, a much broader class of distance measures becomes available. For example, derivatives can be incorporated into the distance metrics. For the height data, we might be interested in patterns of growth curves related to the growth rate (the velocity, i.e., first derivative) or the rate of change in the growth rate (the acceleration, i.e., second derivative). If, for example, our main interest is the growth rate, then we could define the distance between the growth curve functions  $f_1(t)$  and  $f_2(t)$  for two individuals as the square root of the integrated squared distance between the first derivatives of the two height curves:

$$d(f_1(t), f_2(t)) = \sqrt{\int_t \left( \frac{df_1(t)}{dt} - \frac{df_2(t)}{dt} \right)^2 dt}$$

This distance measurement is based on the rate of change of growth, as opposed to the final height achieved.

### Computing a Distance Matrix

For the clustering example, we consider only the height data starting from age 3. The reason for this is that the data in infancy are unstable, and the transition to standing height around age 2 introduces a significant perturbation. We recompute the smoothed fHgt from age 3:

```
> ageRange <- heightData$age >= 3
> heightBasis <- bsplineBasis(fDomain
                             =range(heightData$age[ageRange]),
                             nbasis=16, norder=6)
> fHgt3 <- fVector(object=heightBasis,
```

```
y=heightData[ageRange,2:94],
fArgs=heightData$age[ageRange],
penalty=list(lambda=0.001, linDop=fDop(2)),
fNames=list(age=heightData$age[ageRange],
            child=names(heightData)[2:94],
            height='cm'))
```

The choice of  $\lambda=0.001$  is determined by a procedure described in section Generalized Cross Validation on page 82 .

The S+FDA function `fDist` computes distance matrices from functional data. The following S-PLUS code computes a distance matrix whose  $(i, j)$  element contains the square root of the integrated squared distance between the first derivatives of growth functions  $(i)$  and  $(j)$  for the height data:

```
> distHgt <- sqrt(fDist(fHgt3, linDop=fDop(1)))
```

Now we can apply any clustering method based on distance matrices. For example, the S-PLUS function `hclust` computes clusters for a variety of hierarchical clustering methods from a distance matrix. Here we use average-linkage clustering:

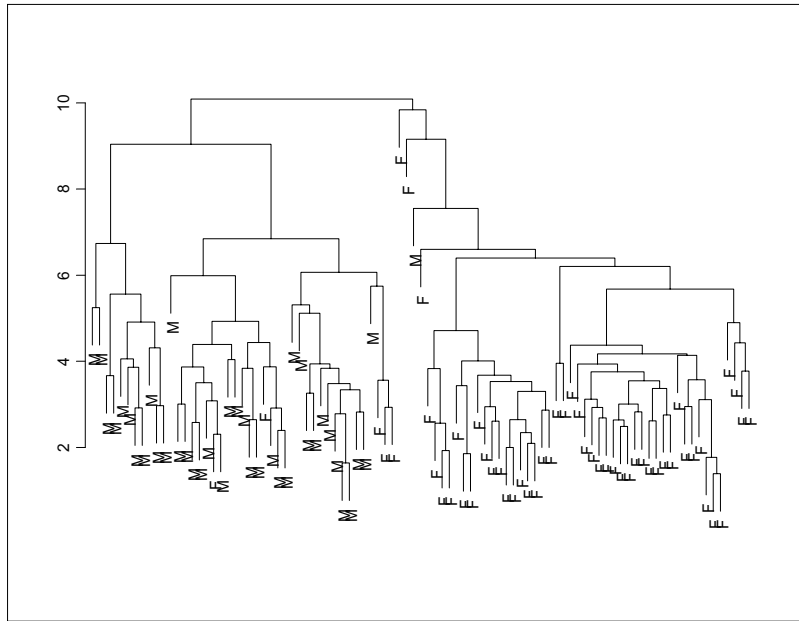
```
> clustHgt <- hclust(distHgt, method="average")
```

A plot of the cluster tree label according to sex is obtained as follows:

```
> sex <- as.factor(c(rep("F", 54), rep("M", 39)))
> plclust(clustHgt, labels=as.character(sex))
```



The result is displayed in Figure 1.8:



**Figure 1.8:** Complete linkage cluster tree labeled according to sex.

## Displaying the Cluster Mean Functions

Since we have heights for both males and females, it would seem natural to group the data by sex. The labeling in Figure 1.8 shows that this grouping is supported by the cluster analysis, indicating that males and females generally have different growth patterns. Only one male appears in the female subtree, and relatively few females appear in the male subtree. To investigate this further, we apply the S-PLUS function `cutree` to obtain the two-group solution:

```
> g <- 2
> groupsHgt <- cutree(clustHgt, k=g)
```

The clusters are as defined by a horizontal line at about distance 9.5 in Figure 1.8. The frequency of males and females in each of the groups is easily obtained using the S-PLUS function `crossstabs`:

```
> crossstabs(~groupsHgt+sex)
```

for which an abbreviated output is shown below:

	F	M	RowTot1
1	49	1	50
2	5	38	43
ColTot1	54	39	93

Group 1 is predominantly female and group 2 predominantly male. We split the data into a list grouped by cluster, and plot the function and derivative means for each group:

```

> splitGroups <- split(fHgt3, groupsHgt)
> par(mfrow=c(2,1))
> plot(1, 20, type="n", xlab="age", ylab="height",
      main="Group Mean Function Heights",
      xlim=c(0, 19), ylim=c(60,200))
> temp <- lapply(1:g, function(i)
  lines(mean(splitGroups[[i]]), lty=i, col=i))
> legend(1, 190, paste(1:g), col=1:g, lty=1:g)
> plot(15, 20, type="n", xlab="age", ylab="height",
      main="Group Mean Derivative Heights",
      xlim=c(0, 19), ylim=c(0,30))
> temp <- lapply(1:g, function(i)
  lines(mean(fVector(splitGroups[[i]]),
    linDop=fDop(1))), lty=i))
> legend(15, 29, paste(1:g), col=1:g,
  lty=1:g, background=0)

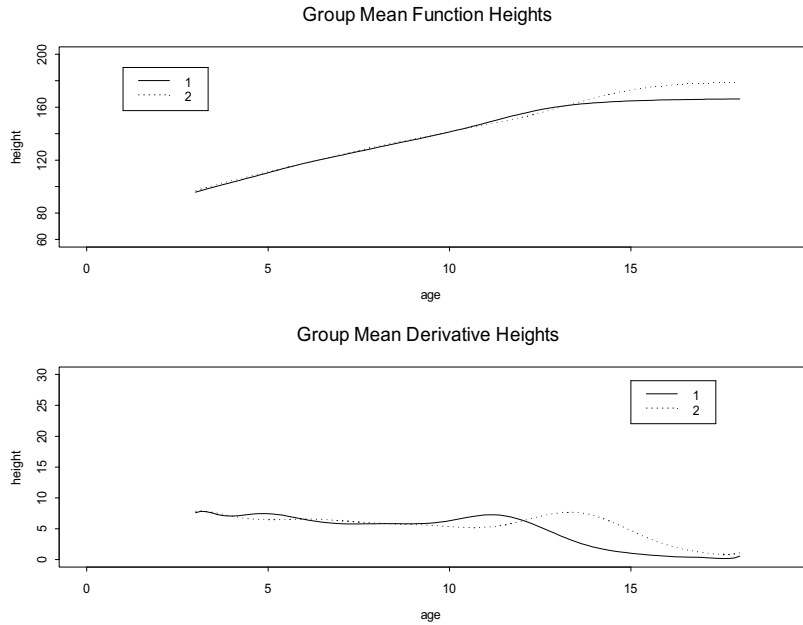
```

The results are shown in Figure 1.9. Since derivatives were used to define the distances, one would expect cluster differences to be reflected in their means, shown in the lower half of Figure 1.9. The display shows that the behavior of the clusters differs with respect to the time and duration of the growth spurt around puberty. There is also a difference in the groups around age 5 where group 1 (mostly females) tends to have a minor growth spurt that is not present in mostly-male group 2.

## Between Cluster Distances

Hierarchical clustering methods produce a grouping for a given number of clusters, but do not include a mechanism for selecting the correct number of clusters. The choice of two groups was based on informal inspection of the clustering tree (Figure 1.8). Because of the

small number of cross-overs from males to females, the two-group solution (males versus females) would seem satisfactory. However, if the labeling according to sex were not available, we would be unlikely to reach this conclusion.



**Figure 1.9:** *The mean function (top) and its first derivative for the two groups in Figure 1.8*

## Summary

This clustering example illustrates the flexibility of functional data analysis methods - when the data are thought of as functions, distance measures based on derivatives are possible, and derivatives can be used to analyze group structure.

## Multidimensional Scaling

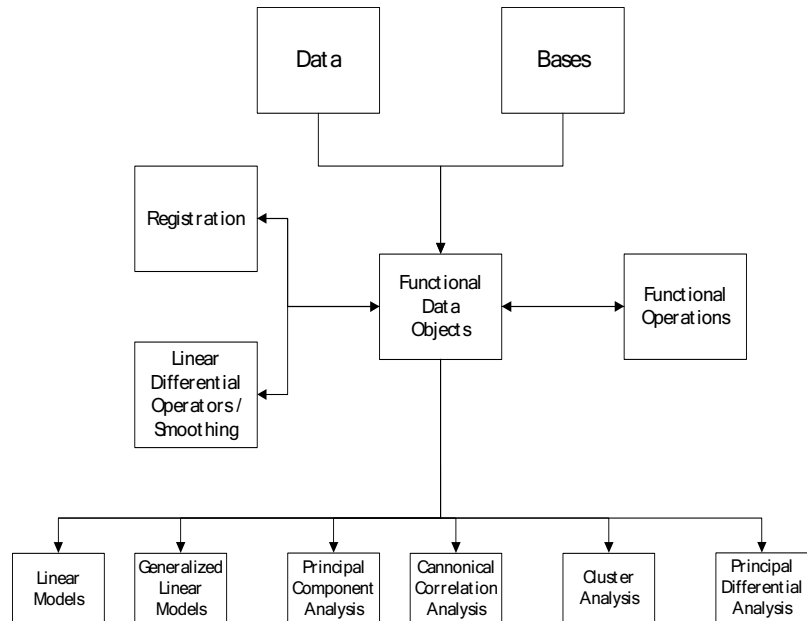
Multidimensional scaling is also possible once a distance matrix is available. We applied the S-PLUS function `cmdscale` to the distance matrix (using the command `cmdscale(distHgt)`) to do a simple multidimensional scaling analysis. In a plot of the (two dimensional) solution (not shown), the males and females are well separated.

## *Chapter 1 Introduction*

The S+FDA library offers many other methods for functional data analysis. These are discussed more fully in subsequent chapters, as well as in Ramsay and Silverman (1997, 2002).

## FDA FLOW CHART

The flowchart in Figure 1.10 represents the organization of this manual.



**Figure 1.10:** FDA flowchart.

Each box in the flowchart represents a chapter. Functional data analysis begins by selecting a basis to represent discrete data in functional form. The data typically correspond to a sample of functions, so that *registration* to remove unimportant differences (e.g. phase and/or amplitude variations) between samples may be necessary. Although the basis representation usually provides some smoothing, it is often desirable to apply one or more smoothing operations before analysis. This smoothing may be accomplished via a penalty on a linear differential operator applied to the functions.

Once a functional data object has been created, it can be analyzed and transformed in ways that are not possible for discrete data. You may perform various arithmetic operations, including differentiation and integration. In addition, a variety of analyses from discrete data analysis have functional analogs: linear and generalized linear

## *Chapter 1 Introduction*

generalized linear modeling, principal component and canonical correlation analysis, and cluster analysis. Principal differential analysis, which has no analog for discrete data, is another option for functional data.