

Curve Alignment by Moments

GARETH M. JAMES*

Abstract

A significant problem with most functional data analyses is that of misaligned curves. Without adjustment, even an analysis as simple as estimation of the mean will fail. One common method to synchronize a set of curves involves equating “landmarks” such as peaks or troughs. The landmarks method can work well but will fail if marker events can not be identified or are missing from some curves. An alternative approach, the “continuous monotone registration” method, works by transforming the curves so that they are as close as possible to a target function. This method can also perform well but is highly dependent on identifying an accurate target function. We develop an alignment method based on equating the “moments” of a given set of curves. These moments are intended to capture the locations of important features which may represent local behavior, such as maximums and minimums, or more global characteristics, such as the slope of the curve averaged over time. Our method works by equating the moments of the curves while also shrinking towards a common shape. This allows us to capture the advantages of both the landmark and continuous monotone registration approaches. The method is illustrated on several data sets and a simulation study is performed.

1 Introduction

The fundamental paradigm of functional data analysis (FDA) involves treating the entire curve or function as the unit of observation rather than individual measurements from the curve (Ramsay and Silverman, 2005). As FDA has become more common many statistical analysis techniques have been adapted to the paradigm. The analysis of functional data possess a number of problems not generally encountered with more classical data. One of the most important is that of misaligned curves. Figure 1 provides a real world illustration of this difficulty using the acceleration curves of ten boys from the Berkeley growth curve study (Tuddenham and Snyder, 1954) where the heights of individuals were recorded at regular intervals until age 18. Figure 1a), which plots smoothed versions of the observed acceleration curves, shows a clear trend of positive and then negative acceleration during the teenage years. However, the onset times, and spread, of these growth spurts can differ by several years so the curves can be considered to be misaligned or “unsynchronized”. The dashed line, which represents the cross-sectional mean based on the observed curves, clearly fails to capture the height of the peaks and troughs and underestimates the rate of change in the acceleration curve during the puberty years. Figure 1b) plots the corresponding curves after synchronization using the approach developed in this paper. Now one can much more clearly discern the general shape of the curves and the gray line, which represents the mean from the synchronized curves, shows that the true

*Associate Professor of Statistics, Marshall School of Business, University of Southern California.

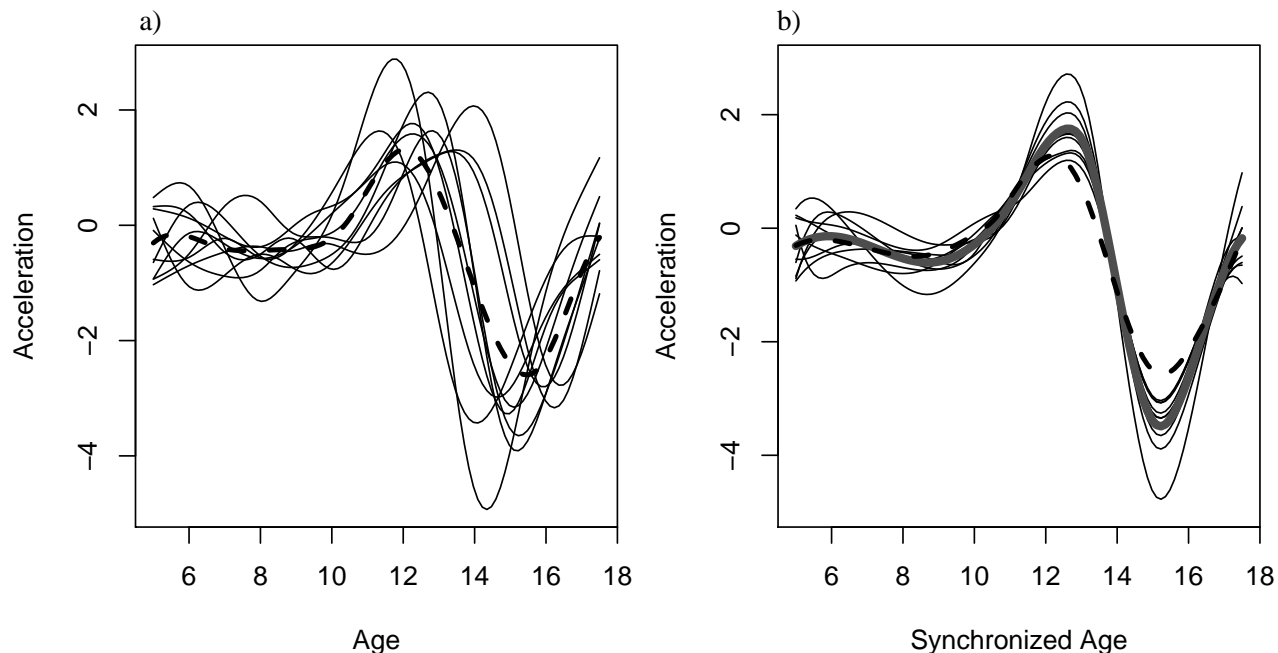


Figure 1: Ten acceleration curves from the Berkeley growth curve study, a) unsynchronized curves and b) after alignment. The dashed lines represent the cross-sectional mean based on the observed curves while the gray solid line corresponds to the mean from the synchronized data.

peaks and troughs are considerably more extreme than was previously apparent. Computing the mean of a set of curves is only one example of the many applications for which proper alignment of the curves is an essential component. For example, functional principal components analysis (James *et al.*, 2000; Rice and Wu, 2001), regression with both functional responses (Zeger and Diggle, 1994) and functional predictors (Ferraty and Vieu, 2002; James and Silverman, 2005), functional linear discriminant analysis (James and Hastie, 2001; Ferraty and Vieu, 2003) and functional clustering (James and Sugar, 2003; Bar-Joseph *et al.*, 2003) all assume that the starting curves are correctly aligned on the x -axis.

The problem of realigning such curves has been studied under different names in several fields. In the statistics literature it is referred to as curve registration (Silverman, 1995; Ramsay and Li, 1998) or, in the context of computing an average curve, structural averaging (Kneip and Gasser, 1992). It is also called curve alignment in biology and time warping in engineering (Sakoe and Chiba, 1978). Any set of curves can be decomposed into “amplitude” functions, which measure differences in the y -axis, and “warping” functions, which measure differences in location on the x -axis. Synchronization requires estimation of the warping functions.

A number of approaches have been proposed for this problem. Marker, or landmark, registration (Kneip and Gasser, 1992) involves selecting common features in the data, such as peaks or troughs, and transforming time so that these features occur together. This method can work well when such features can be easily identified but tends to perform poorly if no obvious and consistent landmarks exist. In addition the landmarks often need to be manually identified, preventing the implementation of a fully automatic approach. An alternative method involves aligning curves using a target function. Silverman (1995) proposed registering

curves using a simple shift in time such that the average squared distance between each curve and a target function is minimized. This idea was extended in Ramsay and Li (1998) using a Procrustes type fitting procedure on a general nonlinear class of time transformations to provide maximal alignment to the target function subject to suitable smoothness of the transformations. This approach, called “continuous monotone registration”, is often very effective but depends heavily on the target function. Generally the cross-sectional mean is used which can provide misleading results if the curves are significantly misaligned. Other recent work in this area includes Kneip *et al.* (2000), Rønn (2001) and Gervini and Gasser (2005).

The aim of this paper is to develop an automatic synchronization method that incorporates the best properties of both the landmark and continuous monotone registration approaches. We start by calculating “moments” for each curve. These moments are intended to capture the locations of important features which may represent local behavior, such as maximums and minimums, or more global characteristics, such as the slope of the curve over time. We then synchronize the curves by both, equating the moments for each curve, which has the effect of aligning common features, and simultaneously shrinking towards a common shape. In situations where obvious marker events are present our approach has the same advantages as landmark registration. Additionally, when there are no events but an accurate target function can be estimated we will achieve similar performance to the continuous monotone registration method. However, we show through the use of theory, simulations and real world examples, that even in situations where the landmark and continuous monotone registration procedures fail i.e. where obvious markers do not exist and an accurate target function can not be computed, our moments based method can still perform well.

General definitions for the moments of an arbitrary function are developed in Section 2. These moments are defined in terms of “feature” functions which can be designed to detect both local and global characteristics of the curves. In Section 3 we provide a model for the observed or unsynchronized curves. The moments from Section 2 are included as a fundamental part of the model. We also discuss alternative types of warping functions, both linear and non-linear. A synchronization procedure to fit our model is presented in Section 4. Our procedure attempts to a) equate the moments among the curves, and hence align the common “features” in analogy with landmark registration, and b) shrink the curves towards a common shape in analogy to the continuous monotone registration approach. We provide an algorithm for implementing our method and demonstrate that the estimates are consistent. Our method is illustrated on the Berkeley growth curve data in Section 5. The results from several simulation studies, comparing the performance of our approach with other synchronization methods, are reported in Section 6. Finally, Section 7 discusses the relationship of our approach to other common methods and suggests some further extensions.

2 Defining the Moments of a Function

In this section we develop definitions for the moments of an arbitrary function, g , in analogy with the moments of a random variable. The fundamental idea is that, just as one can define the distribution of a random variable through its moments and equate two different distributions by transforming to equate the moments, we can also define the shape of a function through its moments and synchronize two curves by equating their moments. We first introduce the concept of a “feature function”, $I_g(t)$, for g and impose the constraints

$$I_g(t) \geq 0 \quad \text{and} \quad \int I_g(t)dt = 1$$

which ensure that I_g is a weighting function. There are various possible choices for $I_g(t)$. Depending on the properties of our data, we may wish to utilize a function that places high weight on the time points corresponding to local features, such as maximums or minimums, or alternatively use a function that places weight according to more global characteristics such as the slope at a given time.

First, we discuss local approaches where most of the weight is concentrated around the time points corresponding to a specific feature in the data. For example, as $r \rightarrow \infty$, $I_g^{\max}(t) \propto (g(t) - \min\{g(t)\})^r$ and $I_g^{\min}(t) \propto (\max\{g(t)\} - g(t))^r$ will respectively concentrate their weight on the global maximum and minimum of $g(t)$. We may wish to search for local, as well as global, maximums and minimums. In this case one could utilize

$$I_g^{\text{local}}(t) \propto \begin{cases} \exp\left(-r \frac{|g^{(1)}(t)|}{\sqrt{|g^{(2)}(t)|}}\right) & g^{(2)}(t) \neq 0 \\ 0 & g^{(2)}(t) = 0. \end{cases}$$

This function places maximum weight on points where the first derivative is zero. However, $I_g^{\text{local}}(t)$ is also high for points with a low first derivative but a high second derivative. Thus, the function effectively searches for local maximums or minimums where g is changing most rapidly. As $r \rightarrow \infty$, $I_g^{\text{local}}(t)$ will place all its weight on the regions around local turning points. Finally, we examine $I_g^{(m)}(t)$ which places weights according to the absolute m th derivative of the curve, $g^{(m)}$, i.e. $I_g^{(m)}(t) \propto |g^{(m)}(t)|$. With $m = 0$ this function puts highest weight on large absolute values of g . With $m = 1$ most weight is placed on time points where g has a large slope and would be used when we are most interested in regions where a curve is changing rapidly. Setting $m = 2$ searches for points with greatest curvature etc. $I_g^{(m)}(t)$ can be considered to be searching for global characteristics of a curve because it is likely to spread its mass over all time points.

Then, for a given choice of I_g , we define the first moment of g by

$$\mu_g^{(1)} = \int t I_g(t) dt$$

and the k th central moment by

$$\mu_g^{(k)} = \int (t - \mu_g^{(1)})^k I_g(t) dt, \quad k \geq 2.$$

$\mu_g^{(1)}$ provides a measure of the center of g on the time axis while $\mu_g^{(2)}$ measures variability in g . Note that the variability is measured in relation to the time axis and not the y , or amplitude, axis. A curve could vary significantly in the y -axis, but still have a low value for $\mu_g^{(2)}$. In general $\mu_g^{(1)}$ will be more useful than the higher order moments when using feature functions such as I_g^{\max} or I_g^{\min} that concentrate on local features. The higher order moments, i.e. $\mu_g^{(k)}$ for $k \geq 2$, increase in importance when using more global feature functions such as $I_g^{(m)}$.

To better understand the properties of $\mu^{(k)}$ we examine the relationship between the moments of a function $h(s)$ and those of the shape invariant function $h(\frac{s-a}{b})$. In this formulation, $h(s)$ is stretched, about $s = 0$, by a factor b and shifted to the right by a . Hence, since $\mu^{(1)}$ is a measure of the center of a function and $\mu^{(k)}$ is a measure of variability about the center, stretching by a factor b should multiply the first moment by b and the k th moment by b^k . For example, one would expect that $\mu_{h(\frac{s-a}{b})}^{(2)}$, which measures the variability of

the transformed function, would equal $b^2\mu_{h(s)}^{(k)}$. Similarly a shift of a should add a to the first moment and leave the higher order moments, which are centered around the first moment, unchanged. We express this mathematically as

$$\mu_{h(\frac{s-a}{b})}^{(1)} = b\mu_{h(s)}^{(1)} + a \quad \text{and} \quad \mu_{h(\frac{s-a}{b})}^{(k)} = b^k\mu_{h(s)}^{(k)}, \quad k \geq 2. \quad (1)$$

Theorem 1 shows that, provided we utilize a certain family of feature functions, these properties will hold.

Theorem 1 *Suppose that $I_g(t)$ is chosen such that*

$$I_{g(\frac{s-a}{b})}(t) \propto I_{g(s)}\left(\frac{t-a}{b}\right), \quad -\infty < t < \infty, \quad (2)$$

for all a, g and $b > 0$. Then (1) will hold for any function $h(s)$.

Condition (2) holds for many large classes of feature functions. In particular, the previously mentioned feature functions all satisfy (2) and hence their corresponding moments all possess the desirable properties given by (1).

Corollary 1 *When utilizing $I_g^{(m)}$, I_g^{\max} , I_g^{\min} or I_g^{local} condition (2) is satisfied, and hence (1) holds. In addition (2) is satisfied for any $I_g^\phi(t) \propto \phi(g(t))$ where $\phi(t)$ is an arbitrary function.*

The feature functions we have utilized represent only a few of the possible choices one could utilize. In fact one of the strengths of our approach is the ability to design functions which best suit one's particular data.

3 The Synchronization Model

Let $Y_1(t), Y_2(t), \dots, Y_N(t)$ represent the unsynchronized functions or curves with Y_i observed at t_1, \dots, t_n where $t_j \in [0, T]$. Suppose we select L feature functions, I_g^1, \dots, I_g^L , and associated moments, $\mu_g^{(1,k)}, \mu_g^{(2,k)}, \dots, \mu_g^{(L,k)}$. Then our synchronization model is given by,

$$Y_i(t_j) = Z_i(W_i(t_j)) + \varepsilon_{ij}, \quad i = 1, \dots, N, \quad (3)$$

$$\mu_{Z_1}^{(l,k)} = \mu_{Z_2}^{(l,k)} = \dots = \mu_{Z_N}^{(l,k)} = \mu_{\bar{Y}}^{(l,k)}, \quad l = 1, \dots, L, \quad \text{and} \quad k = 1, \dots, K_l \quad (4)$$

where $\mu_{\bar{Y}}^{(l,k)} = \frac{1}{N} \sum_i \mu_{Y_i}^{(l,k)}$ and $Z_i(t)$ represents an ‘‘amplitude function’’, which is stretched on the time axis according to a strictly increasing ‘‘warping function’’, $W_i(t)$. In addition ε_{ij} represents *iid* random measurement errors with $E\varepsilon_{ij} = 0$ and $Var(\varepsilon_{ij}) = \sigma^2 < \infty$. Note that we have assumed that the curves are all observed at a common set of points simply for notational convenience. There is nothing in our approach that will prevent it working on curves observed at differing time points.

As with all curve registration methods, (3) has an identifiability problem between Z_i and W_i . Landmark registration achieves identifiable results by assuming certain markers align for every curve. We generalize this approach using the moments condition given by (4) which forces the Z_i 's to have a common ‘‘shape’’. For example, if $I_g^{\max}(t)$, which searches for global maximums, is chosen as the feature function then (4) states that the Z_i 's have a common shape in as much as their global peaks occur at the same time point and that point is equal to the average of the peaks in the observed curves, Y_i . As more feature functions

are chosen (4) forces more alignment in the Z_i 's. Landmark registration can be seen as a special case of (4) because $\mu_{Z_i}^{(l,k)}$ can be used to identify specific marker events in each curve, such as peaks or troughs, in which case (4) simply forces an alignment of landmarks. However, $\mu_{Z_i}^{(l,k)}$ can also be used to measure more general and more global curve characteristics such as the m th derivative as discussed in Section 2. Note that by equating the moments for each curve to $\mu_Y^{(l,k)}$ we are assuming that positive and negative warping cancels out, in terms of the moments, when averaged over all curves. Without this assumption Z_i and W_i will not be identifiable.

We model Z_i and W_i using finite dimensional basis functions. The amplitude function is modeled as $Z_i(t) = \mathbf{z}(t)^T \theta_i$ where $\mathbf{z}(t)$ is a p -dimensional basis function and θ_i represents the corresponding basis coefficients. In the case of the warping functions, since they are restricted to be increasing, we can, without loss of generality, reparameterize them using

$$W_i(t) = \gamma_{i0} + \int_0^t \exp(f_i(s)) ds \quad (5)$$

where γ_{i0} and f_i are unconstrained. As with the amplitude functions we model f using a finite dimensional basis, $f_i(s) = \mathbf{w}(s)^T \gamma_i$, where $\mathbf{w}(s)$ is a q -dimensional basis and γ_i the corresponding coefficients. Several special cases of (5) can be achieved by appropriately restricting the γ_i coefficients. We shall explore two in this paper. The first is the linear warping function $W_i(t) = \alpha_i + \beta_i t$ which is achieved by setting f_i equal to a constant. The second is

$$W_i(t) = \frac{T \int_0^t \exp(f_i(s)) ds}{\int_0^T \exp(f_i(s)) ds}. \quad (6)$$

Equation (6) has the often desirable property that $W_i(0) = 0$ and $W_i(T) = T$ which means that time is taken to run over a consistent time period for all curves. We utilize b-spline bases for both \mathbf{z} and \mathbf{w} but in principal any finite dimensional basis will suffice.

4 Curve Alignment

In this section we detail our curve alignment approach for fitting the model from Section 3.

4.1 A Moments Based Alignment Approach

The aim in fitting our model is to produce estimated curves, $\hat{Y}_{ij} = \mathbf{z}(W_i(t_j))^T \theta_i$, that accurately approximate the observed curves, $Y_{ij} = Y_i(t_j)$, subject to two constraints. First, the shape of the synchronized curves, $Z_i(t)$, should be as close as possible to that of the original curves. Notice that if $W_i'(t) = 1$ for all values of t then $Z_i(t)$ will have an identical shape to $Y_i(t)$. Therefore, we measure the change in shape by examining the departure of $W_i'(t)$ from 1 using $P(W_i) = \left(\int \left\{ [W_i'(t)]^{-1} - 1 \right\} dt \right)^2$ and hence choose a fit such that $P(W_i)$ is small. We penalize the inverse of $W_i'(t)$ to ensure slopes close to zero, which would imply an extremely high level of warping, are strongly discouraged. Second, the shapes of the $Z_i(t)$'s should be as similar as possible to each other. Differences in the shapes can be measured either by examining variability in the θ_i 's from a target μ_θ , $P(\theta_i) = \|\theta_i - \mu_\theta\|^2$, or by concentrating on the spread of the moments, $P(\mu_{Z_i}) =$

$\sum_l \sum_{k=1}^{K_l} \left(\mu_{Z_i}^{(l,k)} - \mu_{\hat{Y}}^{(l,k)} \right)^2$. Hence, we find the θ_i 's, γ_i 's and the μ_θ that minimize

$$Q = \frac{1}{N} \sum_{i=1}^N \left\{ \|\mathbf{Y}_i - \hat{\mathbf{Y}}_i\|^2 + \lambda_{sync} P(\theta_i) + \lambda_{mom} P(\mu_{Z_i}) + \lambda_W P(W_i) \right\} \quad (7)$$

where λ_{sync} , λ_{mom} and λ_W are tuning parameters that determine the impact of each term on the fit. λ_{mom} and λ_{sync} control the balance between the continuous monotone registration and landmark registration methods. Conceptually, setting $\lambda_{mom} = 0$ and minimizing Q is very similar to the continuous monotone registration method of Ramsay and Li (1998). Alternatively, setting $\lambda_{sync} = 0$ and minimizing Q provides a type of generalized landmark registration. Note that including $\|\mathbf{Y}_i - \hat{\mathbf{Y}}_i\|^2$ and $P(\mu_{Z_i})$ ensures that both (3) and (4) from our synchronization model will hold.

For fixed μ_θ , minimizing (7) is relatively simple because we only need minimize Q individually over γ_i and θ_i . This suggests the following iterative algorithm.

1. For fixed μ_θ , minimize Q over γ_i and θ_i for $i = 1$. Repeat for $i = 2, \dots, N$.
2. Set $\mu_\theta = \frac{1}{N} \sum_{i=1}^N \theta_i$.
3. Repeat 1. and 2. until convergence.

Step 1 involves a non-linear optimization but can be achieved with relative ease because we only need optimize over each curve individually and the derivatives of Q can be computed analytically. Note that we optimize over the $\mu_{Z_i}^{(l,k)}$ as part of step 1 i.e. we do not fix $\mu_{Z_i}^{(l,k)}$ at the previous value of θ_i .

Figure 2 uses a simulated data set to illustrate the need for all four terms in (7). Figure 2a) plots ten curves, each generated from the solid grey curve in the center and then ‘‘warped’’ by distorting the time axis by differing amounts. Figure 2b) illustrates the corresponding ten estimates for the Z_i 's, representing the ‘‘synchronized’’ curves, obtained by minimizing (7) with $\lambda_{sync} = \lambda_{mom} = \lambda_W = 0$. The fit is very good, with the estimated standard deviation of the ε_{ij} 's only 0.006, but this approach has clearly done a poor job of synchronizing the data. Alternatively, Figure 2c) shows the results using $\lambda_{sync} = 10$, a small value for λ_W and $\lambda_{mom} = 0$. A high level of synchronization has resulted from the use of $P(\theta_i)$ but the curves bear little relationship to the original ones. In addition, the Z_i 's have been shrunk towards zero resulting in a ten fold increase in the standard deviation of the estimates. As λ_{sync} is reduced and λ_W increased the fit shifts towards that shown in Figure 2b) but at no stage do we get strong synchronization, the correct shape and a good fit to the data. Finally, Figure 2d) provides a plot of the ten estimated Z_i 's after setting $\lambda_{mom} > 0$ and using two moments corresponding to I^{\max} and I^{\min} . Notice that the addition of $P(\mu_{Z_i})$ has enabled us to not only synchronize the data but to also reproduce the original shape of the curves. In addition, the estimated standard deviation is almost identical to that from the fit illustrated in Figure 2b), indicating that the synchronization has not been at the expense of an accurate fit to the data.

There are two reasons for the inadequate fit in Figure 2c). First, because of the significant distortion of the observed Y_i 's the cross-sectional mean, which is used to compute μ_θ , is a poor estimate for the true shape so the curves have been synchronized towards the wrong ‘‘target’’. This is the same problem that one would encounter when using the continuous monotone registration approach on this data. Second, the act of shrinking has resulted in a very poor fit to the original curves. Utilizing $P(\mu_{Z_i})$ has three advantages

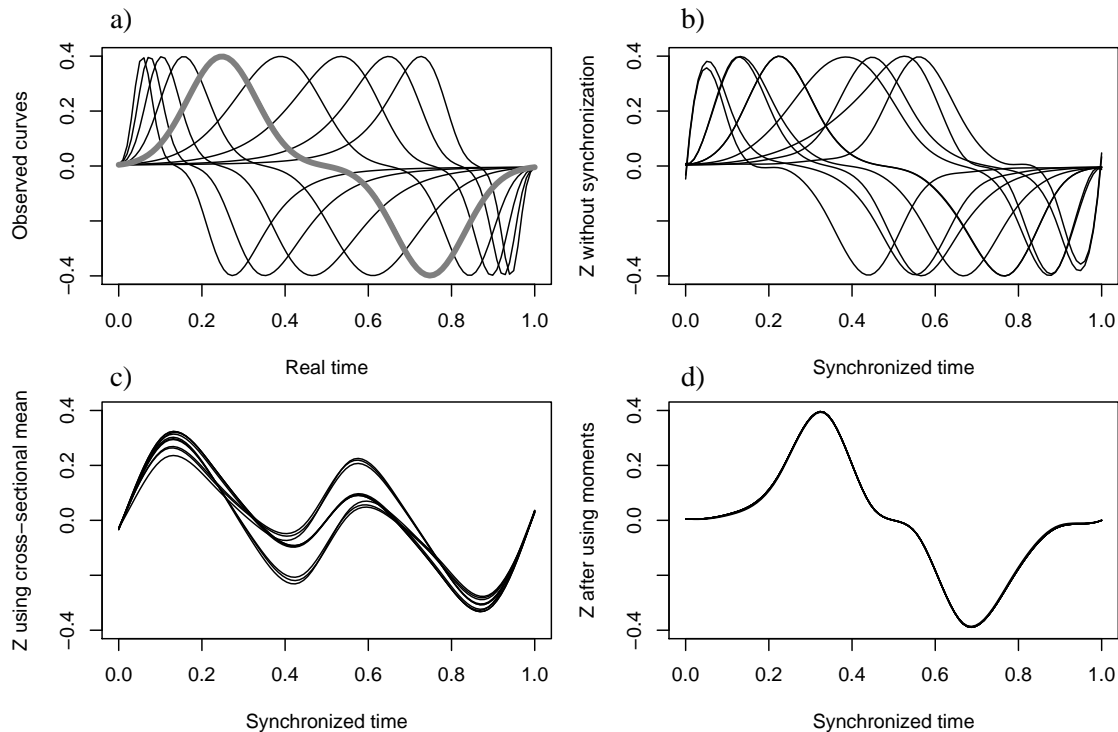


Figure 2: *a) A simulated set of ten curves that have been “warped” in time with the grey line indicating the original shape. b) The estimates for $Z_i(t)$ with $\lambda_{sync} = \lambda_{mom} = \lambda_W = 0$. c) Estimates for $Z_i(t)$ with $\lambda_{mom} = 0$. d) Estimates including all four terms.*

which allows us to address both these problems. First, since the moments are measures of shifts in the time axis, forcing the Z_i 's to have similar moments has no effect on their amplitude and hence does not cause the shrinkage problem observed in Figure 2c). Second, the moments are only a summary of each curve so can often be much more accurately estimated than the entire curve. For example, the cross-sectional mean of the curves in Figure 2a) is a poor estimate of the overall shape of the curves. However, $\mu_{\bar{Y}}^{max}$ and $\mu_{\bar{Y}}^{min}$ still provide good estimates for the maximum and minimum of the original curve that the data was generated from, so the problem of aligning the curves to the wrong shape can be eliminated. Finally, one can choose among a wide range of feature functions when producing the moments. Hence, one can identify specific characteristics or features in the curves and design feature functions accordingly. Since feature functions can theoretically be designed to identify, and hence synchronize towards, any consistent marker events, the landmark registration approach can be seen as a special case of the moments method.

4.2 Asymptotic Theory

Section 5 illustrates the moments method's practical performance on the Berkeley growth curve data and Section 6 provides a comprehensive comparison to other methods on several simulated data sets. However, we can also show that, under general regularity conditions, the method exhibits good large sample properties in terms of asymptotic consistency of the estimators. Let η_0 represent the set of parameters for our model, i.e. $\gamma_1, \dots, \gamma_N$ and $\theta_1, \dots, \theta_N$, and $\hat{\eta}_n$ the corresponding estimates from minimizing (7). Then we first introduce

four assumptions.

A-1 $\mu_{Z_i}^{(l,k)}$ is a continuous function of θ_i for all l and k . Also $\mathbf{z}(W_i(t_j))$ is a continuous function of γ_i .

A-2 $\mathbf{z}(W(t))^T \theta$ is a uniformly continuous function of t i.e. for all $\delta_1 > 0$ there exists $\delta_2 > 0$ such that for all t_1, t_2 , where $|t_1 - t_2| < \delta_2$, it is the case that $|\mathbf{z}(W(t_1))^T \theta - \mathbf{z}(W(t_2))^T \theta| < \delta_1$ for any θ and γ .

A-3 We choose feature functions and corresponding moments such that the synchronization model given by (3) and (4) is identifiable when the curves are observed over a finite set of time points, \mathbf{t} .

A-4 The parameter space is bounded i.e. $\|\eta\|^2 < M$ for some finite M .

We can not hope to have consistent estimators without (A-1) because that would imply that estimating $\mu_{Z_i}^{(l,k)}$ and $\mathbf{z}(W_i(t_j))$ well did not necessarily correspond to estimating the true parameters well. (A-2) places a restriction on the lack of smoothness of the fit. Some level of smoothness must always be imposed on such fits or a line that interpolated the observed values of Y would minimize the criterion. (A-3) is obviously necessary because if the model is unidentifiable we could not select the correct parameters even if we had complete information. (A-4) assumes that the estimators are not allowed to diverge off to infinity. Subject to these four assumptions, we provide the following consistency result.

Theorem 2 *Let $\lambda_{sync,n}$, $\lambda_{W,n}$ and $\lambda_{mom,n}$ represent the tuning parameters as a function of n . Suppose that (A-1) through (A-4) hold, that $\lambda_{sync,n}$ and $\lambda_{W,n}$ are $o(n)$ and that $\lambda_{mom,n}$ is $O(n)$. Then $\hat{\eta}_n$ will be a consistent estimator for η_0 i.e. $\hat{\eta}_n \rightarrow \eta_0$ a.s.*

4.3 Selection of Tuning Parameters

A key component of our synchronization approach is the choice of the tuning parameters λ_{sync} , λ_W and λ_{mom} . The choice of these parameters is governed by a tradeoff between quality of fit i.e. how well the estimated curves fit the observed data, the level of synchronization achieved and the amount of distortion to the shape in performing the synchronization. In general, improving performance in one of these characteristics will cause a deterioration in the other two. An analogy would be choosing between small probabilities of type 1 and type 2 errors in hypothesis tests. Of course the standard approach in that setting is to minimize the probability of a type 2 error subject to an upper bound constraint on the probability of a type 1 error. We take a similar approach here by selecting the tuning parameters to produce the best possible synchronization subject to constraints on the lack of fit and the distortion of the shape.

We measure the level of synchronization, $Sync$, using the average squared deviation of the synchronized curves from their mean curve as a percentage of the same quantity for the unsynchronized curves. Hence, a value of zero would indicate an identical shape for all synchronized curves while one corresponds to no improvement in the synchronization. The lack of fit, σ , is quantified using the average standard deviation between the observed curves, $Y_i(t_j)$, and their “estimates”, $\hat{Y}_i(t_j)$. Finally, the distortion to the shape of the curves is measured using $P(W)$. We then select the tuning parameters so as to minimize $Sync$ subject to σ and $P(W)$ being less than certain upper bounds. Performing this optimization over three parameters is a potentially difficult computational task. Fortunately, the fit turns out to be fairly stable for wide ranges of possible values for λ_W and λ_{mom} , while λ_{sync} has a considerably stronger influence. In the case of λ_{mom}

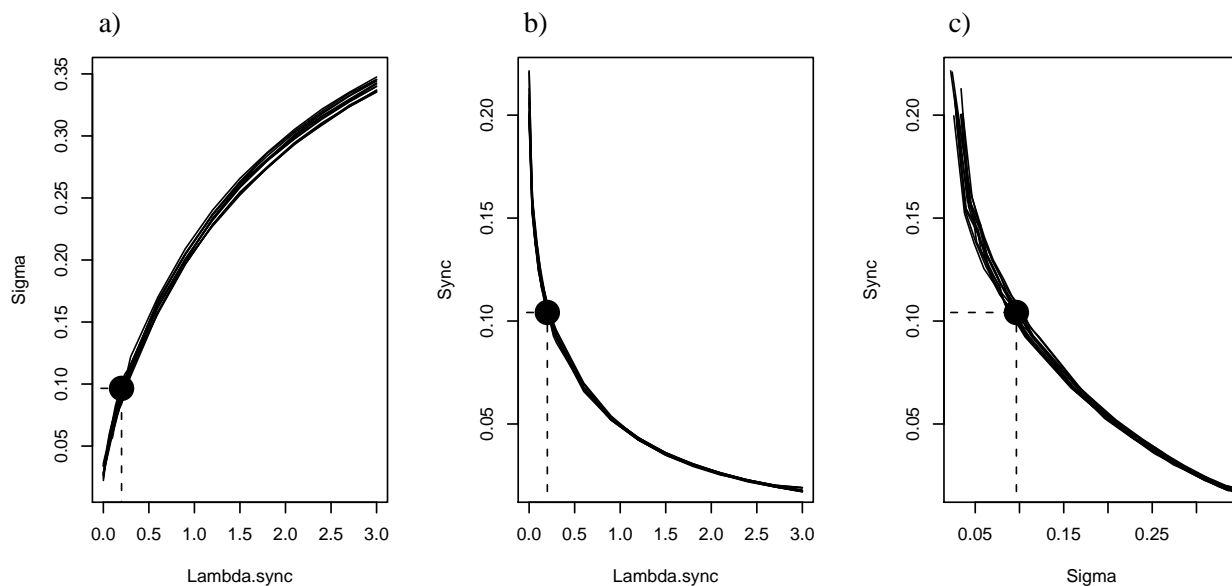


Figure 3: Plots of a) σ versus λ_{sync} , b) Sync versus λ_{sync} and c) Sync versus σ for four different values of λ_{mom} and three different values for λ_W .

it makes intuitive sense that its exact value is not important because the moments are acting to produce an identifiable result so any reasonable weight will make the model identifiable and hence produce a good fit. Hence, it is feasible to implement a grid search over the three parameters where the grid for λ_W and λ_{mom} is very coarse while the grid for λ_{sync} needs to be considerably finer. For the growth curve data, illustrated in the following section, we use values of $10^3, 10^4, 10^5$ and 10^6 for λ_{mom} and values of $10^{-1}, 10^0$ and 10^1 for λ_W . We have found these grids to work well for the problems we have examined. This is consistent with Ramsay and Li (1998) who also found that a small grid of tuning parameters worked over a wide range of applications. In theory cross-validation could be used to select the dimensions of the basis functions \mathbf{z} and \mathbf{w} . However, in practice we have found that, given the flexibility provided by the three tuning parameters, any dimension that provides a reasonably flexible basis will suffice.

5 An Application to the Berkeley Growth Curve Data

In this section we demonstrate the moments based method on the Berkeley growth curve data, discussed in Section 1, utilizing the non-linear warping functions, W_i , given by (6). The data were obtained by fitting a smoothing spline to the second differences of the observed heights for each of ten boys. The smoothing was performed to aid visualizing the resulting curves. We have also performed registration on the raw data with similar results. The first step in implementing our approach involves the choice of the feature functions. This data exhibits clear global maximums and minimums so we elected to utilize I_g^{\max} and I_g^{\min} with $r = 100$. For both feature functions we concentrated on the first moment but one could also have used additional higher

order moments. Next we selected the tuning parameters using the approach from Section 4.3. Figure 3 provides an illustration of this method. Each plot contains 12 separate lines corresponding to four different values for λ_{mom} ($10^3, 10^4, 10^5, 10^6$) and three different values for λ_W ($10^{-1}, 10^0, 10^1$). The 12 lines are almost indistinguishable from each other, emphasizing the insensitivity of the result to the exact choice of λ_{mom} and λ_W . Figure 3a) plots σ as a function of λ_{sync} . Similarly, Figure 3b) plots $Sync$ as a function of λ_{sync} . Finally, Figure 3c) plots $Sync$ as a function of σ . All three plots show a smooth tradeoff between σ and $Sync$ with little effect from the other two tuning parameters. We opted to use tuning parameters that produced the optimal synchronization subject to σ being no larger than 0.1 and $P(W)$ no greater than 0.5. These cutoffs were chosen because they seemed to produce a high level of synchronization with a relatively low increase in σ . The dots on Figure 3 correspond to this fit ($\lambda_{sync} = 0.2, \lambda_W = 10, \lambda_{mom} = 10^5$). We can see that attempting to further synchronize the curves past this point will result in a large increase in σ .

Figure 1b) in Section 1 provides a plot of the synchronized curves, Z_i , from the resulting fit. Notice that the synchronized mean curve not only appears to estimate the correct height for the peaks and troughs but also shifts the peak to a later age from that of the cross-sectional mean. To help judge the accuracy of our procedure Figure 4 provides a comparison to other potential methods. Here we have plotted the estimated mean acceleration curve using five different approaches. In particular we applied our moments method using the above tuning parameters, the moments method with $\lambda_{mom} = 0$, landmark registration (aligning on the peak and the trough of each curve), the continuous monotone registration method, and the cross-sectional mean from the unaligned curves. The cross-sectional mean is well known to be inadequate for this data set (Gasser *et al.*, 1984). However, the landmark method provides a natural gold standard for this problem because it is known to work extremely well in situations such as this one where each curve exhibits a very similar structure (Kneip and Gasser, 1992). All four methods give considerable improvements over the cross-sectional mean but the moments method with $\lambda_{mom} = 10^5$ gives the most similar fit to the landmark approach. The continuous monotone registration method gives the worst performance of the four because it does not take advantage of the specific shape information in the data. Finally, the moments method with $\lambda_{mom} = 0$ gives somewhat intermediate performance. While it does a good job correctly estimating the trough, it fails to identify the correct location of the peak. Again, this is because it fails to make full use of the structure that is present. This illustrates that, while the results are relatively insensitive to the choice of λ_{mom} , this term is still a vital part of the fit.

6 Simulation Study

In this section we compare the performance of our moments based synchronization approach with the continuous monotone registration and landmark methods over four sets of simulations. For each simulation 100 data sets, each consisting of ten curves sampled at 100 equally spaced time points, were generated from a given distribution. Six different synchronization methods were then applied to each data set corresponding to the moments, continuous monotone registration and landmark procedures using both linear and standardized, (6), warping functions. For the moments method we used $K = 1$ moment for each feature function. For each set of simulations the λ parameters were chosen by selecting the values that provided maximum alignment on a preliminary data set subject to constraints on σ and $P(W)$ as discussed in Section 4.3. The simulation results are summarized in Table 1. Two numbers are provided for each simulation-method pair

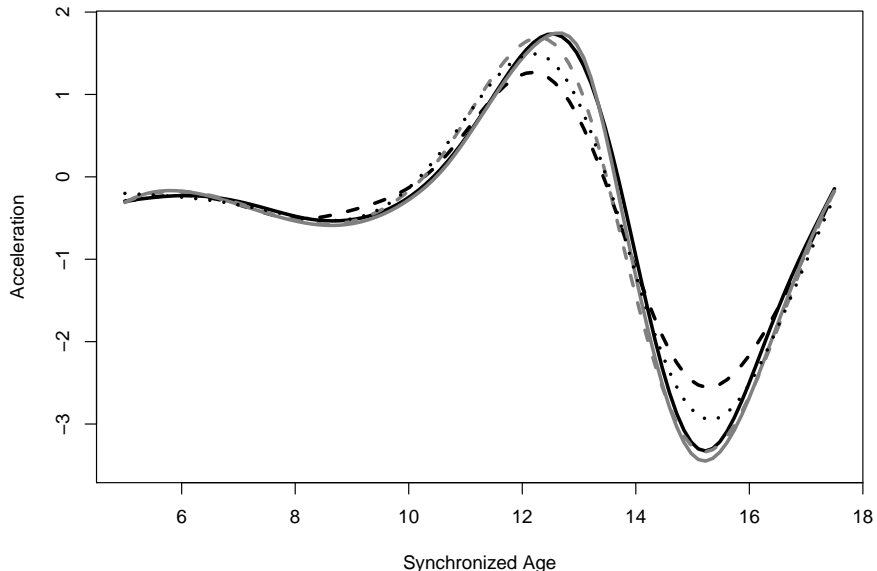


Figure 4: Plots of mean curves on the Berkeley growth curve data using cross-sectional mean (dashed black), continuous monotone registration (dotted), moments method with $\lambda_{mom} = 0$ (dashed grey) and $\lambda_{mom} = 10^5$ (solid grey), and landmark registration (solid black).

corresponding to $Sync$ and σ as defined in Section 4.3. For the moments method σ was produced using $Z_i(W_i(t))$, while for the other two methods it was computed using a smoothing of the curves performed via a smoothing spline prior to synchronization.

Simulation one consisted of curves generated from a standard Gaussian density which were then stretched and shifted in the X or time axis. Figure 5a) illustrates a typical set of curves. We used the peak of each curve as the marker event for the landmark methods and $I^{\max}(t)$ ($r = 100$) for the moments methods. For this simulation the continuous monotone registration and moments methods both worked very well. In particular the continuous monotone registration method produced good results because the cross-sectional mean of the observed curves, used to produce the target function, still had an approximate bell shape. There was little difference between the linear and non-linear warping functions because the true warping was in fact linear. The landmark method, while still providing a considerable level of synchronization, performed relatively less well because, with only one marker, it could not adequately correct for differences in the spread of the curves.

Simulation two had a similar set up to the previous simulation except that half the curves were centered close to 0.7 while the others were centered close to 0.3 (see Figure 5b)). As a result, the cross-sectional mean was bimodal which significantly adversely affected the continuous monotone registration method. The landmark method performed relatively better on this data because shifts in the curve, which it could correct for, formed a larger portion of the lack of synchronization. The moments method was only marginally affected by the bimodal shape of the data.

For simulation three we generated curves using the distribution illustrated in Figure 2a). These curves

$W(t)$	Method	Simulation							
		One		Two		Three		Four	
		<i>Sync</i>	σ	<i>Sync</i>	σ	<i>Sync</i>	σ	<i>Sync</i>	σ
Linear	Cont. Mono. Reg.	0.02	0.0005	76.50	0.0003	78.63*	0.004	75.37*	0.009
	Landmark	11.86*	0.0005	1.15	0.0003	9.39	0.004	15.64	0.009
	Moments	0.07	0.0013	0.50	0.0020	8.33*	0.028	13.71*	0.039
Non-linear	Cont. Mono. Reg.	0.06	0.0005	39.47	0.0003	21.55*	0.004	21.18	0.009
	Landmark	12.32*	0.0005	6.31	0.0003	2.86	0.004	1.42	0.009
	Moments	< 0.01	0.0013	0.59	0.0008	0.76	0.009	1.20	0.014

Table 1: Results from four simulations on six different alignment methods. *Sync* is measured as a percentage so 100 corresponds to no improvement in synchronization. The standard errors on the *Sync* were between 0.15 and 0.45 for those results marked with an * and were less than 0.15 for all others.

were produced using a non-linear warping function and presented a more challenging problem. We utilized both the maximum and minimum points as markers for the landmark methods and $I^{\max}(t)$ and $I^{\min}(t)$ ($r = 100$) for the moments methods. Again, the continuous monotone registration method performed poorly because the cross-sectional mean did not adequately reflect the shape of the curves. The landmark and moments methods both gave good results. For all three procedures the non-linear warping functions worked considerably better than their linear counterparts. Finally, the fourth simulation tested out the effect of noise in the observed curves by adding Gaussian errors with standard deviation of 0.01 to the data from simulation three (see Figure 5c)). We also added a linear drift in the curves to ensure that the moments method still performed well when the curves started and ended at differing values on the Y -axis. In general these changes caused a moderate deterioration in the linear versions of the landmark and moments procedures, presumably because the drift in the curves made it harder for a linear warping function to accurately realign the curves. However, the non-linear versions gave fairly similar performance to those of simulation three. Note that some improvement in the moments method results may have been possible if we had smoothed the curves before applying our approach. However, given the small deterioration from simulation three it is doubtful that any significant gains would have been achieved.

These simulation results may be somewhat unfair to the landmark method because it is difficult to implement this approach in a truly automatic fashion. For example, by manually identifying additional landmarks in the simulated curves one may have been able to produce fits closer to that from the moments approach. However, our attempt here is not necessarily to show that our approach will outperform landmark registration where multiple marker events can be manually identified, since landmark registration is considered the benchmark in this case. Rather, we want to show that the moments method can give comparable results, without the need for manual intervention, when marker events are present, but can also provide accurate results even in the absence of such markers.

Notice that because of the way that the moments method works its σ was somewhat higher on all four simulations than for the continuous monotone registration or landmark methods. This is one of the tradeoffs for a higher level of synchronization. However, the increase is relatively small, particularly for the non-linear warping functions, so the tradeoff clearly seems worthwhile. Simulations two and three

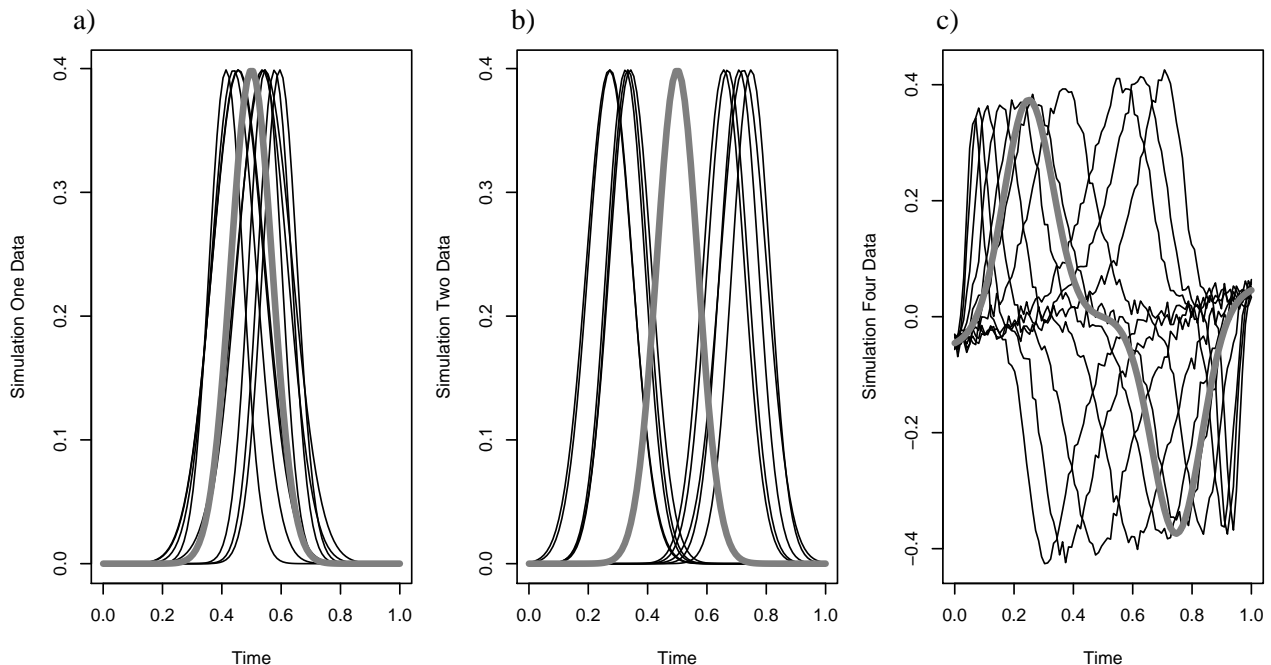


Figure 5: *a) A simulated set of ten curves that have been “warped” in time. This is one of 100 data sets from simulation one. b) One of the 100 data sets from simulation two. c) One of the 100 data sets from simulation four. For each plot, the thicker grey line indicates the original shape.*

illustrate the advantage of combining landmark and continuous monotone registration criteria together. By first synchronizing based on landmarks, such as turning points we can achieve a partial synchronization and then estimate μ_θ well enough to produce a very accurate final alignment. In such situations we have found that the best results are obtained by using a relatively higher value for λ_{mom} in the first few iterations and then reducing λ_{mom} while increasing λ_{sync} in the remaining iterations. This is the approach we took for these simulations.

7 Discussion

In this article we have developed a general moments based approach to the problem of synchronization of functional or curve data. The generally accepted benchmark for such problems is landmark registration which aligns curves by identifying marker events. This approach can be very effective but has two, potentially significant, disadvantages. First, it assumes all curves have consistent marker events and, second, even if the marker events exist one often must manually identify them which is not feasible for large data sets. Alternatively, the continuous monotone registration method works well when an adequate target function, $T(t)$, can be identified but fails when the data is poorly enough aligned that $T(t)$ does not match the shape of the curves. The moments based approach builds on the strengths of both methods and reduces or eliminates their deficiencies. As with the landmark approach, for those curves with marker events, feature functions, such as $I^{\max}(t)$ or $I^{\min}(t)$, can be implemented to synchronize based on these events. However, for curves, or data sets, that do not exhibit such markers more global feature functions, such as $I^{(m)}(t)$, can be utilized.

In this sense our method is an extension of landmark registration. When comparing to the continuous monotone registration approach notice that $\bar{Z}(t) = \mathbf{z}(t)^T \mu_\theta$ can be considered to be the analog of $T(t)$ in that we, at least partially, synchronize the curves towards it. However, even in situations where the cross-sectional mean provides a poor estimate for $T(t)$ and hence the continuous monotone registration method fails, the moments will often induce an accurate enough initial synchronization that $\bar{Z}(t)$ will represent the correct shape. Hence, as the method iterates through the fitting algorithm the synchronization becomes better as opposed to the continuous monotone registration fit where no improvement may be possible. The data from simulations two and three provide a good illustration of this effect. Hence, our approach can also be considered as an extension of continuous monotone registration.

This method could be generalized in several directions. Although, in this article, we have only discussed one-dimensional curves, the moments approach could potentially be extended to multidimensional data. The definition of the feature function, $I_g(t)$, could easily be expanded to such data and hence the moments also. Equating the lower order moments could then be achieved in a similar fashion to the one-dimensional case. The most significant challenge would seem to be dealing with higher order moments on high dimensional data where the number of cross product terms could become unmanageable. Another possible extension is to attempt to model the covariance of the θ_i 's, $Var(\theta_i) = \Theta$. For example, $P(\theta_i)$ could be altered to include Θ using, $P^*(\theta_i) = (\theta_i - \mu_\theta)^T \Theta^{-1} (\theta_i - \mu_\theta)$. There are several possible ways to model Θ . The first, which we have effectively used in $P(\theta_i)$, is to take Θ equal to a multiple of the identity matrix. One could also estimate Θ at each iteration via the sample covariance, $\hat{\Theta} = \frac{1}{N} \sum_i (\theta_i - \mu_\theta)(\theta_i - \mu_\theta)^T$. However, such an unconstrained estimate may be impractical if the dimension of the θ_i 's is large. One solution would be to constrain the rank of Θ and hence significantly reduce the number of parameters to estimate (James *et al.*, 2000). Another appealing alternative would be to design Θ such that $P^*(\theta_i)$ placed no penalty on values of θ_i corresponding to constant vertical shifts of $\mathbf{z}(t)^T \theta_i$. This would mean that two curves that differed only by a constant vertical shift would be considered to be perfectly synchronized and would likely significantly reduce the undesirable shrinkage towards the mean that, for example, is evident in Figure 2c).

Acknowledgements

I would like to thank the Editor, Associate Editor and referees for many helpful suggestions that improved the paper. This work was partially supported by NSF Grant DMS-0705312.

A Appendix

A.1 Proof of Theorem 1

First note that (2) implies that

$$I_g\left(\frac{s-a}{b}\right)(t) = \frac{I_g\left(\frac{t-a}{b}\right)}{\int I_g\left(\frac{t-a}{b}\right) dt} = \frac{1}{b} I_g\left(\frac{t-a}{b}\right)$$

Hence

$$\mu_h^{(1)}\left(\frac{s-a}{b}\right) = \int t I_h\left(\frac{s-a}{b}\right)(t) dt = \int t \frac{1}{b} I_h\left(\frac{t-a}{b}\right) dt = \int (sb+a) I_h(s) ds = b \int I_h(s) ds + a \int I_h(s) ds = b\mu_h^{(1)} + a$$

where $s = \frac{t-a}{b}$. Similarly,

$$\mu_h^{(k)}\left(\frac{s-a}{b}\right) = \int (t - b\mu_h^{(1)} - a)^k I_h\left(\frac{s-a}{b}\right)(t) dt = \int \frac{1}{b} (t - b\mu_h^{(1)} - a)^k I_h\left(\frac{t-a}{b}\right) dt = \int (sb - b\mu_h^{(1)})^k I_h(s) ds = b^k \mu_h^{(k)}.$$

A.2 Proof of Corollary 1

First note that if $I_g^\phi(t) \propto \phi(g(t))$ then $I_{g\left(\frac{s-a}{b}\right)}^\phi(t) \propto \phi\left(g\left(\frac{t-a}{b}\right)\right) \propto I_g^\phi\left(\frac{t-a}{b}\right)$. Next note that

$$\frac{d^m g\left(\frac{t-a}{b}\right)}{dt^m} = \frac{1}{b^m} g^{(m)}\left(\frac{t-a}{b}\right) \quad (8)$$

so $I_{g\left(\frac{s-a}{b}\right)}^{(m)}(t) \propto |g^{(m)}\left(\frac{t-a}{b}\right)| \propto I_g^{(m)}\left(\frac{t-a}{b}\right)$. To show the result for I_g^{\max} note that

$I_{g\left(\frac{s-a}{b}\right)}^{\max}(t) \propto (g\left(\frac{t-a}{b}\right) - \min\{g\left(\frac{t-a}{b}\right)\})^\delta = (g\left(\frac{t-a}{b}\right) - \min\{g(t)\})^\delta \propto I_g^{\max}\left(\frac{t-a}{b}\right)$ and similarly for I_g^{\min} . Finally, by (8)

$$\frac{\frac{dg\left(\frac{t-a}{b}\right)}{dt}}{\sqrt{\frac{d^2 g\left(\frac{t-a}{b}\right)}{dt^2}}} = \frac{g^{(1)}\left(\frac{t-a}{b}\right)/b}{\sqrt{g^{(2)}\left(\frac{t-a}{b}\right)/b^2}} = \frac{g^{(1)}\left(\frac{t-a}{b}\right)}{\sqrt{g^{(2)}\left(\frac{t-a}{b}\right)}}$$

so

$$I_{g\left(\frac{s-a}{b}\right)}^{\text{local}}(t) \propto \exp\left(-\delta \frac{\frac{dg\left(\frac{t-a}{b}\right)}{dt}}{\sqrt{\frac{d^2 g\left(\frac{t-a}{b}\right)}{dt^2}}}\right) \propto \exp\left(-\delta \frac{g^{(1)}\left(\frac{t-a}{b}\right)}{\sqrt{g^{(2)}\left(\frac{t-a}{b}\right)}}\right) \propto I_g^{\text{local}}\left(\frac{t-a}{b}\right).$$

A.3 Proof of Theorem 2

First we state and prove a lemma.

Lemma 1 *Suppose*

$$\sup_t |\mathbf{z}(\hat{W}_n(t))^T \hat{\theta}_n - \mathbf{z}(W_0(t))^T \theta_0| \rightarrow 0 \quad a.s. \quad (9)$$

and

$$\mu_{\hat{Z}_n}^{(l,k)} \rightarrow \mu_{\hat{Y}}^{(l,k)} \quad a.s. \quad \text{for } l = 1, \dots, L \text{ and } k = 1, \dots, K_l \quad (10)$$

where $\hat{Z}_n(t) = \mathbf{z}(t)^T \hat{\theta}_n$ and \hat{W}_n and W_0 respectively represent the warping functions evaluated at $\hat{\gamma}_n$ and γ_0 . Then, provided (A-1), (A-3) and (A-4), hold $\hat{\theta}_n \rightarrow \theta_0$ a.s. and $\hat{\gamma}_n \rightarrow \gamma_0$ a.s.

A.3.1 Proof of Lemma 1

Note we treat each curve individually so we drop the subscript i and let $\eta_0 = \begin{pmatrix} \gamma_0 \\ \theta_0 \end{pmatrix}$ and $\hat{\eta}_n = \begin{pmatrix} \hat{\gamma}_n \\ \hat{\theta}_n \end{pmatrix}$. To reduce notation let

$$f(\eta, t) = \mathbf{z}(W(t))^T \theta.$$

First, note that (9) and (10) imply that there exists Ω^* with $P(\Omega^*) = 1$ s.t. $\forall \omega^* \in \Omega^*$,

$$f(\hat{\eta}_n(\omega^*), t) \rightarrow f(\eta_0, t) \quad \forall t \quad (11)$$

and

$$\mu_{Z_n}^{(l,k)}(\omega^*) \rightarrow \mu_{\hat{Y}}^{(l,k)} \quad \text{for } l = 1, \dots, L \text{ and } k = 1, \dots, K_l. \quad (12)$$

Now, suppose that $\hat{\eta}_n$ does not converge a.s. to η_0 . This implies there exists Ω with $P(\Omega) > 0$ s.t. $\forall \omega \in \Omega$, $\hat{\eta}_n(\omega)$ does not converge to η_0 . Since the intersection of Ω^* and Ω must be nonempty we take a particular $\omega \in \Omega^* \cap \Omega$. Then there exists an infinite subsequence $n'(\omega)$ and $\delta(\omega) > 0$ such that

$$\|\hat{\eta}_{n'(\omega)}(\omega) - \eta_0\| > \delta(\omega) \quad (13)$$

for all $n'(\omega)$. But recall that any bounded sequence must have a convergent subsequence. Hence, by boundedness of γ and θ , (A-4), there must be a subsequence, $n''(\omega)$, of $n'(\omega)$, and a $\eta^*(\omega)$, such that

$$\hat{\eta}_{n''(\omega)}(\omega) \rightarrow \eta^*(\omega). \quad (14)$$

Let W^* represent the warping function evaluated at γ^* . Then, since $\mathbf{z}(\hat{W}_n)$ is a continuous function of $\hat{\gamma}_n$ and $\mu_{Z_n}^{(l,k)}$ is a continuous function of $\hat{\theta}_n$ (by (A-1)), f is continuous and hence (14) implies that

$$f(\hat{\eta}_{n''(\omega)}(\omega), t) \rightarrow f(\eta^*(\omega), t) \quad \forall t \quad (15)$$

and

$$\mu_{Z_n}^{(l,k)}(\omega) \rightarrow \mu_{Z^*}^{(l,k)} \quad \text{for } l = 1, \dots, L \text{ and } k = 1, \dots, K_l. \quad (16)$$

Now, (11) and (15) imply that $f(\eta_0, t) = f(\eta^*(\omega), t)$ for all t while (12) and (16) imply that $\mu_{\hat{Y}}^{(l,k)} = \mu_{Z^*}^{(l,k)}$ for $l = 1, \dots, L$ and $k = 1, \dots, K_l$. By moments identifiability of the model, (A-3), this implies, $\eta^*(\omega) = \eta_0$. But by (13) and (14), $\|\eta^*(\omega) - \eta_0\| > 0$ which is a contradiction. Hence $\hat{\eta}_n \rightarrow \eta_0$ a.s.

A.3.2 Proof of the theorem

Let η represent the set of parameters for our model i.e. $\gamma_1, \dots, \gamma_N, \theta_1, \dots, \theta_N$, and θ_μ . Each curve is evaluated at n time points, t_1, \dots, t_n . Let

$$\begin{aligned} a_n(\eta) &= \frac{1}{n} \sum_{i=1}^N \sum_{j=1}^n \left(Y_{ij} - \mathbf{z}(W_{ij})^T \theta_i \right)^2, \\ b(\eta) &= \lambda_{mom} \frac{1}{n} \sum_{i=1}^N \sum_{l=1}^L \sum_{k=1}^{K_l} \left(\mu_{Z_i}^{(l,k)} - \mu_{\hat{Y}}^{(l,k)} \right)^2, \end{aligned}$$

$c(\boldsymbol{\eta}) = \sum_{i=1}^N \|\boldsymbol{\theta}_i - \boldsymbol{\mu}_\theta\|^2$ and $d(\boldsymbol{\eta}) = \sum_{i=1}^N P(W_i)$ where $W_{ij} = W_{\gamma_i}(t_j)$. So $Q_n(\boldsymbol{\eta}) = a_n(\boldsymbol{\eta}) + b(\boldsymbol{\eta}) + \frac{\lambda_{\text{sync},n}}{n}c(\boldsymbol{\eta}) + \frac{\lambda_{W,n}}{n}d(\boldsymbol{\eta})$ represents (7) using n time points. Let $\boldsymbol{\eta}_0$ represent the true parameters and $\hat{\boldsymbol{\eta}}_n$ the estimators resulting from minimizing Q_n . Then $Q_n(\boldsymbol{\eta}_0) = a_n(\boldsymbol{\eta}_0) + \frac{\lambda_{\text{sync},n}}{n}c(\boldsymbol{\eta}_0) + \frac{\lambda_{W,n}}{n}d(\boldsymbol{\eta}_0)$ where $c(\boldsymbol{\eta}_0)$ and $d(\boldsymbol{\eta}_0)$ are both finite. Note $c(\boldsymbol{\eta}_0)$ is finite since $\boldsymbol{\theta}$ is bounded and $d(\boldsymbol{\eta}_0)$ is finite because, by (5), $0 < W'_{0_i}(t) < \infty$ for $t \in [0, T]$, provided $f_{0_i}(t)$ is bounded and this is the case because γ_{0_i} is bounded. Also, $b(\boldsymbol{\eta}_0) = 0$ because $\mu_{Z_{0_1}}^{(l,k)} = \mu_{Z_{0_2}}^{(l,k)} = \dots = \mu_{Z_{0_N}}^{(l,k)} = \mu_{\bar{Y}}^{(l,k)}$ for all l and k where $Z_{0_i} = \mathbf{z}^T \boldsymbol{\theta}_{0_i}$. Clearly $Q_n(\hat{\boldsymbol{\eta}}_n) \leq Q_n(\boldsymbol{\eta}_0)$ because $\hat{\boldsymbol{\eta}}_n$ is optimized over all $\boldsymbol{\eta}$. Also, $Q_n(\hat{\boldsymbol{\eta}}_n) \geq a_n(\hat{\boldsymbol{\eta}}_n) + b(\hat{\boldsymbol{\eta}}_n)$ because c and d are positive. Hence

$$a_n(\hat{\boldsymbol{\eta}}_n) + b(\hat{\boldsymbol{\eta}}_n) \leq a_n(\boldsymbol{\eta}_0) + \frac{\lambda_{\text{sync},n}}{n}c(\boldsymbol{\eta}_0) + \frac{\lambda_{W,n}}{n}d(\boldsymbol{\eta}_0) \quad (17)$$

Let $\phi_{nij} = (\mathbf{z}(W_{0_{ij}})^T \boldsymbol{\theta}_{0_i} - \mathbf{z}(\hat{W}_{nij})^T \hat{\boldsymbol{\theta}}_{ni})$ and $\varepsilon_{ij} = (Y_{ij} - \mathbf{z}(W_{0_{ij}})^T \boldsymbol{\theta}_{0_i})$ where $W_{0_{ij}} = W(t_j)$ evaluated using the true γ_i and $\hat{W}_{nij} = W(t_j)$ using γ_i from $\hat{\boldsymbol{\eta}}_n$. Then

$$\begin{aligned} a_n(\hat{\boldsymbol{\eta}}_n) &= \frac{1}{n} \sum_{i=1}^N \sum_{j=1}^n \left(Y_{ij} - \mathbf{z}(\hat{W}_{nij})^T \hat{\boldsymbol{\theta}}_{ni} \right)^2 \\ &= \frac{1}{n} \sum_{i=1}^N \sum_{j=1}^n \left(Y_{ij} - \mathbf{z}(W_{0_{ij}})^T \boldsymbol{\theta}_{0_i} + \mathbf{z}(W_{0_{ij}})^T \boldsymbol{\theta}_{0_i} - \mathbf{z}(\hat{W}_{nij})^T \hat{\boldsymbol{\theta}}_{ni} \right)^2 \\ &= \frac{1}{n} \sum_{i=1}^N \sum_{j=1}^n (\varepsilon_{ij} + \phi_{nij})^2 \\ &= \frac{1}{n} \sum_{i=1}^N \sum_{j=1}^n \varepsilon_{ij}^2 + \frac{1}{n} \sum_{i=1}^N \sum_{j=1}^n \phi_{nij}^2 + 2 \frac{1}{n} \sum_{i=1}^N \sum_{j=1}^n \varepsilon_{ij} \phi_{nij} \\ &= a_n(\boldsymbol{\eta}_0) + \frac{1}{n} \sum_{i=1}^N \sum_{j=1}^n \phi_{nij}^2 + 2 \frac{1}{n} \sum_{i=1}^N \sum_{j=1}^n \varepsilon_{ij} \phi_{nij} \end{aligned} \quad (18)$$

Therefore by (17) and (18)

$$\frac{1}{n} \sum_{i=1}^N \sum_{j=1}^n \phi_{nij}^2 + 2 \frac{1}{n} \sum_{i=1}^N \sum_{j=1}^n \varepsilon_{ij} \phi_{nij} + b(\hat{\boldsymbol{\eta}}_n) \leq \frac{\lambda_{\text{sync},n}}{n}c(\boldsymbol{\eta}_0) + \frac{\lambda_{W,n}}{n}d(\boldsymbol{\eta}_0) \quad (19)$$

But notice that the ε_{ij} 's are iid mean zero random variables. Also ϕ_{nij} is a difference of two bounded uniformly continuous functions so is also bounded and uniformly continuous. (Note $\mathbf{z}(W)^T \boldsymbol{\theta}$ is uniformly continuous by (A-2) and is bounded because it is a continuous function of bounded parameters, γ and $\boldsymbol{\theta}$, by (A-1) and (A-4)). Hence, by a standard application of the SLLN, $\frac{1}{n} \sum_{j=1}^n \varepsilon_{ij} \phi_{nij} \rightarrow 0$ a.s. as $n \rightarrow \infty$. (See Theorem 1.13 (ii) in Shao (2003) for a proof of this result.) In addition $\lambda_{\text{sync},n}$ and $\lambda_{W,n}$ are $o(n)$ so the right hand side of (19) also converges to 0. Therefore it must be the case that

$$\frac{1}{n} \sum_{i=1}^N \sum_{j=1}^n \phi_{nij}^2 \rightarrow 0 \quad \text{a.s. for all } i \quad (20)$$

and

$$b_n(\hat{\eta}_n) \rightarrow 0 \quad a.s. \quad (21)$$

Since λ_{mom} is $O(n)$, (21) implies that (10) in Lemma 1 must hold for each curve i . Finally, to show that (9) holds we divide the time interval $[0, T]$ into H equal sized regions R_1, \dots, R_H . Let $n_h = n/H$ equal the number of time points in region h . Then by (20) it must be the case that for every $\omega > 0$, for large enough n ,

$$\frac{1}{n_h} \sum_{j \in R_h} |\phi_{ij}| < \omega \quad a.s. \quad (22)$$

But by uniform continuity, (A-2), there must be a $\delta_H > 0$ such that

$$|(\mathbf{z}(W_{0i}(t))^T \theta_{0i} - \mathbf{z}(\hat{W}_{n_i}(t))^T \hat{\theta}_{n_i}) - \phi_{ij}| < \delta_H \quad (23)$$

for any t and t_j in R_h . Combining (22) and (23) we see that

$$|(\mathbf{z}(W_{0i}(t))^T \theta_{0i} - \mathbf{z}(\hat{W}_{n_i}(t))^T \hat{\theta}_{n_i})| < \delta_H + \omega \quad (24)$$

for any $t \in R_h$ and large enough n . But by making n large enough, this will apply simultaneously for all regions so (24) will hold for all t . Now send $n \rightarrow \infty, H \rightarrow \infty$ and $n/H \rightarrow \infty$. Then $n_h \rightarrow \infty$ so ω can be made arbitrarily small but also $H \rightarrow \infty$ so δ_H can also be made arbitrarily small. Hence (9) holds for each curve i . Therefore the two conditions for Lemma 1 (9 and 10) have been proved and therefore by Lemma 1 the theorem has been proved.

References

- Bar-Joseph, Z., Gerber, G. K., Gifford, D. K., Jaakkola, T. S., and Simon, I. (2003). Continuous representations of time-series gene expression data. *Journal of Computational Biology* **10**, 341–356.
- Ferraty, F. and Vieu, P. (2002). The functional nonparametric model and applications to spectrometric data. *Computational Statistics* **17**, 545–564.
- Ferraty, F. and Vieu, P. (2003). Curves discrimination: a nonparametric functional approach. *Computational Statistics and Data Analysis* **44**, 161–173.
- Gasser, T., Müller, H. G., Khler, W., Molinari, L., and Prader, A. (1984). Nonparametric regression analysis of growth curves (Corr: V12 p1588). *The Annals of Statistics* **12**, 210–229.
- Gervini, D. and Gasser, T. (2005). Nonparametric maximum likelihood estimation of the structural mean of a sample of curves. *Biometrika* **92**, 801–820.
- James, G. M. and Hastie, T. J. (2001). Functional linear discriminant analysis for irregularly sampled curves. *Journal of the Royal Statistical Society, Series B* **63**, 533–550.
- James, G. M., Hastie, T. J., and Sugar, C. A. (2000). Principal component models for sparse functional data. *Biometrika* **87**, 587–602.

- James, G. M. and Silverman, B. W. (2005). Functional adaptive model estimation. *Journal of the American Statistical Association* **100**, 565–576.
- James, G. M. and Sugar, C. A. (2003). Clustering for sparsely sampled functional data. *Journal of the American Statistical Association* **98**, 397–408.
- Kneip, A. and Gasser, T. (1992). Statistical tools to analyze data representing a sample of curves. *Annals of Statistics* **20**, 1266–1305.
- Kneip, A., Li, X., MacGibbon, K. B., and Ramsay, J. O. (2000). Curve registration by local regression. *The Canadian Journal of Statistics* **28**, 1, 19–29.
- Ramsay, J. O. and Li, X. (1998). Curve registration. *Journal of the Royal Statistical Society, B.* **60**, 351–363.
- Ramsay, J. O. and Silverman, B. W. (2005). *Functional Data Analysis*. Springer, 2nd edn.
- Rice, J. A. and Silverman, B. W. (1991). Estimating the mean and covariance structure nonparametrically when the data are curves. *Journal of the Royal Statistical Society, Ser. B* **53**, 233–243.
- Rice, J. A. and Wu, C. O. (2001). Nonparametric mixed effects models for unequally sampled noisy curves. *Biometrics* **57**, 253–259.
- Rønn, B. B. (2001). Nonparametric maximum likelihood estimation for shifted curves. *Journal of the Royal Statistical Society, Series B, Methodological* **63**, 2, 243–259.
- Sakoe, H. and Chiba, S. (1978). Dynamic programming algorithm optimization for spoken word recognition. *IEEE Transactions on Acoustics, Speech, and Signal Processing* **26**, 43–49.
- Shao, J. (2003). *Mathematical Statistics*. Springer-Verlag New York, Inc., 2nd edn.
- Silverman, B. W. (1995). Incorporating parametric effects into functional principal components analysis. *Journal of the Royal Statistical Society, Sec. B* **57**, 673–689.
- Tuddenham, R. D. and Snyder, M. M. (1954). Physical growth of California boys and girls from birth to eighteen years. *University of California Publications in Child Development* **1**, 183–364.
- Zeger, S. L. and Diggle, P. J. (1994). Semiparametric models for longitudinal data with applications to CD4 cell numbers in HIV seroconverters. *Biometrics* **50**, 689–699.