## ARTICLE IN PRESS

Q1 # A Bayesian method for simultaneous registration and clustering of functional observations

Q2 Zizhen Wu, David B. Hitchcock *

*Department of Statistics, University of South Carolina, Columbia, SC 29208, United States*

A R T I C L E   I N F O

A B S T R A C T

We develop a Bayesian method that simultaneously registers and clusters functional data of interest. Unlike other existing methods, which often assume a simple translation in the time domain, our method uses a discrete approximation generated from the family of Dirichlet distributions to allow warping functions of great flexibility. Under this Bayesian framework, a MCMC algorithm is proposed for posterior sampling. We demonstrate this method via simulation studies and applications to growth curve data and cell cycle regulated yeast genes.

© 2016 Elsevier B.V. All rights reserved.

## 1. Introduction

An important example of exploratory data analysis, cluster analysis involves grouping observations that share similar characteristics. In many clustering methods, similarities or dissimilarities between pairs of observations are measured by some relevant distance metric. Methods based on dissimilarity measures include hierarchical clustering and the K-medoids method (Everitt et al., 2011). Another major clustering technique, model-based clustering, requires statistical assumptions about the observations. A popular model-based method (Fraley and Raftery, 2002) assumes a multivariate normal distribution for the measurements and assigns objects to clusters by comparing the posterior group probabilities given the observations.

If the observations to be clustered are functional data, i.e., repeated measures over time or some other domain, one may consider fitting a function to each observational unit using some basis function expansion (Ramsay and Silverman, 2005). Throughout this paper, we use the B-spline basis (De Boor, 2001). One advantage of functional data analysis over traditional multivariate analysis is the ability to examine higher-order derivatives of fitted functions. For example, the first-order derivative of a fitted monotone smoothing function (Ramsay and Silverman, 2005) measuring children's height over a given period represents the estimated growth velocity, and the second-order derivative is the estimated growth acceleration, etc.

Recently, several methods have been developed for clustering functional data. Luan and Li (2003) use mixed-effect models for time-course gene expression and cluster the curves by calculating their posterior cluster probabilities via the EM algorithm. This mixed-effect model is a special case of the model proposed by James and Sugar (2003).

A more challenging yet often-encountered problem is clustering the observations in the presence of time distortions, also known as phase variation. The time distortion is usually modeled by a warping function $h(\cdot)$ (Ramsay and Silverman, 2005), which is a non-decreasing continuous function defined on the time domain $\mathcal{T}$ satisfying the endpoint conditions $h(a) = a$, and $h(b) = b$, where $a$ and $b$ are two endpoints of the time domain. Fig. 1 shows eight warping functions, while

---

* Corresponding author. Tel.: +1 803 777 5346; fax: +1 803 777 4048.
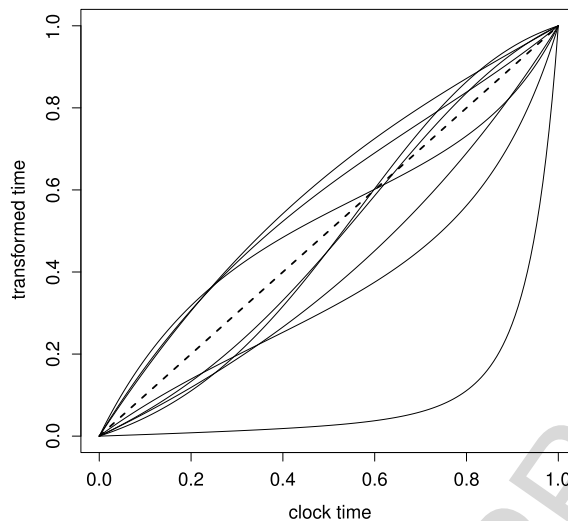   *E-mail address:* hitchcock@stat.sc.edu (D.B. Hitchcock).

**Fig. 1.** Examples of warping functions.

1 the bold dashed line is the 45° reference line representing an identity warp. The cluster structure is blurred by the effect of
2 the time distortions, which should be eliminated for the purpose of clustering. However, it is not feasible to estimate the
3 warping functions without knowing the cluster memberships.

4 Liu and Yang (2009) propose the SACK model, which is capable of clustering functional data when a simple time
5 translation is presented. They translate the shift in the time domain into variation in the measurement space by a first-
6 order Taylor expansion on the B-spline basis functions. The conditional cluster probabilities are calculated via the EM
7 algorithm. Also assuming a simple time translation, Sangalli et al. (2010) propose an iterative method based on a dissimilarity
8 measure called $k$-mean alignment, which iterates among a template identification step, alignment and cluster step, and a
9 normalization step until convergence.

10 To handle more realistic scenarios under arbitrary time warpings, Tang and Müller (2009) propose a method based on
11 pairwise warping. However, this method assumes the mean curves in different clusters are well separated vertically to some
12 degree, a condition potentially too strong for some applications. Zhang and Telesca (2014) propose a hierarchical model for
13 joint curve clustering and registration. They use a reproducing kernel representation of phase variability for registration.

14 In this paper, we develop a Bayesian method for simultaneous clustering and registration of functional data when
15 both arbitrary time distortions and vertical shifts are possible. We model the curves by B-spline basis functions and we
16 approximate the time warping functions by the cumulative sum of realizations from a Dirichlet distribution (following
17 Cheng et al., in press). The posterior cluster memberships are based on a multinomial distribution. Details of our Markov
18 chain Monte Carlo algorithm are introduced in Section 3. Our method of choosing the number of clusters based on a log
19 likelihood is discussed in Section 4. Various simulation studies indicate that our method is capable of clustering curves
20 accurately. We apply our method to the well-known Berkeley growth data and compare our result with that of the SACK
21 model by Liu and Yang (2009) and the kCFC model of Chiou and Li (2007). We also apply our method to the cell-cycle data
22 collected by Alter et al. (2000).

## 2. Model assumption

24 In a functional dataset, we assume that there are $N$ objects, on which we take $K$ measurements over time. Given a certain
25 number of repeated measurements, we may model the response trajectory as a function of time using some basis (such as
26 splines) in the context of functional data analysis.

27 We assume that each observation is composed of a signal function and random error terms, that is,

28 $$\mathbf{Y} = af(\mathbf{t}) + \boldsymbol{\epsilon},$$

29 where $a \in \mathbb{R}^+$ is a stretching/shrinking factor (Zhang and Telesca, 2014), $f(\mathbf{t})$ is the set of underlying responses at the vector
30 of time points $\mathbf{t}$, and $\boldsymbol{\epsilon}$ is an i.i.d. $N(0, \sigma^2)$ error vector.

31 When our observed data must be aligned, we model the effect of the warping function associated with $\mathbf{Y}$ as $\mathbf{Y} = f[h(\mathbf{t})] + \boldsymbol{\epsilon}$,
32 where $h$ is the underlying warping function, and therefore,

33 $$\mathbf{Y}|\boldsymbol{\beta}, \gamma, \sigma^2, a \sim \text{MVN}(af[h(\mathbf{t})], \sigma^2\mathbf{I}).$$

34 For the purpose of clustering, we introduce notation for different groups. For a fixed number of clusters $C$, we use the vector
35 $\mathbf{z}_i = (z_{i1}, z_{i2}, \ldots, z_{iC})$ to denote the cluster membership for the $i$th observation. Note that only one element of $\mathbf{z}_i$ equals