

Linear Transformations and the k -Means Clustering Algorithm: Applications to Clustering Curves

Thaddeus TARPEY

Functional data can be clustered by plugging estimated regression coefficients from individual curves into the k -means algorithm. Clustering results can differ depending on how the curves are fit to the data. Estimating curves using different sets of basis functions corresponds to different linear transformations of the data. k -means clustering is not invariant to linear transformations of the data. The optimal linear transformation for clustering will stretch the distribution so that the primary direction of variability aligns with actual differences in the clusters. It is shown that clustering the raw data will often give results similar to clustering regression coefficients obtained using an orthogonal design matrix. Clustering functional data using an L^2 metric on function space can be achieved by clustering a suitable linear transformation of the regression coefficients. An example where depressed individuals are treated with an antidepressant is used for illustration.

KEY WORDS: Allometric extension; Canonical discriminant analysis; Orthogonal design matrix; Principal component analysis.

1. INTRODUCTION

Functional data applications, where each data point corresponds to a curve, have come to play a prominent role in statistical practice (e.g. Ramsay and Silverman 1997, 2002). The curves in a functional dataset often have a variety of distinctive shapes that can have important interpretations. Representative curve shapes can be found by clustering the curves (e.g. Heckman and Zamar 2000; Abraham et al. 2003; James and Sugar 2003; Luschgy and Pagés 2002; Tarpey and Kinateder 2003). The k -means clustering algorithm (e.g., Forgy 1965; Hartigan and Wong 1979; MacQueen 1967) has been and remains one of the most popular tools for clustering data. When applied to functional data, k -means clustering results vary depending on how the curves are fit to the data. Ultimately, the problem of k -means clustering of functional data boils down to the behavior of the k -means algorithm for different linear transformations of the data which is the focus of this article.

Let $y_1(t), y_2(t), \dots, y_n(t)$ denote a sample of functional responses. In most applications the functions are only observed at

a finite number of time points along with a random error. Thus, a regression model can be used to estimate the function:

$$\mathbf{y}_i = \mathbf{X}\mathbf{b}_i + \boldsymbol{\epsilon}_i, \quad (1)$$

where $\mathbf{y}_i = (y_i(t_1) + \epsilon_{i1}, y_i(t_2) + \epsilon_{i2}, \dots, y_i(t_{m_i}) + \epsilon_{im_i})'$, $\boldsymbol{\epsilon}_i$ is a vector of random errors, \mathbf{b}_i is the $p \times 1$ vector of regression coefficients for the i th function and \mathbf{X} is a design matrix determined by the choice of basis functions used to represent the functions (e.g. Ramsay and Silverman 1997, sec. 3.2). The estimated regression coefficients can be obtained using least-squares

$$\hat{\mathbf{b}}_i = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}_i. \quad (2)$$

A natural way to cluster the functions is to apply the k -means algorithm to the estimated regression coefficients $\hat{\mathbf{b}}_i, i = 1, \dots, n$.

Figure 1 shows fitted curves using a B -spline basis (de Boor 1978) for $n = 414$ depressed subjects treated with Prozac for 12 weeks (McGrath et al. 2000). The functions are the estimated Hamilton Depression (HAM-D) responses as a function of time (scaled to take values between 0 and 1) for each subject where lower HAM-D scores correspond to lower levels of depression. The shapes of the curves are important indicators of the strength of placebo responses and drug responses for individual subjects. However, due to the large number of curves in Figure 1, it is difficult to pick out distinct and representative curve shapes.

Figure 2 shows the $k = 3$ cluster mean curves obtained from the k -means algorithm for four different representations of this functional data. In Figure 2(a), the functional nature of the data is ignored and the raw data (the \mathbf{y}_i 's) were plugged into the k -means algorithm. In panels (b), (c), and (d), the estimated regression coefficients (2) using, respectively, a B -spline basis, a Fourier basis, and a power basis (with an intercept and exponents $-1, 1, 2$, and 3) were clustered. For the power series basis in panel (d), a functional L^2 metric was used in the k -means algorithm instead of the usual Euclidean metric (see Section 3). The B -spline representation used a single knot at $t = 1/2$ which was nearly optimal in terms of a cross-validation prediction error. The resulting design matrix for the B -spline basis was based on $p = 5$ cubic B -spline basis functions. Five basis functions were used for the Fourier and power series as well so that each basis representations corresponds to the same dimension reduction.

The results from clustering shown in panels (a), (b), and (d) in Figure 2 are somewhat similar: (1) the lower curve corresponds to a very strong immediate improvement and then a leveling off indicating an initial placebo response before the drug has an effect; (2) the middle curve shows a steady improvement; (3) the top curve corresponds to subjects experiencing stronger improvement later in the trial, perhaps after the drug has had time to take effect. The cluster mean curves in panel (c) for the Fourier fits are considerably more bunched together and have different

Thaddeus Tarpey is Professor, Department of Mathematics and Statistics, Wright State University, Dayton, OH (E-mail: Thaddeus.tarpey@wright.edu). I am grateful to Minwei Li for programming assistance and to Eva Petkova for helpful discussions related to this work. I wish to thank the referees, an associate editor, and the editor for their comments and suggestions which have improved this article. This work was supported by NIMH grant R01 MH68401.

Hamilton Depression Curves

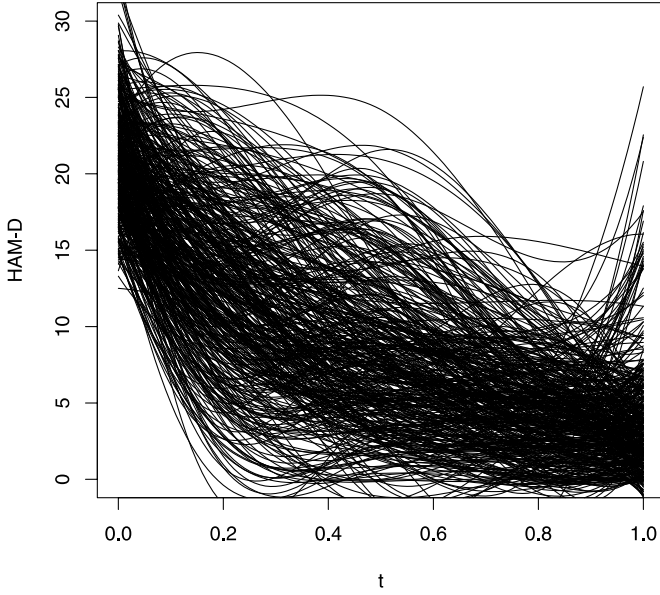


Figure 1. Estimated B -spline curves for HAM-D response over a 12-week period for $n = 414$ depressed subjects taking Prozac.

shapes. Incidentally, the cluster results for the power basis are very similar to the Fourier results (panel (c)) if a Euclidean metric had been used instead of the L^2 metric.

Figure 2 highlights a very important point: The fitted curves for individual subjects using the B -spline, Fourier, and power bases are almost identical and fit the data quite well but the cluster results for the different fits can differ considerably. The reason is because different methods of fitting curves correspond to different linear transformations of the raw data. Therefore the covariance structures for the regression coefficients differ causing the clustering results to differ as well. For instance, the first principal component of the estimated Fourier coefficient distribution accounts for 99.8% of the variability and consequently, the estimated cluster means from the Fourier fits lie approximately along the first principal component axis of the five-dimensional coefficient space. The cluster means from clustering the B -spline coefficients on the other hand lie approximately in the two-dimensional subspace spanned by the first two principal component axes (which explains only 51% and 28% of the variability, respectively).

If \mathbf{A} denotes an arbitrary square invertible matrix, the regression model (1) can be expressed as

$$\mathbf{y}_i = (\mathbf{X}\mathbf{A})(\mathbf{A}^{-1}\mathbf{b}_i) + \boldsymbol{\epsilon}_i = \mathbf{Z}\mathbf{a}_i + \boldsymbol{\epsilon}_i, \quad (3)$$

where $\mathbf{Z} = \mathbf{X}\mathbf{A}$ and $\mathbf{a}_i = \mathbf{A}^{-1}\mathbf{b}_i$. Although the fitted values from (1) and (3) are identical, k -means clustering of the $\hat{\mathbf{b}}_i$'s, and the $\hat{\mathbf{a}}_i$'s generally yield different results. Standardized and weighted regression as well as fitting orthogonal polynomials use transformations similar to (3). It is well known that results from the k -means algorithm depends on how the data is weighted (e.g. Milligan and Cooper 1988; Gnanadesikan, Kettenring, and Tsao 1995; Green, Carmona, and Kim 1990; Milligan 1989). The primary question of interest when clustering functional data

is not necessarily how to choose individual weights, but more generally how best to linearly transform the data prior to clustering.

It is interesting to note that in the Prozac example above that clustering an appropriate linear transformation of the Fourier coefficients produces results almost identical to clustering the B -spline coefficients shown in Figure 2(b). The required linear transformation can be found by letting $\mathbf{P}_F = \mathbf{X}_F(\mathbf{X}_F'\mathbf{X}_F)^{-1}\mathbf{X}_F'$ denote the “hat” matrix for projections onto the column space of \mathbf{X}_F , the Fourier basis design matrix. Let $\hat{\mathbf{X}}_B = \mathbf{P}_F\mathbf{X}_B$ denote the projection of the B -spline basis design matrix \mathbf{X}_B onto the column space of \mathbf{X}_F . Then using $\hat{\mathbf{X}}_B$ in place of \mathbf{X}_B , the estimated B -spline coefficients $\hat{\mathbf{b}}_B$ can be approximated by

$$\begin{aligned} \hat{\mathbf{b}}_B &\approx (\hat{\mathbf{X}}_B'\hat{\mathbf{X}}_B)^{-1}\hat{\mathbf{X}}_B'\mathbf{y} \\ &= (\mathbf{X}_B'\mathbf{P}_F'\mathbf{P}_F\mathbf{X}_B)^{-1}\mathbf{X}_B'\mathbf{P}_B\mathbf{y} \\ &= (\mathbf{X}_B'\mathbf{P}_F\mathbf{X}_B)^{-1}\mathbf{X}_B'\mathbf{X}_F(\mathbf{X}_F'\mathbf{X}_F)^{-1}\mathbf{X}_F'\mathbf{y} \\ &= \mathbf{T}\hat{\mathbf{b}}_F, \end{aligned}$$

where the transformation matrix $\mathbf{T} = (\mathbf{X}_B'\mathbf{P}_F\mathbf{X}_B)^{-1}\mathbf{X}_B'\mathbf{X}_F$, and $\hat{\mathbf{b}}_F$ equals the estimated coefficient vector using the Fourier design matrix. In the Prozac example, the cluster mean curves from clustering the transformed Fourier coefficients $\mathbf{T}\hat{\mathbf{b}}_F$ are essentially indistinguishable from the cluster mean curves obtained by clustering the estimated B -spline coefficients shown in Figure 2(b). It should be noted that transformations of coefficients from one basis to approximate coefficients from another will not always produce nearly identical clustering results.

The remainder of the article is organized as follows: Section 2

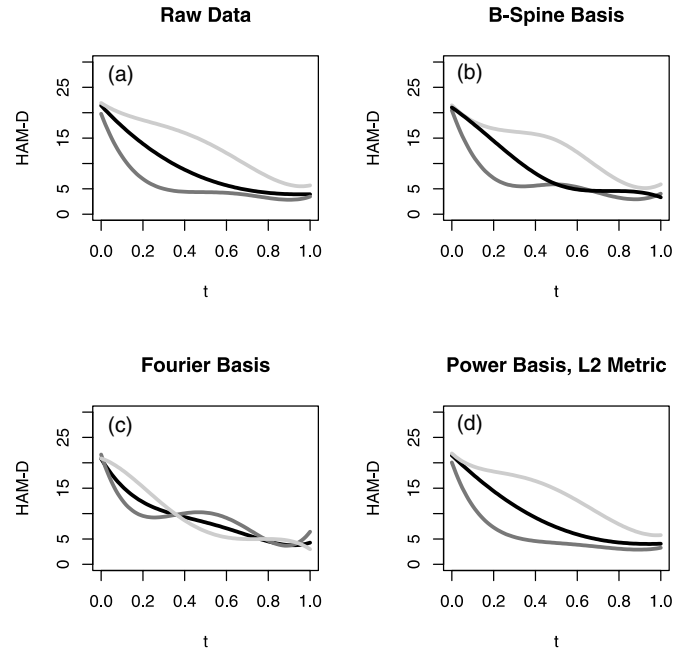


Figure 2. $k = 3$ cluster mean curves for the Prozac data from four different approaches (a) Results from clustering the raw data \mathbf{y}_i 's. (b) Results from clustering estimated B -spline coefficients. (c) Results from clustering estimated Fourier coefficients. (d) Results from clustering estimated coefficients from a power basis fit to the curves.

shows that clustering the raw data will often produce results very similar to clustering estimated regression coefficients from an orthogonal regression. Section 3 shows that the k -means algorithm using an L^2 function space metric is equivalent to clustering regression coefficients after an appropriate linear transformation. Section 4 discusses optimal linear transformations of the data for k -means clustering and provides a simple illustration for clustering linear functions. Section 5 concludes the article.

2. CLUSTERING THE RAW DATA

In the Prozac example of Section 1, it turns out that clustering regression coefficients obtained using an orthogonal design matrix produces cluster mean curves that are essentially indistinguishable from those produced by clustering the raw data shown in Figure 2(a). That is, reducing the dimensionality of the data using the regression model appears to give no advantage over clustering the raw data. This section explains why.

First, certain linear transformations have no effect on k -means clustering. In particular, clustering p -dimensional observations, \mathbf{y}_i 's, and the transformed data

$$\tilde{\mathbf{y}}_i = \boldsymbol{\mu} + c\mathbf{H}\mathbf{y}_i, \quad (4)$$

where $\boldsymbol{\mu} \in \mathbb{R}^p$, $c \in \mathbb{R}^1$, and \mathbf{H} is an orthogonal $p \times p$ matrix, yield identical results. If $\hat{\boldsymbol{\xi}}_j$, $j = 1, \dots, k$, are the cluster means found from clustering the $\tilde{\mathbf{y}}_i$, then $\mathbf{H}'(\hat{\boldsymbol{\xi}}_j - \boldsymbol{\mu})/c$ are the cluster means from clustering the original data.

Let $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}'$ denote the singular value decomposition of the design matrix \mathbf{X} in (1). Then $\mathbf{X}_o = \mathbf{X}\mathbf{V}\mathbf{D}^{-1} = \mathbf{U}$ is an orthogonal design matrix and $\mathbf{b}_o = \mathbf{D}\mathbf{V}\mathbf{b}$ is the vector of associated regression coefficients. The least squares estimator of \mathbf{b}_o is $\hat{\mathbf{b}}_o = (\mathbf{X}_o'\mathbf{X}_o)^{-1}\mathbf{X}_o'\mathbf{y} = \mathbf{U}'(\mathbf{U}\mathbf{b}_o + \boldsymbol{\epsilon}) = \mathbf{b}_o + \mathbf{U}'\boldsymbol{\epsilon}$. Define \mathbf{V} to be a $m \times (m - p)$ matrix with orthonormal columns such that $\mathbf{H} = [\mathbf{U} : \mathbf{V}]$ is an orthogonal $m \times m$ matrix. Since $\mathbf{H}'\mathbf{y}$ is just a rotation of \mathbf{y} , clustering the raw data will yield identical classifications as clustering $\mathbf{H}'\mathbf{y}$ by (4). Now

$$\begin{aligned} \mathbf{H}'\mathbf{y} &= \begin{pmatrix} \mathbf{U}'\mathbf{y} \\ \mathbf{V}'\mathbf{y} \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{U}'(\mathbf{U}\mathbf{b}_o + \boldsymbol{\epsilon}) \\ \mathbf{V}'(\mathbf{U}\mathbf{b}_o + \boldsymbol{\epsilon}) \end{pmatrix} \\ &= \begin{pmatrix} \hat{\mathbf{b}}_o \\ \mathbf{V}'\boldsymbol{\epsilon} \end{pmatrix}. \end{aligned}$$

The raw data, after rotating by \mathbf{H}' , has two orthogonal parts: the estimated orthogonal regression coefficients $\hat{\mathbf{b}}_o$ and a pure error part $\mathbf{V}'\boldsymbol{\epsilon}$. Thus, if the error variance is zero, clustering the raw data is exactly equivalent to clustering the regression coefficients from an orthogonal design matrix. That is, both methods will produce identical clusters. Since the (rotated) raw data has a pure error component $\mathbf{V}'\boldsymbol{\epsilon}$ that presumably contains no information on the true clusters, one would expect that clustering the estimated coefficients $\hat{\mathbf{b}}_o$ from an orthogonal design matrix to do a better job at recovering true clusters than clustering the raw data when the error variance is large. If the error variance is small relative to the variability of the regression coefficients, then clustering the raw data should yield essentially the same

results as clustering estimated coefficients using an orthogonal design matrix which is what happens in the Prozac example of Section 1 when orthogonal design matrices are used.

3. CLUSTERING FUNCTIONAL DATA WITH AN L^2 METRIC

In standard applications of k -means clustering, data points in \mathbb{R}^p are assigned to clusters based minimal Euclidean distance to the cluster centers. If the data are functions, then an L^2 metric in function space may be a more appropriate metric to use for clustering. If $y(t)$ is a functional observation and $\xi(t)$ is a functional cluster mean, then the squared L^2 distance between these two functions on an interval $[T_1, T_2]$ is

$$\|y - \xi\|^2 = \int_{T_1}^{T_2} (y(t) - \xi(t))^2 dt. \quad (5)$$

Suppose functions $y(t)$ are represented using a regression relation

$$y(t) = \beta_0 u_0(t) + \beta_1 u_1(t) + \dots + \beta_p u_p(t) + \epsilon(t).$$

For instance, in a quadratic regression we would have $u_0(t) = 1$, $u_1(t) = t$, $u_2(t) = t^2$. Alternatively, the $u_l(t)$ could be orthogonal polynomials or, in the case of a Fourier expansion, trigonometric functions. Denote the expansion of a cluster mean $\xi(t)$ by $\xi(t) = \sum_{l=0}^p \gamma_l u_l(t)$. The squared L^2 distance between a function $y(t)$ and a cluster mean $\xi(t)$ in (5) can be expressed as

$$\begin{aligned} \|y - \xi\|^2 &= \int_{T_1}^{T_2} (y(t) - \xi(t))^2 dt \\ &= \int_{T_1}^{T_2} \left\{ \sum_{l=0}^p (\beta_l - \gamma_l) u_l(t) \right\}^2 dt \\ &= \sum_{l_1=0}^p \sum_{l_2=0}^p (\beta_{l_1} - \gamma_{l_1})(\beta_{l_2} - \gamma_{l_2}) \int_{T_1}^{T_2} u_{l_1}(t) u_{l_2}(t) dt \\ &= (\boldsymbol{\beta} - \boldsymbol{\gamma})' \mathbf{W} (\boldsymbol{\beta} - \boldsymbol{\gamma}) \\ &= (\mathbf{W}^{1/2} (\boldsymbol{\beta} - \boldsymbol{\gamma}))' (\mathbf{W}^{1/2} (\boldsymbol{\beta} - \boldsymbol{\gamma})), \end{aligned}$$

where $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$ are the vector of regression coefficients for $y(t)$ and $\xi(t)$, respectively, \mathbf{W} is the symmetric $(p + 1) \times (p + 1)$ matrix with elements

$$\int_{T_1}^{T_2} u_{l_1}(t) u_{l_2}(t) dt, \quad (6)$$

and $\mathbf{W}^{1/2}$ is the symmetric square root of \mathbf{W} . Thus, if one wishes to cluster functional data using an L^2 metric, then one can simply plug in the transformed regression coefficients $\mathbf{W}^{1/2}\boldsymbol{\beta}$ into a standard k -means algorithm. This transformation was used in the Prozac example of Section 1 to obtain an L^2 metric clustering of the estimated power series basis coefficients (see Figure 2(d)). As Tarpey and Kinader (2003) noted, clustering regression coefficients is not appropriate if an L^2 metric is desired unless the functions $u_l(t)$ are orthonormal over (T_1, T_2) in which case $\mathbf{W} = \mathbf{I}$ and the L^2 metric in function space is equivalent to a Euclidean metric on the regression coefficients.

4. A CANONICAL TRANSFORMATION FOR CLUSTERING

The k -means algorithm may fail to find true clusters in a dataset if there is substantial variability in the data unrelated to differences in clusters. In fact, there is nothing inherent in the k -means algorithm that guarantees that true clusters will be discovered. Instead the k -means algorithm tends to place sample cluster means where maximal variation occurs in the data. Thus, clustering functional data using the k -means algorithm will perform best if the linear transformations used to fit the curves stretch the data in a direction that corresponds to true cluster differences.

Basically the k -means algorithm begins with an initial set of k cluster means and then assigns individual data points to clusters depending on which cluster center the individual points are nearest. The cluster means are then updated based on the assignment of points to clusters and the algorithm continues to iterate until no more points are reassigned to clusters. Because the algorithm iterates by assigning points to the cluster whose center is closest, the optimization achieved by the algorithm is to find groupings that minimize the within group sum-of-squares, or equivalently, to maximize the between group sum-of-squares.

We will assume that differences between clusters lie in the random regression coefficients \mathbf{b} in (1) and not in the random error ϵ . Let μ_j and Ψ_j denote the mean and covariance matrix, respectively, of the random regression coefficient \mathbf{b} for the j th cluster and let π_j denote the proportion of the population in cluster $j = 1, 2, \dots, k$. The covariance matrix for the \mathbf{b} can be decomposed as

$$\text{cov}(\mathbf{b}) = \mathbf{W} + \mathbf{B}, \quad (7)$$

where

$$\mathbf{W} = \sum_{j=1}^k \pi_j \Psi_j \quad \text{and} \quad \mathbf{B} = \sum_{j=1}^k \pi_j (\mu_j - \mu)(\mu_j - \mu)',$$

are the within cluster and the between cluster covariance matrices, respectively, and where $\mu = \sum_{j=1}^k \pi_j \mu_j$. From (7) one can see that in order to optimize the k -means clustering, a transformation should be used that minimizes the contribution of the within cluster variability while maximizing the between cluster variability. A canonical discriminant function is defined as “linear combinations of variables that best separate the mean vectors of two or more groups of multivariate observations relative to the within-group variance” (Rencher 1993). In canonical discriminant analysis, transformations based on vectors \mathbf{a}_j that successively maximize $(\mathbf{a}_j' \mathbf{B} \mathbf{a}_j) / (\mathbf{a}_j' \mathbf{W} \mathbf{a}_j)$ are used. The solution is to choose the \mathbf{a}_j as the eigenvectors of $\mathbf{W}^{-1} \mathbf{B}$. A canonical transformation for clustering is now defined by first linearly transforming the regression coefficient vector into Fisher’s canonical variates followed by a stretching of the coefficient distribution to accent the between cluster variability and minimize the within cluster variability. In particular, consider a linear transformation that simultaneously diagonalizes \mathbf{W} and \mathbf{B} . Denote the spectral decomposition of $\mathbf{W}^{-1/2} \mathbf{B} \mathbf{W}^{-1/2}$ by $\mathbf{H} \mathbf{D} \mathbf{H}'$ where \mathbf{H} is an orthogonal $p \times p$ matrix and $\mathbf{W}^{1/2}$ is the symmetric square root of \mathbf{W} . Let $\mathbf{\Gamma} = \mathbf{W}^{-1/2} \mathbf{H}$. Assume the eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$ in

\mathbf{D} are arranged from largest to smallest down the diagonal. Then from (7), the covariance matrix of $\mathbf{\Gamma}' \mathbf{b}$ will be

$$\mathbf{I} + \mathbf{D}. \quad (8)$$

In order to accent the between cluster variability and diminish the contribution of the within cluster variability, one can further transform using a *canonical* transformation for clustering

$$\mathbf{C} \mathbf{\Gamma}' \mathbf{b}, \quad (9)$$

where $\mathbf{C} = \text{diag}(c_1, c_2, \dots, c_p)$ and the $c_j \geq 0$ are appropriately chosen constants. From (8), the covariance matrix for the canonically transformed coefficients in (9) is $\mathbf{C}^2 + \mathbf{C}^2 \mathbf{D}$. Thus, choosing large values of c_j corresponding to eigenvalues in \mathbf{D} greater than one inflates the between cluster variability relative to the within cluster variability of the canonically transformed coefficients and setting $c_j = 0$ for eigenvalues between zero and one minimizes the contribution of the within cluster variability. For instance, suppose the cluster means lie on a line. Then multiplying the positive eigenvalue λ_1 in \mathbf{D} by a large value of c_1 transforms the coefficient distribution by stretching it in the direction of the line containing the cluster means. Consequently, the k -means algorithm will place cluster means along this line for large values of c_1 . If the cluster means lie approximately in a q -dimensional plane, then one would choose c_1, \dots, c_q to be large and the remaining c_j to be small. An interesting problem is to determine the optimal settings for the c_j in order to optimize the k -means algorithm according to minimizing a mean squared error or a classification error rate.

The canonical transformation of the regression coefficients in (9) can be adjusted for the random error in a regression model. Letting $\hat{\mathbf{b}} = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y}$ denote the least-squares estimator of \mathbf{b} , it follows from (1) that

$$\text{cov}(\hat{\mathbf{b}}) = \mathbf{W} + \mathbf{B} + \sigma^2 (\mathbf{X}' \mathbf{X})^{-1}, \quad (10)$$

where σ^2 is the error variance and we have assumed the error components are independent. The canonical transformation for $\hat{\mathbf{b}}$ is the same as (9) except \mathbf{W} is replaced by $\mathbf{W} + \hat{\sigma}^2 (\mathbf{X}' \mathbf{X})^{-1}$. The following example illustrates the canonical transformation.

Example: An Simple Illustration of a Canonical Transformation. Consider $k = 3$ clusters of random linear functions $y(t) = b_0 + b_1 t + \epsilon$. The y -intercepts and slopes were simulated from a three-component normal mixture with mean values of y -intercept and slopes equal to $(0, 1)$, $(2, 1)$, and $(3, 3)$ in the three clusters and a common within cluster covariance matrix equal to

$$\Psi = \begin{pmatrix} 2 & -1 \\ -1 & 6 \end{pmatrix}.$$

The proportion of the population in each of the three clusters is taken to be $\pi_1 = \pi_2 = \pi_3 = 1/3$. The error variance is $\sigma^2 = 0.25$. Regression coefficients were estimated via least-squares. In addition, regression coefficient estimates were also estimated using an orthogonal design matrix \mathbf{X}_o where $\mathbf{X}_o' \mathbf{X}_o = \mathbf{I}$. Finally, a canonical transformation for clustering (9) was also used. If c_1 in \mathbf{C} of (9) is too large relative to c_2 , the sample cluster means from the k -means algorithm for the canonically

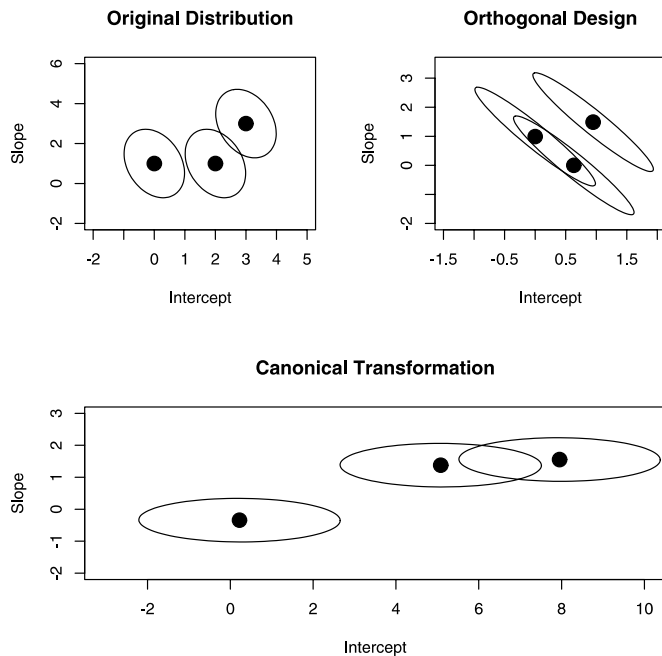


Figure 3. Contours of equal density for a $k = 3$ coefficient distribution. Top-left panel: the original distribution. Top-right panel: coefficient distribution using an orthogonal design matrix. Bottom panel: coefficient distribution using a canonical transformation.

transformed data will lie along a line which will not be optimal because the true cluster means defined above are not colinear. Testing different values of c_1 with simulated data indicated that setting $C = \text{diag}(3.5, 1)$ appears to be a nearly optimal canonical transformation in terms of minimizing the average squared difference between the estimated cluster means and the true cluster means. Figure 3 shows contours of equal density for the three cluster components with the horizontal axis corresponding to the y-intercept and the vertical axis corresponding to the slope. The top-left panel of Figure 3 shows the distribution for the original coefficient distribution; the top-right panel shows the coefficient distribution using the orthogonal design matrix; and the bottom panel shows the coefficients for the canonical transformation distribution.

In order to compare the performance of the three cases shown in Figure 3, 1,000 datasets of size $n = 100$ were simulated for $m_i = 10$ equally spaced time points between $t = 0$ to 1. The k -means algorithm was run for each transformation using the R software (R Development Core Team 2003). The estimated cluster means for the orthogonal design matrix and the canonical coefficient distributions were transformed back to the original scale. For each simulated dataset, a squared error was computed by summing the squared distances between each estimated cluster mean and the nearest true cluster mean. The squared error distributions, plotted in Figure 4, show that clustering the original coefficients (without transforming) does very badly because the primary direction of within-cluster variability is in a direction different than the between-cluster variability. Clustering the coefficients from an orthogonal design matrix does much better in terms of the squared error because the transformation for the orthogonal design matrix provides more separation between

clusters as shown in the top-right panel of Figure 3. The canonical transformation stretches the distributions in the direction with actual cluster differences and therefore performs the best in terms of squared error. With other parameter configurations, it is possible that clustering the original coefficients will perform better than clustering coefficients from an orthogonal design matrix.

The main point of this section is that one should not blindly throw regression coefficients into a clustering algorithm and expect the results to coincide with actual clustering in the data. In particular, as the simulation example above illustrates, the performance of the k -means algorithm for clustering functional data can vary considerably depending how the functional data is transformed prior to clustering.

The optimal canonical transformation for clustering (9) requires knowing the true within- and between-covariance matrices which in practice are unknown. Unfortunately, the sample between covariance matrix $\hat{\mathbf{B}}$ obtained from running the k -means algorithm will often reflect the major variability in the coefficient distribution regardless of whether or not this variability corresponds to true differences in cluster means. In the Prozac example of Section 1, the first eigenvector of the sample between covariance matrix $\hat{\mathbf{B}}$ for the Fourier coefficients is approximately equal to the first eigenvector of the Fourier coefficient covariance matrix, that is, the $k = 3$ estimated cluster means from the k -means algorithm lie approximately on the first principal component axis.

In situations where cluster means lie in a common hyperplane, Bock (1987) proposed a *projection pursuit clustering* algorithm. This algorithm iterates by estimating the common hyperplane using the subspace spanned by the largest eigenvectors from the

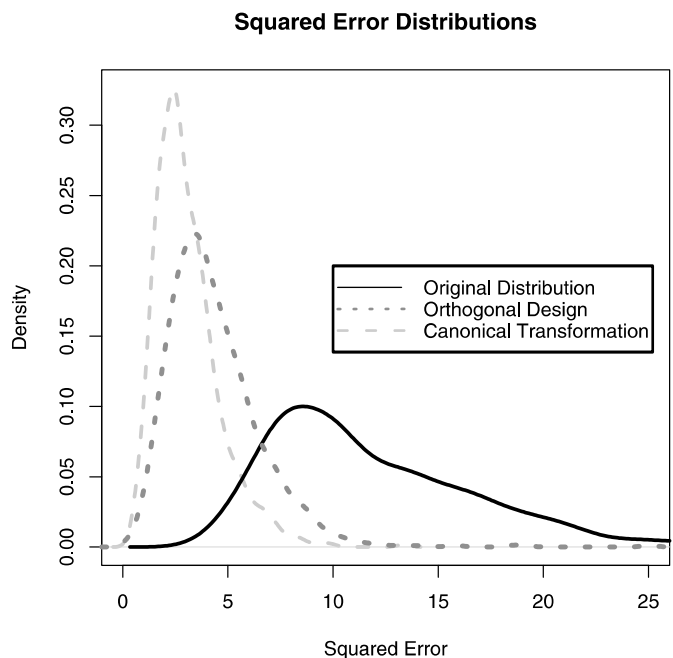


Figure 4. Density distributions for the squared error for clustering the estimated coefficients using (i) the original design matrix (solid curve), (ii) an orthogonal design matrix (dotted curve), and (iii) canonically transformed coefficients (dashed curve).

between group sums-of-squares-and-products matrix and then applying the k -means algorithm to the data projected onto this hyperplane. Bolton and Krzanowski (2003) noted that Bock's algorithm tends to find groups in the direction of the data corresponding to the largest variance and they propose a slightly different projection pursuit index to avoid this problem.

When actual cluster means do indeed lie along the major axis of variation (i.e., the first principal component), the k -means algorithm should perform quite well. This phenomenon occurs frequently in morphometric studies of growth and is called allometric extension (Hills 1982; Bartoletti, Flury, and Nel 1999; Tarpey and Ivey 2006). Let μ_1 and μ_2 denote the means of the two populations and suppose that the eigenvector associated with the largest eigenvalue of the covariance matrices in both populations is the same, call it β_1 . Then the allometric extension model states that $\mu_2 - \mu_1 = \delta\beta_1$ where δ is a constant (Flury 1997, p. 630). The allometric extension model may be reasonable in cases where two (or more) closely related species follow a common growth pattern where one species evolved to a larger overall size. If the first principal component accounts for a large proportion of the overall variance, then the k -means algorithm will tend to place estimated cluster means along the first principal component axis where the true cluster means reside. Thus, one would not want to automatically standardized the data before clustering in these cases because it may hurt the k -means algorithm ability to correctly determine groupings along the primary axis.

In a functional data analysis context, suppose the results of clustering curves produces roughly "parallel" cluster mean curves with the same shape. Parallel cluster mean curves occur quite often in practice when the variability in the intercepts of the curves overwhelms other modes of variation. In these cases, the first principal component variable will tend to coincide with the intercept approximately. Consequently all the cluster mean curves have basically the same shape as the overall mean curve and differ only in their intercepts. This is fine if the actual clusters differ in terms of their intercepts only. However, if curve shapes differ among groups, then the data need to be transformed to minimize the variability of the intercept and allow the k -means algorithm to find distinct curve shapes. A couple possible solutions are to either drop the intercept term when clustering, or to cluster the derivatives of the estimated functions, see Tarpey and Kinateder (2003).

5. DISCUSSION

An appealing aspect of functional data is that the observations are not just ordinary points in Euclidean space, but they are curves with distinct shapes. Clustering functional data is a useful way of determining representative curve shapes in a functional dataset. However, the results from clustering curves depend on how the curves are fit to the data. The k -means clustering algorithm will perform best if the linear transformation used to fit the curves stretches the data in the direction corresponding to true cluster differences. Unfortunately, optimal transformations required for clustering require knowing the true cluster means. A promising approach to solving this problem is to use projection pursuit clustering (Bolton and Krzanowski 2003).

It has been assumed that the error in the regression model

contained no information on the underlying clusters. This assumption may not always hold if the error variances in different clusters differ. In addition, if the wrong model is fit to the data producing nonrandom structure in the residuals, then this structure could contain information on clusters.

Clustering functional data by applying the k -means algorithm to the estimated coefficients is very easy and fast. There are two substantial disadvantages of the k -means algorithm: (i) the algorithm chops up the data into nonoverlapping clusters, whereas in practice distinct groups in the data will often overlap; (ii) the k -means algorithm is completely nonparametric and does not take advantage of any valid parametric assumptions. Finite mixture models do not suffer from these two weaknesses and provide a useful alternative to the k -means algorithm. A simple approach is to plug the estimated coefficients into the EM algorithm for estimating parameters of a finite mixture. A computationally more complicated but highly flexible approach is to express the cluster/mixture model as a random effects model with a latent categorical variable for cluster membership and then estimate the parameters using maximum likelihood via the EM algorithm (James and Sugar 2003; Muthén and Shedden 1999).

[Received April 2006. Revised September 2006.]

REFERENCES

- Abraham, C., Cornillon, P. A., Matzner-Lober, E., and Molinari, N. (2003), "Unsupervised Curve Clustering Using b -Splines," *Scandinavian Journal of Statistics*, 30, 1–15.
- Bartoletti, S., Flury, B., and Nel, D. (1999), "Allometric Extension," *Biometrika*, 55, 1210–1214.
- Bock, H. H. (1987), *On the Interface Between Cluster Analysis, Principal Component Analysis, and Multidimensional Scaling*, Boston: Reidel, pp. 17–34.
- Bolton, R. J., and Krzanowski, W. J. (2003), "Projection Pursuit Clustering for Exploratory Data Analysis," *Journal of Computational and Graphical Statistics*, 12, 121–142.
- de Boor, C. (1978), *A Practical Guide to Splines*, New York: Springer-Verlag.
- Flury, B. (1997), *A First Course in Multivariate Statistics*, New York: Springer.
- Forgy, E. W. (1965), "Cluster Analysis of Multivariate Data: Efficiency vs Interpretability of Classifications," *Biometrics*, 21, 768–769.
- Gnanadesikan, R., Kettenring, J. R., and Tsao, S. L. (1995), "Weighting and Selection of Variables for Cluster Analysis," *Journal of Classification*, 12, 113–136.
- Green, P. E., Carmone, R. J., and Kim, J. (1990), "A Preliminary Study of Optimal Variable Weighting in k -means Clustering," *Journal of Classification*, 7, 271–285.
- Hartigan, J. A., and Wong, M. A. (1979), "A k -Means Clustering Algorithm," *Applied Statistics*, 28, 100–108.
- Heckman, N. E., and Zamar, R. H. (2000), "Comparing the Shapes of Regression Functions," *Biometrika*, 87, 135–144.
- Hills, M. (1982), "Allometry," in *Encyclopedia of Statistical Sciences* (vol. 1), eds. S. Kotz, C. B. Read, and D. L. Banks, New York: Wiley, pp. 48–54.
- James, G., and Sugar, C. (2003), "Clustering for Sparsely Sampled Functional Data," *Journal of the American Statistical Association*, 98, 397–408.
- Luschgy, H., and Pagés, G. (2002), "Functional Quantization of Gaussian Processes," *Journal of Functional Analysis*, 196, 486–531.
- MacQueen, J. (1967), "Some Methods for Classification and Analysis of Multivariate Observations," *Proceedings 5th Berkeley Symposium on Mathematics, Statistics and Probability*, 3, 281–297.
- McGrath, P. J., Stewart, J. W., Petkova, E., Quitkin, F. M., Amsterdam, J. D., Fawcett, J., Reimherr, F. W., Rosenbaum, J. F., and Beasley, C. M. (2000),

- "Predictors of Relapse During Fluoxetine Continuation or Maintenance Treatment for Major Depression," *Journal of Clinical Psychiatry*, 61, 518–524.
- Milligan, G. W. (1989), "A Validation Study of a Variable Weighting Algorithm for Cluster Analysis," *Journal of Classification*, 6, 53–71.
- Milligan, G., and Cooper, M. (1988), "A Study of Standardization of Variables in Cluster Analysis," *Journal of Classification*, 5, 181–204.
- Muthén, B., and Shedden, K. (1999), "Finite Mixture Modeling With Mixture Outcomes Using the EM Algorithm," *Biometrics*, 55, 463–469.
- R Development Core Team (2003), *R: A Language and Environment for Statistical Computing*, Vienna: R Foundation for Statistical Computing.
- Ramsay, J. O., and Silverman, B. W. (1997), *Functional Data Analysis*, New York: Springer.
- (2002), *Applied Functional Data Analysis*, New York: Springer.
- Rencher, A. C. (1993), "Interpretation of Canonical Discriminant Functions, Canonical Variates, and Principal Components," *The American Statistician*, 46, 217.
- Tarpey, T., and Ivey, C. T. (2006), "Allometric Extension for Multivariate Regression Models," *Journal of Data Science*, 4, 479–495.
- Tarpey, T., and Kinateder, K. J. (2003), "Clustering Functional Data," *Journal of Classification*, 20, 93–114.