

KmL: k-means for longitudinal data

Christophe Genolini · Bruno Falissard

Received: 5 May 2009 / Accepted: 12 November 2009 / Published online: 28 November 2009
© Springer-Verlag 2009

Abstract Cohort studies are becoming essential tools in epidemiological research. In these studies, measurements are not restricted to single variables but can be seen as trajectories. Statistical methods used to determine homogeneous patient trajectories can be separated into two families: model-based methods (like Proc Traj) and partition-al clustering (non-parametric algorithms like k-means). KmL is a new implementation of k-means designed to work specifically on longitudinal data. It provides scope for dealing with missing values and runs the algorithm several times, varying the starting conditions and/or the number of clusters sought; its graphical interface helps the user to choose the appropriate number of clusters when the classic criterion is not efficient. To check KmL efficiency, we compare its performances to Proc Traj both on artificial and real data. The two techniques give very close clustering when trajectories follow polynomial curves. KmL gives much better results on non-polynomial trajectories.

Keywords Functional analysis · Longitudinal data · k-means · Cluster analysis · Non-parametric algorithm

C. Genolini (✉) · B. Falissard
Inserm, U669, Paris, France
e-mail: genolini@u-paris10.fr

C. Genolini
Modal'X, Univ Paris Ouest Nanterre La Défense, Paris, France

B. Falissard
Univ Paris-Sud and Univ Paris Descartes, UMR-S0669, Paris, France

B. Falissard
Département de santé publique, AP-HP, Hôpital Paul Brousse, Villejuif, France

1 Introduction

Cohort studies are becoming essential tools in epidemiological research. In these studies, measurements are not restricted to single variables but can be seen as trajectories. As for regular variables, statistical methods can be used to determine homogeneous patient trajectories (Tarpey and Kinader 2003; Rossi et al. 2004; Abraham et al. 2003; James and Sugar 2003). The field of functional cluster analysis can be separated into two families. The first comprises model-based methods. These are related to mixture modelling techniques or latent class analysis (Boik et al. 2008; Atienza et al. 2008). The second family relates to the more classical algorithmic approaches to cluster analysis, such as hierarchical or partitional clustering (Goldstein 1995; Everitt et al. 2001; Ryan 2008). The pros and cons of both approaches are regularly discussed (Magidson and Vermunt 2002; Everitt et al. 2001), even if there is at present little data to show which method is preferable in which situation. In favour of mixture modelling or model-based methods more generally: (1) formal tests can be used to check the validity of the partitioning; (2) results are invariant in linear transformation, so there is no need to standardize variables (this will not be an issue on longitudinal data since all measurements are performed on the same scale), (3) if the model is realistic, inferences about the data-generating process may be possible. On the other hand, traditional algorithmic methods can also have some potential advantages: (1) they do not require any normality or parametric assumptions within clusters (they might be more efficient under a given assumption, but they do not require one; this can be of great interest when the task is to cluster data on which no prior information is available); (2) they are likely to be more robust as regards numerical convergence; (3) in the particular context of longitudinal data, they do not require any assumption regarding the shape of the trajectory (this is likely to be an important point: clustering of longitudinal data is basically an exploratory approach), (4) also in the longitudinal context, they are independent from time-scaling. Even if both methods have been extensively studied, they still present considerable weaknesses, and first of all the difficulty in finding the exact number of clusters. Akaike (1974), Bezdek and Pal (1998), Schwarz (1978), and Sugar and James (2003) provide examples of criteria used to solve this problem. Milligan and Cooper (1985), Shim et al. (2005), Maulik and Bandyopadhyay (2002), Košmelj and Batagelj (1990) compare them using artificial data. Even if criteria perform unequally, all of them fail on a significant proportion of data. Moreover, no study compares criteria specifically on longitudinal data. The problem of cluster selection is indeed an important issue for longitudinal data. More information about clustering longitudinal data can be found in Warren-Liao (2005). Regarding software, longitudinal mixture modeling analysis has been implemented by Jones et al. (2001), Jones and Nagin (2007), Nagin and Tremblay (2001), Jones (2001) in a procedure called Proc Traj on the SAS platform. It has already been extensively used in research on various topics (Jones and Nagin 2007; Clark et al. 2006; Conklin et al. 2005; Nagin 2005). On the R platform (R Development Core Team 2009), S. G. Buyske has proposed the `mmlcr` package, but the statistical background of this routine is not fully documented. `Mplus` (Muthén and Muthén 1998) is also statistical software that provides a general framework that can deal with mixture modeling on longitudinal data. It can be noted that these three procedures are model-based. For the non-parametric solutions, numerous

versions of k-means exist, whether strict (Kaufman and Rousseeuw 1990; Celeux and Govaert 1992) or with variation (Tokushige et al. 2007; Tarpey 2007; García-Escudero and Gordaliza 2005; Vlachos et al. 2003; D'Urso 2004; Lu et al. 2004), but they have considerable drawbacks: 1/ they are not able to deal with missing values; 2/ since the determination of the number of clusters is still an open issue, they require the user to manually re-run k-means several times. In simulation, numerous authors use k-means to compare the different criteria used to find the best cluster number. But the performance of k-means has never been compared to parametric algorithms on longitudinal data.

The rest of this paper is organized as follows: Sect. 2 presents KmL, a package implementing k-means (Lloyd version, 1982) Our package is designed for R platform and is available at (Genolini 2008). It is able to deal with missing values; it also provides an easy way to run the algorithm several times, varying the starting conditions and/or the number of clusters looked for; its graphical interface helps the user to choose the appropriate number of clusters when the classic criterion is not efficient. Section 3 presents simulations on both artificial and real data. Performances of k-means on longitudinal data are compared to Proc Traj results (this appears as the fully dedicated statistical tool that is the most widely used in the literature). Section 4 is the discussion.

2 Algorithm

2.1 Introduction to k-means

K-means is a hill-climbing algorithm (Everitt et al. 2001) belonging to the EM class (Expectation-Maximization) (Celeux and Govaert 1992). Expectation-maximization algorithms work as follows: initially, each observation is assigned to a cluster. Then the optimal clustering is reached by alternating two phases. During the *Expectation* phase, the centers of each cluster (called seeds) are computed. Then the *Maximisation* phase consists in assigning each observation to its “nearest cluster”. The alternation of the two phases is repeated until no further changes occur in the clusters.

More precisely, consider a set S of n subjects. For each subject, an outcome variable Y at t different times is measured. The value of Y for subject i at time k is noted as y_{ik} . For subject i , the sequence y_{ik} is called a trajectory, it is noted $y_i = (y_{i1}, y_{i2}, \dots, y_{it})$. The aim of the clustering is to divide S into g homogeneous sub-groups. Traditionally, k-means can be run using several distances. KmL can use the Euclidean distance $\text{Dist}(y_i, y_j) = \sqrt{\frac{1}{t} \sum_{k=1}^t (y_{ik} - y_{jk})^2}$ and the Manhattan distance $\text{Dist}_M(y_i, y_j) = \frac{1}{t} \sum_{k=1}^t |y_{ik} - y_{jk}|$ (more robust towards outliers (Kaufman and Rousseeuw 1990)).

2.2 Choosing an optimal number of clusters

To choose the optimal number of clusters, KmL uses the Calinski and Harabasz criterion $C(g)$ (Calinski and Harabasz 1974). It has interesting properties, as shown by several authors (Milligan and Cooper 1985; Shim et al. 2005). Let n_m be the number

of trajectories in cluster m ; \bar{y}_m the mean trajectory of cluster m ; \bar{y} the mean trajectory of the whole set S and v' denotes the transposition of vector v . Then the between-variance matrix is $\mathbf{B} = \sum_{m=1}^g n_m (\bar{y}_m - \bar{y})(\bar{y}_m - \bar{y})'$; the trace of the between-variance is the sum of its diagonal coefficients. High between-variance denotes well separated clusters, low between-variance means groups close to each other. The within-variance is $\mathbf{W} = \sum_{m=1}^g \sum_{k=1}^{n_m} (y_{mk} - \bar{y}_m)(y_{mk} - \bar{y}_m)'$. Low within-variance denotes compact groups, high within-variance denotes heterogeneous groups [more details on between and within variance in [Everitt et al. \(2001\)](#)]. The Calinski and Harabasz criterion combines the within and between matrices to evaluate clustering quality. The optimal number of clusters corresponds to the value of g that maximizes $C(g) = \frac{\text{Trace}(\mathbf{B})}{\text{Trace}(\mathbf{W})} \cdot \frac{n-g}{g-1}$ where \mathbf{B} is the between-matrix and \mathbf{W} the within-matrix.

2.3 Avoiding local maxima

One major weakness of hill-climbing algorithms is that they may converge to a local maximum that does not correspond to the best possible clustering in terms of homogeneity. To overcome this problem, different solutions have been proposed. [Hartigan \(1975\)](#) and [Tou and Gonzalez \(1980\)](#) suggest choosing the initial clusters. [Vlachos et al. \(2003\)](#) run a “wavelet” k-means process, modifying the result of a computation and using it as the starting point for the next computation. [Sugar and James \(2003\)](#) and [Hand and Krzanowski \(2005\)](#) suggest running the algorithm several times, and retaining the best solution. It is this approach that has been chosen here. As for the cluster number, the “best” solution is the one that maximizes the between-matrix variance and minimizes the within-variance. Once more, we use the Calinski and Harabasz criterion.

2.4 Dealing with missing value

There are very few studies that try to cluster data assuming missing values ([Hunt and Jorgensen 2003](#)). The simplest way to handle missing data is to exclude trajectories for which certain data are missing. This can severely reduce the sample size, and longitudinal data are especially concerned and subject to missing values (missing values are more likely when an individual is asked to complete certain variables every week than when subjects are asked to complete data only once). In addition, having missing values can be a characteristic that defines a particular cluster, for example an “early drop-out” group.

A different approach has been used here. There is a need to deal with missing data at two different stages. First, during clustering, it is necessary to calculate the distance between two trajectories. Instead of using classic distances as defined in Sect. 2.1, we use distances with Gower adjustment ([Gower 1971](#)): Given y_i and y_j , let w_{ijk} be 0 if y_{ik} or y_{jk} or both are missing, and 1 otherwise; the Euclidian distance with Gower adjustment between y_i and y_j is $\text{Dist}_{\text{Gower}}(y_i, y_j) = \sqrt{\frac{1}{\sum w_{ijk}} \sum_{k=1}^t (y_{ik} - y_{jk})^2 \cdot w_{ijk}}$.

The second problematic step is the calculation of $C(g)$ which helps in the determination of the optimal clustering. At this stage, missing values need to be imputed. We use

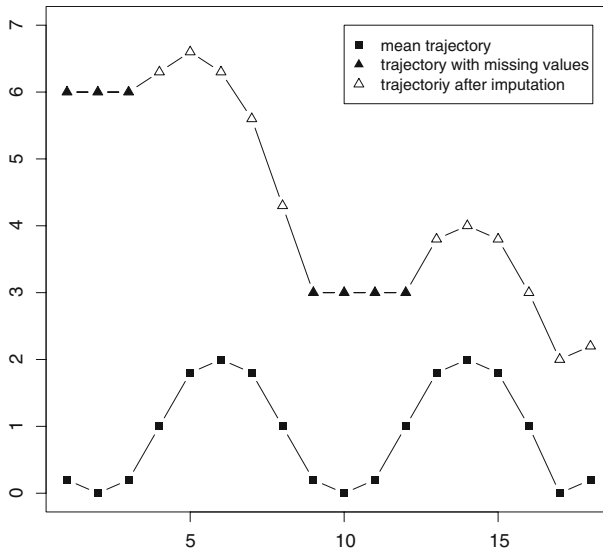


Fig. 1 Example of mean shape copying imputation

the following rules (called *mean shape copying*): if y_{ik} is missing, let y_{ia} and y_{ib} be the closest preceding and following non-missing values of y_{ik} ; let $\overline{y_m} = (\overline{y_{m1}}, \dots, \overline{y_{mt}})$ denote the mean trajectory of y_i cluster. Then $y_{ik} = y_{ia} + (\overline{y_{mk}} - \overline{y_{ma}}) \times \frac{y_{ib} - y_{ia}}{\overline{y_{mb}} - \overline{y_{ma}}}$. If first values are missing, let y_{ib} be the first non-missing value. Then $y_{ik} = y_{ib} + (\overline{y_{mk}} - \overline{y_{mb}})$. If last values are missing, let y_{ia} be the last non-missing value. Then $y_{ik} = y_{ia} + (\overline{y_{mk}} - \overline{y_{ma}})$. Figure 1 gives an example of mean shape copying imputation.

2.5 Implementation of the package

The k-means algorithm used is the Lloyd version (Lloyd 1982). Most of KmL code is written in R using S4 objects (Genolini 2009). The critical part of the programme, clustering, is implemented in two different ways. The first, written in R, provides several options: it can display a graphical representation of the cluster during the convergence of the algorithm; it also lets the user define a distance function that KmL can use to cluster the data. The second, in C (compiled), does not offer any option but is optimized: the C procedure is around 20 times faster than the R procedure. Note that the user does not have to choose between the two functions: KmL automatically selects the fast one when possible, otherwise the slow one.

3 Simulations and applications to real data

3.1 Construction of artificial data sets

To compare the efficiency of Proc Traj and KmL, simulated data were used. We worked on 5,600 data sets defined as follow: a data set is the mixture of several sub-groups.

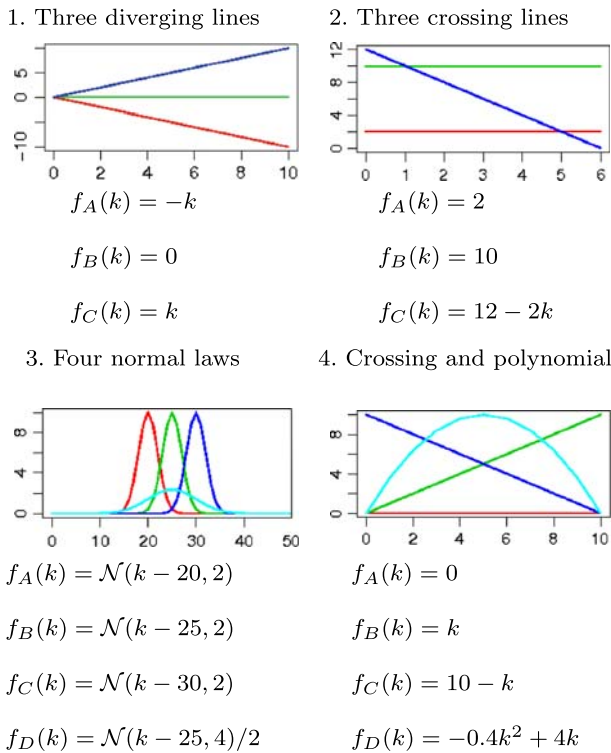


Fig. 2 Trajectory shapes

A subgroup m is defined by a function $f_m(k)$ called the *theoretical trajectory*. Each subject i of a sub-group follows the theoretical trajectory of its subgroup plus a personal variation $\epsilon_i(k)$. The mixture of the different theoretical trajectories is called the *data set shape*. The 5600 data sets were formed varying the data set shape, the number of subjects in each cluster and the personal variations. We defined four data set shapes (presented Fig. 2).

1. “Three diverging lines” is defined by $f_A(k) = -k$; $f_B(k) = 0$; $f_C(k) = k$ with k in $[0 : 10]$.
2. “Three crossing lines” is defined by $f_A(k) = 2$; $f_B(k) = 10$; $f_C(k) = 12 - 2k$ with k in $[0 : 6]$.
3. “Four normal laws” is defined by $f_A(k) = \mathcal{N}(k - 20, 2)$; $f_B(k) = \mathcal{N}(k - 25, 2)$; $f_C(k) = \mathcal{N}(k - 30, 2)$; $f_D(k) = \mathcal{N}(k - 25, 4)/2$ with k in $[0 : 50]$ and $\mathcal{N}(m, \sigma)$ denote the normal law with a mean of m and a standard deviation of σ .
4. “Crossing and polynomial” is defined by $f_A(k) = 0$; $f_B(k) = k$; $f_C(k) = 10 - k$; $f_D(k) = -0.4k^2 + 4k$ with k in $[0 : 10]$.

They were chosen either to correspond to three clearly identifiable clusters (set 1), to present a complex structure (every trajectory intersecting all the others, set 4) or to copy real data (Tremblay (2008) and data presented in Sect. 3.3, sets 2 and 3).

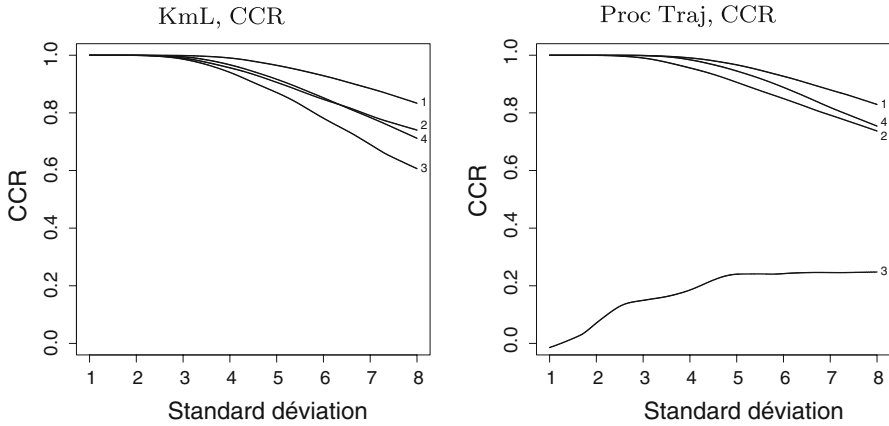


Fig. 3 Comparison of *Correct Classification Rate* between KmL and Proc Traj

Personal variations $\epsilon_i(k)$ are randomised and follow the normal law $\mathcal{N}(0, \sigma)$. Standard deviations increase from $\sigma = 1$ to $\sigma = 8$ (by steps of 0.01). Since the distance between two theoretical trajectories is around 10, $\sigma = 1$ provides “easily identifiable and distinct clusters” whereas $\sigma = 8$ gives “markedly overlapping groups”. The number of subjects in each cluster is set at either 50 or 200. Overall, 4 (data set shape) \times 700 (variance) \times 2 (number of subjects) = 5,600 data sets were created. In a specific data set, the trajectories y_{ik} of an individual belonging to group g is defined by $y_{ik} = fg(k) + \epsilon_i(k)$, with $\epsilon_i(k) \mathcal{N}(0, \sigma^2)$. For the analyses using Proc Traj and KmL, the appropriate number of groups was entered. In addition, the analyses using Proc Traj required the degrees of polynomials that best fitted the trajectories.

3.2 Comparison of KmL and Proc Traj on artificial data sets

Evaluation of KmL and Proc Traj efficiency was performed by measuring two criteria on each clustering C that they found. Firstly, on the artificial data set, the real clustering R is known (the clusters in which each subject should be). The *Correct Classification Rate* (CCR) is the percentage of trajectories that are in the same cluster in C and R (Beauchaine and Beauchaine 2002), that is the percentage of subjects for whom an algorithm makes the right decision. Secondly, working on C , it is possible to evaluate the mean trajectory of each cluster (called the observed trajectory of a cluster). Observed trajectories are an estimation of the theoretical trajectory $f_A(k)$, $f_B(k)$, $f_C(k)$ and $f_D(k)$. An efficient algorithm will find observed trajectories close to the theoretical trajectories. Thus the second criterion, DOT, is the average *Distance between Observed and Theoretical trajectories*. Figures 3 and 4 present the results of the simulations. The graphs present the CCR (resp. the DOT) according to the standard deviation. Table 1 shows the average CCR (resp. the average DOT) for each data set shape.

On dataset shape for 1, 2 and 4, KmL and Proc Traj give very close results whether on CCR or on DOT. In example 3: “Four normal laws”, Proc Traj does not converge,

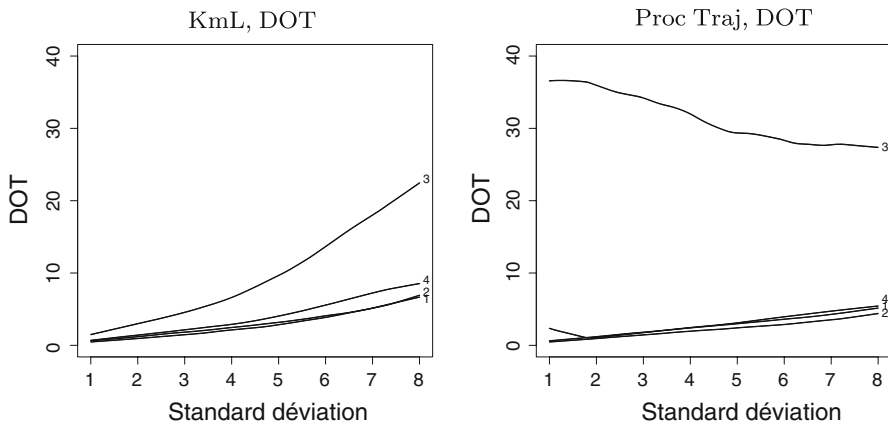


Fig. 4 Comparison of *Distance Observed–Theoretical trajectories* between KmL and Proc Traj

Table 1 Comparison of average DOT and average CCR between KmL and Proc Traj

Data set	KmL	Proc Traj
Average CCR		
1	0.95	0.95
2	0.91	0.91
3	0.86	0.20
4	0.91	0.91
Average DOT		
1	3.17	3.02
2	3.04	2.48
3	9.66	34.28
4	4.24	3.79

or finds results very far removed from the real clusters. KmL performances are as relevant as those obtained on examples 1, 2 and 4.

3.3 Application to real data

The first real example is derived from (Touchette et al. 2007). This study was conducted as part of the Quebec Longitudinal Study of Child Development (Canada) initiated by the Quebec Institute of Statistics. The aim of the study was to investigate the associations between longitudinal sleep duration patterns and behavioral/cognitive functioning at school entry. About 1,492 families participated in the study until the children were 6 years old. Nocturnal sleep duration was measured at 2.5, 3.5, 4, 5, and 6 years of age by an open question on the Self-Administered Questionnaire for the Mother (SAQM). In the original article, a semiparametric model was used to identify subgroups of children who followed different developmental trajectories. They obtained four sleep duration patterns, as illustrated in Fig. 5: a persistent short pattern

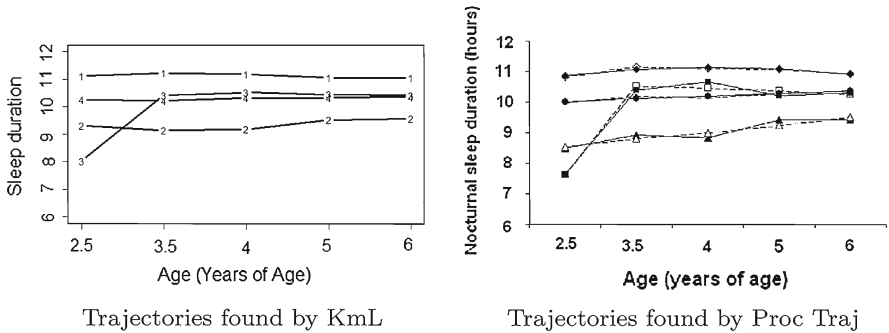


Fig. 5 Sleep duration, means trajectories found by KmL and Proc Traj

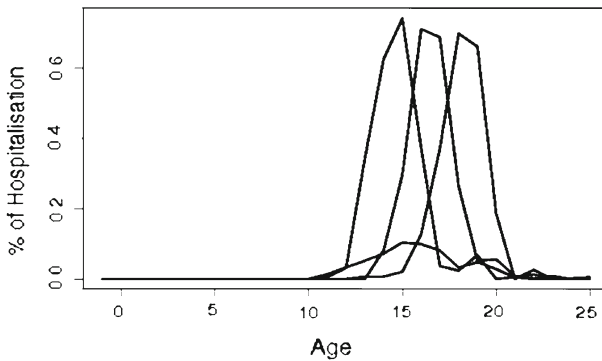


Fig. 6 Hospitalisation length, mean trajectories found by KmL

composed of children sleeping less than 10h per night until age six; a increasing short pattern composed of children who slept fewer hours in early childhood but whose sleep duration increased around 41 months of age, a 10h persistent pattern composed of children who slept persistently approximately 10h per night; and an 11h persistent pattern composed of children who slept persistently around 11h per night.

On this data, KmL finds an optimal solution for a partition into four clusters (as does PROC TRAJ). The trajectories found by both methods are very close (see Fig. 5). The average distance between observed trajectories found by Proc Traj and by KmL is 0.31, which is rather small considering the range of the data (0;12).

The second real example is from a study on the *Trajectories of adolescents hospitalized for Anorexia Nervosa and their social integration in adulthood*, by Hubert, Genolini and Godart (submitted). This study is being conducted at the Institut Mutualiste Montsouris. The authors investigate the relation between adolescent hospitalization for anorexia and their social integration in adulthood. Three hundred and eleven anorexic subjects were included in the study. They were followed from age 0 to 26. The outcome considered here is the annual hospitalisation length, as a percentage. KmL found an optimal solution for a partition into four clusters. The trajectories found by KmL are shown in Fig. 6. Depending on the number of clusters specified in the program, Proc Traj either stated a “false convergence” or gave incoherent results.

4 Discussion

In this article, we present KmL, a new package implementing k-means. The advantage of KmL over the existing procedures (“cluster”, “clusterSim”, “flexclust” or “mclust”) is that it is designed to work specifically on longitudinal data. It provides scope for dealing with missing values; it runs the algorithm several times, varying the starting conditions and/or the number of clusters sought; its graphical interface helps the user to choose the appropriate number of clusters when the classic criterion is not efficient. We also present simulations, and we compare k-means to the latent class model Proc Traj. According to simulations and analysis of real data, k-means seems as efficient as the existing parametric algorithm on polynomial data, and potentially more efficient on non-polynomial data.

4.1 Limitations

The limitations of KmL are inherent in all clustering algorithms. These techniques are mainly exploratory, they cannot statistically test the reality of cluster existence. Moreover, the determination of the optimal cluster number is still an unsettled issue and EM-algorithms can be particularly sensitive to the problem of the local maximum. KmL attempts to deal with these two points by iterating an optimisation process with different initial seeds. Finally, KmL is not model-based, which can be an advantage (non-parametric, more flexible) but also a disadvantage (no scope for testing goodness of fit).

4.2 Advantages

KmL presents some improvement compared to the existing procedures. Since it is a non-parametric algorithm, it does not need any prior information and consequently avoids the issues related to model selection, a frequent concern reported with existing model-based procedures (Nagin 2005, p 65). KmL enables the clustering of trajectories that do not follow polynomial trajectories. Thus, it can deal with a larger set of data (such as Hubert’s hospitalization time in anorexics which follows a normal distribution).

The simulations have shown overall that KmL (like Proc Traj) gives acceptable results for all polynomial examples, even with high levels of noise. A major interest of KmL is that it can work in conjunction with Proc Traj. Finding the number of clusters and the shape of the trajectories (the degree of the polynomial) is still a long and difficult task for Proc Traj users. Running KmL first can give information on both these parameters. In addition, even if Proc Traj has already proved to be an efficient tool in many situations, there is a need to confirm the results, which are mainly of an exploratory nature. When the two algorithms yield similar results, it reinforces confidence in the results.

4.3 Perspectives

A number of unsolved problems need investigation. The optimization of cluster number is a long-standing and important question. Perhaps the particular situation of univariate longitudinal data could yield an efficient solution not yet found in the general context of cluster analysis.

Another interesting point is the generalisation of KmL to problems of higher dimension. At this time, KmL deals only with longitudinal trajectories for a single variable. It would be interesting to develop it for multidimensional trajectories, considering several facets of a patient jointly.

As a last perspective, present algorithms agglomerate trajectories with similar global shape. Thus two trajectories that may be identical in a time translation (one starting early, the other starting late but with the same evolution) will be allocated to two different clusters. One may however consider that the starting time is not really important and that the local shape (the evolution of the trajectory) should be given more emphasis than the overall shape. In this perspective, two individuals with the same development, one starting early and one starting later, would be considered as belonging to the same cluster.

Acknowledgments Thanks to Evelyne Touchette, Tamara Hubert and Nathalie Godart for allowing us to use their data. Thanks to Lionel Riou França, Laurent Orsi and Evelyne Touchette for their helpful advices in programming. *Conflict of interest statement* None

References

- Abraham C, Cornillon P, Matzner-Lober E, Molinari N (2003) Unsupervised curve clustering using B-splines. *Scand J Stat* 30(3):581–595
- Akaike H (1974) A new look at the statistical model identification. *IEEE Trans Autom Control* 19(6):716–723
- Atienza N, García-Heras J, Muñoz-Pichardo J, Villa R (2008) An application of mixture distributions in modelization of length of hospital stay. *Stat Med* 27:1403–1420
- Beauchaine TP, Beauchaine RJ (2002) A Comparison of Maximum Covariance and K-Means Cluster Analysis in Classifying Cases Into Known Taxon Groups. *Psychol Methods* 7(2):245–261
- Bezdek J, Pal N (1998) Some new indexes of cluster validity. In: *IEEE Transactions on Systems, Man and Cybernetics, Part B* 28(3):301–315
- Boik JC, Newman RA, Boik RJ (2008) Quantifying synergism/antagonism using nonlinear mixed-effects modeling: a simulation study. *Stat Med* 27(7):1040–1061
- Calinski T, Harabasz J (1974) A dendrite method for cluster analysis. *Commun Stat* 3(1):1–27
- Celeux G, Govaert G (1992) A classification EM algorithm for clustering and two stochastic versions. *Comput Stat Data Anal* 14(3):315–332
- Clark D, Jones B, Wood D, Cornelius J (2006) Substance use disorder trajectory classes: diachronic integration of onset age, severity, and course. *Addict Behav* 31(6):995–1009
- Conklin C, Perkins K, Sheidow A, Jones B, Levine M, Marcus M (2005) The return to smoking: 1-year relapse trajectories among female smokers. *Nicotine & Tob Res* 7(4):533–540
- D’Urso P (2004) Fuzzy C-means clustering models for multivariate time-varying data: different approaches. *Int J Uncertain Fuzziness Knowl Base Syst* 12(3):287–326
- Everitt BS, Landau S, Leese M (2001) *Cluster analysis*. 4. A Hodder Edwar Arnold Publication, London
- García-Escudero LA, Gordaliza A (2005) A proposal for robust curve clustering. *J Classif* 22(2):185–201
- Genolini C (2008) KmL. <http://christophe.genolini.free.fr/kml/>
- Genolini C (2009) A (Not so) short introduction to S4. <http://cran.r-project.org/>
- Goldstein H (1995) *Multilevel statistical models*. 2. Edwar Arnold, London

- Gower J (1971) A general coefficient of similarity and some of its properties. *Biometrics* 27(4):857–871
- Hand D, Krzanowski W (2005) Optimising k-means clustering results with standard software packages. *Comput Stat Data Anal* 49(4):969–973
- Hartigan J (1975) *Clustering algorithms*. Wiley, New York
- Hunt L, Jorgensen M (2003) Mixture model clustering for mixed data with missing information. *Comput Stat Data Anal* 41(3–4):429–440
- James G, Sugar C (2003) Clustering for sparsely sampled functional data. *J Am Stat Assoc* 98(462):397–408
- Jones BL (2001) Proc traj. <http://www.andrew.cmu.edu/user/bjones/>
- Jones BL, Nagin DS (2007) Advances in group-based trajectory modeling and an SAS procedure for estimating them. *Social Methods & Res* 35(4):542
- Jones BL, Nagin DS, Roeder K (2001) A SAS procedure based on mixture models for estimating developmental trajectories. *Sociological Methods & Research* 29(3):374
- Kaufman L, Rousseeuw PJ (1990) *Finding groups in data: an introduction to cluster Analysis*. Wiley, New York
- Košmelj K, Batagelj V (1990) Cross-sectional approach for clustering time varying data. *J Classif* 7(1):99–109
- Lloyd S (1982) Least squares quantization in PCM. *IEEE Trans Inf Theory* 28(2):129–137
- Lu Y, Lu S, Fotouhi F, Deng Y, Brown SJ (2004) Incremental genetic K-means algorithm and its application in gene expression data analysis. *BMC Bioinformatics* 5:172. <http://www.biomedcentral.com/1471-2105/5/172>
- Magidson J, Vermunt JK (2002) Latent class models for clustering: a comparison with k-means. *Can J Mark Res* 20:37–44
- Maulik U, Bandyopadhyay S (2002) Performance evaluation of some clustering algorithms and validity indices. *IEEE Trans Pattern Anal Mach Intell* 24(12):1650–1654. <http://www.computer.org/portal/web/csdl/doi/10.1109/TPAMI.2002.1114856>
- Milligan GW, Cooper MC (1985) An examination of procedures for determining the number of clusters in a data set. *Psychometrika* 50(2):159–179
- Muthén L, Muthén B (1998) *Mplus user's guide*. Muthén & Muthén 2006, Los Angeles
- Nagin DS (2005) *Group-based modeling of development*. Harvard University Press, Cambridge
- Nagin DS, Tremblay RE (2001) Analyzing developmental trajectories of distinct but related behaviors: a group-based method. *Psychol methods* 6(1):18–34
- R Development Core Team (2009) *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria. <http://www.R-project.org>, ISBN 3-900051-07-0
- Rossi F, Conan-Guez B, Golli AE (2004) Clustering functional data with the SOM algorithm. In: *Proceedings of ESANN*, pp 305–312
- Ryan L (2008) Combining data from multiple sources, with applications to environmental risk assessment. *Stat Med* 27(5):698–710
- Schwarz G (1978) Estimating the dimension of a model. *Ann Stat* 6(2):461–464
- Shim Y, Chung J, Choi I (2005) A comparison study of cluster validity indices using a nonhierarchical clustering algorithm. In: *Proceedings of CIMCA-IAWTIC'05*, IEEE computer society, Washington, vol 1, pp 199–204
- Sugar C, James G (2003) Finding the number of clusters in a Dataset: an information-theoretic approach. *J Am Stat Assoc* 98(463):750–764
- Tarpey T (2007) Linear transformations and the k-means clustering algorithm: applications to clustering curves. *Am Stat* 61(1):34
- Tarpey T, Kinatader K (2003) Clustering functional data. *J classif* 20(1):93–114
- Tokushige S, Yadohisa H, Inada K (2007) Crisp and fuzzy k-means clustering algorithms for multivariate functional data. *Comput Stat* 22(1):1–16
- Tou JTL, Gonzalez RC (1974) *Pattern recognition principles*. Addison-Wesley, Reading
- Touchette E, Petit D, Seguin J, Boivin M, Tremblay R, Montplaisir J (2007) Associations between sleep duration patterns and behavioral/cognitive functioning at school entry. *Sleep* 30(9):1213–1219
- Tremblay RE (2008) *Prévenir la violence dès la petite enfance*. Odile Jacob, Paris
- Vlachos M, Lin J, Keogh E, Gunopulos D (2003) A wavelet-based anytime algorithm for k-means clustering of time series. In: *3rd SIAM international conference on data mining*. San Francisco, May 1–3, 2003, workshop on clustering high dimensionality data and its applications
- Warren-Liao T (2005) Clustering of time series data—a survey. *Pattern Recognit* 38(11):1857–1874