

1 Model based clustering of longitudinal data:
2 application to modeling disease course and gene
3 expression trajectories

4 A. Ciampi¹, H. Campbell, A. Dyachenko, B. Rich, J. McCusker,
5 M. G. Cole

6 **Abstract**

7 We consider the problem of clustering time dependent data. The model is
8 a mixture of regressions, with variance-covariance matrices that are allowed
9 to vary within the extended linear mixed model family. We discuss appli-
10 cations to biomedical data and analyze two longitudinal data sets: one on
11 patients with delirium, and the other on mosquito gene expression following
12 infection.

13 **1 Introduction**

14 Model based clustering (MBC) is increasingly popular as a method for studying
15 multivariate data. Most of the applications rely on the approach proposed by
16 Banfield and Raftery (1993) and implemented in the R-package MCLUST. In
17 this approach data are modeled as a mixture of multivariate normal distributions;
18 an economic parameterization of the variance-covariance matrices is achieved by
19 considering the spectral decomposition of the matrices.

20 New challenges appear when analyzing longitudinal data with non-negligible
21 correlations. We are interested in two broad areas of application: clustering dis-
22 ease trajectories in a clinical setting, and clustering longitudinal data in gene
23 expression at several points in time. The spectral decomposition is of limited help
24 when working with such data, since it does not address the special form that the
25 variance covariance matrices may take. In addition, longitudinal data consist of
26 measurements taken repeatedly on a number of observational units, with the typ-
27 ical feature, especially in clinical settings, that both the number of measurements
28 and the time points may differ across individual units. Analysis of such data is
29 usually performed using the extended linear mixed model (ELMM), see Pinheiro
30 and Bates (2000). However, the ELMM usually assumes a Gaussian distribu-
31 tion for all random effects and error terms. This assumption has been relaxed
32 to include mixtures of Gaussian distributions; see, for instance, Belin and Rubin
33 (1995), Tango (1998), Trottier (1998), Verbeke and Molenberghs (2000), Luan and

¹McGill University & St-Mary's Hospital, Montreal, E-mail:
antonio.ciampi@mcgill.ca

34 Li (2003), Gaffney and Smyth (2003), Celeux et al. (2005), Heard et al. (2006),
 35 Ng et al. (2006), De la Cruz-Mesa (2008). To the best of our knowledge, however,
 36 there is no fully developed method for simultaneously estimating parameters of
 37 both the random effects and the error term.

38 While the primary methods presented in this paper have been developed else-
 39 where, their application is affected by several subjective choices, based on prag-
 40 matic considerations and on assumptions, only valid in specific contexts—these
 41 are the focus of the current paper. We study a general model for representing
 42 mixtures of longitudinal data that generalizes previous attempts: a mixture of
 43 ELLM (Section 2). We develop an EM approach to the estimation of its param-
 44 eters (Section 3) and validate it through simulations (Section 4). In Section 5 is
 45 we apply the model to the study of disease course, analyzing data on delirium
 46 in an elderly population. Section 6 presents an application to longitudinal gene
 47 expression data. Section 7 concludes the paper with a brief discussion.

48 2 The model

Let $Y_i(t_{ij})$ be the observation of the i th individual at time t_{ij} , for $i = 1, \dots, n$, $j = 1, \dots, m_i$, where n is the number of individuals and m_i is the number of time points at which the i th individual has been observed. The ELMM can be written as follows:

$$Y_i = X_i\beta + Z_ib_i + \epsilon_i \quad (1)$$

where X_i and Z_i are design matrices:

$$X_i = \begin{pmatrix} g_1(t_{i1}) & \dots & g_p(t_{i1}) \\ \dots & \dots & \dots \\ g_1(t_{im_i}) & \dots & g_p(t_{im_i}) \end{pmatrix}, \quad Z_i = \begin{pmatrix} h_1(t_{i1}) & \dots & h_q(t_{i1}) \\ \dots & \dots & \dots \\ h_1(t_{im_i}) & \dots & h_q(t_{im_i}) \end{pmatrix},$$

and:

$$\beta = (\beta_1, \dots, \beta_p)', \quad b_i = (b_{i1}, \dots, b_{iq})', \quad \epsilon_i \sim N(0, \sigma^2\Lambda_i)$$

with b_i and ϵ_i assumed independent. Here, Y_i is independent of Y_j for $i \neq j$ and Λ_i is an $m_i \times m_i$ matrix that may depend on i through the time intervals t_{ij} , $j=1, \dots, m_i$ but not otherwise. Typically, Λ_i is parameterized in terms of a relatively small number of variance parameters. Furthermore, the distribution of the random effects, b_i , is assumed to be $N(0, \Psi)$ where Ψ is a symmetric positive definite matrix which may depend on parameters to be estimated. Finally, the g_i 's and the h_i 's denote the elements of a basis in function space. In practice the columns of Z_i are often chosen as a subset of the columns of X_i . We have:

$$Y_i | b_i \sim N(X_i\beta + Z_ib_i, \sigma^2\Lambda_i)$$

and:

$$Y_i \sim N(X_i\beta, \Sigma_i), \quad \Sigma_i = (Z_i\Psi Z_i^T + \sigma^2\Lambda_i).$$

The random effects b_i may be considered as missing data, and maximum likelihood estimation is done by the EM algorithm (see Lindstrom and Bates (1988)).

Since the joint likelihood of $(y_i^T, b_i^T)^T$ is equal to that of $((y_i | b_i)^T, b_i^T)^T$, we have:

$$\begin{bmatrix} y_i | b_i \\ b_i \end{bmatrix} \sim N \left(\begin{pmatrix} X_i \beta + Z_i b_i \\ 0 \end{pmatrix}, \begin{bmatrix} \sigma^2 \Lambda_i & 0 \\ 0 & \Psi \end{bmatrix} \right) \quad (2)$$

Then, as is done for missing data, we can write the ‘‘complete data’’ log-likelihood (Trottier 1998):

$$\begin{aligned} l(\beta, \sigma^2, \Psi, \Lambda | y, b) = & \\ & - \frac{1}{2} \sum_{i=1}^n \left(m_i \log(2\pi) + m_i \log(\sigma^2) + \log(|(\Psi)|) + \log(|(\Lambda_i)|) \right. \\ & \left. + b_i^T (\Psi^{-1}) b_i + \frac{(Y_i - X_i \beta - Z_i b_i)^T (\Lambda_i)^{-1} (Y_i - X_i \beta - Z_i b_i)}{\sigma^2} \right). \end{aligned} \quad (3)$$

Under the assumption that the n individuals are sampled from K distinct component distributions, we can write

$$Y_i = \sum_{k=1}^K \alpha_k (X_i \beta_k + Z_i b_{i(k)} + \epsilon_{i(k)}) \quad (4)$$

49 where the α_k 's are the mixing coefficients:

Under this model formulation, each component, k , is distinct, uniquely defined by β_k , Ψ_k , σ_k^2 , and Λ_k . The log-likelihood of the mixture model can be defined as:

$$\sum_{i=1}^n \log \left(\sum_{k=1}^K \alpha_k \exp \{ l_k(\beta_k, \sigma_k^2, \Psi_k, \Lambda_k | y_i, b_{i(k)}) \} \right) \quad (5)$$

where, $l_k(\cdot)$ is given in Equation (3). Direct maximization of the log-likelihood can be quite difficult due to the sum of terms inside the logarithm. However, we can again complete the data by considering the unobserved latent indicator variables $\delta_{i(k)}$, which is equal to 1 if observation i belongs to cluster k and 0 otherwise, and write the complete data log-likelihood (Celeux et al., 2005):

$$l = \sum_{i=1}^n \sum_{k=1}^K \left\{ \delta_{i(k)} \log(\alpha_k) + \delta_{i(k)} l_k(\beta_k, \sigma_k^2, \Psi_k, \Lambda_k | y_i, b_{i(k)}) \right\} \quad (6)$$

50 where l_k is as in equation (3). Thus, with this ‘‘double completion’’ of the data,
51 maximum likelihood estimates of the parameter vector $\theta = (\alpha, \beta, \sigma^2, \Psi, \Lambda)$ can be
52 obtained using an EM approach as described in the next section.

53 3 Model estimation and inference

54 3.1 EM algorithm

The EM algorithm consists of iterating until convergence between the following E- and M-steps. At iteration $q > 0$, the E-step consists of computing the expectation

of the ‘complete’ log-likelihood knowing the observed data and a current value for the parameters $\theta^{[q]} = (\alpha^{[q]}, \beta^{[q]}, \Psi^{[q]}, (\sigma^2)^{[q]}, \Lambda^{[q]})$, i.e.:

$$Q(\theta | \theta^{[q]}) = E \left[l(\theta | y, \delta, b) | y, \theta^{[q]} \right] = \sum_{i=1}^n \sum_{k=1}^K \left(\tau_{i(k)}^{[q]} \log(\alpha_k) + \tau_{i(k)}^{[q]} E \left[l_k(\beta_k, \sigma_k^2, \Psi_k, \Lambda_k | y_i, b_{i(k)}) | y, \theta^{[q]} \right] \right)$$

where the complete data log-likelihood l is given in Equation (6), and $\tau_{i(k)}^{[q]} = E[\delta_{i(k)} | y, \theta^{[q]}]$ are the so-called posterior probabilities of component membership, computed by:

$$\tau_{i(k)}^{[q]} = P(i \in C_k | y_i, \theta^{[q]}, \alpha^{[q]}) = \frac{\alpha_k^{[q]} g_{m_i}(y_i | \theta_k^{[q]})}{\sum_{l=1}^K \alpha_l^{[q]} g_{m_i}(y_i | \theta_l^{[q]})}$$

55 where C_k denotes the k -th cluster and $g_{m_i}(y_i | \theta_k^{[q]}) = \exp(l_k(\theta_k^{[q]} | y_i))$ denotes the
 56 density of the k -th mixture component. The M-step consists of setting $\theta^{[q+1]} =$
 57 $\arg \max_{\theta} Q(\theta | \theta^{[q]})$. Details are given in Appendix A. Suffice it to say here that
 58 the M-step uses the same numerical methods for the estimation of the parameters
 59 of the Λ_i matrix as in the R functions `lme` and `gls`; therefore, though our approach
 60 follows essentially Celeux et al. (2005), it also borrows from Pinheiro and Bates
 61 (2000).

62 3.2 Initial values

63 It is well known that the EM algorithm can be quite sensitive to the choice of
 64 starting values. A number of different strategies for choosing starting values have
 65 been proposed (McLachlan and Peel, 2000). Following Celeux et al. (2005), we per-
 66 form a large number of short runs (10 iterations) of the EM from different k -means
 67 results. The starting values which initialized the best “short run” solution (i.e.
 68 the short-run solution to achieve the highest log-likelihood), are then selected as
 69 starting values. When the number of observations is not equal across individuals,
 70 k -means is performed on regression parameters obtained from linear regressions
 71 on each individual.

72 3.3 Standard Errors

73 The asymptotic covariance matrix of the maximum-likelihood estimates, $\hat{\theta}$, is equal
 74 to the inverse of the expected Information matrix, $\mathcal{I}(\theta)$, which can be approxi-
 75 mated by following Louis (1982)’s decomposition of $\mathcal{I}(\hat{\theta})$:

$$\mathcal{I}(\hat{\theta}) = E_{\eta}(B(y, \theta)) - E_{\eta}(S(y, \theta)S^T(y, \theta)) + E_{\eta}(S(y, \theta))E_{\eta}(S^T(y, \theta)) \quad (7)$$

where η represents the *missing data*, y , the *observed data* and:

$$B(y, \theta) = \frac{\partial^2 \log l(\theta | y)}{\partial \theta^2}, \quad S(y, \theta) = \frac{\partial \log l(\theta | y)}{\partial \theta} \quad (8)$$

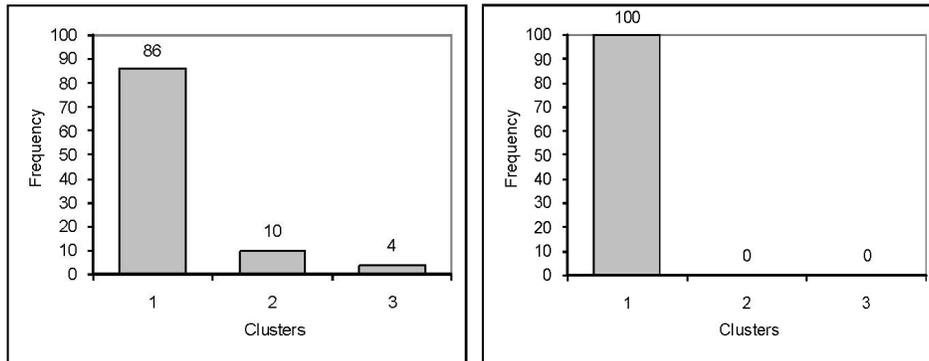


Figure 1: Histograms of number of clusters retrieved by the algorithm when 1 cluster is simulated using AIC(left panel) and BIC(right panel).

76 The standard errors of $\hat{\theta}$ are then given by the diagonal elements of $(\mathcal{I}^{-1}(\hat{\theta}))^{1/2}$.
 77 See Appendix B for details.

78 3.4 Assessing the Number of Clusters

79 Assessing the “correct” number of components or clusters in finite mixture models
 80 is a fundamental and challenging question. The minimum AIC (Akaike Information
 81 Criterion) and BIC (Bayesian Information Criterion) rules are popular choices
 82 and both are presented in our analysis. The performance of the minimum AIC
 83 and BIC rules is investigated in simulation studies presented in the next Section.

84 4 Evaluation through simulation

85 We performed a limited simulation study. We varied K from 1 to 7. For every
 86 fixed K , we generated 100 samples from K multivariate normal distributions with
 87 exchangeable variance-covariance matrices and expectations linear in time; we then
 88 applied our method to estimate the parameters and chose the number of clusters
 89 using both the AIC and the BIC. This was repeated 100 times. The detailed
 90 forms of the distributions were chosen so as to mimic the results of the example
 91 described below. We give in Figures 1, 2 and 3, the results for $K = 1, 3$ and 6
 92 in the form of frequencies of number of retrieved class within the 100 repetitions.
 93 As it can be seen, both criteria perform reasonably well, with a tendency towards
 94 more conservative choices for the BIC and more liberal ones for the AIC. Though
 95 we have not carried out systematic explorations beyond those reported here, our
 96 experience suggests that the behaviour of the AIC and BIC is essentially the same
 97 in many situations.

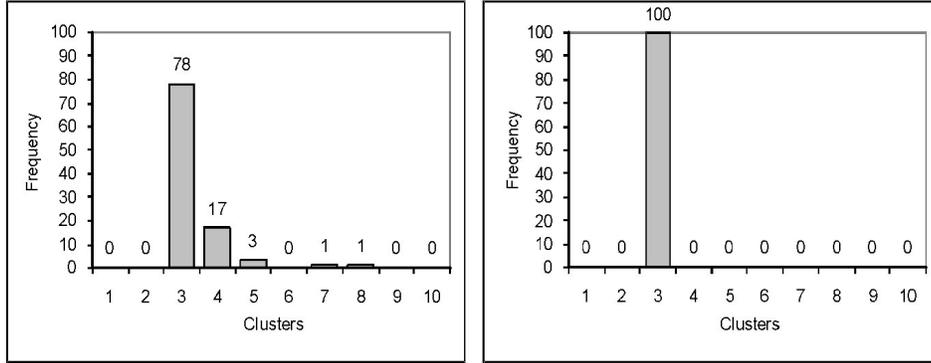


Figure 2: Histograms of number of clusters retrieved by the algorithm when 3 clusters are simulated using AIC(left panel) and BIC(right panel).

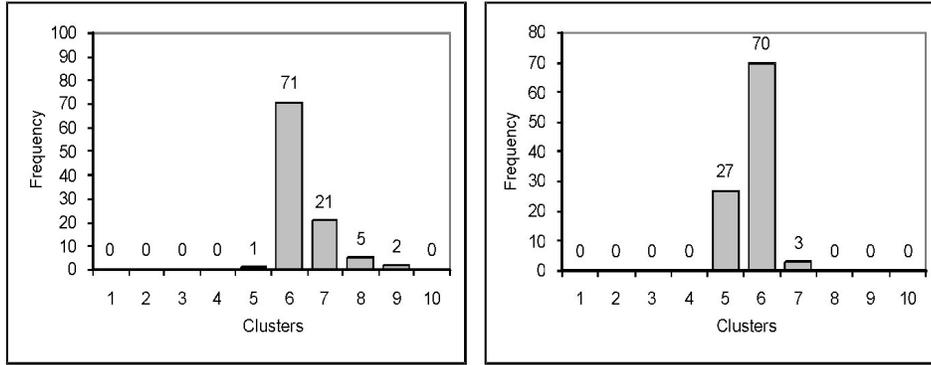


Figure 3: Histograms of number of clusters retrieved by the algorithm when 6 clusters are simulated using AIC(left panel) and BIC(right panel).

98 5 Disease Trajectory Data: Delirium

99 Delirium is a condition often encountered in hospitalized elderly populations. The
 100 Delirium Index (DI), is a validated measure of delirium severity developed at St.
 101 Mary's Hospital, (McCusker et al. 2004). DI scores range from 0 to 21 and
 102 higher scores indicated more severe delirium. We used data from 229 St. Mary's
 103 patients hospitalized between 1996 and 1999. Patients were evaluated with the
 104 DI at enrolment, and several times during the following 15 days. Measurement
 105 times were unequally spaced and differed across individuals. In order to account
 106 for correlation among repeated measurements, four different models were fit:

Independence

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2 Id)$$

Component	α	β	σ^2	Ψ	ϕ
1 – Steady Course	0.153	11.692, -0.015	3.912	3.428	0.438
2 – Fluctuating	0.224	9.956, 0.0078	23.171	2.059	0.387
3 – Worsening	0.204	5.742, 0.032	1.340	2.602	0.001
4 – Recovery	0.139	4.451, -0.343	2.912	0.550	0.615
5 – Fluctuating Recovery	0.279	8.569, -0.316	6.105	2.750	0.215

Table 1: Parameter Estimates for 5 cluster AR(1) and random intercept solution. Except for the slopes of components 1 and 2, all parameters are significantly different from zero at the 0.05 level according to likelihood ratio tests.

Random Intercept

$$y_i = \beta_0 + \beta_1 x_i + b_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2 Id)$$

Autoregressive

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2 R_i), \quad R_i = AR(1)(\phi)$$

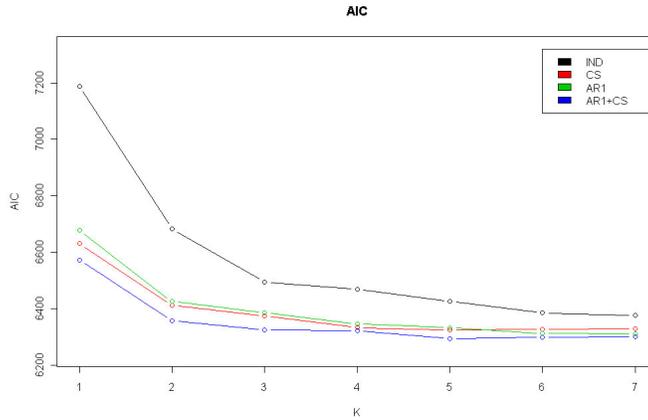
Random intercept and autoregressive

$$y_i = \beta_0 + \beta_1 x_i + b_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2 R_i), \quad R_i = AR(1)(\phi)$$

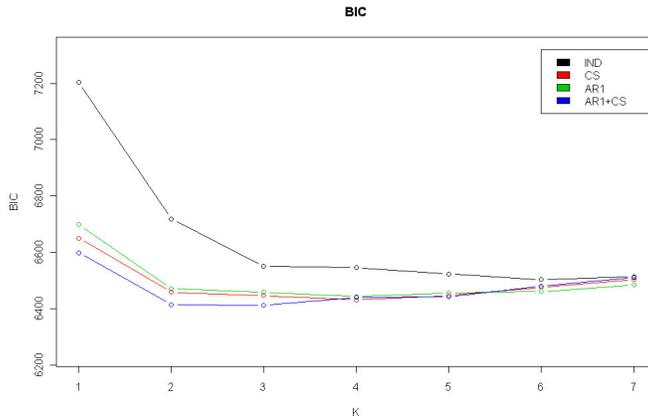
107 The DI curves have a great variety of shapes and by fitting a mixture model to
108 the longitudinal data, it is our goal to reduce these shapes to a few “typical” ones
109 which may be interpreted as distinct courses of the illness. Figure 4 shows AIC
110 and BIC values for the 28 different models fit. A 5-component model with both
111 AR(1) correlation and random intercept is selected for further investigation as it
112 provides good clinical interpretation and is the best model according to the AIC.
113 The mixture model has log-likelihood of -3118.480, AIC = 6284.961 and BIC =
114 6407.968. Parameter estimates appear in Table 1.

115 5.1 Interpretation

116 The **steady course** component represents a course that is quite stable: the slope
117 is negligible (and non-significant at the 0.05 level) and the variance parameter
118 reasonably small. The **fluctuating course** component is similar to the first, but
119 the variance is large, suggesting that patients fluctuate around a stable state. The
120 component named **worsening** has a positive slope, indicating a DI that increases
121 in time, hence a worsening of the delirium. The remaining two components have
122 negative slopes and are therefore named **recovery**: however, in one of them we
123 have a high variance and therefore we qualify the recovery as fluctuating. From a
124 general point of view, differences in the variance parameters may seem uninterest-
125 ing. However in delirium studies it is very important to identify components with
126 large fluctuations since fluctuating severity is considered a fundamental character-
127 istics of ‘true’ delirium.



(a) AIC



(b) BIC

Figure 4: AIC and BIC for different models across K

128 **5.2 Remarks**

129 A few comments are in order. Firstly, the four models considered in this analysis
 130 do not exhaust the possibilities of the ELMM. For example we could have fitted
 131 a model with random intercept and slope; unfortunately further exploration was
 132 limited by computational power. However, we have chosen the four models that
 133 are most currently used in biostatistical practice when analyzing longitudinal data:
 134 they reflect simple and intuitive hypotheses as to how correlation might arise.

135 Secondly, it should be noted that, as the AIC and BIC curves show, the selec-
 136 tion of the number of clusters depends on the model. While this may be seen as
 137 a limitation, it is by no means an uncommon occurrence: for example this depen-

138 dence is commonly observed when working with mixtures of multivariate normal
139 distributions and allowing hypotheses other than homoscedasticity, e.g. using the
140 mclust R package (Banfield and Raftery, 1993).

141 Thirdly, the determination of the number of classes remains a fairly subjective
142 exercise. Indeed neither the AIC nor the BIC provide absolutely objective criteria,
143 as is demonstrated by the numerous alternatives proposed in the literature. In
144 this work we have not attempted to develop new approaches, but have limited
145 ourselves to the most popular ones.

146 Fourthly, though we have shown by limited simulations that the BIC works bet-
147 ter than the AIC in our context, we have actually preferred to retain the 5-cluster
148 solution corresponding to the minimum AIC rather than the 2-cluster solution
149 which minimizes the BIC. This illustrates both the limits and the advantages of
150 using a certain degree of subjectivity. Indeed, the BIC of the 5-cluster solution
151 is not very different from the minimum BIC. On the other hand, a five cluster
152 solution was proposed in a previous work on delirium by Sylvestre et al. (2006),
153 who applied an exploratory approach combining principal component analysis with
154 k -means clustering. The interpretation of our 5-clusters is very similar to the in-
155 terpretation of the five clusters found by these authors, yet it has the advantage
156 of being model based.

157 Finally, though we have not studied robustness systematically, we have found
158 that the point estimates of the fixed effect coefficients are fairly stable regardless
159 of whether or not we include the random effect and/or the AR(1) term. This is
160 encouraging, but further explorations are desirable.

161 6 Time-course gene expression data

162 Microarray analysis is a valuable tool in molecular biology, as it permits to assess
163 the expression levels of a large number of genes simultaneously. In view of the
164 complexity of biological networks, it is useful to study gene expression not only
165 at a specific point in time, as in early microarray experiments, but also longitu-
166 dinally. Expression time profiles can indeed be very useful to find co-regulated
167 and functionally related groups of genes. We analysed a set of longitudinal gene
168 expression data already studied by Heard et al. (2002). The data consists of 2771
169 gene expression time profiles (each with 6 non-equally spaced observations) from
170 mosquitoes which have been infected with a bacterial agent. Visualization of the
171 raw data is not very informative (Figure 5).

172 Heard et al. (2002) proposed a Bayesian model-based hierarchical clustering
173 algorithm to cluster genes having similar expression profiles that led to a 17-cluster
174 solution. From this, interpretable graphs were obtained. Their solution assumed
175 data to be uncorrelated, so that in practice their model is a mixture of ordinary
176 regressions. In contrast, we used the wealth of submodels within the ELMM to find
177 a non-trivial correlation structure that fits the data. To model the trajectories, we
178 used a flexible family of basis functions called the truncated power spline basis, as
179 in Heard et al. (2002):

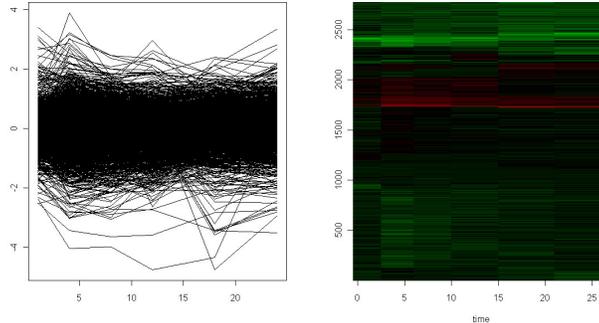


Figure 5: and heat map of the *Salmonella typhi* data presented in Heard et al. (2002)

$$g_1(t_{is}) = 1 \quad g_j(t_{is}) = (t_{is} - t_{i,j-1})_+, s = 1, \dots, n, j = 2, \dots, m_i$$

180 where $(\cdot)_+$ is the positive part function.

181 We fit three different correlation structures: an independence model, an au-
 182 toregressive model as well as a model with random intercept and autoregressive
 183 structure. We varied K from 1 to 26. According to the BIC, see Figure 6, the
 184 best model is a 14-cluster model with both random intercept and autoregressive
 185 structure. A heatmap of the data is presented in Figure 7. Comparing the fitted
 186 values(right) across clusters gives a visual measure of between cluster heteroge-
 187 nity. Comparing the fitted values(right) to the raw observations(left) gives a sense
 188 of within cluster homogeneity. The clustering is shown in more detail in Figure 8.

189 7 Discussion

190 We have presented a straightforward method for modeling heterogeneity in longi-
 191 tudinal data. We have proposed a mixture of regressions with components in the
 192 Extended Linear Mixed Model (ELMM) (Pinheiro and Bates, 2000). The ELMM
 193 consists of a random effect portion (LMM) extended by the addition of an error
 194 term with correlation matrix defined up to a small number of parameters to be
 195 estimated from data. We have limited ourselves to an autoregressive error term.
 196 Our approach to parameter estimation is based on the EM algorithm of Celeux
 197 et al. (2005) for mixtures of LMM, augmented by numerical methods which are
 198 essentially those used in the `lme` and `gls` R functions of Pinheiro and Bates (2000).
 199 The theoretical novelty of this approach is modest: it permits, on the one hand,
 200 to deal with correlated errors of a type that is important in applied research, and,
 201 on the other, suggests further extension to a catalog of possible error correlation
 202 structures such as those contained in the `lme` and `gls` R functions. Although other
 203 authors have considered mixtures of regressions for longitudinal data, no one has

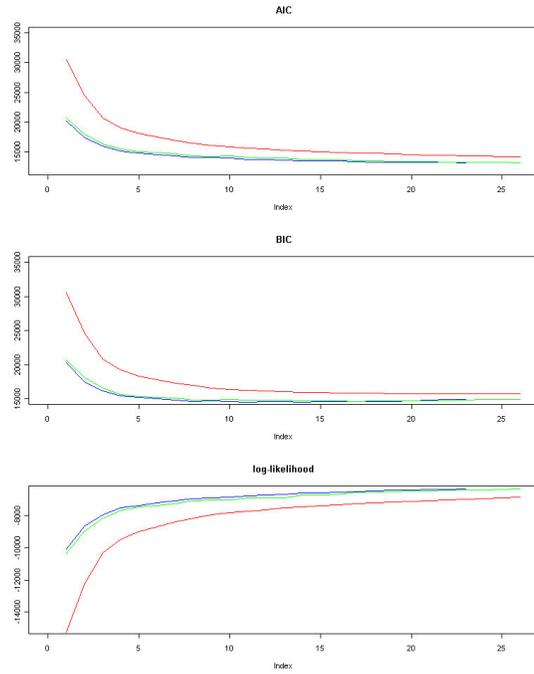


Figure 6: log-likelihood, AIC and BIC plots for different models(independence(red), autoregressive(green) and autoregressive with random intercept(blue)) fit for $K = 1, \dots, 26$

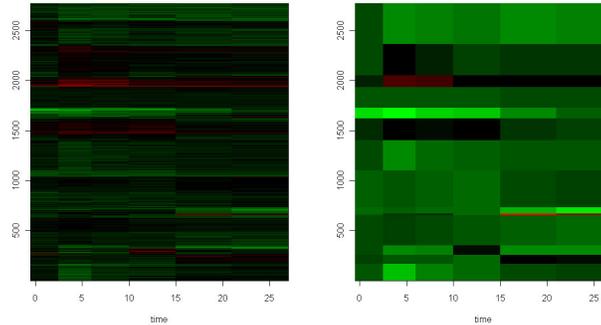


Figure 7: Clustered gene expression profiles form the *Salmonella typhi* data

204 yet achieved the generality that can be achieved with the ELMM. An appropri-
 205 ate modeling of the correlation structure of longitudinal data is important: not
 206 only does it provide useful insight into the dynamical process under study, but it
 207 also leads to fewer clusters, hence to a more economical model of the data. We

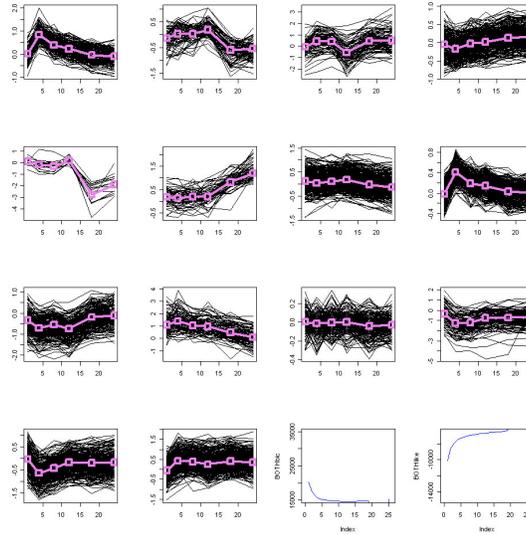


Figure 8: 14 cluster solution suggested by autoregressive AR(1) model with random intercept model. $BIC = 14528.19$

208 have also shown that mixtures of regressions offer an important tool in classifying
 209 course of diseases from clinical data and longitudinal gene expression data,
 210 providing easy-to-interpret analyses.

211 Further research will aim to speed up the EM algorithm, which will also allow
 212 us to study even richer correlation structures. Amelioration of computing efficiency
 213 will allow us to carry out more extensive simulations and to study the robustness
 214 of key features of our models, e.g. fixed effect parameter estimates and selection
 215 of the number of clusters. We plan also to revisit the Bayesian approach of Heard
 216 et al. (2002).

217 References

- 218 [1] Banfield J.D., Raftery A.E. (1993). “Model based Gaussian and non Gaussian
 219 clustering.” *Biometrics*, **49**: 803–821.
- 220 [2] Belin, T.R, Rubin, D.B. (1995). “The analysis of repeated-measures data on
 221 schizophrenic reaction times using mixture models.” *Statistics in Medicine*
 222 **14**: 747–768.
- 223 [3] Celeux, G., Lavergne, C., Martin, O. (2005). “Mixture of linear mixed models
 224 for clustering gene expression profiles from repeated microarray experiments”.
 225 *Statist. Model.* **5**: 243–267.

- 226 [4] De la Cruz-Mesía, R., Quintana, F.A., Marshall, G. (2008) “Model Based
227 Clustering for Longitudinal Data.” *Computational Statistics and Data Anal-*
228 *ysis* **52**: 1441–1457.
- 229 [5] Gaffney, S.J., Smyth, P. (2003). “Curve clustering with random effects regres-
230 sion mixtures.” In: Bishop, C.M., Frey, B.J. (Eds.), *Proceedings of the Ninth*
231 *International Workshop on Artificial Intelligence and Statistics*, KeyWest,
232 FL.
- 233 [6] Heard, Holmes, Stephens (2006). “A Quantitative Study of Gene Regulation
234 Involved in the Immune Response of Anopheline Mosquitoes: An Applica-
235 tion of Bayesian Hierarchical Clustering of Curves.” *Journal of the American*
236 *Statistical Association* **101**(473): 18–29.
- 237 [7] Jennrich and Schluchter (1986) *Unbalanced Repeated-Measures Models with*
238 *Structured Covariance Matrices*
- 239 [8] Luan Y., Li H. (2003). “Clustering of time course gene expression data using
240 a mixed-effect model with B-splines.” *Bioinformatics*, **19**, 474–482.
- 241 [9] McCusker J, Cole M, Dendukuri N, Belzile E. (2004). “The Delirium Index,
242 a Measure of the Severity of Delirium: New Findings on Reliability, Validity,
243 and Responsiveness”. *Journal of the American Geriatric Society*, **52**:1744-
244 1749.
- 245 [10] Ng, S. K., McLachlan, G. J., Wang, K., Ben-Tovim Jones, L., Ng, S. W. (2006)
246 “A Mixture model with random-effects components for clustering correlated
247 gene-expression profiles.” *Bioinformatics* **22**: 1745–1752.
- 248 [11] Pinheiro J.C., Bates D. (2000). *Mixed-Effects Models in S and S-PLUS*.
249 Springer, New York.
- 250 [12] Tango T. (1998). “A mixture model to classify individual profiles of repeated
251 measurements.” In: *Data Science, Classification and Related Topics* (eds by
252 C. Hayashi, et al.), Springer-Verlag, Tokyo, 247–254.
- 253 [13] Sylvestre, M. P., McCusker, L., Cole, M., Regeasse, A. Belzile, E., and Abra-
254 hamowicz, M. (2006). “Classification of patterns of delirium severity scores
255 over time in an elderly population.” *International Psychogeriatrics*, **18**: 667–
256 680.

257 A Details on the M-step of the EM algorithm

We have for the $q+1$ th step:

$$\alpha_k^{[q+1]} = \sum_{i=1}^n \frac{\tau_{i(k)}^{[q]}}{n}$$

$$\beta_k^{[q+1]} = \left[\left(\sum_{i=1}^n \tau_{i(k)}^{[q]} \left((y_i - Z_i b_{i(k)}^{[q]}) \right)^T \left((\Lambda_{ik}^{[q]})^{-1} (-X_i) \right) \right) \left(\sum_{i=1}^n \tau_{i(k)}^{[q]} (X_i)^T \left((\Lambda_{ik}^{[q]})^{-1} (-X_i) \right) \right)^{-1} \right]^T$$

$$(\sigma_k^2)^{[q+1]} = \frac{1}{\sum_{i=1}^n \tau_{i(k)}^{[q]} m_i} \sum_{i=1}^n \tau_{i(k)}^{[q]} \left[E(e_{ik}^T \Lambda_{ik}^{-1} e_{ik} \mid y_i, \theta^{[q-1]}) \right]$$

where

$$E(e_{ik}^T \Lambda_{ik}^{-1} e_{ik} \mid y_i, \theta^{[q-1]}) =$$

$$\sigma_k^{4[q-1]} (y_i - X_i \beta_k^{[q-1]})^T \Sigma_{ik}^{-1[q-1]} \Lambda_{ik}^{[q-1]} \Sigma_{ik}^{-1[q-1]} (y_i - X_i \beta_k^{[q-1]}) +$$

$$m_i \sigma_k^{2[q-1]} - \sigma_k^{4[q-1]} \text{tr}(\Sigma_{ik}^{-1[q-1]} \Lambda_{ik}^{[q-1]})$$

and

$$\Psi_k^{[q+1]} = \frac{1}{\sum_{i=1}^n \tau_{i(k)}^{[q]} m_i} \sum_{i=1}^n \tau_{i(k)}^{[q]} E(b_{i(k)} b_{i(k)}^T \mid y_i, \theta^{[q-1]})$$

where

$$E(b_{i(k)} b_{i(k)}^T \mid y_i, \theta^{[q-1]}) = E(b_{i(k)} \mid y_i, \theta^{[q-1]}) E(b_{i(k)} \mid y_i, \theta^{[q-1]})^T +$$

$$m_i \Psi_k^{[q-1]} - \Psi_k^{[q-1]} Z_i^T \Sigma_{ik}^{-1[q-1]} Z_i \Psi_k^{[q-1]}$$

and

$$E(b_{i(k)} \mid y_i, \theta^{[q-1]}) = \Psi_k^{[q-1]} Z_i^T \Sigma_{ik}^{-1[q-1]} (y_i - X_i \beta_k^{[q-1]})$$

258 Finally, consider the positive-definite matrices Λ_{ik} (there are nK of these). There
 259 are different ways that such matrices may be parametrized, depending on as-
 260 sumptions regarding the intra-individual covariance structure (Pineiro and Bates,
 261 2000). Let ϕ_k denote the set of parameters used in the parametrization of $\{\Lambda_{ik}\}_{i=1, \dots, n}$.
 262 To estimate these parameters at iteration $[q]$, we use numerical maximization meth-
 263 ods (e.g the R function `nlminb()`).

264 B Details on computing the standard error of $\hat{\theta}$

265 Jennrich and Schluchter(1986) provide equations for the required score vector
 266 statistics and Hessian matrix in the homogeneous model. The required first and
 267 second derivatives for the mixture model are presented bellow. Derivatives with
 268 respect to the parameters that define Λ , (ϕ) , must be calculated by numerical
 269 methods.

$$\frac{\partial l(\theta \mid Y)}{\partial \alpha_k} = \sum_{i=1}^n \frac{\delta_{i(k)}}{\alpha_k} - \frac{\delta_{i(K)}}{\alpha_K}, \quad k=1, \dots, K-1 \quad (9)$$

$$\frac{\partial l(\theta \mid Y)}{\partial \beta_k} = \sum_{i=1}^n \delta_{i(k)} (X_i^T (\Sigma_{ik}^{-1})) (y_i - X_i \beta_k), \quad k=1, \dots, K \quad (10)$$

$$\frac{\partial l(\theta | Y)}{\partial \sigma_k^2} = \sum_{i=1}^n \frac{\delta_{i(k)}}{2} \text{tr}(\Sigma_{ik}^{-1}(y_i - X_i \beta_k)(y_i - X_i \beta_k) \Sigma_{ik}^{-1} \Lambda^{-1}), \quad k=1, \dots, K \quad (11)$$

$$\frac{\partial l(\theta | Y)}{\partial \Psi_k} = \sum_{i=1}^n \frac{\delta_{i(k)}}{2} \text{tr}(\Sigma_{ik}^{-1}(y_i - X_i \beta_k)(y_i - X_i \beta_k) \Sigma_{ik}^{-1} Z_i Z_i^T), \quad k=1, \dots, K \quad (12)$$

$$\frac{\partial^2 l(\theta | Y)}{\partial \alpha_k^2} = \sum_{i=1}^n \frac{-\delta_{i(k)}}{\alpha_k^2} - \frac{\delta_{i(K)}}{\alpha_K^2}, \quad k=1, \dots, K-1 \quad (13)$$

$$\frac{\partial^2 l(\theta | Y)}{\partial \beta_k^2} = \sum_{i=1}^n -\delta_{i(k)} (X_i^T \Sigma_{ik}^{-1} X_i), \quad k=1, \dots, K \quad (14)$$

$$\frac{\partial^2 l(\theta | Y)}{\partial (\sigma_k^2)^2} = \sum_{i=1}^n \frac{-\delta_{i(k)}}{2} \text{tr}(\Sigma_{ik}^{-1} \Lambda_{ik} \Sigma_{ik}^{-1} (2(y_i - X_i \beta_k)(y_i - X_i \beta_k)^T - \Sigma_{ik}) \Sigma_{ik}^{-1} \Lambda_{ik}), \quad k=1, \dots, K \quad (15)$$

$$\frac{\partial^2 l(\theta | Y)}{\partial \sigma_k^2 \partial \beta_{jk}} = \sum_{i=1}^n -\delta_{i(k)} X_{jk} (\Sigma_{ik}^{-1} \Lambda_{ik} \Sigma_{ik}^{-1} (y_i - X_i \beta_k)), \quad k=1, \dots, K \quad j=1, \dots, p \quad (16)$$

$$\frac{\partial^2 l(\theta | Y)}{\partial \Psi_k^2} = \sum_{i=1}^n \frac{-\delta_{i(k)}}{2} \text{tr}(\Sigma_{ik}^{-1} Z_i Z_i^T \Sigma_{ik}^{-1} (2(y_i - X_i \beta_k)(y_i - X_i \beta_k)^T - \Sigma_{ik}) \Sigma_{ik}^{-1} Z_i Z_i^T), \quad k=1, \dots, K \quad (17)$$

$$\frac{\partial^2 l(\theta | Y)}{\partial \Psi_k \partial \beta_{jk}} = \sum_{i=1}^n -\delta_{i(k)} X_{jk} (\Sigma_{ik}^{-1} Z_i Z_i^T \Sigma_{ik}^{-1} (y_i - X_i \beta_k)), \quad k=1, \dots, K \quad j=1, \dots, p \quad (18)$$