

"This is the pre print version of the following article: "Item selection by latent class-based methods: an application to nursing home evaluation", which has been published in final form at DOI 10.1007/s11634-016-0232-3. This article may be used for non-commercial purposes and may not be enhanced, enriched or otherwise transformed into a derivative work, without express permission.

URL <https://link.springer.com/content/pdf/10.1007/s11634-016-0232-3.pdf>

Advances in Data Analysis and Classification

Item selection by Latent Class-based methods: an application to nursing home evaluation

--Manuscript Draft--

Manuscript Number:	ADAC-D-14-00063
Full Title:	Item selection by Latent Class-based methods: an application to nursing home evaluation
Article Type:	S.I. : Advances in Latent Variables: Methods, Models and Applications
Corresponding Author:	Silvia Pandolfi, Ph.D University of Perugia Perugia, ITALY
Corresponding Author Secondary Information:	
Corresponding Author's Institution:	University of Perugia
Corresponding Author's Secondary Institution:	
First Author:	Silvia Pandolfi, Ph.D
First Author Secondary Information:	
Order of Authors:	Silvia Pandolfi, Ph.D Francesco Bartolucci, Ph.D Giorgio Eduardo Montanari, Ph.D
Order of Authors Secondary Information:	
Abstract:	<p>The evaluation of nursing homes is usually based on the administration of questionnaires made of a large number of polytomous items. In such a context, the Latent Class (LC) model represents a useful tool for clustering subjects in homogenous groups corresponding to different degrees of impairment of the health conditions. It is known that the performance of model-based clustering and the accuracy of the choice of the number of latent classes may be affected by the presence of irrelevant or noise variables. In this paper, we show the application of an item selection algorithm to real data collected within a project, named ULISSE, on the quality-of-life of elderly patients hosted in italian nursing homes. This algorithm, which is closely related to that proposed by Dean and Raftery in 2010, is aimed at finding the subset of items which provides the best clustering according to the Bayesian Information Criterion. At the same time, it allows us to select the optimal number of latent classes. Given the complexity of the ULISSE study, we perform a validation of the results by means of a sensitivity analysis to different specifications of the initial subset of items and of a resampling procedure.</p>
Suggested Reviewers:	

Noname manuscript No.
(will be inserted by the editor)

Item selection by Latent Class-based methods: an application to nursing home evaluation

Francesco Bartolucci · Giorgio E.
Montanari · Silvia Pandolfi

Received: date / Accepted: date

F. Bartolucci
Department of Economics, University of Perugia, Perugia, ITALY
E-mail: bart@stat.unipg.it

S. Pandolfi
Department of Economics, University of Perugia, Perugia, ITALY
Tel.: +39-075-5855236
Fax: +39-075-5855950
E-mail: pandolfi@stat.unipg.it

G.E. Montanari
Department of Political Sciences, University of Perugia, Perugia, ITALY
E-mail: giorgio.montanari@unipg.it

Noname manuscript No. (will be inserted by the editor)
--

Item selection by Latent Class-based methods: an application to nursing home evaluation

Received: date / Accepted: date

Abstract The evaluation of nursing homes is usually based on the administration of questionnaires made of a large number of polytomous items. In such a context, the Latent Class (LC) model represents a useful tool for clustering subjects in homogenous groups corresponding to different degrees of impairment of the health conditions. It is known that the performance of model-based clustering and the accuracy of the choice of the number of latent classes may be affected by the presence of irrelevant or noise variables. In this paper, we show the application of an item selection algorithm to real data collected within a project, named ULISSE, on the quality-of-life of elderly patients hosted in Italian nursing homes. This algorithm, which is closely related to that proposed by Dean and Raftery in 2010, is aimed at finding the subset of items which provides the best clustering according to the Bayesian Information Criterion. At the same time, it allows us to select the optimal number of latent classes. Given the complexity of the ULISSE study, we perform a validation of the results by means of a sensitivity analysis to different specifications of the initial subset of items and of a resampling procedure.

Keywords Bayesian Information Criterion · Expectation-Maximization algorithm · Polytomous items · Quality-of-life · ULISSE project

1 Introduction

The evaluation of appropriateness of long-term care facilities is assuming a role of increasing relevance due to the rapid growth of demand for long-term care services for elderly people. The main cause is represented by the rapid aging of the population and also by the changes in the family structure and in the socio-economic context. Furthermore, the debate on population aging focuses on the effects that such a phenomenon has on the welfare and on the

1 health care system in various countries (Galasso and Profeta, 2007; Breyer
2 et al, 2010). In this regard, health care quality measurement and performance
3 evaluation of nursing homes represent a challenging issue to assure the quality
4 of services and to allocate resources efficiently.
5

6 Issues related to population aging are particularly relevant in Italy, which
7 is one of the European countries with the highest proportion of elderly people,
8 where this proportion is expected to increase over the next decades (Kohler
9 et al, 2002). In this country, the ULISSE project (Lattanzio et al, 2010) has
10 been carried out to obtain relevant data for health care planning. The purpose
11 of the project is to document the change in elderly patients' health status
12 and the ability of the health care system to satisfy their needs. The dataset
13 obtained from this project was collected by the administration of a question-
14 naire to patients hosted in a sample of Italian nursing homes. The question-
15 naire is made of a large number of polytomous items about different aspects
16 of the quality-of-life and health status of these patients. In such a context, a
17 model-based clustering procedure may be applied to evaluate the performance
18 of nursing homes. In particular, patients may be clustered in homogeneous
19 groups according to their health conditions in order to describe the case-mix
20 of the nursing homes. The resulting clustering may have important applica-
21 tions in the context of the nursing home evaluation, and therefore affecting the
22 system of financial support, when the clusters correspond to different degrees
23 of impairment of the patients' health conditions. However, the performance
24 of model-based clustering procedure may be degraded by the presence of ir-
25 relevant items. Moreover, the administration of a questionnaire made of a
26 large number of items may be lengthy and expensive. Due to tiring effects,
27 using several items may also induce the respondent to provide inaccurate re-
28 sponses. This is particularly relevant when the questionnaire is periodically
29 administered. Therefore, methods which allow us to select the smallest subset
30 of items useful for clustering are of interest. These methods may lead to a
31 reduction of the costs of the data collection process, and a better quality of
32 the collected data. Moreover, reducing the dimension of the dataset implies
33 that it may be more easily analyzed by complex statistical models.
34

35 Motivated by availability of the ULISSE dataset, in this paper we adopt
36 the item selection algorithm proposed by Dean and Raftery (2010), which may
37 be applied when a large number of items is included in a questionnaire. The
38 adopted algorithm is based on the latent class (LC) model (Lazarsfeld, 1950;
39 Lazarsfeld and Henry, 1968; Goodman, 1974) which represents an important
40 tool of analysis of data collected by questionnaires made of polytomous items.
41 As it is well known, the model relies on a discrete latent variable, which defines
42 a certain number of latent classes, and it assumes the independence of the
43 responses to the items given this variable. Therefore, the use of the LC model
44 is justified when the items measure one or more latent traits, such as the
45 quality-of-life or the tendency toward a certain behavior. In geriatrics, the LC
46 model is used to measure mobility disability (Bandeem-Roche et al, 1997), to
47 study behavioral syndromes in Alzheimer' patients (Moran et al, 2004), and
48 to test the validity of certain physical frailty measures (Bandeem-Roche et al,
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 2006). Moreover, Lafortune et al (2009) uses the LC analysis to model the
2 heterogeneity in elderly individuals' health status.
3

4 It is important to recall that the LC model produces a *model-based cluster-*
5 *ing* (Fraley and Raftery, 2002), with clusters corresponding to the latent
6 classes. Once the model is fitted, a subject is assigned to the latent class cor-
7 responding to the highest posterior probability, that is, the conditional prob-
8 ability of the latent class given the observed data. It is also worth noting that
9 latent variable models for ordinal variables are also present in the literature,
10 such as the graded response model (Samejima, 1969, 1996). In this model, the
11 response probability is expressed as a function of one or more latent variables
12 through a specific link function (e.g., the cumulative logit link); see Bacci
13 et al (2014) and the reference therein. Given the ordinal nature of the items
14 composing the questionnaire considered in our application, this model could
15 make sense in the present context. However, we prefer to avoid such a para-
16 metrization and to rely on a standard LC model. The main advantages of this
17 choice are the simplicity of the resulting approach and the reduced number of
18 parametric assumptions.
19

20 The item selection algorithm we adopt aims at finding the set of items
21 which provides the best value of the Bayesian Information Criterion (BIC)
22 index (Schwarz, 1978; Kass and Raftery, 1995). As motivated by Dean and
23 Raftery (2010), this leads to selecting the subset of items which are indeed
24 useful for clustering. It is worth noting that, at the same time, this method
25 allows us to choose the number of latent classes. As usual in the context
26 of LC models, the parameters estimation is performed by the Expectation-
27 Maximization algorithm (EM; Goodman, 1974; Dempster et al, 1977).
28

29 More in detail, the implementation of the algorithm is based on a stepwise
30 scheme that, starting from an initial set of items, at each iteration performs
31 both inclusion and exclusion steps till a certain optimal criterion is satisfied.
32 We also extend the version proposed in Dean and Raftery (2010) by includ-
33 ing an additional step aimed at initializing, with a large number of starting
34 values, the estimation algorithm, so as to face the problem of local maxima
35 of the model log-likelihood, without being computationally expensive. This
36 problem is particularly evident in applications involving a very large number
37 of statistical units and response variables, as in the ULISSE project. In such an
38 application, we also assess the performance of the item selection algorithm by
39 implementing a sensitivity analysis of the final results with respect to different
40 specifications of the initial subset of items and by validating the solution on
41 the basis of resampling procedures.
42

43 The remainder of this paper is structured as follows. The ULISSE dataset
44 is described in the following section. In Section 3 we briefly illustrate the LC
45 model on which the item selection algorithm is based together with maximum
46 likelihood estimation of the model, via the EM algorithm. In Section 4 we illus-
47 trate the item selection algorithm based on the implementation proposed by
48 Dean and Raftery (2010). In Section 5 we present the results of this approach
49 applied to the ULISSE dataset, whereas Section 6 provides a final discussion.
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

The approach proposed in this paper has been implemented in a series of R functions which rely on a Fortran code to make the execution faster. These functions are available to the reader upon request.

2 The ULISSE dataset

The ULISSE project (“Un Link Informatico sui Servizi Sanitari Esistenti per l’Anziano” - “A Computerized Network on Health Care Services for Older People”) is aimed at studying the health status of elderly patients who currently receive health care assistance in Italy (Lattanzio et al, 2010). The main purpose of the study is to improve the knowledge of the characteristics and the quality of health care services provided to elderly in Italy. The project was carried out by a Research Group established by the Italian Ministry of Health. The study considers three different levels of health care assistance: that provided by acute care facilities, that provided by nursing homes, and that provided at home. Overall, 23 acute geriatric or internal medicine hospital units, 31 nursing homes, and 11 home care services have been involved in the project. In the analysis here presented, we consider only data collected in the nursing homes.

The ULISSE project is based on a longitudinal survey; in the nursing homes recruited by the project, all residents, or a maximum of 100 randomly selected residents for bigger nursing homes, were evaluated at admission and then re-evaluated at 6 and 12 months after the admission. Only long stay residents (i.e., permanently admitted to the nursing home) aged at least 65 years were included in the study. For the analysis here presented, we consider only the first wave of interviews, which covers 1739 patients.

The detailed patients information were collected using the classification system VAOR (“Valutazione dell’Anziano Ospite di Residenza”) that represents the Italian version of the interRAI Minimum Data Set (MDS) for nursing home care (Morris et al, 1991; Hawes et al, 1997). The questionnaire, which was filled up by the nursing assistants, is made of 75 items covering several aspects of the health status of the elderly patients. The items are polytomous, with categories generally ordered according to increasing difficulty levels in accomplishing a certain task or severeness of a specific aspect of the health conditions. The complete list of items, with the corresponding number of response categories, is reported in Appendix.

The 75 items are grouped into eight different sections (indexed by d) of the questionnaire, concerning:

1. Cognitive Conditions (CC);
2. Auditory and View Fields (AVF);
3. Humor and Behavioral Disorders (HBD);
4. Activities of Daily Living (ADL);
5. Incontinence (I);
6. Nutritional Field (NF);
7. Dental Disorders (DD);

8. Skin Conditions (SC).

Table 1 shows the average of the percentage of missing values computed with respect to the items composing each section of the questionnaire. We observe that the percentage of missing responses is considerably different between sections, going from 0.54% for section AVF to 9% for section ADL.

Table 1 Average percentage of missing values for each section of the questionnaire.

d	section	% missing
1	CC	0.78
2	AVF	0.54
3	HBD	1.96
4	ADL	9.00
5	I	1.61
6	NF	5.58
7	DD	2.18
8	SC	6.48

3 The Latent Class model

Let J denote the number of items in the questionnaire of interest and, for a sample of n respondents, let Y_{ij} denote the response variable for subject i and item j , with $i = 1, \dots, n$ and $j = 1, \dots, J$. Let l_j denote the number of response categories of item j , labeled from 0 to l_j-1 . Some of the responses may not be observed for different reasons. We rely on the standard assumption of *missing at random* (MAR; Rubin, 1976; Little and Rubin, 2002) to deal with these missing responses, as described in the following. Finally, let $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{iJ})$ be the vector of all response variables for subject i .

In order to explain the dependence structure between the response variables, the LC model assumes the existence of a discrete latent variable U_i which has, *a priori*, the same distribution for every subject i . This distribution is based on k support points, labeled from 1 to k . Each support point corresponds to a latent class in the population and has a specific weight (or *a priori* probability); these weights are denoted by π_1, \dots, π_k . Moreover, the conditional probability that individual i in class u provides response y to item j is

$$\lambda_{j|u}(y) = p(Y_{ij} = y | U_i = u), \quad j = 1, \dots, J, u = 1, \dots, k, y = 0, \dots, l_j - 1.$$

Overall, the number of non-redundant parameters of the model is

$$g_k = (k - 1) + k \sum_j (l_j - 1).$$

The basic assumption underlying the LC model is that of *local independence*. This assumption is formulated by requiring that, for $i = 1, \dots, n$, the variables

Y_{i1}, \dots, Y_{iJ} are conditionally independent given the latent variable U_i . This assumption implies that

$$p(\mathbf{y}_i|u) = P(\mathbf{Y}_i = \mathbf{y}_i|U_i = u) = \prod_{j=1}^J \lambda_{j|u}(y_{ij}),$$

with $\mathbf{y}_i = (y_{i1}, \dots, y_{iJ})$. Moreover, the *manifest probability* of the response pattern \mathbf{y}_i for subject i is denoted by

$$p(\mathbf{y}_i) = \sum_{u=1}^k p(\mathbf{y}_i|u)\pi_u.$$

Another quantity of interest is the posterior probability that a subject with observed response configuration \mathbf{y}_i belongs to latent class u . Using standard rules, this probability is equal to

$$p(u|\mathbf{y}_i) = \frac{p(\mathbf{y}_i|u)\pi_u}{p(\mathbf{y}_i)}, \quad u = 1, \dots, k. \quad (1)$$

These posterior probabilities are used to allocate subjects in the different latent classes, as will be clarified in the sequel.

As defined above, in the presence of missing responses to the items, we rely on the MAR hypothesis. This assumption states that the probability of the observed missingness pattern, given the observed and the unobserved data, does not depend on the unobserved data (see, among others, Lu and Copas, 2004). Therefore, provided that the model for the missing data mechanism is separated from that of the LC model, the missing responses are *ignorable* for likelihood-based inference. The corresponding LC model may be formulated by introducing the missing data indicator M_{ij} which is equal to 1 when subject i does not respond to item j and to 0 otherwise. Then, for a given subject i , we have $\mathbf{M}_i = (M_{i1}, \dots, M_{iJ})$. The corresponding response pattern is given by $(\mathbf{m}_i, \mathbf{y}_{i,obs})$ in which \mathbf{m}_i is a configuration of \mathbf{M}_i and the subvector $\mathbf{y}_{i,obs}$ contains the observed components of \mathbf{Y}_i .

The MAR assumption implies that the parameters may be estimated on the basis of the log-likelihood of the vectors of observed responses $\mathbf{y}_{i,obs}$ only, without worrying about the model for missingness (see, among others, Harel and Schafer, 2009). In particular, the conditional probability of the observed response configuration, given the latent class, is simply

$$p(\mathbf{y}_{i,obs}|u) = \int_{\mathbf{y}_{i,mis}} p(\mathbf{y}_i|u) d\mathbf{y}_{i,mis} = \prod_{j:m_{ij}=0} \lambda_{j|u}(y_{ij}),$$

and then

$$p(\mathbf{y}_{i,obs}) = \sum_u p(\mathbf{y}_{i,obs}|u)\pi_u. \quad (2)$$

Note that under this assumption the number of parameters to be estimated remains the same as in the standard LC model.

3.1 Maximum likelihood estimation

Given the assumption of independence between the sample units, the general formulation of the log-likelihood function of the proposed model is

$$\ell(\boldsymbol{\theta}) = \sum_i \log p(\mathbf{y}_i).$$

In the presence of missing values considered as MAR, the above expression becomes

$$\ell(\boldsymbol{\theta}) = \sum_i \log p(\mathbf{y}_{i,obs}),$$

where the manifest probability $p(\mathbf{y}_{i,obs})$ is computed as in (2). In the above expressions, $\boldsymbol{\theta}$ is a short-hand notation for all model parameters. In order to estimate these parameters, we maximize $\ell(\boldsymbol{\theta})$ by the EM algorithm (Dempster et al, 1977).

3.2 Expectation-Maximization algorithm

The EM algorithm is based on the *complete-data likelihood* that we could compute if we knew the value of the latent variable U_i for every unit i in the sample. This is equivalent to the knowledge of the latent class to which every subject belongs. As usual, we represent such an information by the set of dummy variables z_{iu} , $i = 1, \dots, n$, $u = 1, \dots, k$, where z_{iu} is equal to 1 if respondent i belongs to latent class u and to 0 otherwise. Then, under the MAR assumption, we can write the complete-data log-likelihood as

$$\begin{aligned} \ell^*(\boldsymbol{\theta}) &= \sum_i \sum_u z_{iu} \log [p(\mathbf{y}_{i,obs}|u)\pi_u] \\ &= \sum_i \sum_u z_{iu} \sum_{j:m_{ij}=0} \log \lambda_{j|u}(y_{ij}) + \sum_u z_{+u} \log \pi_u, \end{aligned} \quad (3)$$

where $z_{+u} = \sum_i z_{iu}$ is the number of subjects in latent class u . We have an explicit solution for the maximum of $\ell^*(\boldsymbol{\theta})$ with respect to the model parameters, which, for $u = 1, \dots, k$ is

$$\tilde{\pi}_u = \frac{z_{+u}}{n}, \quad (4)$$

$$\tilde{\lambda}_{j|u}(y) = \frac{\sum_i I(y_{ij} = y)(1 - m_{ij})z_{iu}}{\sum_i z_{iu}(1 - m_{ij})}, \quad j = 1, \dots, J, y = 0, \dots, l_j - 1, \quad (5)$$

where $I(\cdot)$ is the indicator function equal to 1 if its argument is true and to 0 otherwise.

In order to maximize the model log-likelihood, the EM algorithm alternates the following two steps until convergence, starting from an initial guess of the model parameters in $\boldsymbol{\theta}$:

- 1 – **E-step**: compute the conditional expected value of the complete-data log-likelihood $\ell^*(\boldsymbol{\theta})$ given the observed data and the current value of the parameters;
- 2
- 3
- 4 – **M-step**: update the model parameters by maximizing the expected value
- 5 obtained at the E-step.
- 6

7 Both steps are simple to implement. In practice, the E-step consists of
8 obtaining the posterior expected value of every dummy variable z_{iu} , that is,

$$9 \hat{z}_{iu} = p(u|\mathbf{y}_i), \quad i = 1, \dots, n, u = 1, \dots, k,$$

10 which may be computed according to (1). At the M-step we maximize the
11 expected value of the complete-data log-likelihood, which is obtained by substituting
12 every dummy variable z_{iu} in (3) with \hat{z}_{iu} , and in this way we update
13 the parameter vector $\boldsymbol{\theta}$. For this maximization we use formulae (4) and (5),
14 with \hat{z}_{iu} instead of z_{iu} .

15 As mentioned in Section 3, the posterior probabilities \hat{z}_{iu} can be used for
16 clustering, that is, to allocate subjects in the k latent classes. In particular,
17 on the basis of the output of the EM algorithm, we assign subject i to latent
18 class u when $\hat{z}_{iu} = \hat{z}_i^*$, where \hat{z}_i^* is the maximum of $\hat{z}_{i1}, \dots, \hat{z}_{ik}$. For this reason,
19 Magidson and Vermunt (2001) and Vermunt and Magidson (2002) refer to this
20 kind of model as an *LC cluster model*.

21 3.2.1 Initialization of the algorithm

22 A typical problem of the latent variable and finite mixture models is the mul-
23 timodality of the likelihood. Obviously, in the presence of multiple local maxima,
24 the EM algorithm converges to one of them, which is not ensured to be
25 the global maximum. In this case, it is advisable to use a random initialization
26 strategy which consists of repeatedly initializing the algorithm from a
27 large number of randomly chosen starting values for the parameters. When
28 more starting values are used, the final estimate is the one corresponding to
29 the largest likelihood value that has been found at convergence of the EM
30 algorithm (Biernacki et al, 2003; Karlis and Xekalaki, 2003). This solution is
31 not guaranteed to correspond to the global maximum; however, it is rather
32 obvious that the chance of reaching the global maximum increases with the
33 number of the starting values that are tried. The problem of likelihood mul-
34 timodality is particularly severe for the LC model used in our analysis, due
35 to the large number of items and to the fact that these items generally have
36 more than two response categories.

37 In our approach we adopt a random initialization which is based on drawing
38 each latent class weight π_u from a uniform distribution between 0 and 1
39 and then normalizing these random draws so that they sum to 1. In a similar
40 way we randomly choose the conditional response probabilities $\lambda_{j|u}(y)$,
41 $y = 0, \dots, l_j - 1$. Moreover, we use a number of random initializations which
42 increases with the number of latent classes, because the latter affects the number
43 of parameters and then the expected number of local maxima.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

4 Item selection procedure

In this section we describe the procedure for item selection based on the method proposed by Dean and Raftery (2010). This method follows a stepwise scheme that, starting from an initial set of items, performs a series of inclusion and exclusion steps until a suitable stopping rule is satisfied, also allowing the selection of the proper number of latent classes.

4.1 Inclusion-exclusion algorithm with random check

The inclusion-exclusion algorithm for item selection is based on assessing the importance of a certain item by comparing two LC models. In the first model, the item is assumed to provide additional information about clustering allocation, beyond that contained in the already selected items; in the second model, this item does not provide additional information useful for clustering and then it is independent of the latent variable defining the latent classes. The two models are compared via BIC index (Schwarz, 1978), which is seen as an approximation of the Bayes Factor (Kass and Raftery, 1995).

In more detail, let $\mathcal{A}^{(0)}$ denote the initial set of items and let $k^{(0)}$ the corresponding number of latent classes. At the h th iteration, the item selection algorithm performs an inclusion and an exclusion step, so as to obtain $\mathcal{A}^{(h)}$ and $k^{(h)}$, as follows:

- *Inclusion step:* each item j in $\bar{\mathcal{A}}^{(h-1)}$, the complement of $\mathcal{A}^{(h-1)}$ with respect to the full set of items, is singly proposed for inclusion in $\mathcal{A}^{(h)}$. The item to be included is selected on the basis of the difference between BIC of the LC model for the items in $\mathcal{A}^{(h-1)} \cup j$ (optimized over the number of classes from 2 to k_{\max} , with k_{\max} *a priori* fixed) and BIC of the LC model in which item j is assumed to be independent of the latent class. This difference is as follow

$$BIC_{diff}(\mathcal{A}^{(h-1)}, j) = \min_{2 \leq k \leq k_{\max}} BIC_k(\mathcal{A}^{(h-1)} \cup j) - [BIC_1(j) + BIC_{k^{(h-1)}}(\mathcal{A}^{(h-1)})],$$

where

$$BIC_k(\mathcal{A}) = -2 \hat{\ell}_k(\mathcal{A}) + g_k(\mathcal{A}) \log(n), \quad (6)$$

with $\hat{\ell}_k(\mathcal{A})$ being the maximum of the log-likelihood of the LC model applied to the data referred to the items in \mathcal{A} , and $g_k(\mathcal{A})$ being the corresponding number of free parameters. Obviously, BIC_1 is the BIC index of the single class LC model, which corresponds to the model of independence. Note that each step of the algorithm also selects the number of classes k , since BIC_k is minimized over k from 2 to k_{\max} . The item included is the one with the smallest negative $BIC_{diff}(\mathcal{A}^{(h-1)}, j)$, and $\mathcal{A}^{(h)} = \mathcal{A}^{(h-1)} \cup j$ (with $k^{(h)}$ updated). If no item yields a negative BIC_{diff} , then we set $\mathcal{A}^{(h)} = \mathcal{A}^{(h-1)}$.

1 – *Exclusion step*: each item j in $\mathcal{A}^{(h)}$ is singly proposed for exclusion. The
 2 item to be removed from $\mathcal{A}^{(h)}$ is selected on the basis of the same criterion
 3 as above, that is,
 4

$$5 \quad BIC_{diff}(\mathcal{A}^{(h)} \setminus j, j) = BIC_{k^{(h)}}(\mathcal{A}^{(h)}) - [BIC_1(j) + \min_{2 \leq k \leq k_{\max}} BIC_k(\mathcal{A}^{(h)} \setminus j)].$$

6
 7
 8 The item with the highest positive value of $BIC_{diff}(\mathcal{A}^{(h)} \setminus j, j)$ is re-
 9 moved from $\mathcal{A}^{(h)}$ and $k^{(h)}$ is updated. If no item is found with a positive
 10 $BIC_{diff}(\mathcal{A}^{(h)} \setminus j, j)$, $\mathcal{A}^{(h)}$ is left unchanged.
 11

12 The algorithm ends when no item is added to $\mathcal{A}^{(h)}$ and no item is removed from
 13 $\mathcal{A}^{(h)}$. It has to be clear that different LC models are estimated at each step of
 14 this algorithm. These models are different in the set of items and in the number
 15 of latent classes. Obviously, an initialization of the EM algorithm is required
 16 for each of these models. In particular, Dean and Raftery (2010) propose to use
 17 the parameter estimates, available at the end of the previous step of the item
 18 selection algorithm, to obtain these starting values. In more detail, they use the
 19 estimated posterior probabilities \hat{z}_{iu} as reasonable starting values for models
 20 involving the updated dataset, with one more or one less item. When a different
 21 number of latent classes is considered, the new latent classes are obtained by
 22 collapsing two or more closest classes, in terms of Euclidean distance between
 23 the corresponding conditional response probabilities, or by splitting one or
 24 more classes into new classes.
 25

26 As the above initialization does not prevent the problem of the likelihood
 27 multimodality that is particularly severe for the LC model, we propose to in-
 28 clude an additional step, after each inclusion and exclusion step aimed at per-
 29 forming a check based on several random initializations of the EM algorithm.
 30 In particular, as outlined in Section 3.2.1, we initialize the EM algorithm by
 31 a large number of random starting values, proportional to the current num-
 32 ber of latent states, and we take the estimates corresponding to the highest
 33 log-likelihood at convergence of the algorithm. This random check, which is
 34 performed once an item has been included or removed and the corresponding
 35 number of latent classes has been selected, allows us to assess the convergence
 36 to the global maximum of the model likelihood, without being too computa-
 37 tionally expensive.
 38

39 The item selection algorithm at issue also requires to properly choose the
 40 initial set $\mathcal{A}^{(0)}$ of items. For this aim, Dean and Raftery (2010) propose to
 41 estimate an LC model with at least 2 classes for all the items and to order
 42 these items in terms of variability of the corresponding estimated response
 43 probabilities across classes. Then, they choose the smallest set of items that,
 44 among the items having the highest variability, allows at least a 2-class model
 45 to be identified. In the application which follows, we also perform a sensitivity
 46 analysis of the final solution to different initial set of items.
 47

48 It is worth noting that $BIC_{diff}(\mathcal{A}, j)$ corresponds to the difference between
 49 the BIC index for an LC model based on the assumption that the items in
 50 $\mathcal{A} \cup j$ depend on the latent variable (and the items in $\mathcal{A} \setminus j$ do not depend on
 51
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

the latent variable) and that for the corresponding model in which only the items in \mathcal{A} depend on the latent variable (and the items in $\bar{\mathcal{A}}$ do not depend on the latent variable). Consequently, the item selection algorithm above is aimed at finding the set of items and the number of classes which minimize the index:

$$BIC_{tot,k}(\mathcal{A}) = BIC_k(\mathcal{A}) + BIC_1(\bar{\mathcal{A}}). \quad (7)$$

Through the last index it is then possible to compare the solutions obtained by the item selection algorithm, even with a different number of latent classes or a different number of items included in the set \mathcal{A} .

5 Application to the ULISSE dataset

In this section, we illustrate the results obtained from the application of the item selection algorithm to the ULISSE dataset described in Section 2.

5.1 Item selection

First of all, in order to apply the inclusion-exclusion algorithm, we fit the LC model described in Section 3, considering the full set of 75 items, denoted by \mathcal{J} , for a number of latent classes from 2 to k_{\max} , with $k_{\max} = 10$. For each k , we initialize the EM algorithm by means of $100 \times (k - 1)$ random starting values. The number of latent classes corresponding to the minimum of $BIC_k(\mathcal{J})$ defined in (6) is $k = 8$.

In order to select the initial set of clustering items, for each category of each item, we then calculate the variance of its estimated conditional probability across classes. For each item, we finally add up these variances and we order the items according to this sum. This because items with high values of this sum have high between-class variability, and therefore may be more useful for clustering. We consider different sizes of the initial set $\mathcal{A}^{(0)}$, equal to 3, 10, 20, 30, 75, to perform the inclusion-exclusion algorithm. Using different choices allows us to study the sensitivity of the final solution with respect to different starting sets. For each initial set of items, we select the initial number of latent classes $k^{(0)}$, always on the basis of $BIC_k(\mathcal{A}^{(0)})$, for $k = 2, \dots, k_{\max}$, and we apply the item selection procedure.

Table 2 reports the results obtained, in terms of number of items, number of classes, and corresponding $BIC_k(\hat{\mathcal{A}})$ and $BIC_{tot,k}(\hat{\mathcal{A}})$ index, defined in (7), of the final model. The output of the inclusion-exclusion algorithm shows a moderate dependence on the choice of the initial set of clustering items. Moreover, the random check allows us to increase the chance of reaching a global maximum of the model log-likelihood, with always better results in terms of $BIC_{tot,k}(\hat{\mathcal{A}})$ with respect to the algorithm without the random check. From the results, we also note that the selected items are mostly equivalent, apart from four items, and the number of classes varies from $\hat{k} = 8$ and $\hat{k} = 10$, leading to different values of the $BIC_{tot,k}(\hat{\mathcal{A}})$ index. The best result,

in terms of $BIC_{tot,k}(\hat{\mathcal{A}})$, is obtained with an initial set of 30 items, that leads to selecting 50 items and $\hat{k} = 9$ latent classes.

Table 2 Comparison between the results of the inclusion-exclusion algorithm for item selection with respect to different sizes of the initial set of clustering items (in boldface are the quantities corresponding to the best solution in terms of $BIC_{tot,k}(\hat{\mathcal{A}})$)

size of $\mathcal{A}^{(0)}$	$k^{(0)}$	<i>random check</i>	# items	\hat{k}	$BIC_k(\hat{\mathcal{A}})$	$BIC_{tot,k}(\hat{\mathcal{A}})$
3	2	<i>yes</i>	53	8	129,344.90	165,353.30
3	2	<i>no</i>	49	10	123,023.90	165,573.40
10	10	<i>yes</i>	50	9	124,815.80	165,320.90
10	10	<i>no</i>	53	8	129,388.40	165,396.90
20	10	<i>yes</i>	50	9	124,816.50	165,321.50
20	10	<i>no</i>	50	9	124,835.10	165,340.10
30	9	<i>yes</i>	50	9	124,799.10	165,304.20
30	9	<i>no</i>	51	9	127,332.90	165,486.10
75	8	<i>yes</i>	53	8	129,351.60	165,360.00
75	8	<i>no</i>	53	8	129,365.90	165,374.40

5.2 Validation by resampling

Given the nature of the search algorithm and the complexity of the study we are dealing with, the selected number of latent classes and the final set of clustering items may also be sensitive to the specific data used for the analysis. In order to address this issue, we validate the results obtained above by sampling with replacement, from the original dataset, $B = 99$ random samples of the same size n of the original one.

For each sample, we then select the optimal number of latent classes, denoted by k_b , corresponding to the minimum $BIC_{k_b}(\mathcal{I})$, $k_b = 2, \dots, k_{\max}$, $b = 1, \dots, B$. As done for the original data, we order the full set of items on the basis of the variability of their estimated conditional response probabilities across classes with the aim of selecting the initial subset of items, $\mathcal{A}_b^{(0)}$, and we apply the item selection algorithm with random check. For each sample, we consider a size of the initial set equal to 30 items, which give the best results, in terms of $BIC_{tot,k}$, in the application to the original data.

In Table 3 we summarize the results of the above validation procedure. In particular, for each item in the full set, we report the number of times that it has been selected with respect to the different starting set in the original data, if it has been included in the best solution (obtained with a starting set of 30 items as illustrated in Table 2), and the number of times that the item has been selected with respect to the different random samples obtained in the validation procedure. From the table, we observe that 45 items are always included

1 in the different final solutions, both with respect to different specifications of
 2 the initial set of items and with respect to the different random samples in
 3 the validation procedure. Moreover, items 14 and 68 are always included apart
 4 from one and three random samples. On the other hand, we note that 7 items
 5 are never included in the final solutions, whereas 7 items are included only in
 6 very few solutions provided by the random samples (see items 19, 22, 28, 36,
 7 39, 72, and 73) . The remaining 14 items are in intermediated situations. In
 8 conclusion, more than three-quarter of the random samples confirm the results
 9 obtained by the item selection algorithm.

10 With respect to the sections of the questionnaire, we observe that all items
 11 referred to sections CC, AVF, and ADL and two items of section I are retained
 12 in $\hat{\mathcal{A}}$. On the contrary, most of the excluded items belongs to sections DD,
 13 HBD, NF and SC.
 14
 15

16 5.3 Parameter estimates

17 With the aim of evaluating the performance of the nursing homes, the adopted
 18 approach also allows us to cluster their patients into the different latent classes
 19 according to their health conditions on the basis of the parameter estimates.
 20 This may be useful to describe the case-mix of the nursing homes and to es-
 21 timate their ability of retaining patients in the groups corresponding to better
 22 health conditions.
 23

24 In this section, we report the estimation results based on the best solu-
 25 tion, in terms of $BIC_{tot,k}(\hat{\mathcal{A}})$, provided by the inclusion-exclusion algorithm,
 26 which selects 50 items with $\hat{k} = 9$. Since the items are categorical, with a
 27 different number of categories, we report the estimated conditional response
 28 probabilities, $\hat{\lambda}_{j|u}(y)$, by assigning an equally-spaced score between 0 and 1 to
 29 the different response categories. Then, we compute the average of the scores,
 30 weighted with the corresponding response probabilities. This amounts to com-
 31 puting the following *item mean score*
 32

$$33 \hat{\mu}_{j|u} = \frac{1}{l_j - 1} \sum_y (y - 1) \hat{\lambda}_{j|u}(y), \quad j \in \hat{\mathcal{A}}, u = 1, \dots, \hat{k}, y = 0, \dots, l_j - 1.$$

34 In particular, a value of $\hat{\mu}_{j|u}$ close to 0 corresponds to a low probability of
 35 suffering from a certain pathology, whereas a value close to 1 corresponds to
 36 a high probability of suffering from the same pathology. To summarize these
 37 results, we also compute the *section mean score* $\hat{\mu}_{d|u}$ as the average of $\hat{\mu}_{j|u}$ for
 38 the items in $\hat{\mathcal{A}}$ composing each section d of the questionnaire.
 39

40 In order to have a clearer interpretation of the results, we order the latent
 41 classes on the basis of the values of $\hat{\mu}_{d|u}$ assumed in the section denoted by
 42 ADL (Activity of Daily Living) of the questionnaire. This is the section with
 43 the highest difference between the maximum and the minimum value of the
 44 section mean score $\hat{\mu}_{d|u}$ across classes.
 45
 46
 47
 48
 49
 50
 51
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

Table 3 Final results of the item selection algorithm (with random check) and of the validation by resampling. j is the item index, #sel. is the number of times that item j has been selected with respect to the different starting sets, best indicates if item included in the best solution, #resamp. is the number of times that item j has been selected with respect to the different samples.

section	j	#sel.	best	#resamp	section	j	#sel.	best	#resamp
CC	1	5	×	99	HBD	39			4
CC	2	5	×	99	ADL	40	5	×	99
CC	3	5	×	99	ADL	41	5	×	99
CC	4	5	×	99	ADL	42	5	×	99
CC	5	5	×	99	ADL	43	5	×	99
CC	6	5	×	99	ADL	44	5	×	99
CC	7	5	×	99	ADL	45	5	×	99
CC	8	5	×	99	ADL	46	5	×	99
CC	9	5	×	99	ADL	47	5	×	99
CC	10	5	×	99	ADL	48	5	×	99
CC	11	5	×	99	ADL	49	5	×	99
CC	12	5	×	99	ADL	50	5	×	99
CC	13	5	×	99	ADL	51	5	×	99
AVF	14	5	×	98	ADL	52	5	×	99
AVF	15	5	×	99	ADL	53	5	×	99
AVF	16	5	×	99	ADL	54	5	×	99
AVF	17	5	×	99	ADL	55	5	×	99
AVF	18	5	×	99	ADL	56	5	×	99
HBD	19			5	ADL	57	5	×	99
HBD	20	5	×	96	I	58	5	×	99
HBD	21			12	I	59	5	×	99
HBD	22			2	I	60	5	×	79
HBD	23			10	NF	61	5	×	99
HBD	24			25	NF	62	5	×	99
HBD	25				NF	63			
HBD	26	2		59	NF	64	2		42
HBD	27			31	NF	65			
HBD	28			1	NF	66	2		30
HBD	29			20	DD	67			36
HBD	30			23	DD	68	5	×	96
HBD	31				DD	69	5	×	75
HBD	32	5	×	99	DD	70			
HBD	33	5	×	99	DD	71			
HBD	34	5	×	99	DD	72			9
HBD	35	5	×	99	SC	73			3
HBD	36			1	SC	74			
HBD	37			17	SC	75	5	×	99
HBD	38	5	×	99					

For each latent class, Table 4 shows the values of $\hat{\mu}_{d|u}$ and the estimated class weights $\hat{\pi}_u$, together with the difference between the maximum and the minimum value of $\hat{\mu}_{d|u}$ for each section of the questionnaire. As we can note, the latter is high for sections ADL, CC, I, and AVF, and low for the remaining sections. The smallest among these differences is observed for section DD, which, consequently, tends to discriminate less between subjects with respect to the other sections. The first latent class, that includes around 17% of subjects, corresponds to the best health conditions with respect to all the

1 pathologies measured by the sections of the questionnaire, apart from sec-
 2 tion NF, DD and SC. On the other hand, the 9th latent class, which includes
 3 about 11% of patients, corresponds to cases with the worst health conditions
 4 for almost all the pathologies. Intermediate classes show a different case-mix
 5 depending on the section mean score pattern. For instance, in the 3rd class are
 6 included patients with severe cognitive conditions (CC) and consistent impair-
 7 ment referred to sections AVF, HBD, and I. Moreover, in the same class we
 8 also register a moderate impairment of the activities of daily living (ADL). In
 9 the 7th class are instead included patients with worse conditions than those
 10 assigned to the 3rd class; in particular, in addition to section CC, we also
 11 register a severe impairment of the incontinence conditions (section I), and a
 12 worsening of the pathologies measured by sections ADL, NF and SC.
 13
 14

15
 16 **Table 4** Estimated section mean score, $\hat{\mu}_{d|u}$, for each latent class u and each section d of
 17 the questionnaire, together with the estimated weights $\hat{\pi}_u$ and the difference between the
 18 largest and the smallest estimated section mean score for each section, under the latent
 19 ignorability assumption.

u	d								$\hat{\pi}_u$
	1 (CC)	2 (AVF)	3 (HBD)	4 (ADL)	5 (I)	6 (NF)	7 (DD)	8 (SC)	
1	0.041	0.098	0.066	0.090	0.230	0.085	0.392	0.047	0.169
2	0.361	0.231	0.205	0.131	0.347	0.078	0.390	0.031	0.105
3	0.694	0.444	0.464	0.238	0.688	0.168	0.398	0.056	0.084
4	0.143	0.178	0.089	0.315	0.407	0.148	0.377	0.100	0.096
5	0.596	0.360	0.292	0.527	0.776	0.215	0.338	0.067	0.098
6	0.079	0.141	0.078	0.628	0.609	0.125	0.406	0.145	0.103
7	0.757	0.593	0.357	0.680	0.894	0.337	0.371	0.172	0.102
8	0.479	0.328	0.180	0.739	0.844	0.240	0.401	0.200	0.129
9	0.735	0.666	0.259	0.895	0.888	0.439	0.380	0.315	0.114
$\max_u(\hat{\mu}_{d u}) -$	0.716	0.568	0.398	0.805	0.664	0.360	0.068	0.284	
$\min_u(\hat{\mu}_{d u})$									

37 6 Conclusions

38 In this paper, we illustrate the application of an algorithm for item selection,
 39 when items are used for clustering purposes, which is based on the latent
 40 class (LC) model (Lazarsfeld, 1950; Lazarsfeld and Henry, 1968; Goodman,
 41 1974). This algorithm closely follows the one proposed by Dean and Raftery
 42 (2010), and aims at finding the optimal subset of items useful for clustering
 43 searching for the best result in terms of the Bayesian Information Criterion
 44 (BIC, Schwarz, 1978).
 45

46 More in detail, we illustrate an application based on a dataset collected
 47 within the Italian project named ULISSE (Lattanzio et al, 2010), regarding
 48 the quality-of-life of elderly hosted in nursing homes. As typically happens in
 49
 50
 51
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

1 such a context, the questionnaire used to collect the data is made of a large
2 number of polytomous items. This may lead to a lengthy and expensive ad-
3 ministration of the questionnaire and may induce the respondents to provide
4 inaccurate responses. In this respect, the algorithm for item selection we illus-
5 trate may lead to a sensible reduction of the number of items for clustering
6 purposes. Moreover, by removing irrelevant or noise items, it may improve
7 the performance of the model-based clustering procedure and the accuracy of
8 the choice of the number of latent classes. The adopted algorithm extends the
9 inclusion-exclusion algorithm proposed by Dean and Raftery (2010), by in-
10 cluding an additional step, which we call random check, aimed at initializing,
11 with a large number of random starting values, the estimation algorithm, so
12 as to prevent the problem of the multimodality of the likelihood.

13
14 In the present application to the ULISSE dataset, we also perform a sensi-
15 tivity analysis of the final solution with respect to different specifications of
16 the initial set of clustering items and a validation of the results by means of
17 a resampling procedure. The results confirm that the random check allows us
18 to increase the chance of reaching the global maximum of the log-likelihood,
19 especially in the presence of complex models characterized by a large number
20 of items and estimated latent classes. Moreover, the validation procedure may
21 be useful in applications concerning complex phenomena, where the results
22 may be sensitive to the specific data used in the analysis and may be affected
23 by potential outliers in the respondents.

24
25 The best result, in terms of BIC, leads to selecting 50 items, out of the
26 75 considered, and 9 latent classes. This reduction implies clear advantages in
27 terms of setting up a questionnaire which may be more easily administered,
28 especially in a longitudinal context in which we have repeated measurements.
29 Most of the selected items belong to sections of the questionnaire referred to
30 cognitive conditions, auditory and view fields, activities of daily living and
31 incontinence. The remaining sections are the ones that tend to discriminate
32 less between subjects in the estimated latent classes.

33
34 Once the optimal subset of items has been selected together with the cor-
35 responding number of latent classes, the estimation results may be used to
36 assign subjects to homogenous classes, which is one of the main aim of the LC
37 model. This may have important implications in the context of the ULISSE
38 project, where patients are assigned to different latent classes corresponding to
39 different levels of impairment. This is useful for evaluating the long-term nurs-
40 ing homes performance with respect to their ability in improving the patients
41 health conditions or in delaying their worsening.

42 43 44 45 **References**

46
47 Bacci S, Bartolucci F, Gnaldi M (2014) A class of multidimensional latent
48 class IRT models for ordinal polytomous item responses. *Communications*
49 *in Statistics - Theory and Methods* 43:787–800
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1 Bandeen-Roche K, Miglioretti DL, Zeger SL, Rathouz PJ (1997) Latent vari-
2 able regression for multiple discrete outcomes. *Journal of the American Stat-*
3 *istical Association* 92:1375–1386
- 4 Bandeen-Roche K, Xue QL, Ferrucci L, Walston J, Guralnik JM, Chaves P,
5 Zeger SL, Fried LP (2006) Phenotype of frailty: characterization in the wo-
6 men’s health and aging studies. *The Journals of Gerontology Series A: Bio-*
7 *logical Sciences and Medical Sciences* 61:262–266
- 8 Biernacki C, Celeux G, Govaert G (2003) Choosing starting values for the
9 EM algorithm for getting the highest likelihood in multivariate Gaussian
10 mixture models. *Computational Statistics & Data Analysis* 41:561–575
- 11 Breyer F, Costa-Font J, Felder S (2010) Ageing, health, and health care. *Ox-*
12 *ford Review of Economic Policy* 26:674–690
- 13 Dean N, Raftery A (2010) Latent class analysis variable selection. *Annals of*
14 *the Institute of Statistical Mathematics* 62:11–35
- 15 Dempster AP, Laird NM, Rubin DB (1977) Maximum Likelihood from incom-
16 plete data via the EM algorithm. *Journal of the Royal Statistical Society,*
17 *Series B* 39:1–38
- 18 Fraley C, Raftery AE (2002) Model-based clustering, discriminant analysis,
19 and density estimation. *Journal of the American Statistical Association*
20 97:611–631
- 21 Galasso V, Profeta P (2007) How does ageing affect the welfare state?
22 *European Journal of Political Economy* 23:554–563
- 23 Goodman L (1974) Exploratory latent structure analysis using both identifi-
24 able and unidentifiable models. *Biometrika* 61:215–231
- 25 Harel O, Schafer JL (2009) Partial and latent ignorability in missing-data
26 problems. *Biometrika* 96:37–50
- 27 Hawes C, Morris JN, Phillips CD, Fries BE, Murphy K, Mor V (1997) De-
28 velopment of the nursing home resident assessment instrument in the USA.
29 *Age and Ageing* 26:19–25
- 30 Karlis D, Xekalaki E (2003) Choosing initial values for the EM algorithm for
31 finite mixtures. *Computational Statistics & Data Analysis* 41:577–590
- 32 Kass R, Raftery A (1995) Bayes factors. *Journal of the American Statistical*
33 *Association* 90:773–795
- 34 Kohler H, Billardi FC, Ortega J (2002) The emergence of lowest-low fertility in
35 Europe during the 1990s. *Population and Development review* 28:641–680
- 36 Lafortune L, Beland F, Bergman H, Ankri J (2009) Health status transitions in
37 community-living elderly with complex care needs: a latent class approach.
38 *BMC Geriatrics* 9:6
- 39 Lattanzio F, Mussi C, Scafato E, Ruggiero C, Dell’Aquila G, Pedone C, Mam-
40 marella F, Galluzzo L, Salvioli G, Senin U, Carbonin PU, Bernabei R, Cher-
41 ubini A (2010) Health care for older people in Italy: The U.L.I.S.S.E. project
42 (un link informatico sui servizi sanitari esistenti per l’anziano - a computer-
43 ized network on health care services for older people). *J Nutr Health Aging*
44 14:238–42
- 45 Lazarsfeld PF (1950) The logical and mathematical foundation of latent struc-
46 ture analysis. In: S A Stouffer EAS L Guttman (ed) *Measurement and Pre-*
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

-
- 1 diction, Princeton University Press, New York
2 Lazarsfeld PF, Henry NW (1968) Latent Structure Analysis. Houghton Mifflin,
3 Boston
4 Little RJA, Rubin DB (2002) Statistical Analysis with Missing Data, 2nd edn.
5 Wiley Series in Probability and Statistics, Wiley
6 Lu G, Copas JB (2004) Missing at random, likelihood ignorability and model
7 completeness. *The Annals of Statistics* 32:754–765
8 Magidson J, Vermunt JK (2001) Latent class factor and cluster models, bi-
9 plots and related graphical displays. *Sociological Methodology* 31:223–264
10 Moran M, Walsh C, Lynch A, Coen RF, Coakley D, Lawlor BA (2004) Syn-
11 dromes of behavioural and psychological symptoms in mild Alzheimer’s dis-
12 ease. *International Journal of Geriatric Psychiatry* 19:359–364
13 Morris J, Hawes C, Murphy K, et al (1991) Resident Assessment Instrument
14 Training Manual and Resource Guide. Eliot Press, Natick, MA
15 Rubin DB (1976) Inference and missing data. *Biometrika* 63:581–592
16 Samejima F (1969) Estimation of ability using a response pattern of graded
17 scores. *Psychometrika Monograph* 17
18 Samejima F (1996) Evaluation of mathematical models for ordered polychot-
19 omous responses. *Behaviormetrika* 23:17–35
20 Schwarz G (1978) Estimating the dimension of a model. *The Annals of Stat-*
21 istics 6:461–464
22 Vermunt JK, Magidson J (2002) Latent class cluster analysis. In: Hagenaars
23 JA, McCutcheon AL (eds) *Applied latent class analysis*, Cambridge Univer-
24 sity Press, Cambridge, UK
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Appendix

<i>j</i>	#	item description
cat.		
Section CC		
01	2	Short-term memory (0 = “recalls what recently happened (5 minutes)”, 1 = “does not recall”)
02	2	Long-term memory (0 = “keeps some past memories green”, 1 = “does not keep some past memories green”)
03	2	Memory status (0 = “recalls the actual season”, 1 = “does not recall the actual season”)
04	2	Memory status (0 = “recalls where is his room”, 1 = “does not recall where is his room”)
05	2	Memory status (0 = “recalls the names and faces of the staff”, 1 = “does not recall the names and faces of the staff”)
06	2	Memory status (0 = “recalls where he is”, 1 = “does not recall where he is”)
07	4	Decision about his daily activities (from 0 = “independent decisions” to 3 = “unable to decide”)
08	3	Easily sidetracked (from 0 = “problems absent” to 2 = “problems worsened in the last week”)
09	3	Altered perception or awareness of surrounding (from 0 = “problems absent” to 2 = “problems worsened in the last week”)
10	3	Disorganized speech (from 0 = “problems absent” to 2 = “problems worsened in the last week”)
11	3	Restlessness movements (from 0 = “problems absent” to 2 = “problems worsened in the last week”)
12	3	Lethargic spans (from 0 = “problems absent” to 2 = “problems worsened in the last week”)
13	3	Change in the cognitive conditions during the day (from 0 = “problems absent” to 2 = “problems worsened in the last week”)
Section AVF		
14	4	Hearing (from 0 = “no hearing impairment” to 3 = “severe hearing impairment”)
15	4	Ability to make itself understood (from 0 = “understood” to 3 = “seldom/never understood”)
16	3	Clear language (from 0 = “clear language” to 2 = “no language”)
17	4	Ability to understand others (from 0 = “understands” to 3 = “seldom/never understands”)
18	5	Sight in conditions of adequate lighting (from 0 = “no sight impairment” to 4 = “severe sight impairment”)
Section HDB		
19	3	Negative statements (from 0 = “symptom not showed” to 2 = “symptom daily showed”)
20	3	Repetitive questions (from 0 = “symptom not showed” to 2 = “symptom daily showed”)
21	3	Repetitive verbalizations (from 0 = “symptom not showed” to 2 = “symptom daily showed”)
22	3	Persistent anger with himself or others (from 0 = “symptom not showed” to 2 = “symptom daily showed”)
23	3	Self deprecation disesteem (from 0 = “symptom not showed” to 2 = “symptom daily showed”)
24	3	Fears that are not real (from 0 = “symptom not showed” to 2 = “symptom daily showed”)
25	3	To believe himself to be dying (from 0 = “symptom not showed” to 2 = “symptom daily showed”)
26	3	To complain about his health (from 0 = “symptom not showed” to 2 = “symptom daily showed”)
27	3	Repeated events anxiety (from 0 = “symptom not showed” to 2 = “symptom daily showed”)
28	3	Unpleasant mood in morning (from 0 = “symptom not showed” to 2 = “symptom daily showed”)
29	3	Insomnia/problems with sleep (from 0 = “symptom not showed” to 2 = “symptom daily showed”)
30	3	Expressions of sad-faced (from 0 = “symptom not showed” to 2 = “symptom daily showed”)
31	3	Easily tears (from 0 = “symptom not showed” to 2 = “symptom daily showed”)
32	3	Repetitive movements (from 0 = “symptom not showed” to 2 = “symptom daily showed”)
33	3	Abstention from activities of interest (from 0 = “symptom not showed” to 2 = “symptom daily showed”)
34	3	Reduced local interactions (from 0 = “symptom not showed” to 2 = “symptom daily showed”)
35	4	To wander aimlessly (from 0 = “problem absent” to 3 = “problem daily encountered”)
36	4	Offensive language (from 0 = “problem absent” to 3 = “problem daily encountered”)
37	4	Physically aggressive (from 0 = “problem absent” to 3 = “problem daily encountered”)
38	4	Socially inappropriate behavior (from 0 = “problem absent” to 3 = “problem daily encountered”)
39	4	To refuse assistance (from 0 = “problem absent” to 3 = “problem daily encountered”)

Table 5 Description of the full set of items.

<i>j</i>	#	item description	cat.
Section ADL			
40	5	Moving to/from lying position (from 0 = "independent" to 4 = "totally dependent")	
41	5	Moving to/from bed, chair, wheelchair (from 0 = "independent" to 4 = "totally dependent")	
42	5	Walking between different points within the room (from 0 = "independent" to 4 = "totally dependent")	
43	5	Walking in the corridor (from 0 = "independent" to 4 = "totally dependent")	
44	5	Walking into the nursing home ward (from 0 = "independent" to 4 = "totally dependent")	
45	5	Walking outside the nursing home ward (from 0 = "independent" to 4 = "totally dependent")	
46	5	Dressing (from 0 = "independent" to 4 = "totally dependent")	
47	5	Eating (from 0 = "independent" to 4 = "totally dependent")	
48	5	Using the toilet room (from 0 = "independent" to 4 = "totally dependent")	
49	5	Personal hygiene (from 0 = "independent" to 4 = "totally dependent")	
50	5	Taking full-body bath/shower (from 0 = "independent" to 4 = "totally dependent")	
51	4	Balance problems (from 0 = "does not have balance problems" to 3 = "needs physical assistance")	
52	3	Mobility in the neck (0 = "no limitation", 1 = "unilateral limitation", 2 = "bilateral limitation")	
53	3	Mobility in the arm including shoulder or elbow (0 = "no limitation", 1 = "unilateral limitation", 2 = "bilateral limitation")	
54	3	Movements of the hand including wrist or finger (0 = "no limitation", 1 = "unilateral limitation", 2 = "bilateral limitation")	
55	3	Mobility in the leg and hip (0 = "no limitation", 1 = "unilateral limitation", 2 = "bilateral limitation")	
56	3	Mobility in the foot and ankle (0 = "no limitation", 1 = "unilateral limitation", 2 = "bilateral limitation")	
57	3	Other movements (0 = "no limitation", 1 = "unilateral limitation", 2 = "bilateral limitation")	
Section I			
58	5	Fecal incontinence (from 0 = "continence" to 4 = "incontinence")	
59	5	Urinary incontinence (from 0 = "continence" to 4 = "incontinence")	
60	2	Elimination of feces (0 = "adequate", 1 = "not adequate")	
Section NF			
61	2	Chewing problem (0 = "no problem", 1 = "problems")	
62	2	Swallowing problem (0 = "no problem", 1 = "problems")	
63	2	Mouth pain (0 = "no problem", 1 = "problems")	
64	2	Taste of many foods (0 = "does not complain", 1 = "complains")	
65	2	Hungry (0 = "does not complain", 1 = "complains")	
66	2	Food on his plate (0 = "does not leave it", 1 = "leaves it")	
Section DD			
67	2	Debris present in mouth prior to going to bed at night (0 = "problem absent", 1 = "problem present")	
68	2	Dentures/removable bridge (0 = "absent", 1 = "present")	
69	2	Some/all natural teeth lost and does not have/does not use dentures (or partial plates) (0 = "problem absent", 1 = "problem present")	
70	2	Broken, loose, or carious teeth (0 = "problem absent" 1 = "problem present")	
71	2	Inflamed gums, swollen or bleeding gums, oral abscesses, ulcers or rashes (0 = "problem absent", 1 = "problem present")	
72	2	Dentures or removable bridge daily cleaned by resident or staff (0 = "absent", 1 = "present")	
Section SC			
73	5	Pressure ulcer (from 0 = "no pressure ulcer" to 4 = "stage 4")	
74	5	Stasis ulcers (from 0 = "no pressure ulcer" to 4 = "stage 4")	
75	2	Resolved or cured ulcer (0 = "absent", 1 = "present")	

Table 6 Description of the full set of items (continued).