

A note on the mixture transition distribution and hidden Markov models

Francesco Bartolucci^{a,*} and Alessio Farcomeni^b

We discuss an interpretation of the mixture transition distribution (MTD) for discrete-valued time series which is based on a sequence of independent latent variables which are occasion-specific. We show that, by assuming that this latent process follows a first order Markov Chain, MTD can be generalized in a sensible way. A class of models results which also includes the hidden Markov model (HMM). For these models we outline an EM algorithm for the maximum likelihood estimation which exploits recursions developed within the HMM literature. As an illustration, we provide an example based on the analysis of stock market data referred to different American countries.

Keywords: Backward–forward recursions; discrete-valued time series; EM-algorithm; state-space models.

1. INTRODUCTION

Let $X_t, t = 1, \dots, T$, be a sequence of random variables having support $\{1, \dots, k\}$ and let x_t denote a realization of X_t . This sequence is said to follow a mixture transition distribution (MTD) of order l , MTD_l for short, when

$$p(x_t | x_1, \dots, x_{t-1}) = p(x_t | x_{t-l}, \dots, x_{t-1}) = \sum_h \lambda_h \pi_{x_{t-h}, x_t}, \quad t > l, \quad (1)$$

where $\lambda_h, h = 1, \dots, l$, are *weights* and $\pi_{j_1, j_2}, j_1, j_2 = 1, \dots, k$, are *transition probabilities*. The former ones are subject to the constraints $\lambda_h \geq 0, h = 1, \dots, l$, and $\sum_h \lambda_h = 1$. Assumption (1) implies that the joint probability of the entire sequence of random variables is given by

$$p(x_1, \dots, x_T) = p(x_1, \dots, x_l) \prod_{t>l} \sum_h \lambda_h \pi_{x_{t-h}, x_t}, \quad (2)$$

where $p(x_1, \dots, x_l)$ denotes the joint probability of the first l observations which may be arbitrarily defined. This model was introduced by Raftery (1985a); see also Raftery and Tavaré (1994) who discussed more general constraints on the parameters λ_h . For an exhaustive review on MTD, see Berchtold and Raftery (2002). Obviously, a similar model may be assumed for a sequence of continuous random variables $X_t, t = 1, \dots, T$. In this case, we write

$$f(x_t | x_1, \dots, x_{t-1}) = f(x_t | x_{t-l}, \dots, x_{t-1}) = \sum_h \lambda_h \tau(x_t | x_{t-h}), \quad t > l,$$

with $f(\cdot)$ standing for conditional density function and $\tau(\cdot)$ denoting a suitable transition kernel. A similar extension is possible for the model proposed in this article.

With respect to a Markov Chain model of order l , an MTD model with the same order has the advantage of being much more parsimonious because it is based on $(l - 1) + k(k - 1)$ parameters and this number increases linearly with l . We recall that a Markov Chain of order l is instead based on $k^l(k - 1)$ parameters; this number increases exponentially with l . In both cases we do not consider the parameters used to define the initial probability $p(x_1, \dots, x_l)$. Also note that these transition probabilities can be lag-specific, so that

$$p(x_t | x_1, \dots, x_{t-1}) = p(x_t | x_{t-l}, \dots, x_{t-1}) = \sum_h \lambda_h \pi_{x_{t-h}, x_t}^{(h)}, \quad t > l, \quad (3)$$

and then a generalization of the MTD model results. This generalized model is indicated by $gMTD_l$ (Raftery, 1985b; Berchtold, 1998). Even in this case, the number of parameters increases linearly with l , since it is given by $(l - 1) + kl(k - 1)$.

Although the MTD model is generally justified by claiming its parsimony and good adaptation when fitting discrete-valued time series, there are different interpretations and justifications that can be additionally put forward. First of all (1) implies that

^aUniversity of Perugia

^bUniversity of Rome "La Sapienza"

*Correspondence to: Department of Economics, Finance and Statistics, University of Perugia, 06123 Perugia, Italy

[†]E-mail: bart@stat.unipg.it

$$p(X_t = j_2 | X_{t-l} = \dots = X_{t-1} = j_1) = \pi_{j_1, j_2}, \quad t > l,$$

so that π_{j_1, j_2} is the probability that the chain moves to state j_2 after it persisted in state j_1 for a period of length l . On the other hand, the weights λ_h can be directly interpreted as the relative influence of each of the previous l occasions on the present.

A more interesting interpretation of the MTD model may be obtained by introducing the occasion-specific latent variables $Z_t, t = l + 1, \dots, T$, which are independent and identically distributed and are also independent of X_1, \dots, X_T . Each variable Z_t has a discrete distribution with support $\{1, \dots, l\}$ and mass probabilities $\lambda_1, \dots, \lambda_l$. In particular, we can easily show that the MTD model, formulated in its generalized version based on (3), is equivalent to a model based on the assumption

$$p(x_t | x_1, \dots, x_{t-1}, z_{l+1}, \dots, z_t) = p(x_t | x_{t-l}, \dots, x_{t-1}, z_t) = \sum_h I(z_t = h) \pi_{x_{t-h}, x_t}^{(h)}, \quad t > l, \tag{4}$$

where $I(\cdot)$ is the indicator function. According to (4), the response variable X_t depends only on the lagged variable X_{t-h} , where the lag h is chosen by a random mechanism which is not directly observable. Then, for $t > l$ and given $Z_t = h$ and $X_{t-h} = j_1, \pi_{j_1, j_2}^{(h)}$ is the conditional probability of $X_t = j_2$, i.e. $p(X_t = j_2 | X_{t-h} = j_1, Z_t = h) = \pi_{j_1, j_2}^{(h)}$. This latent variable interpretation of the gMTD model motivates the use of the EM algorithm (Dempster *et al.*, 1977) for parameter estimation; see also Le *et al.* (1996).

The assumption that Z_{l+1}, \dots, Z_T is a sequence of independent random variables implies that, at each time occasion, the lag on which to rely is independent of the lags previously adopted. In several contexts, this is far to be realistic. Then, we propose a generalization of the MTD and gMTD models based on the assumption that the sequence Z_{l+1}, \dots, Z_T follows a hidden Markov Chain. Further generalizations are possible, but are easily seen to lead to models in which the number of parameters can be high and whose fit involves computationally intensive algorithms. The proposed generalization is illustrated in Section 2, where we show that the resulting model also generalizes the hidden Markov model (HMM); see MacDonald and Zucchini (1997). Likelihood inference for the proposed model is discussed in Section 3.

2. HIDDEN MARKOV EXTENSION OF THE MIXTURE TRANSITION DISTRIBUTION

The proposed generalization is based on assumption (4) with $Z_t, t = l + 1, \dots, T$, that follows a homogenous first-order Markov Chain with initial probabilities $\rho_h = p(Z_{l+1} = h), h = 1, \dots, l$, and transition probabilities $\phi_{h_1, h_2} = p(Z_t = h_2 | Z_{t-1} = h_1), h_1, h_2 = 1, \dots, l$, for $t > l + 1$.

So, to compute the conditional probability $p(x_{l+1}, \dots, x_T)$, and then $p(x_1, \dots, x_T)$ as in (2), we can exploit a forward recursion which recalls a well-known recursion in the HMM literature. First of all consider that

$$p(x_{l+1}, z_{l+1} | x_1, \dots, x_l) = \rho_h \pi_{x_{l+1-h}, x_{l+1}} \tag{5}$$

and that, for any $t > l + 1$, we have

$$p(x_{l+1}, \dots, x_t, z_t | x_1, \dots, x_l) = \sum_h p(x_{l+1}, \dots, x_{t-1}, z_{t-1} = h | x_1, \dots, x_l) \phi_{h, z_t} \pi_{x_{t-z_t}, x_t}^{(z_t)}. \tag{6}$$

By computing (5) and then (6) for $t = l+2, \dots, T$, we obtain $p(x_{l+1}, \dots, x_T, z_T | x_1, \dots, x_l)$ and consequently the conditional probability of the last $T - l$ observations given the first l observations as

$$p(x_{l+1}, \dots, x_T | x_1, \dots, x_l) = \sum_h p(x_{l+1}, \dots, x_T, z_T = h | x_1, \dots, x_l).$$

Moreover, we have

$$p(x_{l+1} | x_1, \dots, x_l) = \sum_h \rho_h \pi_{x_{l+1-h}, x_{l+1}}^{(h)}$$

which is the same as (3), whereas, for $t > l + 1$, the above assumptions imply that

$$p(x_t | x_1, \dots, x_{t-1}) = \sum_h \lambda_h^{(t)}(x_1, \dots, x_{t-1}) \pi_{x_{t-h}, x_t}^{(h)}, \tag{7}$$

with $\lambda_h^{(t)}(x_1, \dots, x_{t-1})$ denoting the conditional probability of $Z_t = h$ given all the previous observations, which may be computed as

$$\lambda_h^{(t)}(x_1, \dots, x_{t-1}) = \frac{\sum_m p(x_{l+1}, \dots, x_{t-1}, z_{t-1} = m | x_1, \dots, x_l) \phi_{m, h}}{\sum_m p(x_{l+1}, \dots, x_{t-1}, z_{t-1} = m | x_1, \dots, x_l)}.$$

Clearly, expression (7) is a generalization of (3) in which the mixing weights are time-varying and depend on the previous observations. The way in which each weight varies according to t and the previous observations depends on the latent transition probabilities. It is also clear that the above model generalizes not only the MTD and gMTD models, but also the HMM; then we will indicate it by HM-gMTD_l, where l is the lag order.

It is worth noting that the HM-gMTD_l model specializes into the gMTD_l model when $\phi_{h_1, h_2} = \rho_{h_2}$, $h_1, h_2 = 1, \dots, l$, and then the latent variables Z_t are independent of each other and have the same distribution with mass probabilities ρ_1, \dots, ρ_l . On the other hand, the HM-gMTD_l model specializes into the HMM when $\pi_{j_1, j_2}^{(h)} = \pi_{j_2}^{(h)}$, $j_1, j_2 = 1, \dots, k, h = 1, \dots, l$, so that the distribution of each observation does not depend on the previous observations, but only on the corresponding latent variable. Note that when such an assumption is made on the manifest probabilities, the latent process can be considered to start at $t = 1$. Other different models can arise according to the constraints which are put on the parameters of the HM-gMTD_l model.

The above points are summarized in Table 1, where we also indicate how to compute the number of parameters of the HM-gMTD_l model and the most important submodels; see also Table 2 for numerical examples about the application of these rules.

It can be appreciated that the HM-gMTD class is flexible enough to contain many models commonly used for discrete-value time series. The HM-MTD specialization provides a generalization of MTD which is still quite parsimonious while providing interesting insights into persistency phenomena of the series. Also note that the number of parameters of the HM-gMTD_l model is $l^2 - 1 + kl(k - 1)$ which increases quadratically, rather than linearly, in l . In any case, this number is usually much smaller than that of an ordinary Markov Chain model with the same lag, especially when the manifest transition probabilities $\pi_{j_1, j_2}^{(h)}$ are assumed to be constant in h , and then the HM-MTD_l model results. A further reduction in the number of parameters can be achieved by assuming a specific structure for the latent transition matrix with elements ϕ_{h_1, h_2} . For instance, we can assume this matrix to be symmetric, tridiagonal, or even with off-diagonal elements equal to each other. For an illustration of constraints on this type in a similar context see Bartolucci (2006).

Raftery (1985a) showed that the MTD model has the same equilibrium distribution as the first-order Markov Chain with the same transition probabilities, no matter the MTD order. In parallel with that result, we prove below that for any finite l , the stationary distribution of the HM-MTD_l model coincides with that of the corresponding first order Markov Chain with transition probabilities π_{j_1, j_2} , $j_1, j_2 = 1, \dots, k$. It is then straightforward to see that any HM-gMTD_l model has stationary distribution given by a suitable mixture of the stationary distributions associated to each matrix of transition probabilities with elements $\pi_{j_1, j_2}^{(h)}$, $h = 1, \dots, l$.

THEOREM 1. *Let X_1, X_2, \dots be distributed according to the HM-MTD_l model, with l finite, and let π_1, \dots, π_k denote the probability masses of the stationary distribution associated to the transition probabilities π_{j_1, j_2} , $j_1, j_2 = 1, \dots, k$. Then, as t goes to infinity, $p(X_t = j) \rightarrow \pi_j$, $j = 1, \dots, k$.*

Table 1. List of models nested into the HM-gMTD model with the corresponding number of parameters

Model	Nested models	Constraints	#parameters
HM-gMTD _l	HM-MTD _l , gMTD _l , HMM, MTD _l	-	$l^2 - 1 + kl(k - 1)$
HM-MTD _l	MTD _l	$\pi_{j_1, j_2}^{(h)} = \pi_{j_1, j_2}$	$l^2 - 1 + k(k - 1)$
gMTD _l	MTD _l	$\phi_{h_1, h_2} = \rho_{h_2}$	$l - 1 + kl(k - 1)$
HMM		$\pi_{j_1, j_2}^{(h)} = \pi_{j_2}^{(h)}$	$l^2 - 1 + l(k - 1)$
MTD _l		$\phi_{h_1, h_2} = \rho_{h_2}, \pi_{j_1, j_2}^{(h)} = \pi_{j_1, j_2}$	$l - 1 + k(k - 1)$

Table 2. Comparison between the models listed in Table 1 and a Markov Chain model of order l in terms of number of parameters

Model	$l(k = 2)$				
	1	2	3	4	5
HM-gMTD _l	2	7	14	23	34
HM-MTD _l	2	5	10	17	26
gMTD _l	2	5	8	11	14
HMM	1	5	11	19	29
MTD _l	2	3	4	5	6
MarkovChain	2	4	8	16	32
Model	$l(k = 3)$				
	1	2	3	4	5
HM-gMTD _l	6	15	26	39	54
HM-MTD _l	6	9	14	21	30
gMTD _l	6	13	20	27	34
HMM	2	7	14	23	34
MTD _l	6	7	8	9	10
MarkovChain	6	18	54	162	486
Model	$l(k = 4)$				
	1	2	3	4	5
HM-gMTD _l	12	27	44	63	84
HM-MTD _l	12	15	20	27	36
gMTD _l	12	25	38	51	64
HMM	3	9	17	27	39
MTD _l	12	13	14	15	16
MarkovChain	12	48	192	768	3072

PROOF. First of all consider that

$$p(X_t = j) = \sum_h p(X_t = j | Z_t = h)p(Z_t = h), \quad t > l.$$

For any h and j , $p(X_t = j | Z_t = h) \rightarrow \pi_j$ as t goes to infinity. Then the result obviously holds because $\sum_h p(Z_t = h) = 1$. \square

3. LIKELIHOOD INFERENCE

In the following, we outline an EM algorithm (Dempster *et al.*, 1977) which may be used for the maximum likelihood estimation of the parameters of the HM-gMTD, model and then of each nested model listed in Table 1. The algorithm is formulated for the case in which we observe a single time series x_1, \dots, x_T , but it can be easily adapted to the case of panel data in which we observe short sequences of observations for a sample of n statistical units.

When we observe a single time series, the *log-likelihood* to be maximized is

$$\ell(\theta) = \log p(x_{l+1}, \dots, x_T | x_1, \dots, x_l) + \log p(x_1, \dots, x_l),$$

where θ is the vector of all model parameters and the first component at rhs may be computed by the recursion illustrated in Section 2. The second component at rhs, i.e. $\log p(x_1, \dots, x_l)$, is not of direct interest and it is treated as a constant term.

The EM algorithm is based on the maximization of a suitable expectation of the log-likelihood of the *complete data* which are represented by z_{l+1}, \dots, z_T further to the observations x_1, \dots, x_T . This log-likelihood has expression

$$\begin{aligned} \ell^*(\theta) &= \log p(x_{l+1}, \dots, x_T, z_{l+1}, \dots, z_T | x_1, \dots, x_l) \\ &= \sum_{t>l} \sum_h d_{t,h} \log(\pi_{x_t-h, x_t}^{(h)}) + \sum_h d_{l+1,h} \log(\rho_h) + \sum_{h_1} \sum_{h_2} \log(\phi_{h_1, h_2}) \sum_{t>l+1} d_{t-1, h_1} d_{t, h_2}, \end{aligned}$$

where $d_{t,h} = I(z_t = h)$ is a dummy variable equal to 1 if the latent process is in state h at occasion t and to 0 otherwise. Consequently, $\sum_{t>l+1} d_{t-1, h_1} d_{t, h_2}$ is equal to the number of transitions from state h_1 to state h_2 .

At the E-step, the algorithm computes the conditional expected value of each $d_{t,h}$ and $d_{t-1, h_1} d_{t, h_2}$ given the observed data. Note that

$$\begin{aligned} E(d_{t,h} | x_1, \dots, x_T) &= p(Z_t = h | x_1, \dots, x_T), \\ E(d_{t-1, h_1} d_{t, h_2} | x_1, \dots, x_T) &= p(Z_{t-1} = h_1, Z_t = h_2 | x_1, \dots, x_T); \end{aligned}$$

these *posterior probabilities* may be obtained by recursions taken from the HMM literature which we describe below. See MacDonald and Zucchini (1997) for a general description and Bartolucci (2006) for an efficient implementation based on the matrix notation. Also see Bartolucci and Besag (2002) for alternative recursions.

For $t > l$, let

$$\begin{aligned} \alpha_t(h) &= p(x_{l+1}, \dots, x_t, Z_t = h | x_1, \dots, x_l), \\ \beta_t(h) &= p(x_{t+1}, \dots, x_T | x_1, \dots, x_t, Z_t = h), \end{aligned}$$

so that $p(x_{l+1}, \dots, x_T | x_1, \dots, x_l) = \sum_h \alpha_T(h)$. The first quantity corresponds to (5) when $t = l + 1$ and, because of (7), may be recursively computed as

$$\alpha_t(h) = \sum_m \alpha_{t-1}(m) \phi_{m,h} \pi_{x_{t-h}, x_t}^{(h)}$$

for $t > l + 1$. Similarly, $\beta_t(h)$ may be computed by the backward recursion

$$\beta_t(h) = \sum_m \beta_{t+1}(m) \phi_{h,m} \pi_{x_{t+1-m}, x_{t+1}}^{(m)},$$

initialized with $\beta_T(h) = 1$ for $h = 1, \dots, l$. It is straightforward to see that

$$p(Z_t = h | x_1, \dots, x_T) = \frac{\alpha_t(h) \beta_t(h)}{p(x_{l+1}, \dots, x_T | x_1, \dots, x_l)},$$

and

$$p(Z_{t-1} = h_1, Z_t = h_2 | x_1, \dots, x_T) = \frac{\alpha_{t-1}(h_1) \phi_{h_1, h_2} \pi_{x_{t-h_2}, x_t}^{(h_2)} \beta_t(h_2)}{p(x_{l+1}, \dots, x_T | x_1, \dots, x_l)}.$$

At the M-step, the algorithm updates the parameter estimates by maximizing the expected value of $\ell^*(\theta)$, obtained by substituting to each $d_{t,h}$ and $d_{t-1, h_1} d_{t, h_2}$ the expected values computed as above. Under the largest model, HM-gMTD_l, explicit solutions are available, i.e.

$$\pi_{j_1, j_2}^{(h)} = \frac{\sum_{t>l} p(Z_t = h | x_1, \dots, x_T) I(x_{t-h} = j_1, x_t = j_2)}{\sum_{t>l} p(Z_t = h | x_1, \dots, x_T) I(x_{t-h} = j_1)}, \quad h = 1, \dots, l, \quad j_1, j_2 = 1, \dots, k, \quad (8)$$

for what concerns the manifest transition probabilities,

$$\rho_h = p(Z_{l+1} = h | x_1, \dots, x_T), \quad h = 1, \dots, l,$$

for the initial probabilities of the latent process, and

$$\phi_{h_1, h_2} = \frac{\sum_{t>l+1} p(Z_{t-1} = h_1, Z_t = h_2 | x_1, \dots, x_T)}{\sum_{t>l+1} p(Z_{t-1} = h_1 | x_1, \dots, x_T)}, \quad h_1, h_2 = 1, \dots, l,$$

for its transition probabilities.

Note that in case the HM-MTD_l model is assumed, the manifest transition probabilities are updated as

$$\pi_{j_1, j_2} = \frac{\sum_h \sum_{t>l} p(Z_t = h | x_1, \dots, x_T) I(x_{t-h} = j_1, x_t = j_2)}{\sum_h \sum_{t>l} p(Z_t = h | x_1, \dots, x_T) I(x_{t-h} = j_1)}, \quad j_1, j_2 = 1, \dots, k,$$

instead of by (8). Moreover, when the MTD_l model is assumed, the initial probabilities of the latent process are updated as

$$\rho_h = \frac{\sum_{t>l} p(Z_t = h | x_1, \dots, x_T)}{T - l}, \quad h = 1, \dots, l,$$

and we let $\phi_{h_1, h_2} = \rho_{h_2}$, $h_1, h_2 = 1, \dots, l$, since in this case the latent transition matrix is assumed to have each row equal to ρ_1, \dots, ρ_l . In case the HMM is assumed, the algorithm reduces to a standard EM algorithm to fit this model. Finally, under more elaborated constraints on the latent transition matrix, e.g. this matrix is assumed tridiagonal, updating the estimates of its elements requires more sophisticated rules which may be taken from Bartolucci (2006).

The EM algorithm described above is guaranteed to lead to a local maximum of the likelihood. To increase the chance of catching the global maximum, common strategies involve multistart and/or initialization from opportune starting values (for instance obtained from maximum likelihood estimation of models nested in the assumed one).

Once the maximum likelihood estimate has been obtained, we can predict the most likely sequence of latent states through a Viterbi algorithm (Viterbi, 1967) along the same lines as Bartolucci and Farcomeni (2009). We also refer to Bartolucci and Farcomeni (2009) for a method to compute the standard errors for the parameter estimates which is based on the numerical derivative of the score vector; the latter is directly obtained from the EM algorithm. These standard errors may be used to construct confidence intervals and testing statistical hypotheses on the parameters. A more general way to test such hypotheses is by the likelihood ratio statistic. Note, however, that the null asymptotic distribution of this statistic is not ensured to be a standard chi-squared distribution when the hypothesis of interest is that certain elements of the latent transition matrix are equal to 0. This happens, for instance, when we assume that this matrix is tridiagonal. In this case, the asymptotic distribution is of chi-bar-squared type (Bartolucci, 2006), i.e. a mixture of chi-squared distributions with suitable weights; for a general description of this distribution see Shapiro (1988).

Finally, a fundamental point concerns model choice with respect to both the order l of the lag and possible constraints on the parameters; see Table 1. In the MTD literature, the Bayesian Information Criterion (BIC) of Schwarz (1978) seems to be preferred among the available selection criteria. This criterion is based on the minimization of the index $BIC = -2\ell(\hat{\theta}) + g \log(T - l)$, where $\hat{\theta}$ is the vector of parameter estimates obtained at convergence of the EM algorithm and g is the number of non-redundant parameters. Modifications of the penalization terms are required with panel data so as to take into account the sample size also. In the HMM literature, BIC is known to perform well in choosing the order of the model even if its theoretical properties are not so clear; see Celeux and Durand (2008) and the references therein. These reasons lead us to suggest BIC as an adequate selection criterion for the proposed model, as an alternative to other criteria such as the Akaike Information Criterion (Akaike, 1973).

4. AN EXAMPLE

For illustration we use the Stock Market Data analyzed in Dias *et al.* (2008). For the markets of Argentina, Brazil, Canada, Chile, Mexico, Peru and United States we drew from the Datastream database the daily closing price from July 4, 1994, to September 27, 2007. All series are denominated in US dollars, and for each of them we model the daily rates of returns

$$x_t = \log(P_t/P_{t-1}), \quad t = 1, \dots, 3454,$$

where P_t denotes the closing price on day t .

We fit a HM-MTD_l model to the time series of each country separately. Being in the continuous case, the model is based on a kernel transition function $\tau(\cdot | \cdot)$ such that, given $X_{t-h} = x_{t-h}$ and $Z_t = h$, X_t follows a Normal distribution with mean x_{t-h} and standard deviation σ_h which is specific to each time series. The order l has been selected on the basis of the Bayesian Information Criterion outlined in Section 3. The results are displayed in Table 3 for each country.

The order l for each chain, with the only exception of Peru, shows that the memory of the process very likely lasts for an entire week (6–8 days).

Table 3. Results from fitting the HM-MTD_l model to the time series of the daily rates of return of Argentina, Brazil, Canada, Chile, Mexico, Peru and United States referred to the period from July 4, 1994, to September 27, 2007

Stock market	<i>l</i>	log-likelihood	#parameters	BIC
AR	7	-5093.26	55	10641.27
BR	6	-5290.49	41	10919.98
CN	7	-5093.07	55	10640.90
CL	6	-5243.97	41	10824.93
MX	7	-5122.85	55	10700.46
PE	4	-5320.54	19	10798.18
US	8	-5094.95	71	10776.94

As an illustration, we show in detail the results for the US stock market. The eight Normal transition kernels have standard deviations $\sigma_h, h = 1, \dots, 8$, equal to

$$0.60, 0.59, 0.65, 0.47, 0.44, 2.67, 0.78, 0.48.$$

Note the very high standard deviation of the sixth latent state. For the latent transition matrix Φ with elements $\phi_{h_1, h_2}, h_1, h_2 = 1, \dots, 8$, we have the following estimate, where the largest element in each row is in bold

$$\hat{\Phi} = \begin{pmatrix} 0.00 & \mathbf{0.33} & 0.04 & 0.09 & 0.20 & 0.00 & 0.13 & 0.23 \\ 0.00 & 0.19 & 0.06 & \mathbf{0.30} & 0.24 & 0.01 & 0.19 & 0.01 \\ 0.00 & 0.00 & 0.34 & 0.00 & 0.00 & 0.02 & \mathbf{0.38} & 0.26 \\ 0.00 & 0.11 & 0.00 & 0.23 & \mathbf{0.37} & 0.00 & 0.13 & 0.15 \\ \mathbf{0.39} & 0.11 & 0.03 & 0.21 & 0.04 & 0.03 & 0.00 & 0.20 \\ 0.00 & 0.00 & 0.01 & 0.00 & 0.01 & \mathbf{0.61} & 0.34 & 0.03 \\ 0.08 & 0.12 & \mathbf{0.31} & 0.09 & 0.10 & 0.06 & 0.10 & 0.14 \\ \mathbf{0.34} & 0.29 & 0.04 & 0.05 & 0.15 & 0.00 & 0.00 & 0.13 \end{pmatrix}.$$

The corresponding stationary distribution is

$$0.12, 0.16, 0.10, 0.14, 0.16, 0.04, 0.13, 0.15.$$

Moreover, the estimated initial probabilities are all equal to 0, with the exception of the fifth latent state that has probability equal to 1.

So as to illustrate the features of the model also with discrete response variables, we discretized the US log-returns in five categories using the following cut-points: -0.02, -0.01, 0.01, 0.02. On the resulting data, the HM-MTD_l model is fit for an increasing number of latent states until BIC does not decrease. The results are in Table 4.

When $l = 4$, the manifest transition matrix Π with elements $\pi_{j_1, j_2}, j_1, j_2 = 1, \dots, 5$, is estimated as:

$$\hat{\Pi} = \begin{pmatrix} 0.25 & 0.19 & 0.05 & 0.14 & \mathbf{0.37} \\ 0.08 & 0.21 & \mathbf{0.38} & 0.23 & 0.09 \\ 0.01 & 0.07 & \mathbf{0.84} & 0.07 & 0.00 \\ 0.03 & 0.08 & \mathbf{0.71} & 0.15 & 0.03 \\ 0.08 & 0.21 & \mathbf{0.46} & 0.14 & 0.10 \end{pmatrix},$$

which corresponds to a strong persistency ($\pi_{33} = 0.84$) in the third state of returns around 0. The estimates for the first row are particularly interesting, indicating that either extremely negative returns in the class $(-\infty, -0.02]$ or extremely positive returns in the class $(0.02, \infty)$ are expected when the process relies on a time occasion with return in the extremely negative class $(-\infty, -0.02]$. On the other hand, when the process relies on a time occasion with return in the extremely positive class, a situation of much lower volatility is likely to follow, with most likely return in the class $(-0.01, 0.01]$.

The estimated transition probabilities for the latent process Φ are:

$$\hat{\Phi} = \begin{pmatrix} 0.00 & 0.27 & 0.36 & \mathbf{0.37} \\ 0.18 & 0.40 & 0.00 & \mathbf{0.42} \\ 0.00 & \mathbf{0.68} & 0.32 & 0.00 \\ \mathbf{0.66} & 0.02 & 0.17 & 0.15 \end{pmatrix},$$

Table 4. Summary of HM-MTD model fit for different choices on the number of latent states for the US stock market data divided in five classes

<i>l</i>	log-likelihood	#parameters	BIC
1	-29838	20	61305
2	-29326	23	60525
3	-29047	28	60374
4	-28729	35	60308
5	-28538	44	60660

Boldface are the data referred to the model with the smallest BIC.

while the estimates of the initial probabilities place all mass on the fourth latent state.

For the chosen model HM-MTD₄, we also tested the hypothesis of independence of the latent transitions, which would result in an ordinary MTD₄. This hypothesis holds when all the rows of Φ are equal each other. The likelihood ratio test statistic for this hypothesis is equal to 266, which leads us to reject it.

Acknowledgements

We acknowledge the financial support of the 'Einaudi Institute for Economics and Finance', Rome – IT. Francesco Bartolucci also acknowledges the financial support of the Italian Government (PRIN – 2007).

REFERENCES

- Akaike, H. (1973) Information theory and an extension of the maximum likelihood principle. In *Second International Symposium on Information Theory* (eds B. N. Petrov and F. Csaki). Budapest: Akademiai Kiado, pp. 267–81.
- Bartolucci, F. (2006) Likelihood inference for a class of latent Markov models under linear hypotheses on the transition probabilities. *Journal of the Royal Statistical Society, Series B* **68**, 155–78.
- Bartolucci, F. and Besag, J. (2002) A recursive algorithm for Markov random fields. *Biometrika* **89**, 724–30.
- Bartolucci, F. and Farcomeni, A. (2009) A multivariate extension of the dynamic logit model for longitudinal data based on a latent Markov heterogeneity structure. *Journal of the American Statistical Association* **104**, 816–31.
- Berchtold, A. (1998) *Chaînes de Markov et Modèles de Transition: Applications aux Sciences Sociales*. Paris: Hermes.
- Berchtold, A. and Raftery, A. E. (2002) The mixture transition distribution model for high-order Markov chains and non-Gaussian time series. *Statistical Science* **17**, 328–56.
- Celeux, G. and Durand, J.-B. (2008) Selecting hidden Markov model state number with cross-validated likelihood. *Computational Statistics* **23**, 541–64.
- Dempster, A. P., Laird, N. M. and Rubin, D. B. (1977) Maximum likelihood from incomplete data via the EM algorithm (with discussion). *Journal of the Royal Statistical Society, Series B* **39**, 1–38.
- Dias, J. G., Vermunt, J. K., and Ramos, S. (2008) Heterogeneous hidden Markov models. *Proceedings of Compstat 2008*.
- Le, N. D., Martin, D. R. and Raftery, A. E. (1996) Modeling flat stretches, bursts, and outliers in time series using mixture transition distribution models. *Journal of the American Statistical Association* **91**, 1504–1515.
- MacDonald, I. L. and Zucchini, W. (1997) *Hidden Markov and Other Models for Discrete Valued Time Series*. London: Chapman and Hall.
- Raftery, A. E. (1985a) A model for high-order Markov chains. *Journal of the Royal Statistical Society, Series B* **47**, 528–39.
- Raftery, A. E. (1985b) A new model for discrete-valued time series: Autocorrelations and extensions. *Rassegna di Metodi Statistici ed Applicazioni* **3**, 149–162.
- Raftery, A. E. and Tavaré, S. (1994) Estimation and modelling repeated patterns in high order Markov chains with the mixture transition distribution model. *Applied Statistics* **43**, 179–199.
- Schwarz, G. (1978) Estimating the dimension of a model. *Annals of Statistics* **6**, 461–64.
- Shapiro, A. (1988) Towards a unified theory of inequality constrained testing in multivariate analysis. *International Statistical Review* **56**, 49–62.
- Viterbi, A. J. (1967) Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. *Transactions on Information Theory* **13**, 260–69.