



————— Technical Report No. 161 —————

Dirichlet Mixtures of Bayesian Linear Gaussian State-Space Models: a Variational Approach

Silvia Chiappa¹, David Barber²

————— June 2007 —————

¹ Department Schölkopf, email: silvia.chiappa@tuebingen.mpg.de; ² University College London, Department of Computer Science, email: d.barber@cs.ucl.ac.uk

Dirichlet Mixtures of Bayesian Linear Gaussian State-Space Models: a Variational Approach

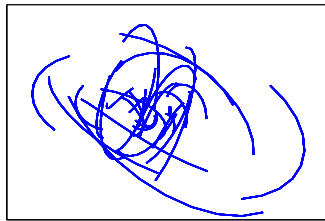
Silvia Chiappa, David Barber

Abstract. We describe two related models to cluster multidimensional time-series under the assumption of an underlying linear Gaussian dynamical process. In the first model, times-series are assigned to the same cluster when they show global similarity in their dynamics, while in the second model times-series are assigned to the same cluster when they show simultaneous similarity. Both models are based on Dirichlet Mixtures of Bayesian Linear Gaussian State-Space models in order to (semi) automatically determine an appropriate number of components in the mixture, and to additionally bias the components to a parsimonious parameterization. The resulting models are formally intractable and to deal with this we describe a deterministic approximation based on a novel implementation of Variational Bayes.

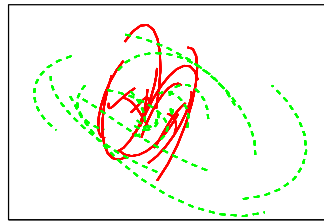
1 Introduction

Clustering is a large topic in machine learning and related areas, and the aim of this report is to provide an additional methodology which may be better suited to applications in which a dynamical system is believed to be underlying the data. This situation is common in time-series based on natural phenomena, since the equations believed to be describing the physical world can often be modeled as Markovian dynamics on an underlying state-system.

Perhaps the most straightforward approach to perform clustering of a set of time-series is to consider each time-series as a ‘static’ vector, thereby transforming the time-series clustering problem into a more standard clustering-of-vectors problem. For clustering vectors, any of a number of methods may be applied, ranging from K-Means [1], to more recent methods based on probabilistic mixture models [2, 3, 4, 5]. However, for long or high dimensional time-series, this approach becomes computationally problematic, and therefore features of the signal are used instead, typically extracted from short windows of the time-series. Each time-series is then represented by either a single, or set of feature vectors, which may be then clustered by any standard static clustering technique[6]. Whilst in feature extraction methods the length and the dimensionality of the time-series no longer constitutes a problem, it is not always clear what the appropriate feature extraction method should be, nor how one should deal with missing data. For example, the times-series generated by the same dynamical system with different initial conditions can look highly dissimilar (see Fig. 1a). In such cases, it is not obvious which features should be used to perform clustering. On the other hand, a method which can explicitly model the dynamics of the time-series would be able to perform clustering without the need of a preliminary feature extraction.



(a) Unclustered Trajectories



(b) Clustered Trajectories

Figure 1: (a) Thirty trajectories length $T = 10$ resulting from the dynamics of two different LGSSMs, both with hidden dimension $H = 4$. Plotted are the points $([v_t]_1, [v_t]_2)$. (b) The trajectories as labeled by our algorithm – all labels are consistent with the generating mechanism.

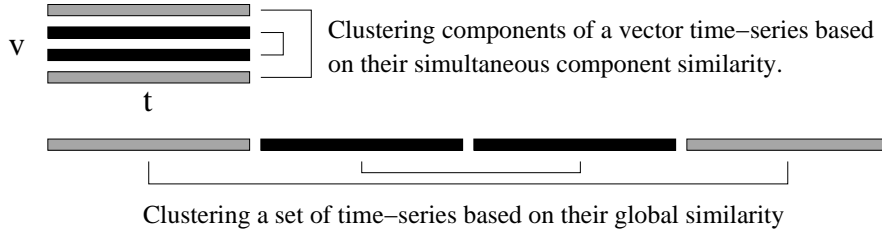


Figure 2: (Top) Time-series clustering based on simultaneous similarity versus (Bottom) time-series clustering based on global similarity.

Our method is therefore motivated by the desiderata to perform clustering by explicitly modeling the dynamics of the time-series. Furthermore, we are interested in the case in which the number of clusters is not known in advance, and therefore in a model which can automatically determine an appropriate number of clusters. To prevent overfitting, we would also like to encourage each cluster to be described by a parsimonious parameterization.

Recently, there has been considerable interest in Dirichlet Process Mixture Models [2, 3, 4, 5] which provide a (semi) automatic way to determine an appropriate number of clusters. The Dirichlet mixture results in a Polya distribution on the cluster assignments, enabling an appropriate number of clusters to be found by the model. In a continued theme of development of this class of techniques, our time-series model will be based on a Dirichlet Mixture of Bayesian Linear Gaussian State-Space models (DMBLGSSMs). The Bayesian approach places a prior on the parameters of each mixture component, encouraging them to have the smallest parameterization consistent with the data.

We will give two main variants of the DMBLGSSM, in order to perform clustering based on either *simultaneous* or *global* similarity of the time-series dynamics. The distinction between these two approaches will be spelled out mathematically in Section 4, whilst an informal sketch is presented in Fig. 2:

Simultaneous Similarity Our simultaneous similarity clustering approach assigns two time-series to the same cluster if they are derived from the *same realization* of a dynamical process.

Global Similarity The global similarity method will assign two time-series to the same cluster if they are generated by *different realizations* of the same dynamical process.

The resulting Bayesian time-series clustering model is formally computationally intractable, and therefore approximations need to be considered. To the best of our knowledge, whilst sampling methods have been applied in a similar more constrained context [7], a Variational Bayesian treatment of this class of models is new, as is our application to clustering based on simultaneous dynamical similarity. Variational Bayes is a deterministic approximation scheme which has the potential advantage of speed over sampling techniques.

We will first describe the BLGSSM in Section 2 and the general procedure of Dirichlet Process mixture models in Section 3. We will then marry the two in Section 4 to form the clustering methods based on simultaneous and global dynamical similarity. In Section 5 we will give an illustrative demonstration of the performance of the two approaches.

2 Bayesian Linear Gaussian State-Space Models

In a Linear Gaussian State-Space Model (LGSSM)¹ [8, 9, 10], a sequence of observations $v_{1:T} \equiv v_1, \dots, v_T$ is generated from an underlying dynamical system on $h_{1:T}$ according to:

$$v_t = B h_t + \eta_t^v, \quad \eta_t^v \sim \mathcal{N}(\mathbf{0}_V, \Sigma_V), \quad h_t = A h_{t-1} + \eta_t^h, \quad \eta_t^h \sim \mathcal{N}(\mathbf{0}_H, \Sigma_H), \quad h_1 \sim \mathcal{N}(\mu, \Sigma),$$

where $\mathcal{N}(m, S)$ denotes a Gaussian with mean m and covariance S , and $\mathbf{0}_X$ denotes an X -dimensional zero vector. The observation v_t has dimension V and the hidden state h_t has dimension H . Probabilistically, the LGSSM is defined by:

$$p(v_{1:T}, h_{1:T} | \Theta) = p(v_1 | h_1) p(h_1) \prod_{t=2}^T p(v_t | h_t) p(h_t | h_{t-1}),$$

¹Also called Kalman Filter/Smother, Linear Dynamical System.

with

$$p(v_t|h_t) = \mathcal{N}(Bh_t, \Sigma_V), \quad p(h_t|h_{t-1}) = \mathcal{N}(Ah_{t-1}, \Sigma_H), \quad p(h_1) = \mathcal{N}(\mu, \Sigma)$$

and where $\Theta = \{A, B, \Sigma_H, \Sigma_V, \mu, \Sigma\}$ denotes the model parameters. Thanks to the simple structure of the model, most quantities of interest, such as the posterior density $p(h_t|v_{1:T}, \Theta)$ and likelihood

$$p(v_{1:T}|\Theta) = \int_{h_{1:T}} p(v_{1:T}, h_{1:T}|\Theta) \quad (1)$$

can be computed efficiently in $O(T)$ operations [11].

In a Bayesian treatment of the model, a parameter prior $p(\Theta|\hat{\Theta})$ is defined, where $\hat{\Theta}$ are the associated hyperparameters, resulting in the marginal likelihood

$$p(v_{1:T}|\hat{\Theta}) = \int_{\Theta, h_{1:T}} p(v_{1:T}, h_{1:T}|\Theta)p(\Theta|\hat{\Theta}). \quad (2)$$

In a full Bayesian treatment we would define additional prior distributions over the hyperparameters $\hat{\Theta}$. Here we take instead the ML-II (‘evidence’) framework, in which the optimal set of hyperparameters is found by maximizing $p(v_{1:T}|\hat{\Theta})$ with respect to $\hat{\Theta}$ [12, 13, 14].

Whilst the integral required to compute the likelihood (Eq. (1)) is tractable, the result of this integral couples the parameters Θ in the integrand of the marginal likelihood (Eq. (2)). An exact implementation of the Bayesian LGSSM is then formally intractable. Recently Variational Bayes (VB) has been applied to this model as a route to a computationally efficient approximate implementation [12, 13, 14, 15, 16, 17].

The most challenging part of implementing the VB method is performing inference over $h_{1:T}$. Some authors [12, 16, 17] have developed their own specialized routines, based on Belief Propagation, since standard LGSSM inference routines appear, at first sight, not to be applicable. However, in [15] is shown how the VB treatment of the LGSSM can be implemented using any standard LGSSM inference routine, including those specifically addressed to improve numerical stability [9, 18, 19]. This approach can also be used for the Dirichlet Mixture case, as we will see below.

For the parameter priors, here we define Gaussians on the elements of A and on the columns of B :

$$p(A|\alpha, \Sigma_H^{-1}) = \prod_{i,j=1}^H \frac{\alpha_{ij}^{1/2}}{\sqrt{2\pi} [\Sigma_H]_{ii}} e^{-\frac{\alpha_{ij}}{2} [\Sigma_H^{-1}]_{ii} (A_{ij} - \hat{A}_{ij})^2},$$

$$p(B|\beta, \Sigma_V^{-1}) = \prod_{j=1}^H \frac{\beta_j^{V/2}}{\sqrt{|2\pi\Sigma_V|}} e^{-\frac{\beta_j}{2} (B_j - \hat{B}_j)^\top \Sigma_V^{-1} (B_j - \hat{B}_j)},$$

where \hat{A} and \hat{B} are hyperparameters defining our preferred values for the transition and emission matrices. The dependency of the priors on Σ_H and Σ_V renders the VB implementation feasible. The choice $\hat{A} \equiv 0, \hat{B} \equiv 0$ creates an Automatic Relevance Determination (ARD) bias towards a simple dynamics and eliminates unnecessary parameters of the model [20]. This form of prior, in which individual elements of the transition A and columns on the emission B are biased to be small, is appropriate for the global similarity clustering of Section 4.1². For the simultaneous similarity clustering of Section 4.2 we modify the emission prior by biasing elements of B to be small.

The conjugate priors for general inverse covariances Σ_H^{-1} and Σ_V^{-1} are Wishart distributions. In the simpler case of diagonal covariances $\Sigma_H^{-1} = dg(\tau)$ ³ and $\Sigma_V^{-1} = dg(\rho)$ these become Gamma distributions (see Appendix A.1). The hyperparameters to optimize are then $\alpha, \beta, \mu, \Sigma$ ⁴ and the parameters of Gamma or Wishart distributions.

²A prior that prefer simultaneously the i^{th} row and column of A to be small would be preferable. However, it is not straightforward to make the VB feasible in this case.

³The notation $dg(\tau)$ indicates diagonal matrix with the elements of the vector τ on the main diagonal.

⁴We do not put priors on μ and Σ , which will be formally considered as hyperparameters.

3 Dirichlet Mixture Models

Given a set of observations $v^{1:N} \equiv v^1, \dots, v^N$, the clustering task is to assign each observation v^n to one of a finite set of cluster labels $k = 1, \dots, K$. To do so, we introduce a cluster indicator variable $z^n \in \{1, \dots, K\}$ for each observation. The joint distribution of all observations is then given by

$$p(v^{1:N} | \Theta^{1:K}) = \sum_{z^{1:N}} \left\{ \prod_n p(v^n | z^n, \Theta^{1:K}) \right\} p(z^{1:N}), \quad (3)$$

where $p(v^n | z^n = k, \Theta^{1:K}) \equiv p(v^n | \Theta^k)$ denotes that the parameters of cluster k are used to determine the probability of observation v^n . To model the joint cluster allocations we define

$$p(z^{1:N}) = \int_{\pi} \left\{ \prod_n p(z^n | \pi) \right\} p(\pi), \quad (4)$$

where $p(z^n = k | \pi) \equiv \pi_k$ is a multinomial distribution. Using the Dirichlet prior $p(\pi) \propto \prod_{k=1}^K \pi_k^{\gamma/K-1}$, the integral in Eq. (4) gives rise to the Polya distribution (see Appendix A.3):

$$p(z^{1:N}) = \frac{\Gamma(\gamma)}{\Gamma(N + \gamma)} \prod_{k=1}^K \frac{\Gamma(N_k + \gamma/K)}{\Gamma(\gamma/K)}, \quad (5)$$

where $N_k \equiv \sum_{n=1}^N I[z^n = k]$ counts the number of times that state k occurs in the indicators⁵. In the limit of infinite K , the prior expected number of clusters for a set of N sequences is [21]

$$\sum_{n=1}^N \frac{\gamma}{\gamma + n - 1} \approx \gamma \log \left(\frac{N}{\gamma} + 1 \right),$$

which remains finite for finite γ . In our experiments γ will typically be treated as a hyperparameter and optimized with respect to the marginal likelihood.

In our work, we will consider K to be finite. This is in contrast to Dirichlet *Process* Mixture Models [3, 4, 5], in which the $K \rightarrow \infty$ limit is formally taken. This can be achieved, for example, by writing down a sampling algorithm for the finite dimensional case, and then taking the limit $K \rightarrow \infty$. If the sampler is initialized with a small number of clusters, the sampling algorithm generates new clusters until sufficiently many are present to explain the data well. In practice, since only a finite number of mixture components is effectively used, we prefer the finite K case. An advantage of this is that we retain an explicit expression for the marginal likelihood which is then amenable to fast deterministic approximation schemes.

Computing the resulting model likelihood in Eq. (3) is intractable, and a useful deterministic approximation is to use a lower bound based on the ‘collapsed’⁶ variational KL divergence $\text{KL} \left(\prod_{n=1}^N q(z^n) || p(z^{1:N} | v^{1:N}) \right)$ [22]⁷.

4 Dirichlet Mixture of Bayesian LGSSMs

Our approach to clustering is to form a Dirichlet mixture of Bayesian LGSSMs (DMBLGSSMs). This has the advantage of determining the number of clusters, where each cluster may also be biased towards a dynamical system of a preferred form. We will start by describing a model for time-series clustering based on global similarity. We will then introduce a clustering method based on simultaneous similarity as a modification of this first model.

4.1 Clustering based on Global Similarity

The graphical representation of Dirichlet Mixture of Bayesian LGSSMs is given in Fig. 3. The likelihood term $p(v^n | z^n, \Theta^{1:K})$ for each temporal sequence $v^n \equiv v_{1,\dots,T}^n$ in Eq. (3) is defined by the likelihood of the LGSSM (Eq.

⁵ $I[a = b] = 1$ if $a = b$ and 0 otherwise.

⁶The ‘uncollapsed’ joint approximation of $p(z^{1:N}, \pi | v^{1:N})$ is seductive since it is more straightforward to compute the KL divergence in the joint $(z^{1:N}, \pi)$ space under the factorized assumption $q(z^{1:N})q(\pi)$. However, the strong implicit coupling between $z^{1:N}$ and π means that the factorized approximation is often insufficiently accurate to be practically useful [22].

⁷Here and in the rest of the report (if this does not cause confusion), we omit the conditioning on the observations for q .

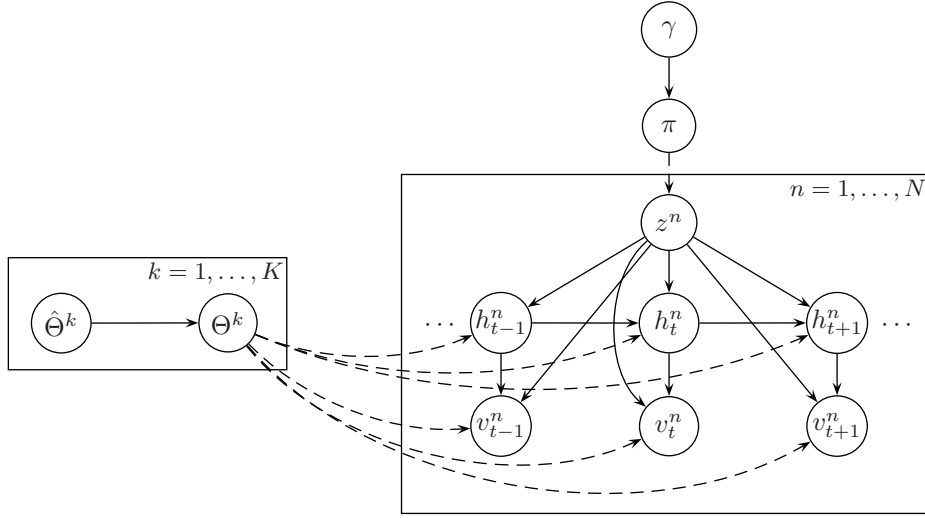


Figure 3: Graphical representation of the Dirichlet Mixture of Bayesian LGSSMs for performing clustering based on global similarity.

(1)). The approach we take to extend this to a Dirichlet mixture of the *Bayesian* LGSSMs is to introduce a distribution q , for which we assume that the following factorizations hold:

$$\begin{aligned} q(h_{1:T}^{1:N} | z^{1:N}, \Theta^{1:K}) &\equiv q(h_{1:T}^{1:N} | z^{1:N}) \\ q(z^{1:N}, \Theta^{1:K}) &\equiv q(z^{1:N}) q(\Theta^{1:K}) \\ q(z^{1:N}) &\equiv \prod_{n=1}^N q(z^n). \end{aligned}$$

We then consider the variational approximation⁸:

$$p(z^{1:N}, h_{1:T}^{1:N}, \Theta^{1:K} | v_{1:T}^{1:N}, \hat{\Theta}^{1:K}) \approx \left\{ \prod_{n=1}^N q(z^n) q(h_{1:T}^n | z^n) \right\} \prod_{k=1}^K q(\Theta^k).$$

Taking the KL divergence between the right and left hand sides of the above then gives the following lower bound on the log-likelihood, $\log p(v_{1:T}^{1:N} | \hat{\Theta}^{1:K}) \geq \mathcal{F}(\hat{\Theta}^{1:K}, q)$, of the Dirichlet mixture of Bayesian LGSSMs:

$$\begin{aligned} \mathcal{F} &\equiv \sum_{k=1}^K H_q(\Theta^k) + \sum_{n=1}^N \sum_{k=1}^K q(z^n = k) H_q(h_{1:T}^n | z^n = k) + \sum_{n=1}^N H_q(z^n) + \sum_{k=1}^K \langle \log p(\Theta^k | \hat{\Theta}^k) \rangle_{q(\Theta^k)} \\ &\quad + \langle \log p(z^{1:N}) \rangle_{\prod_n q(z^n)} + \sum_{n=1}^N \sum_{k=1}^K q(z^n = k) \langle \log p(v_{1:T}^n, h_{1:T}^n | \Theta^k) \rangle_{q(\Theta^k) q(h_{1:T}^n | z^n = k)}, \end{aligned} \quad (6)$$

where $H_q(x)$ denotes the entropy of the distribution $q(x)$, and $\langle \cdot \rangle_q$ denotes expectation with respect to q . Variational Bayes then proceeds by iteratively maximizing the lower bound with respect to the q distributions for fixed hyperparameters $\hat{\Theta}$ and vice-versa until no further improvement is found. The resulting updates for q are given by:

$$\begin{aligned} q(\Theta^k) &\propto p(\Theta^k | \hat{\Theta}^k) e^{\sum_{n=1}^N q(z^n = k) \langle \log p(v_{1:T}^n, h_{1:T}^n | \Theta^k) \rangle_{q(h_{1:T}^n | z^n = k)}} \\ q(z^n = k) &\propto e^{H_q(h_{1:T}^n | z^n = k) + \langle \log p(z^n = k | z^{1:n}) \rangle_{\prod_{m \neq n} q(z^m)} + \langle \log p(v_{1:T}^n, h_{1:T}^n | \Theta^k) \rangle_{q(h_{1:T}^n | z^n = k) q(\Theta^k)}} \\ q(h_{1:T}^n | z^n = k) &\propto e^{\langle \log p(v_{1:T}^n, h_{1:T}^n | \Theta^k) \rangle_{q(\Theta^k)}} \end{aligned}$$

⁸From the assumptions on p and q , it follows that the q that maximizes the lower bound \mathcal{F} (Eq. (6)) satisfies $q(h_{1:T}^{1:N} | z^{1:N}) = \prod_{n=1}^N q(h_{1:T}^n | z^n)$ and $q(\Theta^{1:K}) = \prod_{k=1}^K q(\Theta^k)$ (see Appendix B.1).

where z^{-n} indicates all indicator variables except for z^n . We will discuss each specific update below. Full details are given in Appendix B.2.

Missing Observations

One of the advantages using a LGSSM model for each cluster is the ease with which missing observations can be dealt with. Indeed if the multidimensional vector v_t has some missing components corresponding to unobserved or corrupted measurements, we can integrate these out in the likelihood. This can be formally introduced in our model by the replacement $B \leftarrow W_t^n B$ in the bound \mathcal{F} , where W_t^n is the identity matrix with diagonal elements corresponding to missing observation components replaced by zeros, and by replacing missing components in v_t^n by zeros. We will give the updates for this general framework.

Throughout, we assume that Σ_V is diagonal in the case in which there are missing observations. This assumption makes the formula for the mean update independent of Σ_V , which is a computational convenience.

Updates for $q(B^k, [\Sigma_V^k]^{-1})$

To simplify the notation, in the following (and elsewhere where clear from the context) we will omit the dependency of the model parameters Θ^k and hyperparameters $\hat{\Theta}^k$ on the mixture k . The choice Normal-Wishart(Gamma) prior $p(B, \Sigma_V^{-1})$ give rise to a Normal-Wishart(Gamma) approximated posterior $q(B, \Sigma_V^{-1})$:

$$q(B, \Sigma_V^{-1}) \propto e^{\sum_{n=1}^N q(z^n=k) \sum_{t=1}^T \langle \log p(v_t^n | h_t^n, B, \Sigma_V^{-1}) \rangle_{q(h_t^n | z^n=k)}} p(B | \beta, \Sigma_V^{-1}) p(\Sigma_V^{-1} | \hat{\Theta}).$$

In the following, we decompose the joint $q(B, \Sigma_V^{-1}) \equiv q(B | \Sigma_V^{-1}) q(\Sigma_V^{-1})$.

$q(B | \Sigma_V^{-1})$: In Appendix B.2.1 we show that $q(vc(B) | \Sigma_V^{-1}) = \mathcal{N}(\mu_B, \Sigma_B)$ where⁹:

$$\begin{aligned} \Sigma_B &= (I_H \otimes \Sigma_V) \left(\underbrace{\sum_{n=1}^N q(z^n=k) \sum_{t=1}^T \langle h_t^n (h_t^n)^\top \rangle_{q(h_t^n | z^n=k)} \otimes W_t^n + dg(\beta) \otimes I_V}_{H_{BM}} \right)^{-1} \\ \mu_B &= H_{BM}^{-1} vc \left(\underbrace{\sum_{n=1}^N q(z^n=k) \sum_{t=1}^T W_t^n v_t^n \langle h_t^n \rangle_{q(h_t^n | z^n=k)}^\top + \hat{B} dg(\beta)}_{N_B} \right). \end{aligned} \quad (7)$$

We remind the reader that in the case of missing observations, Σ_V in Eq. (7) is constrained to be diagonal. For the case of no missing observations, $W_t^n = I_V$, Eq. (7) simplifies, as shown in Appendix B.2.1.

$q(\Sigma_V^{-1})$: For the case in which there are no missing observations $W_t^n = I_V$, we may consider the general Wishart prior $p(\Sigma_V^{-1} | \nu_V, S_V) = \mathcal{W}(\nu_V, S_V)$, for which the updates are (see Appendix B.2.2):

$$q(\Sigma_V^{-1}) = \mathcal{W} \left(\nu_V + T \sum_{n=1}^N q(z^n=k), \left(S_V^{-1} + \sum_{n=1}^N q(z^n=k) \sum_{t=1}^T v_t^n (v_t^n)^\top - N_B H_B^{-1} N_B^\top + \hat{B} dg(\beta) \hat{B}^\top \right)^{-1} \right),$$

with

$$H_B \equiv \sum_{n=1}^N q(z^n=k) \sum_{t=1}^T \langle h_t^n (h_t^n)^\top \rangle_{q(h_t^n | z^n=k)} + dg(\beta).$$

Under the simpler diagonal constraint $\Sigma_V^{-1} = dg(\rho)$, where each diagonal element ρ_i follows a Gamma prior $\mathcal{G}(b_1^i, b_2^i)$, the optimal updates are (see Appendix B.2.2):

$$q(\rho_i) = \mathcal{G} \left(b_1^i + \frac{T}{2} \sum_{n=1}^N q(z^n=k), b_2^i + \frac{1}{2} \left(\sum_{n=1}^N q(z^n=k) \sum_{t=1}^T [v_t^n]_i^2 - [G_B]_{ii} + \sum_j \beta_j \hat{B}_{ij}^2 \right) \right), \quad (8)$$

⁹ $vc(B)$ denotes the vector formed by stacking the columns of the matrix B . This column vector formulation simplifies mathematical notations. \otimes indicates the Kronecker product and I_X is the identity matrix of dimension $X \times X$.

where

$$G_B \equiv M_B N_B^T, \quad M_B \equiv N_B H_B^{-1} \quad (9)$$

For the case of missing observations, we are restricted to the above diagonal covariance constraint, in which case we replace M_B in Eq. (9) by the $V \times H$ matrix

$$M_{BM} = mc(H_{BM}^{-1} vc(N_B))$$

where $mc(x)$ denotes reshaping the vector to form a matrix by reverse column-stacking.

Updates for $q(A^k, [\Sigma_H^k]^{-1})$

The optimal $q(A, \Sigma_H^{-1})$ is given by:

$$q(A, \Sigma_H^{-1}) \propto e^{\sum_{n=1}^N q(z^n=k) \sum_{t=2}^T \langle \log p(h_t^n | h_{t-1}^n, A, \Sigma_H^{-1}) \rangle_{q(h_{t-1:t}^n | z^n=k)}} p(A | \alpha, \Sigma_H^{-1}) p(\Sigma_H^{-1} | \hat{\Theta}).$$

As above, without loss of generality, we decompose this as $q(A | \Sigma_H^{-1}) q(\Sigma_H^{-1})$.

$q(vr(A) | \Sigma_H^{-1})$: In Appendix B.2.3 we show that $q(vr(A) | \Sigma_H^{-1}) = \mathcal{N}(\mu_A, \Sigma_A)$ where¹⁰:

$$\begin{aligned} \Sigma_A &= bdg(H_{1A}^{-1} [\Sigma_H]_{11}, \dots, H_{HA}^{-1} [\Sigma_H]_{HH}) \\ \mu_A &= vert\left([N_A]_{1'}, H_{1A}^{-1}{}^T, \dots, [N_A]_{H'}, H_{HA}^{-1}{}^T\right) \end{aligned}$$

where $[N_A]_{i'}$ indicates the i -th row of matrix N_A , defined as

$$[N_A]_{ij} = \sum_{n=1}^N q(z^n=k) \langle [h_{t-1}^n]_j [h_t^n]_i \rangle_{q(h_{t-1:t}^n | z^n=k)} + \alpha_{ij} \hat{A}_{ij}.$$

In addition, we define

$$[H_{iA}]_{jl} \equiv \sum_{n=1}^N q(z^n=k) \sum_{t=2}^T \langle [h_{t-1}^n]_j [h_t^n]_l \rangle_{q(h_{t-1:t}^n | z^n=k)} + \alpha_{ij} \delta_{jl}$$

$q(\Sigma_H^{-1})$: For $\Sigma_H^{-1} = dg(\tau)$, where each diagonal element τ_i follows a Gamma prior $\mathcal{G}(a_1^i, a_2^i)$, the updates are (see Appendix B.2.4):

$$q(\tau_i) = \mathcal{G}\left(a_1^i + \frac{T-1}{2} \sum_{n=1}^N q(z^n=k), a_2^i + \frac{1}{2} \left(\sum_{n=1}^N q(z^n=k) \sum_{t=2}^T \langle [h_t^n]_i^2 \rangle_{q(h_t^n | z^n=k)} - [G_A]_i + \sum_j \alpha_{ij} \hat{A}_{ij}^2 \right) \right),$$

where $[G_A]_i \equiv [N_A]_{i'} H_{iA}^{-1} [N_A]_{i'}^T$.

Updates for $q(z^n)$

The optimal $q(z^n=k)$ is given by:

$$q(z^n=k) \propto e^{H_q(h_{1:T}^n | z^n=k) + \langle \log p(z^n=k | z^{-n}) \rangle_{q(z^{-n})} + \langle \log p(v_{1:T}^n, h_{1:T}^n | \Theta^k) \rangle_{q(h_{1:T}^n | z^n=k)} q(\Theta^k)}$$

The first term in the exponent can be computed exactly as described in the Appendix B.5. The average $\langle \log p(z^n=k | z^{-n}) \rangle_{(z^{-n})}$ needs attention since, naively, $p(z^n=k | z^{-n})$ possesses little structure to enable the average to be tractable. Whilst the naive exponential complexity can be reduced, we employ the second order Taylor expansion approximation of [22], as shown in detail in Appendix B.2.5.

¹⁰ $vr(A)$ denotes the vector formed by stacking the rows of the matrix A . Unlike for B , the choice of a row vector formulation is more appropriate for simplifying mathematical notations. $bdg(x_1, \dots, x_n)$ indicates the block diagonal matrix with blocks x_1, \dots, x_n , while $vert(x_1, \dots, x_n)$ stands for vertically concatenating the arguments x_1, \dots, x_n .

Inference on $q(h_{1:T}^n | z^n = k)$

The optimal $q(h_{1:T}^n | z^n = k)$ is given by:

$$q(h_{1:T}^n | z^n = k) \propto e^{\langle \log p(v_{1:T}^n, h_{1:T}^n | \Theta^k) \rangle_{q(\Theta^k)}}. \quad (10)$$

This term is closely related to a standard VB approximation to a Bayesian LGSSM. Clearly the structure of $q(h_{1:T}^n | z^n = k)$ is a pairwise Markov chain, and inference algorithms such as Belief Propagation can be used [12, 17]. However, we take the approach discussed in [15], which reformulates the problem such that standard LGSSM inference routines can be applied. This both simplifies the development and can be advantageous in regimes of numerical instability. The central idea is to use the following decomposition:

$$\langle (v_t^n - W_t^n B h_t^n)^\top \Sigma_V^{-1} (v_t^n - W_t^n B h_t^n) \rangle_{q(B, \Sigma_V^{-1})} = \underbrace{(v_t^n - W_t^n \langle B \rangle h_t^n)^\top \langle \Sigma_V^{-1} \rangle (v_t^n - W_t^n \langle B \rangle h_t^n)}_{\text{mean}} + \underbrace{(h_t^n)^\top (S_B)_t^n h_t^n}_{\text{fluctuation}},$$

and similarly

$$\langle (h_t^n - A h_{t-1}^n)^\top \Sigma_H^{-1} (h_t^n - A h_{t-1}^n) \rangle_{q(A, \Sigma_H^{-1})} = \underbrace{(h_t^n - \langle A \rangle h_{t-1}^n)^\top \langle \Sigma_H^{-1} \rangle (h_t^n - \langle A \rangle h_{t-1}^n)}_{\text{mean}} + \underbrace{(h_{t-1}^n)^\top S_A h_{t-1}^n}_{\text{fluctuation}},$$

where

$$(S_B)_t^n \equiv \langle B^\top W_t^n \Sigma_V^{-1} W_t^n B \rangle - \langle B \rangle^\top W_t^n \langle \Sigma_V^{-1} \rangle W_t^n \langle B \rangle, \quad S_A \equiv \langle A^\top \Sigma_H^{-1} A \rangle - \langle A \rangle^\top \langle \Sigma_H^{-1} \rangle \langle A \rangle.$$

The analytical expressions for these covariances are given in Appendix B.3. The mean terms represent the contribution of a standard LGSSM with parameters A , B , Σ_H^{-1} and Σ_V^{-1} replaced by their average values.

The key observation is to consider the extra ‘fluctuation’ terms as having been generated from fictitious zero-valued observations. This way, we can see Eq. (10) as the posterior of a standard LGSSM for which any of the standard algorithms in the literature [11] may be applied to perform inference.

More specifically we want to represent Eq. (10) directly as the posterior distribution $\tilde{q}(h_{1:T}^n | \tilde{v}_{1:T}^n)$ of an LGSSM by augmenting v_t^n and B as¹¹:

$$\tilde{v}_t^n \equiv \text{vert}(v_t^n, \mathbf{0}_H, \mathbf{0}_H), \quad \tilde{B}_t^n \equiv \text{vert}(W_t^n \langle B \rangle, U_A, (U_B)_t^n),$$

where U_A is the Cholesky decomposition of S_A , so that $U_A^\top U_A = S_A$ (similarly, $(U_B)_t^n$ is the Cholesky decomposition of $(S_B)_t^n$). The equivalent LGSSM $\tilde{q}(h_{1:T}^n | \tilde{v}_{1:T}^n)$ is then completed by specifying¹²:

$$\tilde{A} \equiv \langle A \rangle, \quad \tilde{\Sigma}_H \equiv \langle \Sigma_H^{-1} \rangle^{-1}, \quad \tilde{\Sigma}_V \equiv \text{bdg}(\langle W_t^n \Sigma_V^{-1} \rangle^{-1}, I_H, I_H), \quad \tilde{\mu} \equiv \mu, \quad \tilde{\Sigma} \equiv \Sigma.$$

The validity of this parameter assignment can be checked by showing that, up to negligible constants, the exponent of this augmented LGSSM has the same form as the exponent in Eq. (10). We can now apply any standard inference routine to compute $q(h_{1:T}^n | v_{1:T}^n) = \tilde{q}(h_{1:T}^n | \tilde{v}_{1:T}^n)$ [8, 9, 19]¹³. In Appendix B.2.6 we describe the standard predictor-corrector form of the Kalman Filter, together with the Rauch-Tung-Striebel Smoother [9] and we show that a slight modification of the predictor-corrector algorithm produces a more efficient procedure obviating the need to consider fictitious outputs explicitly.

4.2 Clustering based on Simultaneous Similarity

In Section 4.1 we described a time-series clustering approach based on global similarity of the dynamics. We now introduce a different method, in which time-series are clustered on the basis of their simultaneous dynamical similarity. We will show that this can be achieved by a modification of the previously described framework.

As a motivating scenario, consider a situation in which at each time t , the stock prices of a set of V companies is known. Making a model of the stock prices of the companies over a long time is a considerable challenge.

¹¹There are several ways of achieving a similar augmentation. We chose this since, in the non-Bayesian limit $U_A = (U_B)_t^n = \mathbf{0}_{HH}$, no numerical instabilities would be introduced ($\mathbf{0}_{HH}$ is a $H \times H$ matrix of zero elements).

¹²For time T , $\tilde{B}_T^n \equiv \text{vert}(W_T^n \langle B \rangle, \mathbf{0}_{HH}, (U_B)_T^n)$.

¹³Note that, since the augmented LGSSM $\tilde{q}(h_{1:T}^n | \tilde{v}_{1:T}^n)$ is designed to match the *fully* clamped distribution $q(h_{1:T}^n | v_{1:T}^n)$, the filtered posterior $\tilde{q}(h_t^n | \tilde{v}_{1:t}^n)$ does not correspond to $q(h_t^n | v_{1:t}^n)$.

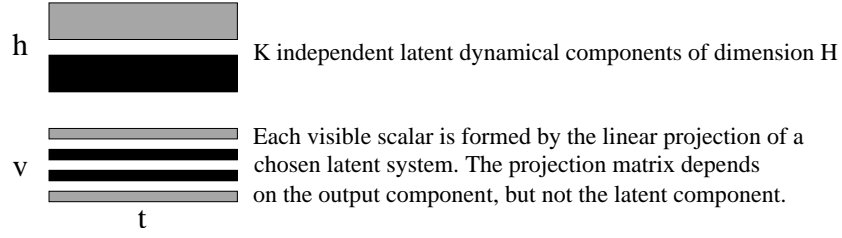


Figure 4: Time-series clustering based on simultaneous similarity.

However, it is also suspected that certain stocks follow a similar underlying temporal profile, at least over a short time-scale. For example, the stocks of high-tech companies might be strongly dependent in the sense that their movements are correlated. Similarly, the stocks of oil companies might be more strongly related to each other than to companies outside the oil group. Our interest is therefore to cluster the companies into groups of ‘correlated’ activities. Whilst making a model for the long-time resulting vector time-series of the companies is very complex, grouping stocks together based on their simultaneous dynamical behavior over a short time-scale may be much simpler since we can compare at each time-point how the movements of the stocks are correlated with each other.

For clarity, consider only the task of clustering a set of V *unidimensional* time-series $v_{1:T}^{1:V}$; the extension to the multidimensional case is straightforward). Unlike the global clustering method, we are now interested in simultaneous dynamical similarity, for which we assume that all time-series have the same length T ¹⁴. It is useful to compactly rewrite this set as one V -dimensional time-series $v_{1:T}$ with components $[v_{1:T}]_i \equiv [v_1]_i, \dots, [v_T]_i, i = 1, \dots, V$. We are thus interested in clustering the components i , that is to assign the i^{th} sequence $[v_{1:T}]_i$ into one of K clusters on the basis of its simultaneous dynamical similarity with other sequences in the same cluster. In order to do so, we consider K independent dynamical systems. Each unidimensional time-series is then a one-dimensional projection from one of the dynamical systems. The assignment of an output to a latent dynamical system is fixed throughout the time-series. See Fig. 4 for an informal sketch of this setup. More precisely, our model assumes that the emission parameters of the LGSSM $\Theta_e \equiv \{B, \Sigma_V\}$ do not depend on the cluster k , while the transition parameters $\Theta_t^k \equiv \{A^k, \Sigma_H^k, \mu^k, \Sigma^k\}$ depend on the cluster k . We may cluster components by using an indicator z^i for each component and define $p([v_t]_i | h_t^{1:K}, z^i = k, \Theta_e) \equiv p([v_t]_i | h_t^k, \Theta_e)$. That is, when z^i is in state k , the i^{th} component is drawn from the dynamics of the k^{th} LGSSM. More precisely, the emission term is given by

$$p([v_t]_i | h_t^{1:K}, z^i = k, \Theta_e) \equiv \mathcal{N}(B_i h_t^k, [\Sigma_V]_{ii}),$$

while the transition term is given by:

$$p(h_t^k | h_{t-1}^k, \Theta_t^k) \equiv \mathcal{N}(A^k h_{t-1}^k, \Sigma_H^k).$$

The emissions and transitions over time give:

$$p([v_{1:T}]_i | h_{1:T}^{1:K}, z^i = k, \Theta_e) = \prod_{t=1}^T p([v_t]_i | h_t^{1:K}, z^i = k, \Theta_e)$$

$$p(h_{1:T}^k | \Theta_t^k) = p(h_1^k | \Theta_t^k) \prod_{t=2}^T p(h_t^k | h_{t-1}^k, \Theta_t^k).$$

Each of the K linear dynamical systems proceeds independently of the rest, giving:

$$p(h_{1:T}^{1:K} | \Theta_t^{1:K}) = \prod_{k=1}^K p(h_{1:T}^k | \Theta_t^k).$$

The covariance Σ_V is constrained to be diagonal, so that

$$p(v_{1:T} | h_{1:T}^{1:K}, z^{1:V}, \Theta_e) = \prod_{i=1}^V p([v_{1:T}]_i | h_{1:T}^{1:K}, z^i, \Theta_e).$$

¹⁴For synchronized time-series of differing lengths, in principle one could treat this as a missing-data problem, along the lines previously described.

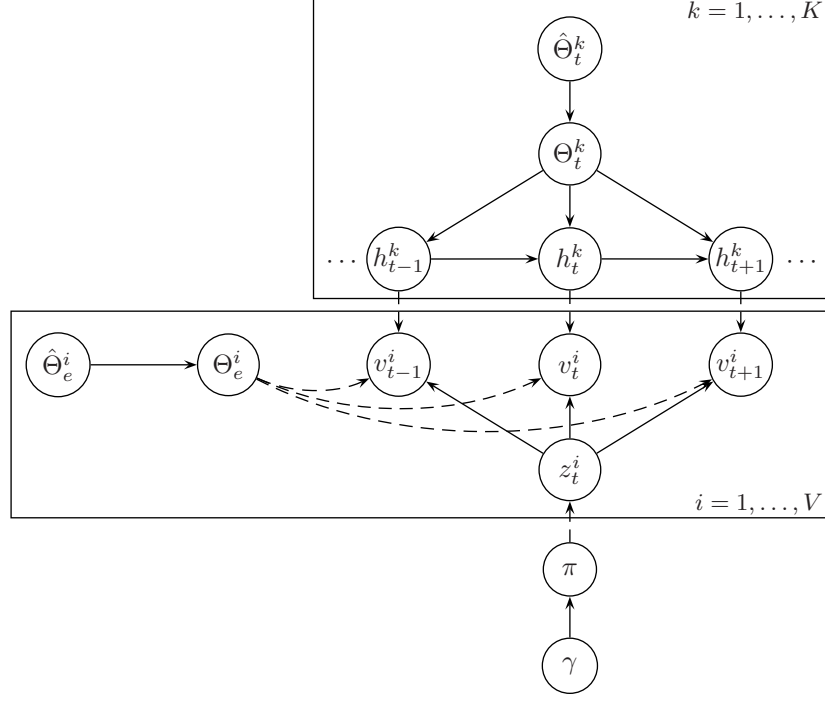


Figure 5: Graphical representation of the Dirichlet Mixture of Bayesian LGSSMs for performing clustering based on simultaneous similarity.

The graphical representation of this model is presented in Fig. 5. In the Bayesian version, we define the joint density

$$p(v_{1:T}, h_{1:T}^{1:K}, z^{1:V}, \Theta_e, \Theta_t^{1:K} | \hat{\Theta}_e, \hat{\Theta}_t^{1:K}) \\ = \prod_t \left\{ \prod_i p([v_t]_i | h_t^{1:K}, z^i, \Theta_e) \prod_k p(h_t^k | h_{t-1}^k, \Theta_t^k) \right\} p(z^{1:V}) p(\Theta_e | \hat{\Theta}_e) p(\Theta_t^{1:K} | \hat{\Theta}_t^{1:K}).$$

As in Section 4, we can automatically learn the number of clusters K by placing a Polya distribution on $p(z^{1:V})$ (Eq. (5)), where now N is replaced with V . For the LGSSM parameters, we define Gaussian-Gamma priors, similar to the one described in Section 2.

To form a tractable marginal log-likelihood bound, we use the variational approximation¹⁵:

$$p(h_{1:T}^{1:K}, z^{1:V}, \Theta_e, \Theta_t^{1:K} | v_{1:T}) \approx q(h_{1:T}^{1:K}) q(z^{1:V}) q(\Theta_e, \Theta_t^{1:K}),$$

and further make the factorization assumption $q(z^{1:V}) \equiv \prod_i q(z^i)$. The independence assumptions on p and q imply $q(h_{1:T}^{1:K}) = \prod_k q(h_{1:T}^k)$ and $q(\Theta_t^{1:K}) = \prod_k q(\Theta_t^k)$. Then the lower bound on $\log p(v_{1:T} | \hat{\Theta})$ is given by:

$$\mathcal{F} \equiv \sum_k H_q(h_{1:T}^k) + \sum_i H_q(z^i) + H_q(\Theta_e) + \sum_k H_q(\Theta_t^k) \\ + \sum_{i,t} \langle \log p([v_t]_i | h_t^{1:K}, z^i, \Theta_e) \rangle_{q(z^i)q(h_{1:T}^k)q(\Theta_e)} + \sum_{t,k} \langle \log p(h_t^k | h_{t-1}^k, \Theta_t^k) \rangle_{q(h_{1:t-1}^k)q(\Theta_t^k)} \\ + \langle \log p(z^{1:V}) \rangle_{\prod_i q(z^i)} + \langle \log p(\Theta_e | \hat{\Theta}_e) \rangle_{q(\Theta_e)} + \sum_k \langle \log p(\Theta_t^k | \hat{\Theta}_t^k) \rangle_{q(\Theta_t^k)}.$$

¹⁵The factorization $q(h_{1:T}^{1:K} | z^{1:V}) \equiv q(h_{1:T}^{1:K})$ is assumed in order to avoid the computational issue of having to consider all possible combinations of $z^{1:V}$.

The resulting optimal q distributions are:

$$\begin{aligned}
q(A^k, (\Sigma_H^k)^{-1}) &\propto p\left(A^k, (\Sigma_H^k)^{-1} | \hat{\Theta}_t^k\right) e^{\sum_t \langle \log p(h_t^k | h_{t-1}^k, \Theta_t^k) \rangle_{q(h_{t-1:t}^k)}} \\
q(B, \Sigma_V^{-1}) &\propto p\left(B, \Sigma_V^{-1} | \hat{\Theta}_e\right) e^{\sum_{i,t} \langle \log p([v_t]_i | h_t^{1:K}, z^i, \Theta_e) \rangle_{q(z^i)q(h_{1:K}^k)}} \\
q(h_{1:T}^k) &\propto e^{\sum_{i,t} q(z^i=k) \langle \log p([v_t]_i | h_t^k, \Theta_e) \rangle_{q(\Theta_e)} + \sum_t \langle \log p(h_t^k | h_{t-1}^k, \Theta_t^k) \rangle_{q(\Theta_t^k)}} \\
q(z^i = k) &\propto e^{\langle \log p(z^i | z^{-i}) \rangle_{\prod_{j \neq i} z^j} + \sum_t \langle \log p(v_t^i | h_t^k, \Theta_e) \rangle_{q(h_t^k)q(\Theta_e)}}
\end{aligned}$$

The specific updates can be derived similarly as for the method described in Section 4.1. We point out that for performing inference on $q(h_{1:T}^k)$ a similar approach as the one described in Section 4.1 can be applied with the replacement $\langle \rho_i \rangle \leftarrow q(z^i = k) \langle \rho_i \rangle$, $S_B \leftarrow \sum_i q(z^i = k) H_{iB}^{-1}$ for the case in which there are not missing observations.

4.3 Relation to Previous Work

Older works based on autoregressive (AR) models include [23], which uses a mixture of ARMA models, with the number of mixtures determined by the BIC criterion. Our method is similar to [7] in which specially constrained LGSSMs were used to form components in a Dirichlet mixture. Similar to our procedure, the authors used a Bayesian prior to encourage simplicity of each LGSSM. In [7], sampling is slow since computing the likelihood of the LGSSM requires $O(T)$ operations, so a single MCMC update is $O(T)$. This can be seen as an extension of [3], which discusses a sampling approach for a Dirichlet Process mixture of Factor Analyzers. In [3], mixing of the MCMC chain is slow, and the method is prohibitive for large datasets and also large observation vectors.

5 Demonstration

We performed several experiments to test the two models presented above on their clustering ability in ‘difficult’ situations. We will give some illustrative examples below.

5.1 Clustering based on Global Similarity

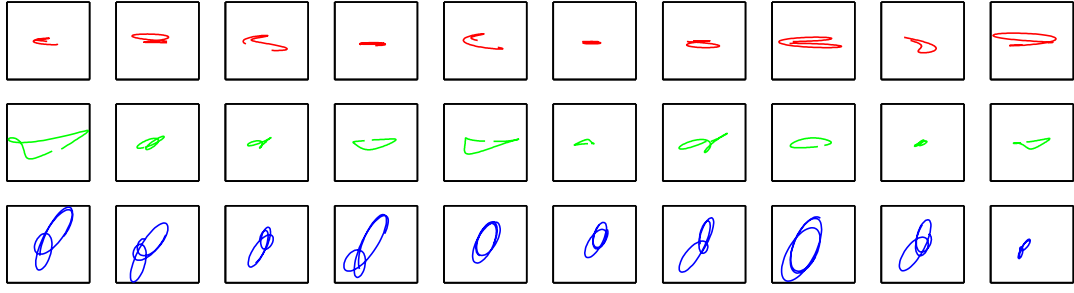
We tested the model on thirty synthetic sequences generated by two LGSSMs with $V = 2$, $H = 5$ and $T = 140$. The parameters were chosen so that all time-series had visually dissimilar dynamical trajectories, see Fig. 6a. Our model with initial $K = 5$ clusters and a $H = 10$ latent dimension perfectly clustered the data into three groups (see Fig. 6b). Thanks to the priors enforcing a low number of clusters, and simplicity of each cluster model, we consistently found the same clustering using different initial K and H , provided they are sufficiently large. This illustrates that the DMBLGSSM is capable of clustering time-series for which a common method based on feature extraction would be more difficult to apply, since it is not clear which type of features can be used to group the time-series.

5.1.1 Missing Observations

We generated fifty synthetic sequences from two LGSSMs with $V = 2$, $H = 5$ and $T = 30$. The dynamics A^k of the two models were set to differ by a small amount (the eigenvalues are close and have the same stability properties). The mixing matrices, the noise covariances and prior means were set to be independent on the LGSSM, that is $B^k = B$, $\Sigma^k = \Sigma$, $\Sigma_H^k = \Sigma_H$ and $\Sigma_V^k = \Sigma_V$. Therefore the two models differ slightly only in the deterministic part of the dynamics A^k . We removed randomly 10% of the data from each channel (for a total 20% of the available time values). In Fig. 7 we plot four samples, two from LGSSM1 (blue) and two from LGSSM2 (red). Notice that, as expected from the setup, it is not possible to visually identify a common structure for two samples from the same cluster or a structure specific to each cluster. We run the DMBLGSSM with initial number of mixture $M = 8$ and hidden dimensionality $H = 7$. The model could correctly learn the appropriate number of mixture components and assigned all samples to the correct cluster. However, when removing 15% of the data from each channel, the model incorrectly assigned all samples to one cluster only.

5.2 Clustering based on Simultaneous Similarity

As a simple illustration of the clustering method based on simultaneous similarity, in Fig. 8 we plot a set of $V = 6$ output sequences of length $T=250$ which were generated by projecting from two independent LGSSM of



(a) Unclustered Trajectories



(b) Clustered Trajectories

Figure 6: (a) Thirty trajectories length $T = 140$ resulting from the dynamics of three different LGSSMs, all with hidden dimension $H = 5$. Plotted are the points $([v_t]_1, [v_t]_2)$. Different colors correspond to different underlying LGSSMs. (b) Our method correctly identifies three clusters, with all cluster labels consistent with the known generating mechanism. The trajectories belonging to the same cluster are plotted in the same subfigure.

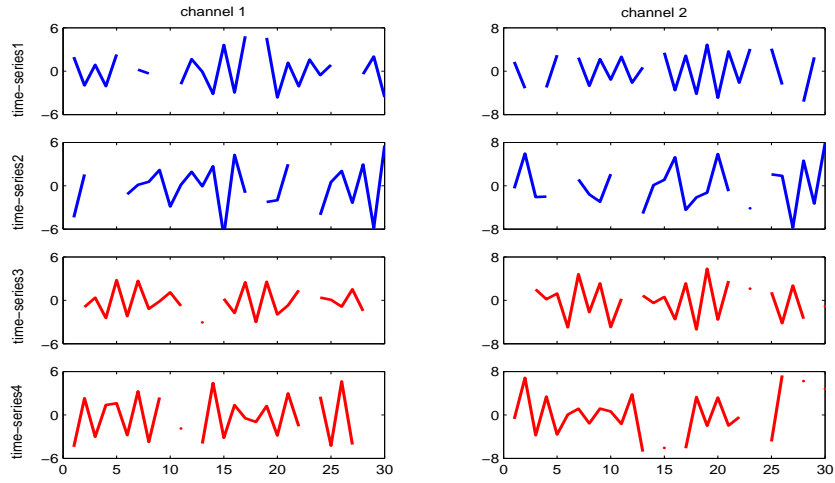


Figure 7: Four synthetic time-series generated from two LGSSMs with $V=2$, $H=5$ and $T=30$. Different colors correspond to different underlying LGSSMs. Plotted in the first and second columns are the first and second component of v_t respectively. No specific structure which identifies each cluster is visible.

dimension $H = 6$ (different colors correspond to different underlying LGSSMs). We trained our model on this data, assuming $K = 4$ latent linear dynamical systems, each of dimension $H = 8$. Pleasingly, the method correctly discarded two of the unneeded clusters, and identified the first three outputs (from top to bottom) as belonging to cluster 1, and the bottom three as belonging to cluster 2, consistent with the way the data was generated.

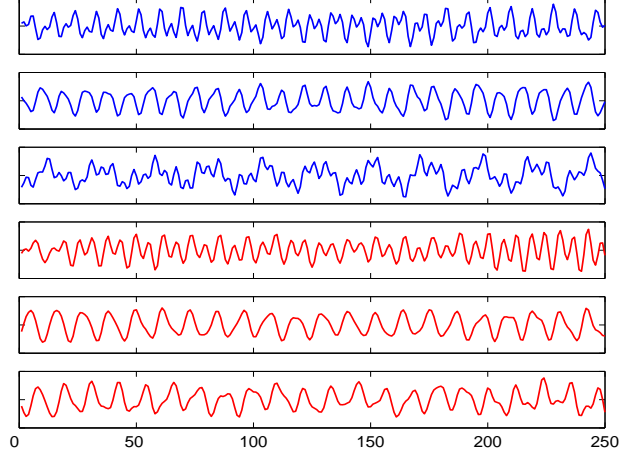


Figure 8: Clustering based on Simultaneous Similarity. Our model correctly identifies two clusters, assigning the top three output sequences to one cluster and the bottom three to another.

6 Discussion

Throughout our experiments, we found that the marginal likelihood bound is generally a reliable measure of clustering quality, provided that we consider models from the same class. However, for two models with different parameterizations (H , K values), typically one could not rely on their corresponding likelihood bounds to determine which model performed best. There are potentially several reasons why this may be the case, bearing in mind that the bound is only an approximation of the true marginal likelihood. Indeed, a point often overlooked in the literature is that the likelihood approximation, since it will typically approximate a single mode of the posterior, will miss $K!$ modes in the equivalence class defined by permuting the cluster labels. However, this correction alone cannot always account for the sometimes poor quality of the bound as a relative performance criterion across models. It might be that in difficult cases the factorization between parameters and latent states explicit in the VB approximation is too crude to accurately capture the mass of the posterior. Similar potential difficulties with the VB method have previously been reported [24].

APPENDIX

A

A.1 Wishart and Gamma Distributions

Wishart Distribution

Let Σ be a $s \times s$ positive definite symmetric matrix of random variables and let S be a positive definite matrix of size $s \times s$. Then, Σ has a Wishart distribution $\mathcal{W}(\nu, S)$ if it has a probability density function given by:

$$p(\Sigma|\nu, S) = \frac{1}{Z} |\Sigma|^{(\nu-s-1)/2} e^{-\frac{1}{2} \text{tr}[S^{-1}\Sigma]},$$

where $Z = 2^{\nu s/2} |S|^{\nu/2} \pi^{s(s-1)/4} \prod_{i=1}^s \Gamma\left(\frac{\nu+1-i}{2}\right)$.

Gamma Distribution

A random variable σ has a Gamma distribution $\mathcal{G}(\nu_1, \nu_2)$ if it has a probability density function given by:

$$p(\sigma|\nu_1, \nu_2) = \frac{\nu_2^{\nu_1}}{\Gamma(\nu_1)} \sigma^{\nu_1-1} e^{-\nu_2 \sigma}.$$

A.2 Kullback-Leibler Divergence of Gaussian, Gamma and Wishart Distributions

The Kullback-Leibler divergence $\text{KL}(q||p) = \langle \log p/q \rangle_q$ between two s -dimensional Gaussian distributions $q(x|\mu_q, \Sigma_q) = \mathcal{N}(\mu_q, \Sigma_q)$ and $p(x|\mu_p, \Sigma_p) = \mathcal{N}(\mu_p, \Sigma_p)$ is given by:

$$\text{KL}(q||p) = \frac{1}{2} \left(\log \left(\frac{\det \Sigma_p}{\det \Sigma_q} \right) + \text{tr} [\Sigma_p^{-1} \Sigma_q] + (\mu_p - \mu_q)^T \Sigma_p^{-1} (\mu_p - \mu_q) - s \right).$$

The KL divergence between two Gamma distributions $q(\sigma|q_1, q_2) = \mathcal{G}(q_1, q_2)$ and $p(\sigma|p_1, p_2) = \mathcal{G}(p_1, p_2)$ is given by:

$$\text{KL}(q||p) = q_1 \log q_2 - p_1 \log p_2 - \log \frac{\Gamma(q_1)}{\Gamma(p_1)} + (q_1 - p_1)(\psi(q_1) - \log q_2) - q_1 \left(1 - \frac{p_2}{q_2} \right),$$

where $\psi(\cdot)$ derivative of the gamma function logarithm. The KL divergence between two s -dimensional Wishart distributions $q(\Sigma|\nu_q, S_q) = \mathcal{W}(\nu_q, S_q)$ and $p(\Sigma|\nu_p, S_p) = \mathcal{W}(\nu_p, S_p)$:

$$\text{KL}(q||p) = \log \frac{Z_{\nu_p S_p}}{Z_{\nu_q S_q}} + \frac{\nu_q - \nu_p}{2} \langle \ln |\Sigma| \rangle_q + \frac{1}{2} \nu_q \text{tr} [S_p^{-1} S_q - I_s],$$

where $\langle \ln |\Sigma| \rangle_q = \sum_{i=1}^s \psi(\frac{\nu_q + 1 - i}{2}) + s \log 2 + \log |S_q|$.

A.3 Dirichlet Distribution

The Dirichlet distribution of order $K \geq 2$ with parameters $\alpha_1, \dots, \alpha_K$ has a probability density function given by:

$$\frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K x_k^{\alpha_k - 1}, \quad \text{with } \sum_{k=1}^K x_k = 1.$$

The factor $\frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)}$ is the normalizing constant, that is $\frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^K \alpha_k)} = \prod_{k=1}^K \int_0^1 x_k^{\alpha_k - 1} dx_k$. In order to see that, we first show that $\frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx$ ¹⁶.

$$\begin{aligned} \Gamma(\alpha)\Gamma(\beta) &= \int_0^\infty t^{\alpha-1} e^{-t} dt \int_0^\infty s^{\beta-1} e^{-s} ds \\ &= \int_0^\infty \int_0^\sigma \tau^{\alpha-1} (\sigma - \tau)^{\beta-1} e^{-\sigma} d\sigma d\tau \\ &= \int_0^\infty \int_0^1 y^{\alpha-1} x^{\alpha-1} y^{\beta-1} (1-x)^{\beta-1} e^{-y} y dy dx \\ &= \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx \Gamma(\alpha + \beta), \end{aligned}$$

with the change of variables $\sigma = t + s, \tau = t$, which implies $\sigma \in [0, \infty], \tau \in [0, \sigma] (s = \sigma - \tau > 0)$ and $|J| = 1$; and the change of variables $x = \frac{\tau}{\sigma}, y = \sigma$, which gives $|J| = y$. Using this result, we find:

$$\begin{aligned} \prod_{k=1}^K \Gamma(\alpha_k) &= \int_0^\infty t_1^{\alpha_1-1} e^{-t_1} dt_1 \dots \int_0^\infty t_K^{\alpha_K-1} e^{-t_K} dt_K \\ &= \int_0^\infty t_1^{\alpha_1-1} e^{-t_1} \dots \int_0^\infty t_{K-2}^{\alpha_{K-2}-1} e^{-t_{K-2}} \int_0^1 u_{K-1}^{\alpha_{K-1}-1} (1-u_{K-1})^{\alpha_K-1} \int_0^\infty \sigma^{\alpha_{K-1}+\alpha_K-1} e^{-\sigma} \\ &= \dots \\ &= \left[\prod_{k=1}^{K-1} \int_0^1 u_k^{\alpha_k-1} (1-u_k)^{\sum_{j=k+1}^K \alpha_j-1} \right] \Gamma \left(\sum_{k=1}^K \alpha_k \right). \end{aligned}$$

¹⁶This is the normalizing constant of the Beta distribution.

Therefore we have to show that $\prod_{k=1}^{K-1} \int_0^1 u_k^{\alpha_k-1} (1-u_k)^{\sum_{j=k+1}^K \alpha_j-1} = \prod_{k=1}^K \int_0^1 x_k^{\alpha_k-1} dx_k$. This can be done by induction. If we assume:

$$\prod_{k=1}^{K-2} \int_0^1 u_k^{\alpha_k-1} (1-u_k)^{\sum_{j=k+1}^{K-1} \alpha_j-1} = \int_0^1 \cdots \int_0^1 \prod_{k=1}^{K-2} x_k^{\alpha_k-1} \left(1 - \sum_{k=1}^{K-2} x_k\right)^{\alpha_{K-1}-1} dx_{1:K-2},$$

then

$$\underbrace{\prod_{k=1}^{K-2} \int_0^1 u_k^{\alpha_k-1} (1-u_k)^{\sum_{j=k+1}^K \alpha_j-1}}_A = \int_0^1 \cdots \int_0^1 \prod_{k=1}^{K-2} x_k^{\alpha_k-1} \left(1 - \sum_{k=1}^{K-2} x_k\right)^{\alpha_{K-1}+\alpha_K-1} dx_{1:K-2},$$

and

$$\begin{aligned} A &= \int_0^1 u_{K-1}^{\alpha_{K-1}-1} (1-u_{K-1})^{\alpha_K-1} \\ &= \int_0^1 \cdots \int_0^1 \prod_{k=1}^{K-2} x_k^{\alpha_k-1} \left(1 - \sum_{k=1}^{K-2} x_k\right)^{\alpha_{K-1}+\alpha_K-1} dx_{1:K-2} \int_0^1 u_{K-1}^{\alpha_{K-1}-1} (1-u_{K-1})^{\alpha_K-1} \\ &= \int_0^1 \cdots \int_0^1 \prod_{k=1}^{K-2} y_k^{\alpha_k-1} \left(1 - \sum_{k=1}^{K-2} y_k\right)^{\alpha_{K-1}-1} \left(1 - \frac{y_K}{1 - \sum_{k=1}^{K-2} y_k}\right)^{\alpha_{K-1}-1} y_K^{\alpha_K-1} \frac{1}{1 - \sum_{k=1}^{K-2} y_k} dy_{1:K-2,K} \\ &= \int_0^1 \cdots \int_0^1 \prod_{k=1}^{K-2,K} y_k^{\alpha_k-1} \left(1 - \sum_{k=1}^{K-2,K} y_k\right)^{\alpha_{K-1}-1} dy_{1:K-2,K}, \end{aligned}$$

with the change of variables $y_K = (1 - \sum_{k=1}^{K-2} x_k)(1 - u_{K-1})$, $y_1 = x_1, \dots, y_{K-2} = x_{K-2}$, which implies $|J| = \frac{1}{1 - \sum_{k=1}^{K-2} y_k}$.

Polya Distribution

Consider a multinomial distribution $p(z^{1:N}|\pi) = \prod_{k=1}^K \pi_k^{N_k}$ and a symmetric Dirichlet distribution $p(\pi) = \frac{\Gamma(\gamma)}{\Gamma(\gamma/K)^K} \prod_{k=1}^K \pi_k^{\gamma/K-1}$. From what we have seen above:

$$p(z^{1:N}) = \frac{\Gamma(\gamma)}{\Gamma(\gamma/K)^K} \int_{\pi} \prod_{k=1}^K \pi_k^{\gamma/K+N_k-1} = \frac{\Gamma(\gamma)}{\Gamma(\gamma/K)^K} \frac{\prod_{k=1}^K \Gamma(N_k + \gamma/K)}{\Gamma(N + \gamma)},$$

and

$$\begin{aligned} p(z^1 = 1 | z^{2:N}) &= \frac{p(z^1 = 1, z^{2:N})}{p(z^{2:N})} \\ &= \frac{\prod_{k=1}^K \Gamma(\gamma/K + N_k)}{\Gamma(\gamma + N)} \frac{\Gamma(\gamma + N - 1)}{\Gamma(\gamma/K + N_1 - 1) \prod_{k=2}^K \Gamma(\gamma/K + N_k)} \\ &= \frac{\Gamma(\gamma/K + N_1)}{\Gamma(\gamma + N)} \frac{\Gamma(\gamma + N - 1)}{\Gamma(\gamma/K + N_1 - 1)} = \frac{\gamma/K + N_1 - 1}{\gamma + N - 1}. \end{aligned}$$

A.4 Kronecker Product and Vectorization

The following properties of Kronecker product and matrix vectorization hold:

$$\begin{aligned}
\text{tr} [A^\top B] &= \text{vc}(A)^\top \text{vc}(B) = \text{vr}(A)^\top \text{vr}(B) \\
(A \otimes B)^\top &= A^\top \otimes B^\top \\
(A \otimes B)^{-1} &= A^{-1} \otimes B^{-1} \\
\text{vc}(ABC) &= (C^\top \otimes A) \text{vc}(B) \\
\text{vc}(ABC)^\top &= \text{vc}(B)^\top (C \otimes A^\top) \\
\text{vr}(ABC)^\top &= \text{vr}(B)^\top (A^\top \otimes C) \\
(A \otimes B)(C \otimes D) &= AC \otimes BD
\end{aligned}$$

B

B.1 Independence Assumptions on the q Distribution

Consider the model described in Section 4.1. The independence assumptions made on the q distribution are the following: $q(h_{1:T}^{1:N} | z^{1:N}, \Theta^{1:K}) \equiv q(h_{1:T}^{1:N} | z^{1:N})$, $q(z^{1:N}, \Theta^{1:K}) \equiv q(z^{1:N})q(\Theta^{1:K})$ and $q(z^{1:N}) \equiv \prod_{n=1}^N q(z^n)$. The first two assumptions on q and the assumptions on p imply that the optimal q satisfies: $q(h_{1:T}^{1:N} | z^{1:N}) = \prod_{n=1}^N q(h_{1:T}^n | z^n)$ and $q(\Theta^{1:K}) = \prod_{k=1}^K q(\Theta^k)$. Indeed the lower bound is given by:

$$\begin{aligned}
\mathcal{F} \equiv & H_q(\Theta^{1:K}) + \sum_{k^{1:N}} q(z^{1:N} = k^{1:N}) H_q(h_{1:T}^{1:N} | z^{1:N} = k^{1:N}) + H_q(z^{1:N}) + \left\langle \sum_{k=1}^K \log p(\Theta^k | \hat{\Theta}^k) \right\rangle_{q(\Theta^{1:K})} \\
& + \langle \log p(z^{1:N}) \rangle_{q(z^{1:N})} + \sum_{k^{1:N}} q(z^{1:N} = k^{1:N}) \left\langle \sum_{n=1}^N \log p(v_{1:T}^n, h_{1:T}^n | \Theta^{k^n}) \right\rangle_{q(\Theta^{1:K}) q(h_{1:T}^{1:N} | z^{1:N} = k^{1:N})}.
\end{aligned}$$

By maximizing \mathcal{F} with respect to $q(h_{1:T}^{1:N})$ and $q(\Theta^{1:K})$, we obtain:

$$\begin{aligned}
q(h_{1:T}^{1:N} | z^{1:N} = k^{1:N}) &\propto e^{\langle \sum_{n=1}^N \log p(v_{1:T}^n, h_{1:T}^n | \Theta^{k^n}) \rangle_{q(\Theta^{1:K})}} \\
&= e^{\sum_{n=1}^N \langle \log p(v_{1:T}^n, h_{1:T}^n | \Theta^{k^n}) \rangle_{q(\Theta^{k^n})}} \\
&= \prod_{n=1}^N q(h_{1:T}^n | z^n = k^n),
\end{aligned}$$

and

$$\begin{aligned}
q(\Theta^{1:K}) &\propto \prod_{k=1}^K p(\Theta^k | \hat{\Theta}^k) e^{\langle \sum_{n=1}^N \log p(v_{1:T}^n, h_{1:T}^n | \Theta^{k^n}) \rangle_{q(h_{1:T}^{1:N}, z^{1:N} = k^{1:N})}} \\
&= \prod_{k=1}^K p(\Theta^k | \hat{\Theta}^k) e^{\sum_{n=1}^N q(z^n = k) \langle \log p(v_{1:T}^n, h_{1:T}^n | \Theta^k) \rangle_{q(h_{1:T}^n | z^n = k)}} \\
&= \prod_{k=1}^K q(\Theta^k).
\end{aligned}$$

Notice that, under these assumptions, the optimal $q(z^{1:N})$ would be:

$$q(z^{1:N} = k^{1:N}) \propto p(z^{1:N} = k^{1:N}) e^{H_q(h_{1:T}^{1:N} | z^{1:N} = k^{1:N}) + \langle \sum_{n=1}^N \log p(v_{1:T}^n, h_{1:T}^n | \Theta^{k^n}) \rangle_{q(\Theta^{1:K}) q(h_{1:T}^{1:N} | z^{1:N} = k^{1:N})}},$$

which does not factorize because $p(z^{1:N})$ does not factorize.

B.2 Parameter Updates

In this Section we describe in details the updates for the q distribution for the model described in Section 4.1. The updates for the model described in Section 4.2 are similar and thus omitted.

To simplify the notation, it is useful to use a column vectorization for B , $vc(B)$, and a row vectorization for A , $vr(A)$ for the respective distributions. For the same reason, we will omit the dependency of the model parameter Θ^k and hyperparameter $\hat{\Theta}^k$ on the mixture k .

Updates for $q(B^k, [\Sigma_V^k]^{-1})$

The optimal $q(B, \Sigma_V^{-1})$ is a Gaussian-Wishart(Gamma) distribution given by:

$$\frac{1}{|2\pi\Sigma_V|^{\frac{T}{2}} \sum_{n=1}^N q(z^n=k)} e^{-\frac{1}{2} \sum_{n=1}^N q(z^n=k) \sum_{t=1}^T \langle (v_t^n - W_t^n B h_t^n)^\top \Sigma_V^{-1} (v_t^n - W_t^n B h_t^n) \rangle_{q(h_t|z^n=k)}} p(B|\beta, \Sigma_V^{-1}) p(\Sigma_V^{-1}|\hat{\Theta}) \quad (11)$$

The exponent (excluding $p(\Sigma_V^{-1})$) is given by $-\frac{1}{2}\mathcal{E}$, where:

$$\mathcal{E} = \sum_{n=1}^N q(z^n=k) \left(\sum_{t=1}^T (v_t^n)^\top \Sigma_V^{-1} v_t^n - 2 \sum_{t=1}^T \langle h_t^n \rangle^\top B^\top W_t^n \Sigma_V^{-1} v_t^n + \sum_{t=1}^T \langle (h_t^n)^\top B^\top W_t^n \Sigma_V^{-1} W_t^n B h_t^n \rangle \right) + \beta_j B_j^\top \Sigma_V^{-1} B_j - 2 \hat{B}_j^\top \Sigma_V^{-1} B_j + \hat{B}_j^\top \Sigma_V^{-1} \hat{B}_j.$$

B.2.1 Determining $q(B^k | [\Sigma_V^k]^{-1})$

Optimally, $q(B|\Sigma_V^{-1})$ is a Gaussian. If we assume that Σ_V is diagonal for the case in which there are missing observations, then $W_t^n \Sigma_V^{-1} W_t^n = W_t^n \Sigma_V^{-1}$ and, using the properties described in Appendix A.4, we can write the quadratic term in B of \mathcal{E} as

$$\begin{aligned} & \sum_{n=1}^N q(z^n=k) \sum_{t=1}^T \langle (h_t^n)^\top B^\top W_t^n \Sigma_V^{-1} B h_t^n \rangle + \sum_j \beta_j B_j^\top \Sigma_V^{-1} B_j \\ &= tr \left[\sum_{n=1}^N q(z^n=k) \sum_{t=1}^T \langle h_t^n (h_t^n)^\top \rangle B^\top W_t^n \Sigma_V^{-1} B + dg(\beta) B^\top \Sigma_V^{-1} B \right] \\ &= vc(B)^\top \left(\underbrace{\sum_{n=1}^N q(z^n=k) \sum_{t=1}^T \langle h_t^n (h_t^n)^\top \rangle \otimes W_t^n + dg(\beta) \otimes I_V}_{H_{BM}} \right) (I_H \otimes \Sigma_V^{-1}) vc(B), \end{aligned}$$

that is the covariance of $q(vc(B) | \Sigma_V^{-1})$ is given by:

$$\Sigma_B = (I_H \otimes \Sigma_V) H_{BM}^{-1}.$$

The linear term in B is given by:

$$\begin{aligned} & tr \left[\sum_{n=1}^N q(z^n=k) \sum_{t=1}^T \langle h_t^n \rangle (v_t^n)^\top \Sigma_V^{-1} W_t^n B + dg(\beta) \hat{B}^\top \Sigma_V^{-1} B \right] \\ &= vc \left(\Sigma_V^{-1} \left(\underbrace{W_t^n \sum_{n=1}^N q(z^n=k) \sum_{t=1}^T v_t^n \langle h_t^n \rangle^\top + \hat{B} dg(\beta)}_{N_B} \right) \right)^\top vc(B), \end{aligned}$$

that is the mean of $q(vc(B) | \Sigma_V^{-1})$ is given by:

$$\mu_B = \Sigma_B vc(\Sigma_V^{-1} N_B) = H_{BM}^{-1} vc(N_B).$$

In the case in which there are not missing observations ($W_t^n = I_V$), the formula for the covariance reduces to:

$$\Sigma_B = \left(\underbrace{\sum_{n=1}^N q(z^n = k) \sum_{t=1}^T \langle h_t^n (h_t^n)^\top \rangle + dg(\beta)}_{H_B} \right)^{-1} \otimes \Sigma_V,$$

and the mean becomes:

$$\mu_B = vc(N_B H_B^{-1}).$$

B.2.2 Determining $q([\Sigma_V^k]^{-1})$

In Section B.2.1, we have shown that $q(B|\Sigma_V^{-1})$ is Gaussian with exponent:

$$-\frac{1}{2} (vc(B) - \mu_B)^\top \Sigma_V^{-1} (vc(B) - \mu_B). \quad (12)$$

The part of Eq. (12) which is not explicitly present in Eq. (11) is given by:

$$-\frac{1}{2} \mu_B^\top \Sigma_V^{-1} \mu_B = -\frac{1}{2} tr[M_B N_B^\top \Sigma_V^{-1}].$$

where $M_B = N_B H_B^{-1}$ for the case in which there are not missing observations, while M_B is the $V \times H$ matrix formed by the vector $H_{BM}^{-1} vc(N_B)$ for the case of missing observations. The negative of this term, together with the part in the exponent of Eq. (11) of which contains a dependency Σ_V , gives as exponent for $q(\Sigma_V^{-1})$:

$$\frac{1}{2} tr[M_B N_B^\top \Sigma_V^{-1}] - \frac{1}{2} \sum_{n=1}^N q(z^n = k) \sum_{t=1}^T (v_t^n)^\top \Sigma_V^{-1} v_t^n - \frac{1}{2} B_j^\top \Sigma_V^{-1} B_j - \frac{1}{2} tr[S_V^{-1} \Sigma_V^{-1}].$$

The terms that contain dependency on $|\Sigma_V|$ are given by:

$$\frac{|\Sigma_B|^{1/2} |dg(\beta)|^{V/2} |\Sigma_V^{-1}|^{(\nu_V - V - 1)/2}}{|\Sigma_V|^{\frac{T}{2} \sum_{n=1}^N q(z^n = k)} |\Sigma_V|^{H/2}} = \frac{|H_B^{-1}|^{V/2} |dg(\beta)|^{V/2} |\Sigma_V^{-1}|^{(\nu_V - V - 1)/2}}{|\Sigma_V|^{\frac{T}{2} \sum_{n=1}^N q(z^n = k)}}.$$

That means that if Σ_V^{-1} follows a Wishart distribution $\mathcal{W}(\nu_V, S_V)$ (which we permit when there are no missing observations) the updates are:

$$q(\Sigma_V^{-1}) = \mathcal{W}\left(\nu_V + T \sum_{n=1}^N q(z^n = k), \left(S_V^{-1} + \sum_{n=1}^N q(z^n = k) \sum_{t=1}^T v_t^n (v_t^n)^\top - G_B + \hat{B} dg(\beta) \hat{B}^\top\right)^{-1}\right),$$

where $G_B \equiv N_B H_B^{-1} N_B^\top$. Instead, for the constraint $\Sigma_V^{-1} = dg(\rho)$, where each diagonal element ρ_i follows a Gamma prior $\mathcal{G}(b_1^i, b_2^i)$, the optimal updates are:

$$q(\rho_i) = \mathcal{G}\left(b_1^i + \frac{T}{2} \sum_{n=1}^N q(z^n = k), b_2^i + \frac{1}{2} \left(\sum_{n=1}^N q(z^n = k) \sum_{t=1}^T [v_t^n]_i^2 - [G_B]_{ii} + \sum_j \beta_j \hat{B}_{ij}^2 \right)\right).$$

Updates for $q(A^k, [\Sigma_H^k]^{-1})$

The optimal $q(A, \Sigma_H^{-1})$ is a Gaussian-Gamma distribution given by:

$$\frac{1}{|2\pi \Sigma_H|^{\frac{T-1}{2} \sum_{n=1}^N q(z^n = k)}} e^{-\frac{1}{2} \sum_{n=1}^N q(z^n = k) \sum_{t=2}^T \langle (h_t^n - A h_{t-1}^n)^\top \Sigma_H^{-1} (h_t^n - A h_{t-1}^n) \rangle_{q(h_{t-1:t}^n | z^n = k)}} p(A | \alpha, \Sigma_H^{-1}) p(\Sigma_H^{-1} | \hat{\Theta}) \quad (13)$$

B.2.3 Determining $q(A|\Sigma_H)$

Optimally, $q(A|\Sigma_H^{-1})$ is a Gaussian. In order to obtain independence of the mean and other quantities from Σ_H^{-1} , here we have to assume that Σ_H^{-1} is diagonal (this will not be the case for other choices of the prior $p(A|\alpha, \Sigma_H^{-1})$).

The quadratic term in A of the exponent of Eq. (13) is given by:

$$\begin{aligned} & \sum_{n=1}^N q(z^n = k) \sum_{t=2}^T \langle (h_{t-1}^n)^\top A^\top \Sigma_H^{-1} A h_{t-1}^n \rangle + \sum_{ij} \alpha_{ij} A_{ij}^\top [\Sigma_H^{-1}]_{ii} A_{ij} \\ &= \text{tr} \left[\sum_{n=1}^N q(z^n = k) \sum_{t=2}^T \langle h_{t-1}^n (h_{t-1}^n)^\top \rangle A^\top \Sigma_H^{-1} A \right] + \underbrace{\text{vr}(A)^\top \text{bdg}([\Sigma_H^{-1}]_{11} \text{dg}(\alpha_1'), \dots, [\Sigma_H^{-1}]_{HH} \text{dg}(\alpha_{H'}))}_{D_A} \text{vr}(A) \\ &= \text{vr}(A)^\top \left(\left[\Sigma_H^{-1} \otimes \sum_{n=1}^N q(z^n = k) \sum_{t=2}^T \langle h_{t-1}^n (h_{t-1}^n)^\top \rangle \right] + D_A \right) \text{vr}(A) \\ &= \text{vr}(A)^\top \text{bdg}([\Sigma_H^{-1}]_{11} H_{1A}, \dots, [\Sigma_H^{-1}]_{HH} H_{HA}) \text{vr}(A), \end{aligned}$$

where $[H_{iA}]_{jl} \equiv \sum_{n=1}^N q(z^n = k) \sum_{t=2}^T \langle [h_{t-1}^n]_j [h_{t-1}^n]_l \rangle_{q(h_{t-1}^n)} + \alpha_{ij} \delta_{jl}$.

This means that the covariance of $q(\text{vr}(A) | \Sigma_H^{-1})$ is given by

$$\Sigma_A = \text{bdg}([\Sigma_H]_{11} H_{1A}^{-1}, \dots, [\Sigma_H]_{HH} H_{HA}^{-1}).$$

The linear term is given by:

$$\begin{aligned} & \sum_{n=1}^N q(z^n = k) \sum_{t=2}^T \langle (h_t^n)^\top \Sigma_H^{-1} A h_{t-1}^n \rangle + \sum_{ij} \alpha_{ij} \hat{A}_{ij}^\top [\Sigma_H^{-1}]_{ii} A_{ij} \\ &= \text{tr} \left[\sum_{n=1}^N q(z^n = k) \sum_{t=2}^T \langle h_{t-1}^n (h_t^n)^\top \rangle \Sigma_H^{-1} A \right] + \text{vr}(\hat{A})^\top D_A \text{vr}(A) \\ &= \text{vr} \left(\Sigma_H^{-1} \sum_{n=1}^N q(z^n = k) \sum_{t=2}^T \langle h_t^n (h_{t-1}^n)^\top \rangle \right)^\top \text{vr}(A) + \text{vr}(\hat{A})^\top D_A \text{vr}(A), \end{aligned}$$

that is the mean of $q(\text{vr}(A) | \Sigma_H^{-1})$ is given by:

$$\begin{aligned} \mu_A &= \Sigma_A \left(\text{vr} \left(\Sigma_H^{-1} \sum_{n=1}^N q(z^n = k) \sum_{t=2}^T \langle h_t^n (h_{t-1}^n)^\top \rangle \right) + D_A \text{vr}(\hat{A}) \right) \\ &= \text{vert} \left(([N_A]_1, H_{1A}^{-1})^\top, \dots, ([N_A]_{H'}, H_{HA}^{-1})^\top \right). \end{aligned}$$

where $[N_A]_{ij} = \sum_{n=1}^N \sum_{t=2}^T q(z^n = k) \langle [h_{t-1}^n]_j [h_t^n]_i \rangle_{q(h_{t-1}^n, h_t^n)} + \alpha_{ij} \hat{A}_{ij}$

B.2.4 Determining $q([\Sigma_H^k]^{-1})$

The missing part of $(\text{vr}(A) - \mu_A) \Sigma_A^{-1} (\text{vr}(A) - \mu_A)$ in the exponent of $q(A, \Sigma_H^{-1})$ (Eq. (13)) is given by:

$$\sum_i [\Sigma_H^{-1}]_{ii} [N_A]_{i'} H_{iA}^{-1} [N_A]_{i'}^\top.$$

For $\Sigma_H^{-1} = \text{dg}(\tau)$, where each element τ_i follows a Gamma prior $\mathcal{G}(a_1^i, a_2^i)$, the updates are:

$$q(\tau_i) = \mathcal{G} \left(a_1^i + \frac{T-1}{2} \sum_{n=1}^N q(z^n = k), a_2^i + \frac{1}{2} \left(\sum_{t=2}^T \langle [h_t]_i^2 \rangle - [G_A]_i + \sum_j \alpha_{ij} \hat{A}_{ij}^2 \right) \right),$$

where $[G_A]_i \equiv [N_A]_{i'} H_{iA}^{-1} [N_A]_{i'}^\top$.

B.2.5 Updates for $q(z)$

The average $\langle \log p(z^n = k | z^{-n}) \rangle_{\prod_{m \neq n} q(z^m)}$ is approximated using a second order Taylor expansion. More specifically:

$$p(z^n = k | z^{-n}) = \frac{N_{k,-n} + \gamma/K}{N - 1 + \gamma} \equiv f(N_{k,-n}), \quad (14)$$

where $N_{k,-n} \equiv N_k - I[z^n = k]$ is the number of times z is in state k , excluding z^n . The quantities $N_{k,-n}$ are sums of Bernoulli variables and may be approximated with a Gaussian with mean and variance given by:

$$M_{k,-n} \equiv \sum_{m=1, m \neq n}^N q(z^m = k), \quad S_{k,-n} \equiv \sum_{m=1, m \neq n}^N q(z^m = k)(1 - q(z^m = k)).$$

We then approximate $\langle f(N_{k,-n}) \rangle$ in Eq. (14) by using a second order Taylor expansion¹⁷:

$$\langle f(N_{k,-n}) \rangle = f(M_{k,-n}) + \frac{1}{2} f''(M_{k,-n}) S_{k,-n}.$$

B.2.6 Inference on $q(h_{1:T})$

In this section we describe a standard algorithm the standard predictor-corrector form of the Kalman Filter, together with the Rauch-Tung-Striebel Smoother from the LGSSM literature, which can be used for performing inference on terms such as:

$$q(h_{1:T}) \propto e^{\langle \log p(v_{1:T}, h_{1:T} | \Theta) \rangle_{q(\Theta)}}.$$

This requires defining a new set of $\tilde{A}, \tilde{B}, \tilde{\Sigma}_H, \tilde{\Sigma}_V, \tilde{\mu}, \tilde{\Sigma}$ parameters of the type described in Section 4.1. We also give a slight modification of the predictor-corrector algorithm which obviates the need to introduce fictitious outputs.

Algorithm 1 describe the standard predictor-corrector form of the Kalman Filter, together with the Rauch-Tung-Striebel Smoother [9] for computing $q(h_t | v_{1:T}) = \tilde{q}(h_t | \tilde{v}_{1:T})$.

There are two variants of the FORWARD pass. Either we may call procedure FORWARD in Algorithm 1 with parameters $\tilde{A}, \tilde{B}, \tilde{\Sigma}_H, \tilde{\Sigma}_V, \tilde{\mu}, \tilde{\Sigma}$ and the augmented visible variables \tilde{v}_t in which we use steps 1a, 2a, 5a and 6a. This is exactly the predictor-corrector form of a Kalman Filter [9]. Otherwise, in order to reduce the computational cost, we may call procedure FORWARD with the parameters $\tilde{A}, \langle B \rangle, \tilde{\Sigma}_H, \langle \Sigma_V^{-1} \rangle^{-1}, \tilde{\mu}, \tilde{\Sigma}$ and the original visible variable v_t in which we use steps 1b (where $U_{AB}^\top U_{AB} \equiv S_A + S_B$ ¹⁸), 2b, 5b and 6b. The two algorithms are mathematically equivalent, as shown below. Computing $q(h_t | v_{1:T}) = \tilde{q}(h_t | \tilde{v}_{1:T})$ is then completed by calling the common BACKWARD pass.

The important point here is that the reader may supply any standard Kalman Filtering/Smoothing routine, and simply call it with the appropriate parameters. In some parameter regimes, or in very long time-series, numerical stability may be a serious concern, for which several stabilized algorithms have been developed over the years, for example the square-root forms [9, 18, 19].

Equivalence of Algorithm 1 a and b

The filtered covariance P_t^t obtained from Algorithms 1 a and b are equivalents. Indeed, let suppose that we have demonstrated the equivalence at time $t - 1$, then by repetitive application of the matrix inversion lemma¹⁹ we obtain:

$$\begin{aligned} P_t^t &= P_t^{t-1} - P_t^{t-1} \tilde{B}^\top (\tilde{B} P_t^{t-1} \tilde{B}^\top + \tilde{\Sigma}_V)^{-1} P_t^{t-1} \\ &= ((P_t^{t-1})^{-1} + \tilde{B}^\top \tilde{\Sigma}_V^{-1} \tilde{B})^{-1} \\ &= \left(\underbrace{(P_t^{t-1})^{-1} + S_A + S_B}_{P^{-1}} + \langle B \rangle^\top \langle \Sigma_V^{-1} \rangle \langle B \rangle \right)^{-1} \\ &= P - P \langle B \rangle^\top (\langle B \rangle P \langle B \rangle^\top + \langle \Sigma_V^{-1} \rangle^{-1})^{-1} \langle B \rangle P \end{aligned}$$

¹⁷The potentially more accurate procedure of using Quadrature fails in this case, since the arguments under Gaussian Quadrature take the function out of defined regions.

¹⁸At time T , we need to define U_{AB} such that $U_{AB}^\top U_{AB} \equiv S_B$.

¹⁹Matrix inversion lemma: if the matrices A, B, C, D satisfy $B^{-1} = A^{-1} + C^\top D^{-1} C$, where all inverses are assumed to exist, then $B = A - A C^\top (C A C^\top + D)^{-1} C A$.

Algorithm 1 LGSSM: Forward and backward recursive updates. The smoothed posterior $p(h_t|v_{1:T})$ is returned in the mean \hat{h}_t^T and covariance P_t^T .

procedure FORWARD

1a: $P \leftarrow \Sigma$
1b: $P \leftarrow D\Sigma$, where $D \equiv I - \Sigma U_{AB}^\top (I + U_{AB} \Sigma U_{AB}^\top)^{-1} U_{AB}$
2a: $\hat{h}_1^0 \leftarrow \mu$
2b: $\hat{h}_1^0 \leftarrow D\mu$
3: $K \leftarrow PB^\top (BPB^\top + \Sigma_V)^{-1}$, $P_1^1 \leftarrow (I - KB)P$, $\hat{h}_1^1 \leftarrow \hat{h}_1^0 + K(v_1 - B\hat{h}_1^0)$
for $t \leftarrow 2, T$ **do**
4: $P_t^{t-1} \leftarrow AP_{t-1}^{t-1}A^\top + \Sigma_H$
5a: $P \leftarrow P_t^{t-1}$
5b: $P \leftarrow D_t P_t^{t-1}$, where $D_t \equiv I - P_t^{t-1} U_{AB}^\top (I + U_{AB} P_t^{t-1} U_{AB}^\top)^{-1} U_{AB}$
6a: $\hat{h}_t^{t-1} \leftarrow A\hat{h}_{t-1}^{t-1}$
6b: $\hat{h}_t^{t-1} \leftarrow D_t A\hat{h}_{t-1}^{t-1}$
7: $K \leftarrow PB^\top (BPB^\top + \Sigma_V)^{-1}$, $P_t^t \leftarrow (I - KB)P$, $\hat{h}_t^t \leftarrow \hat{h}_t^{t-1} + K(v_t - B\hat{h}_t^{t-1})$

end for

end procedure

procedure BACKWARD

for $t \leftarrow T-1, 1$ **do**
 $\overleftarrow{A}_t \leftarrow P_t^t A^\top (P_{t+1}^t)^{-1}$
 $P_t^T \leftarrow P_t^t + \overleftarrow{A}_t (P_{t+1}^T - P_{t+1}^t) \overleftarrow{A}_t^\top$
 $\hat{h}_t^T \leftarrow \hat{h}_t^t + \overleftarrow{A}_t (\hat{h}_{t+1}^T - A\hat{h}_t^t)$

end for

end procedure

where P can be written as $P = P_t^{t-1} - P_t^{t-1} U_{AB}^\top (U_{AB} P_t^{t-1} U_{AB}^\top + I)^{-1} U_{AB} P_t^{t-1}$.

The contribution from the observations to form the mean \hat{h}_t^t is equivalent in the two algorithms, indeed:

$$\begin{aligned}
P_t^{t-1} \tilde{B}^\top (\tilde{B} P_t^{t-1} \tilde{B}^\top + \tilde{\Sigma}_V)^{-1} &= \left(P_t^{t-1} - P_t^{t-1} \tilde{B}^\top \tilde{\Sigma}_V^{-1} \tilde{B} (\tilde{B}^\top \tilde{\Sigma}_V^{-1} \tilde{B} + (P_t^{t-1})^{-1})^{-1} \right) \tilde{B}^\top \tilde{\Sigma}_V^{-1} \\
&= \left(P_t^{t-1} - P_t^{t-1} (P_t^{t-1} + (\tilde{B}^\top \tilde{\Sigma}_V^{-1} \tilde{B})^{-1})^{-1} P_t^{t-1} \right) \tilde{B}^\top \tilde{\Sigma}_V^{-1} \\
&= ((P_t^{t-1})^{-1} + \tilde{B}^\top \tilde{\Sigma}_V^{-1} \tilde{B})^{-1} \tilde{B}^\top \tilde{\Sigma}_V^{-1} \\
&= P_t^t \tilde{B}^\top \tilde{\Sigma}_V^{-1} \\
&= \text{ver} P_t^t \langle B \rangle^\top \langle \Sigma_V^{-1} \rangle^{-1}, P_t^t U_A^\top, P_t^t U_B^\top \\
&= \text{vert} \left(P \langle B \rangle^\top (\langle B \rangle P \langle B \rangle^\top + \langle \Sigma_V^{-1} \rangle^{-1}), P_t^t U_A^\top, P_t^t U_B^\top \right)
\end{aligned}$$

The first element in vert is equivalent to the contribution of Algorithm 1 b , while the second and third do not contribute given that the corresponding elements in \tilde{v}_t are zeros. Finally, the contribution from $\hat{A}h_{t-1}^{t-1}$ to form the mean \hat{h}_t^t is

$$\begin{aligned}
(I - P_t^{t-1} \tilde{B}^\top (\tilde{B} P_t^{t-1} \tilde{B}^\top + \tilde{\Sigma}_V)^{-1} \tilde{B}) P_t^t (P_{t-1}^t)^{-1} &= \left(I - P \langle B \rangle^\top (\langle B \rangle P \langle B \rangle^\top + \langle \Sigma_V^{-1} \rangle^{-1})^{-1} \langle B \rangle \right) P (P_{t-1}^t)^{-1} \\
&= \left(I - P \langle B \rangle^\top (\langle B \rangle P \langle B \rangle^\top + \langle \Sigma_V^{-1} \rangle^{-1})^{-1} \langle B \rangle \right) D_t P_{t-1}^t (P_{t-1}^t)^{-1}
\end{aligned}$$

B.3 Parameter Covariance

The parameter covariance S_B introduced in Section 4.1 is given by:

$$\begin{aligned}
[(S_B)_t^n]_{jl} &= \text{tr} \left[\left\langle \Sigma_V^{-1} W_t^n (B_l - \langle B_l \rangle) (B_j - \langle B_j \rangle)^\top \right\rangle_{q(B, \Sigma_V^{-1})} \right] \\
&= \sum_{i, o, p=1}^V \left[\left\langle [\Sigma_V^{-1}]_{io} W_{op}^n [\Sigma_B]_{p+H(l-1), i+H(j-1)} \right\rangle_{q(\Sigma_V^{-1})} \right] \\
&= \sum_{i, o, p=1}^V \left\langle [\Sigma_V^{-1}]_{io} W_{op}^n \sum_{qr} [I_H]_{lr} [\Sigma_V]_{pq} [H_{BM}^{-1}]_{q+H(r-1), i+H(j-1)} \right\rangle \\
&= \sum_i W_{ii}^n [H_{BM}^{-1}]_{i+H(l-1), i+H(j-1)} .
\end{aligned}$$

In the particular case in which there are not missing observations this reduces to:

$$S_B = V H_B^{-1}.$$

Analogously $S_A = \sum_i H_{iA}$. Indeed:

$$\begin{aligned}
[S_A]_{jl} &= \text{tr} \left[\left\langle \Sigma_H^{-1} (A_j - \langle A_j \rangle) (A_l - \langle A_l \rangle)^\top \right\rangle \right] \\
&= \sum_{i, k} \left\langle [\Sigma_H^{-1}]_{ik} [\Sigma_H]_{ki} [H_{kA}]_{lj} \right\rangle \\
&= \sum_i [H_{iA}]_{jl} .
\end{aligned}$$

B.4 Hyperparameter Updates

Updates for β_j

If we compute the derivative of Eq. (6) with respect to the hyperparameter β_j , $j = 1, \dots, H$ and set it to zero we obtain:

$$\beta_j = \frac{V}{\left\langle (B_j - \hat{B}_j)^\top \Sigma_V^{-1} (B_j - \hat{B}_j) \right\rangle},$$

where

$$\begin{aligned}
\left\langle (B_j - \hat{B}_j)^\top \Sigma_V^{-1} (B_j - \hat{B}_j) \right\rangle &= \underbrace{\left\langle (B_j - \langle B_j \rangle)^\top \Sigma_V^{-1} (B_j - \langle B_j \rangle) \right\rangle}_{[S_B]_{jj}} \\
&\quad + \langle B_j \rangle^\top \langle \Sigma_V^{-1} \rangle \langle B_j \rangle - 2 \langle B_j^\top \rangle \langle \Sigma_V^{-1} \rangle \hat{B}_j + \hat{B}_j^\top \langle \Sigma_V^{-1} \rangle \hat{B}_j.
\end{aligned}$$

Updates for α_{ij}

If we compute the derivative of Eq. (6) with respect to the hyperparameter α_{ij} , $i, j = 1, \dots, H$, and set it to zero we obtain:

$$\alpha_{ij} = \frac{1}{\left\langle [\Sigma_H^{-1}]_{ii} (A_{ij} - \hat{A}_{ij})^2 \right\rangle},$$

where

$$\begin{aligned}
\left\langle [\Sigma_H^{-1}]_{ii} (A_{ij} - \hat{A}_{ij})^2 \right\rangle &= \underbrace{\left\langle [\Sigma_H^{-1}]_{ii} (A_{ij} - \langle A_{ij} \rangle)^2 \right\rangle}_{[H_{iA}^{-1}]_{jj}} \\
&\quad + \langle [\Sigma_H^{-1}]_{ii} \rangle \langle A_{ij} \rangle^2 - 2 \langle [\Sigma_H^{-1}]_{ii} \rangle \langle A_{ij} \rangle \hat{A}_{ij} + \langle [\Sigma_H^{-1}]_{ii} \rangle \hat{A}_{ij}^2.
\end{aligned}$$

Updates for b_1, b_2

For the constraint $\Sigma_V^{-1} = dg(\rho)$, each diagonal element ρ_i follows a Gamma prior $\mathcal{G}(b_1^i, b_2^i)$.

In order to constrain b_1^i to be positive, we set $b_1^i = b^2$. The derivative of Eq. (6) with respect to b is given by:

$$2b \log b_2^i - 2b\psi(b^2) + 2b \langle \log \rho_i \rangle_{q(\rho_i)},$$

where

$$\begin{aligned} \langle \log \rho_i \rangle_{q(\rho_i)} &= \frac{q_2^{q_1}}{\Gamma(q_1)} \int_{\rho_i} \rho_i^{q_1-1} e^{-q_2 \rho_i} \log \rho_i \\ &= \frac{q_2^{q_1}}{\Gamma(q_1)} \frac{\partial}{\partial q_1} \frac{\Gamma(q_1)}{q_2^{q_1}} \\ &= \psi(q_1) - \log q_2. \end{aligned}$$

Given that we cannot obtain a close form update for b_1^i , we have to use some optimization method. Setting to zero the derivative of Eq. (6) with respect to b_2^i , we obtain:

$$b_2^i = \frac{b_1^i}{\langle \rho_i \rangle_{q(\rho_i)}}.$$

Similar updates can be obtained for a_1 and a_2 .

Updates for ν_V, S_V

If Σ_V^{-1} follows a Wishart prior distribution $p(\Sigma_V^{-1} | \nu_V, S_V) = \mathcal{W}(\nu_V, S_V)$, the derivative of Eq. (6) with respect to ν_V is given by:

$$-\frac{V}{2} \log 2 - \frac{1}{2} \log |S_V| - \sum_{i=1}^V \psi \left(\frac{\nu_V + 1 - i}{2} \right) + \langle \log |\Sigma_V^{-1}| \rangle_{q(\Sigma_V^{-1})},$$

where

$$\begin{aligned} \langle \log |\Sigma_V^{-1}| \rangle_{q(\Sigma_V^{-1})} &= \frac{1}{Z_q} \int \log |\Sigma_V^{-1}| |\Sigma_V^{-1}|^{\frac{\nu_q - V - 1}{2}} e^{-\frac{1}{2} \text{tr}[S_q^{-1} \Sigma_V^{-1}]} \\ &= \frac{2}{Z_q} \frac{\partial Z_q}{\partial \nu_q} \\ &= \sum_i \psi \left(\frac{\nu_q + 1 - i}{2} \right) + V \log 2 + \log |S_q|. \end{aligned}$$

Setting to zero the derivative of Eq. (6) with respect to S_V , we obtain:

$$S_V = \frac{1}{\nu_V} \Sigma_V^{-1}.$$

Updates for γ

In order to constrain γ to be positive we set $\gamma = \delta^2$. The derivative of Eq. (6) with respect to δ is given by:

$$-2\delta \left(\psi(\delta^2) - \psi(N + \delta^2) - \psi \left(\frac{\delta^2}{K} \right) \right) + 2 \frac{\delta}{K} \sum_{k=1}^K \langle \psi(N_k + \delta^2/K) \rangle_{\prod_{n=1}^N q(z^n)},$$

which does not give a close form update for γ .

B.5 Computing the Log-likelihood Bound

The log-likelihood bound (Eq. (6)) can be rewritten as:

$$\begin{aligned} \mathcal{F} = & \sum_{n=1}^N \sum_{k=1}^K q(z^n = k) H_q(h_{1:T}^n | z^n = k) + \sum_{n=1}^N H_q(z^n) - \sum_{k=1}^K \left\langle \log \frac{q(A|\Sigma_H^{-1})}{p(A|\Sigma_H^{-1})} \right\rangle_{q(A, \Sigma_H^{-1})} \\ & - \sum_{k=1}^K \left\langle \log \frac{q(B|\Sigma_V^{-1})}{p(B|\Sigma_V^{-1})} \right\rangle_{q(B, \Sigma_V^{-1})} - \sum_{k=1}^K \left\langle \log \frac{q(\Sigma_H^{-1})}{p(\Sigma_H^{-1})} \right\rangle_{q(\Sigma_H^{-1})} - \sum_{k=1}^K \left\langle \log \frac{q(\Sigma_V^{-1})}{p(\Sigma_V^{-1})} \right\rangle_{q(\Sigma_V^{-1})} \\ & + \langle \log p(z^{1:N}) \rangle_{\prod_{n=1}^N q(z^n)} + \sum_{n=1}^N \sum_{k=1}^K q(z^n = k) \langle \log p(v_{1:T}^n, h_{1:T}^n | \Theta^k) \rangle_{q(\Theta)q(h_{1:T}^n | z^n = k)}. \end{aligned}$$

Each entropic element in the first sum can be computed as:

$$\begin{aligned} - \left\langle \sum_{t=1}^{T-1} \log \frac{q(h_t^n, h_{t+1}^n | v_{1:T})}{q(h_{t+1}^n | v_{1:T})} + \log q(h_T^n | v_{1:T}) \right\rangle_{q(h_{1:T}^n | v_{1:T}, z^n = k)} &= \sum_{t=1}^{T-1} \frac{1}{2} \log \det(P_t^T - \overleftarrow{A}_t P_{t+1}^T \overleftarrow{A}_t^\top) \\ &+ \frac{1}{2} \log \det(P_T^T) + \frac{T}{2} H(1 + \log(2\pi)) \end{aligned}$$

where P_t^T is the covariance of $q(h_t^n | v_{1:T})$, $\overleftarrow{A}_t = P_t^t A^\top (A P_t^t A^\top + \Sigma_H)^{-1}$, and where we have used the property

$$\det \begin{pmatrix} A & B \\ C & D \end{pmatrix} = \det(AD - BD^{-1}CD).$$

Terms such as $\langle \log q(A|\Sigma_H^{-1}) / p(A|\Sigma_H^{-1}) \rangle_{q(A, \Sigma_H^{-1})}$ can be computed as

$$\frac{1}{2} \left(\log \left(\frac{\det \Sigma_p}{\det \Sigma_q} \right) + \text{tr} \Sigma_p^{-1} \Sigma_q + (\mu_p - \mu_q)^\top \langle \Sigma_p^{-1} \rangle_{q(\Sigma_H)} (\mu_p - \mu_q) - H \right),$$

where μ_q and Σ_q are the mean and covariance of $q(A|\Sigma_H)$, and μ_p and Σ_p are the mean and covariance of $p(A|\Sigma_H)$. Notice that this simple formula comes from our choice for the prior, which makes Σ_q and Σ_p to have a common dependency on Σ_H and the means μ_q and μ_p not dependent on Σ_H .

Notice also that when $q(z^n = k) = 0$ for a given k and $n = 1, \dots, N$, that is when no sequences are assigned to component k , we have $q(A|\Sigma_H^{-1}) = p(A|\Sigma_H^{-1})$, $q(B|\Sigma_V^{-1}) = p(B|\Sigma_V^{-1})$, $q(\Sigma_H^{-1}) = p(\Sigma_H^{-1})$ and $q(\Sigma_V^{-1}) = p(\Sigma_V^{-1})$ (this can be seen from the updates in Appendix B.2). Therefore all terms coming from component k will give no contribution to the bound. This means that, if we consider only these terms, a model with M mixture components of which only M' are active has the same bound as a model with M' mixture components. However, the term $\langle \log p(z^{1:N}) \rangle_{\prod_{n=1}^N q(z^n)}$ gives a contribution which reflects the fact that our prior on $p(\pi)$ assumes that there are M mixture components. Therefore this term will penalize this model with respect to a model with only M' mixture components.

References

- [1] R. O. Duda, P. E. Hart, and D. G. Stork. *Pattern Classification*. Wiley, New York, 2000.
- [2] Z. Ghahramani and M. J. Beal. Variational inference for Bayesian mixtures of factor analysers. In *Advances in Neural Information Processing Systems 12 (NIPS)*, pages 449–455, 2000.
- [3] D. Görür. *Nonparametric Bayesian Discrete Latent Variable Models for Unsupervised Learning*. Phd thesis, Technischen Universität Berlin, 2007.
- [4] A. Kottas. Dirichlet Process Mixtures of Beta Distributions, with Applications to Density and Intensity Estimation. In *Workshop on Learning with Nonparametric Bayesian Methods, 23rd International Conference on Machine Learning (ICML)*, 2006.
- [5] C. E. Rasmussen. The infinite Gaussian mixture model. In *Advances in Neural Information Processing Systems 12 (NIPS)*, pages 554–560, 2000.
- [6] K. Kalpakis, S. Gada, and V. Puttagunta. Distance measures for effective clustering of ARIMA time-series. pages 273–280, 2001.
- [7] L. Y. Inoue, M. Neira, C. Nelson, M. Gleave, and R. Etzioni. Cluster-based network model for time-course gene expression data. *Biostatistics*, 2007.
- [8] Y. Bar-Shalom and X.-R. Li. *Estimation and Tracking: Principles, Techniques and Software*. Artech House, 1998.
- [9] M. S. Grewal and A. P. Andrews. *Kalman Filtering: Theory and Practice Using MATLAB*. John Wiley and Sons, Inc., 2001.
- [10] R. H. Shumway and D. S. Stoffer. *Time Series Analysis and its Applications*. Springer, 2000.
- [11] J. Durbin and S. J. Koopman. *Time Series Analysis by State Space Methods*. Oxford Univ. Press, 2001.
- [12] M. J. Beal. *Variational Algorithms for Approximate Bayesian Inference*. Ph.D. thesis, Gatsby Computational Neuroscience Unit, University College London, 2003.
- [13] D. J. C. MacKay. Ensemble learning and evidence maximisation. Unpublished manuscript: www.variational-bayes.org, 1995.
- [14] H. Valpola and J. Karhunen. An unsupervised ensemble learning method for nonlinear dynamic state-space models. *Neural Computation*, 14:2647–2692, 2002.
- [15] D. Barber and S. Chiappa. Unified inference for variational Bayesian linear Gaussian state-space models. In *Advances in Neural Information Processing Systems 20 (NIPS)*, 2006.
- [16] M. J. Beal, F. Falciani, Z. Ghahramani, C. Rangel, and D. L. Wild. A Bayesian approach to reconstructing genetic regulatory networks with hidden factors. *Bioinformatics*, 21:349–356, 2005.
- [17] A. T. Cemgil and S. J. Godsill. Probabilistic phase vocoder and its application to interpolation of missing values in audio signals. In *13th European Signal Processing Conference*, 2005.
- [18] M. Morf and T. Kailath. Square-root algorithms for least-squares estimation. *IEEE Transactions on Automatic Control*, 20:487–497, 1975.
- [19] P. Park and T. Kailath. New square-root smoothing algorithms. *IEEE Transactions on Automatic Control*, 41:727–732, 1996.
- [20] D. J. C. MacKay. *Information Theory, Inference and Learning Algorithms*. Cambridge University Press, 2003.
- [21] C. Antoniak. Mixtures of Dirichlet processes with applications to Bayesian nonparametric problems. *The Annals of Statistics*, 2:1152–1174, 1974.

- [22] K. Kurihara, M. Welling, and Y. W. Teh. Collapsed variational Dirichlet process mixture models. In *International Joint Conference on Artificial Intelligence (IJCAI)*, 2007.
- [23] Y. Xiong and D-Y. Yeung. Mixtures of ARMA models for model-based time series clustering. *IEEE International Conference on Data Mining (ICDM)*, pages 717–720, 2002.
- [24] D. J. C. MacKay. Local minima, symmetry-breaking, and model pruning in variational free energy minimization. Inference Group, Cavendish Laboratory, Cambridge, U.K., 2001.

List of Notations

$[x_t]_i$	i -th element of the vector x_t	1
$x_{1:T}$	Shorthand for x_1, \dots, x_T	2
$\mathcal{N}(m, S)$	Gaussian distribution with mean m and covariance S	2
$\mathbf{0}_X$	X -dimensional zero vector	2
x^\top	Transpose of x	3
x_i	i -th element (column) of the vector (matrix) x	3
$[x_t]_{ii}$	ii -th element of the matrix x_t	3
$dg(x)$	Diagonal matrix with the elements of x on the main diagonal	3
$x^{1:N}$	Shorthand for x^1, \dots, x^N	4
$H_q(x)$	Entropy $-\int_x q(x) \log q(x)$	5
$\langle x \rangle_q$	Expectation of x with respect to q	5
$\neg n$	All indices except for n	6
$vc(x)$	Vector formed by vertically stacking the columns of matrix x	6
$x \otimes y$	Kronecker product between x and y	6
I_N	Identity matrix of size $N \times N$	6
$\mathcal{W}(\nu, S)$	Wishart distribution with parameters ν and S	6
$\mathcal{G}(a_1, a_2)$	Gamma distribution with parameters a_1 and a_2	6
$mc(x)$	The matrix formed from stacking the elements of the vector x , columnwise	7
$bdg(x_1, \dots, x_n)$	Block-diagonal matrix with blocks x_1, \dots, x_n	7
$vert(x_1, \dots, x_n)$	Vertical concatenation of x_1, \dots, x_n	7
$x_{i'}$	i -th row of the matrix x	7
$tr[x]$	Trace of x	13