

Variational Learning for Switching State-Space Models

Zoubin Ghahramani

Geoffrey E. Hinton

Gatsby Computational Neuroscience Unit, University College London, London WC1N 3AR, U.K.

We introduce a new statistical model for time series that iteratively segments data into regimes with approximately linear dynamics and learns the parameters of each of these linear regimes. This model combines and generalizes two of the most widely used stochastic time-series models—hidden Markov models and linear dynamical systems—and is closely related to models that are widely used in the control and econometrics literatures. It can also be derived by extending the mixture of experts neural network (Jacobs, Jordan, Nowlan, & Hinton, 1991) to its fully dynamical version, in which both expert and gating networks are recurrent. Inferring the posterior probabilities of the hidden states of this model is computationally intractable, and therefore the exact expectation maximization (EM) algorithm cannot be applied. However, we present a variational approximation that maximizes a lower bound on the log-likelihood and makes use of both the forward and backward recursions for hidden Markov models and the Kalman filter recursions for linear dynamical systems. We tested the algorithm on artificial data sets and a natural data set of respiration force from a patient with sleep apnea. The results suggest that variational approximations are a viable method for inference and learning in switching state-space models.

1 Introduction ---

Most commonly used probabilistic models of time series are descendants of either hidden Markov models (HMM) or stochastic linear dynamical systems, also known as state-space models (SSM). HMMs represent information about the past of a sequence through a single discrete random variable—the hidden state. The prior probability distribution of this state is derived from the previous hidden state using a stochastic transition matrix. Knowing the state at any time makes the past, present, and future observations statistically independent. This is the Markov independence property that gives the model its name.

SSMs represent information about the past through a real-valued hidden state vector. Again, conditioned on this state vector, the past, present, and future observations are statistically independent. The dependency between

the present state vector and the previous state vector is specified through the dynamic equations of the system and the noise model. When these equations are linear and the noise model is gaussian, the SSM is also known as a linear dynamical system or Kalman filter model.

Unfortunately, most real-world processes cannot be characterized by either purely discrete or purely linear-gaussian dynamics. For example, an industrial plant may have multiple discrete modes of behavior, each with approximately linear dynamics. Similarly, the pixel intensities in an image of a translating object vary according to approximately linear dynamics for subpixel translations, but as the image moves over a larger range, the dynamics change significantly and nonlinearly.

This article addresses models of dynamical phenomena that are characterized by a combination of discrete and continuous dynamics. We introduce a probabilistic model called the switching SSM inspired by the divide-and-conquer principle underlying the mixture-of-experts neural network (Jacobs, Jordan, Nowlan, & Hinton, 1991). Switching SSMs are a natural generalization of HMMs and SSMs in which the dynamics can transition in a discrete manner from one linear operating regime to another. There is a large literature on models of this kind in econometrics, signal processing, and other fields (Harrison & Stevens, 1976; Chang & Athans, 1978; Hamilton, 1989; Shumway & Stoffer, 1991; Bar-Shalom & Li, 1993). Here we extend these models to allow for multiple real-valued state vectors, draw connections between these fields and the relevant literature on neural computation and probabilistic graphical models, and derive a learning algorithm for all the parameters of the model based on a structured variational approximation that rigorously maximizes a lower bound on the log-likelihood.

In the following section we review the background material on SSMs, HMMs, and hybrids of the two. In section 3, we describe the generative model—the probability distribution defined over the observation sequences—for switching SSMs. In section 4, we describe a learning algorithm for switching state-space models that is based on a structured variational approximation to the expectation-maximization algorithm. In section 5 we present simulation results in both an artificial domain, to assess the quality of the approximate inference method, and a natural domain. We conclude with section 6.

2 Background

2.1 State-Space Models. An SSM defines a probability density over time series of real-valued observation vectors $\{Y_t\}$ by assuming that the observations were generated from a sequence of hidden state vectors $\{X_t\}$. (Appendix A describes the variables and notation used throughout this article.) In particular, the SSM specifies that given the hidden state vector at one time step, the observation vector at that time step is statistically independent from all other observation vectors, and that the hidden state vectors

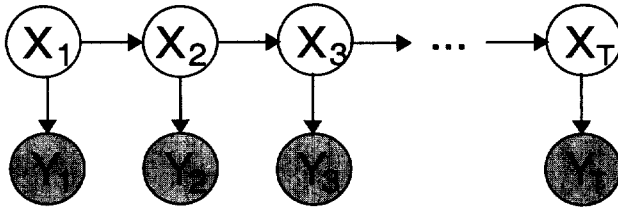


Figure 1: A directed acyclic graph (DAG) specifying conditional independence relations for a state-space model. Each node is conditionally independent from its nondescendants given its parents. The output Y_t is conditionally independent from all other variables given the state X_t ; and X_t is conditionally independent from X_1, \dots, X_{t-2} given X_{t-1} . In this and the following figures, shaded nodes represent observed variables, and unshaded nodes represent hidden variables.

obey the Markov independence property. The joint probability for the sequences of states X_t and observations Y_t can therefore be factored as:

$$P(\{X_t, Y_t\}) = P(X_1)P(Y_1|X_1) \prod_{t=2}^T P(X_t|X_{t-1})P(Y_t|X_t). \quad (2.1)$$

The conditional independencies specified by equation 2.1 can be expressed graphically in the form of Figure 1. The simplest and most commonly used models of this kind assume that the transition and output functions are linear and time invariant and the distributions of the state and observation variables are multivariate gaussian. We use the term *state-space model* to refer to this simple form of the model. For such models, the state transition function is

$$X_t = AX_{t-1} + w_t, \quad (2.2)$$

where A is the state transition matrix and w_t is zero-mean gaussian noise in the dynamics, with covariance matrix Q . $P(X_1)$ is assumed to be gaussian. Equation 2.2 ensures that if $P(X_{t-1})$ is gaussian, so is $P(X_t)$. The output function is

$$Y_t = CX_t + v_t, \quad (2.3)$$

where C is the output matrix and v_t is zero-mean gaussian output noise with covariance matrix R ; $P(Y_t|X_t)$ is therefore also gaussian:

$$P(Y_t|X_t) = (2\pi)^{-D/2} |R|^{-1/2} \exp \left\{ -\frac{1}{2} (Y_t - CX_t)' R^{-1} (Y_t - CX_t) \right\}, \quad (2.4)$$

where D is the dimensionality of the Y vectors.

Often the observation vector can be divided into input (or predictor) variables and output (or response) variables. To model the input-output behavior of such a system—the conditional probability of output sequences given input sequences—the linear gaussian SSM can be modified to have a state-transition function,

$$X_t = AX_{t-1} + BU_t + w_t, \quad (2.5)$$

where U_t is the input observation vector and B is the (fixed) input matrix.¹

The problem of *inference* or *state estimation* for an SSM with known parameters consists of estimating the posterior probabilities of the hidden variables given a sequence of observed variables. Since the local likelihood functions for the observations are gaussian and the priors for the hidden states are gaussian, the resulting posterior is also gaussian. Three special cases of the inference problem are often considered: filtering, smoothing, and prediction (Anderson & Moore, 1979; Goodwin & Sin, 1984). The goal of filtering is to compute the probability of the current hidden state X_t given the sequence of inputs and outputs up to time t — $P(X_t | \{Y\}_1^t, \{U\}_1^t)$.² The recursive algorithm used to perform this computation is known as the *Kalman filter* (Kalman & Bucy, 1961). The goal of *smoothing* is to compute the probability of X_t given the sequence of inputs and outputs up to time T , where $T > t$. The Kalman filter is used in the forward direction to compute the probability of X_t given $\{Y\}_1^t$ and $\{U\}_1^t$. A similar set of *backward* recursions from T to t completes the computation by accounting for the observations after time t (Rauch, 1963). We will refer to the combined forward and backward recursions for smoothing as the Kalman smoothing recursions (also known as the RTS, or Rauch-Tung-Streifel smoother). Finally, the goal of *prediction* is to compute the probability of future states and observations given observations up to time t . Given $P(X_t | \{Y\}_1^t, \{U\}_1^t)$ computed as before, the model is simulated in the forward direction using equations 2.2 (or 2.5 if there are inputs) and 2.3 to compute the probability density of the state or output at future time $t + \tau$.

The problem of learning the parameters of an SSM is known in engineering as the *system identification* problem; in its most general form it assumes access only to sequences of input and output observations. We focus on maximum likelihood learning in which a single (locally optimal) value of the parameters is estimated, rather than Bayesian approaches that treat the parameters as random variables and compute or approximate the posterior distribution of the parameters given the data. One can also distinguish between on-line and off-line approaches to learning. On-line recursive algorithms, favored in real-time adaptive control applications, can be obtained by computing the gradient or the second derivatives of the log-likelihood (Ljung

¹ One can also define the state such that $X_{t+1} = AX_t + BU_t + w_t$.

² The notation $\{Y\}_1^t$ is shorthand for the sequence Y_1, \dots, Y_t .

& Söderström, 1983). Similar gradient-based methods can be obtained for off-line methods. An alternative method for off-line learning makes use of the expectation maximization (EM) algorithm (Dempster, Laird, & Rubin, 1977). This procedure iterates between an E-step that fixes the current parameters and computes posterior probabilities over the hidden states given the observations, and an M-step that maximizes the expected log-likelihood of the parameters using the posterior distribution computed in the E-step. For linear gaussian state-space models, the E-step is exactly the Kalman smoothing problem as defined above, and the M-step simplifies to a linear regression problem (Shumway & Stoffer, 1982; Digalakis, Rohlicek, & Ostendorf, 1993). Details on the EM algorithm for SSMs can be found in Ghahramani and Hinton (1996a), as well as in the original Shumway and Stoffer (1982) article.

2.2 Hidden Markov Models. Hidden Markov models also define probability distributions over sequences of observations $\{Y_t\}$. The distribution over sequences is obtained by specifying a distribution over observations at each time step t given a *discrete* hidden state S_t , and the probability of transitioning from one hidden state to another. Using the Markov property, the joint probability for the sequences of states S_t and observations Y_t , can be factored in exactly the same manner as equation 2.1, with S_t taking the place of X_t :

$$P(\{S_t, Y_t\}) = P(S_1)P(Y_1|S_1) \prod_{t=2}^T P(S_t|S_{t-1})P(Y_t|S_t). \quad (2.6)$$

Similarly, the conditional independencies in an HMM can be expressed graphically in the same form as Figure 1. The state is represented by a single multinomial variable that can take one of K discrete values, $S_t \in \{1, \dots, K\}$. The state transition probabilities, $P(S_t|S_{t-1})$, are specified by a $K \times K$ transition matrix. If the observables are discrete symbols taking on one of L values, the observation probabilities $P(Y_t|S_t)$ can be fully specified as a $K \times L$ observation matrix. For a continuous observation vector, $P(Y_t|S_t)$ can be modeled in many different forms, such as a gaussian, mixture of gaussians, or neural network. HMMs have been applied extensively to problems in speech recognition (Juang & Rabiner, 1991), computational biology (Baldi, Chauvin, Hunkapiller, & McClure, 1994), and fault detection (Smyth, 1994).

Given an HMM with known parameters and a sequence of observations, two algorithms are commonly used to solve two different forms of the inference problem (Rabiner & Juang, 1986). The first computes the posterior probabilities of the hidden states using a recursive algorithm known as the *forward-backward* algorithm. The computations in the forward pass are exactly analogous to the Kalman filter for SSMs, and the computations in the backward pass are analogous to the backward pass of the Kalman smoothing

equations. As noted by Bridle (pers. comm., 1985) and Smyth, Heckerman, and Jordan (1997), the forward-backward algorithm is a special case of exact inference algorithms for more general graphical probabilistic models (Lauritzen & Spiegelhalter, 1988; Pearl, 1988). The same observation holds true for the Kalman smoothing recursions. The other inference problem commonly posed for HMMs is to compute the single most likely sequence of hidden states. The solution to this problem is given by the Viterbi algorithm, which also consists of a forward and backward pass through the model.

To learn maximum likelihood parameters for an HMM given sequences of observations, one can use the well-known Baum-Welch algorithm (Baum, Petrie, Soules, & Weiss, 1970). This algorithm is a special case of EM that uses the forward-backward algorithm to infer the posterior probabilities of the hidden states in the E-step. The M-step uses expected counts of transitions and observations to reestimate the transition and output matrices (or linear regression equations in the case where the observations are gaussian distributed). Like SSMs, HMMs can be augmented to allow for input variables, such that they model the conditional distribution of sequences of output observations given sequences of inputs (Cacciatore & Nowlan, 1994; Bengio & Frasconi, 1995; Meila & Jordan, 1996).

2.3 Hybrids. A burgeoning literature on models that combine the discrete transition structure of HMMs with the linear dynamics of SSMs has developed in fields ranging from econometrics to control engineering (Harrison & Stevens, 1976; Chang & Athans, 1978; Hamilton, 1989; Shumway & Stoffer, 1991; Bar-Shalom & Li, 1993; Deng, 1993; Kadirkamanathan & Kadirkamanathan, 1996; Chaer, Bishop, & Ghosh, 1997). These models are known alternately as hybrid models, SSMs with switching, and jump-linear systems. We briefly review some of this literature, including some related neural network models.³

Shortly after Kalman and Bucy solved the problem of state estimation for linear gaussian SSMs, attention turned to the analogous problem for switching models (Ackerson & Fu, 1970). Chang and Athans (1978) derive the equations for computing the conditional mean and variance of the state when the parameters of a linear SSM switch according to arbitrary and Markovian dynamics. The prior and transition probabilities of the switching process are assumed to be known. They note that for M models (sets of parameters) and an observation length T , the exact conditional distribution of the state is a gaussian mixture with M^T components. The conditional mean and variance, which require far less computation, are therefore only summary statistics.

³ A review of how SSMs and HMMs are related to simpler statistical models such as principal components analysis, factor analysis, mixture of gaussians, vector quantization, and independent components analysis (ICA) can be found in Roweis and Ghahramani (1999).

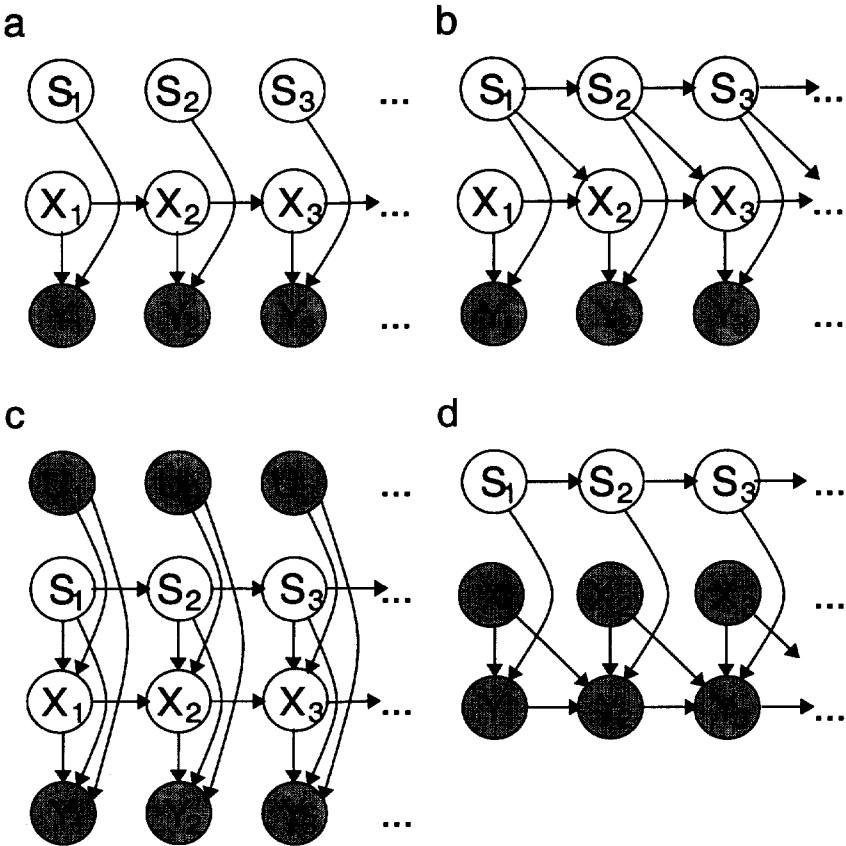


Figure 2: Directed acyclic graphs specifying conditional independence relations for various switching state-space models. (a) Shumway and Stoffer (1991): the output matrix (C in equation 2.3) switches independently between a fixed number of choices at each time step. Its setting is represented by the discrete hidden variable S_t ; (b) Bar-Shalom and Li (1993): both the output equation and the dynamic equation can switch, and the switches are Markov; (c) Kim (1994); (d) Fraser and Dimitriadis (1993): outputs and states are observed. Here we have shown a simple case where the output depends directly on the current state, previous state, and previous output.

Shumway and Stoffer (1991) consider the problem of learning the parameters of SSMs with a single real-valued hidden state vector and switching output matrices. The probability of choosing a particular output matrix is a prespecified time-varying function, independent of previous choices (see Figure 2a). A pseudo-EM algorithm is derived in which the E-step, which

in its exact form would require computing a gaussian mixture with M^T components, is approximated by a single gaussian at each time step.

Bar-Shalom and Li (1993; sec. 11.6) review models in which both the state dynamics and the output matrices switch, and where the switching follows Markovian dynamics (see Figure 2b). They present several methods for approximately solving the state-estimation problem in switching models (they do not discuss parameter estimation for such models). These methods, referred to as generalized pseudo-Bayesian (GPB) and interacting multiple models (IMM), are all based on the idea of collapsing into one gaussian the mixture of M gaussians that results from considering all the settings of the switch state at a given time step. This avoids the exponential growth of mixture components at the cost of providing an approximate solution. More sophisticated but computationally expensive methods that collapse M^2 gaussians into M gaussians are also derived. Kim (1994) derives a similar approximation for a closely related model, which also includes observed input variables (see Figure 2c). Furthermore, Kim discusses parameter estimation for this model, although without making reference to the EM algorithm. Other authors have used Markov chain Monte Carlo methods for state and parameter estimation in switching models (Carter & Kohn, 1994; Athaide, 1995) and in other related dynamic probabilistic networks (Dean & Kanazawa, 1989; Kanazawa, Koller, & Russell, 1995).

Hamilton (1989, 1994, sec. 22.4) describes a class of switching models in which the real-valued observation at time t , Y_t , depends on both the observations at times $t-1$ to $t-r$ and the discrete states at time t to $t-r$. More precisely, Y_t is gaussian with mean that is a linear function of Y_{t-1}, \dots, Y_{t-r} and of binary indicator variables for the discrete states, S_t, \dots, S_{t-r} . The system can therefore be seen as an $(r+1)$ th order HMM driving an r th order autoregressive process, and is tractable for small r and number of discrete states in S .

Hamilton's models are closely related to hidden filter HMM (HFHMM; Fraser & Dimitriadis, 1993). HFHMMs have both discrete and real-valued states. However, the real-valued states are assumed to be either observed or a known, deterministic function of the past observations (i.e., an embedding). The outputs depend on the states and previous outputs, and the form of this dependence can switch randomly (see Figure 2d). Because at any time step the only hidden variable is the switch state, S_t , exact inference in this model can be carried out tractably. The resulting algorithm is a variant of the forward-backward procedure for HMMs. Kehagias and Petridis (1997) and Pawelzik, Kohlmorgen, and Müller (1996) present other variants of this model.

Elliott, Aggoun, and Moore (1995; sec. 12.5) present an inference algorithm for hybrid (Markov switching) systems for which there is a separate observable from which the switch state can be estimated. The true switch states, S_t , are represented as unit vectors in \mathbb{R}^M , and the estimated switch state is a vector in the unit square with elements corresponding to the es-

estimated probability of being in each switch state. The real-valued state, X_t , is approximated as a gaussian given the estimated switch state by forming a linear combination of the transition and observation matrices for the different SSMs weighted by the estimated switch state. Elliott et al. also derive control equations for such hybrid systems and discuss applications of the change-of-measures whitening procedure to a large family of models.

With regard to the literature on neural computation, the model presented in this article is a generalization of both the mixture-of-experts neural network (Jacobs et al., 1991; Jordan & Jacobs, 1994) and the related mixture of factor analyzers (Hinton, Dayan, & Revow, 1997; Ghahramani & Hinton, 1996b). Previous dynamical generalizations of the mixture-of-experts architecture consider the case in which the gating network has Markovian dynamics (Cacciatore & Nowlan, 1994; Kadiramanathan & Kadiramanathan, 1996; Meila & Jordan, 1996). One limitation of this generalization is that the entire past sequence is summarized in the value of a single discrete variable (the gating activation), which for a system with M experts can convey on average at most $\log M$ bits of information about the past. In the models we consider here, both the experts and the gating network have Markovian dynamics. The past is therefore summarized by a state composed of the cross-product of the discrete variable and the combined real-valued state-space of all the experts. This provides a much wider information channel from the past. One advantage of this representation is that the real-valued state can contain componential structure. Thus, attributes such as the position, orientation, and scale of an object in an image, which are most naturally encoded as independent real-valued variables, can be accommodated in the state without the exponential growth required of a discretized HMM-like representation.

It is important to place the work in this article in the context of the literature we have just reviewed. The hybrid models, state-space models with switching, and jump-linear systems we have described all assume a single real-valued state vector. The model considered in this article generalizes this to multiple real-valued state vectors.⁴ Unlike the models described in Hamilton (1994), Fraser and Dimitradis (1993), and the current dynamical extensions of mixtures of experts, in the model we present, the real-valued state vectors are hidden. The inference algorithm we derive, which is based on making a structured variational approximation, is entirely novel in the context of switching SSMs. Specifically, our method is unlike all the approximate methods we have reviewed in that it is not based on fitting a single gaussian to a mixture of gaussians by computing the mean and covariance of the mixture.⁵ We derive a learning algorithm for all of the parameters

⁴ Note that the state vectors could be concatenated into one large state vector with factorized (block-diagonal) transition matrices (cf. factorial hidden Markov model; Ghahramani & Jordan, 1997). However, this obscures the decoupled structure of the model.

⁵ Both classes of methods can be seen as minimizing Kullback-Liebler (KL) diver-

of the model, including the Markov switching parameters. This algorithm maximizes a lower bound on the log-likelihood of the data rather than a heuristically motivated approximation to the likelihood. The algorithm has a simple and intuitive flavor: It decouples into forward-backward recursions on a HMM, and Kalman smoothing recursions on each SSM. The states of the HMM determine the soft assignment of each observation to an SSM; the prediction errors of the SSMs determine the observation probabilities for the HMM.

3 The Generative Model

In switching SSMs, the sequence of observations $\{Y_t\}$ is modeled by specifying a probabilistic relation between the observations and a hidden state-space comprising M real-valued state vectors, $X_t^{(m)}$, and one discrete state vector S_t . The discrete state, S_t , is modeled as a multinomial variable that can take on M values: $S_t \in \{1, \dots, M\}$; for reasons that will become obvious we refer to it as the *switch* variable. The joint probability of observations and hidden states can be factored as

$$\begin{aligned} P\left(\left\{S_t, X_t^{(1)}, \dots, X_t^{(M)}, Y_t\right\}\right) \\ = P(S_1) \prod_{t=2}^T P(S_t | S_{t-1}) \cdot \prod_{m=1}^M P\left(X_1^{(m)}\right) \prod_{t=2}^T P\left(X_t^{(m)} | X_{t-1}^{(m)}\right) \\ \cdot \prod_{t=1}^T P\left(Y_t | X_t^{(1)}, \dots, X_t^{(M)}, S_t\right), \end{aligned} \quad (3.1)$$

which corresponds graphically to the conditional independencies represented by Figure 3. Conditioned on a setting of the switch state, $S_t = m$, the observable is multivariate gaussian with output equation given by state-space model m . Notice that m is used as both an index for the real-valued state variables and a value for the switch state. The probability of the observation vector Y_t is therefore

$$\begin{aligned} P\left(Y_t | X_t^{(1)}, \dots, X_t^{(M)}, S_t = m\right) \\ = |2\pi R|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2} \left(Y_t - C^{(m)} X_t^{(m)}\right)' R^{-1} \left(Y_t - C^{(m)} X_t^{(m)}\right)\right\}, \end{aligned} \quad (3.2)$$

where R is the observation noise covariance matrix and $C^{(m)}$ is the output matrix for SSM m (cf. equation 2.4 for a single linear-gaussian SSM). Each

gences. However, the KL divergence is asymmetrical, and whereas the variational methods minimize it in one direction, the methods that merge gaussians minimize it in the other direction. We return to this point in section 4.2.

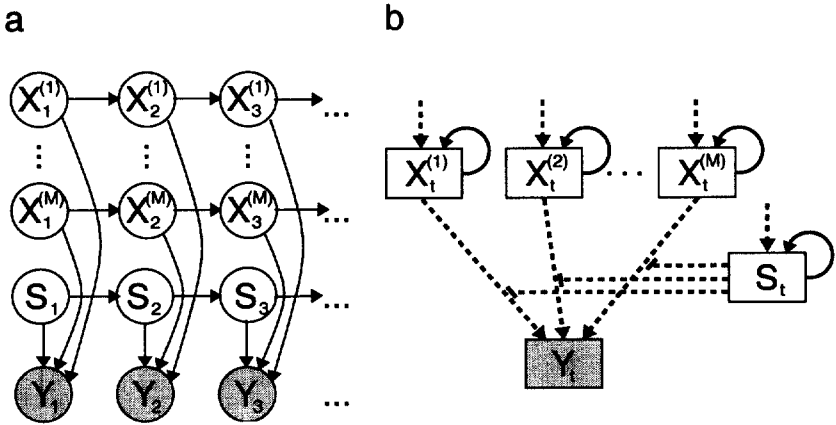


Figure 3: (a) Graphical model representation for switching state-space models. S_t is the discrete switch variable, and $X_t^{(m)}$ are the real-valued state vectors. (b) Switching state-space model depicted as a generalization of the mixture of experts. The dashed arrows correspond to the connections in a mixture of experts. In a switching state-space model, the states of the experts and the gating network also depend on their previous states (solid arrows).

real-valued state vector evolves according to the linear gaussian dynamics of an SSM with differing initial state, transition matrix, and state noise (see equation 2.2). For simplicity we assume that all state vectors have identical dimensionality; the generalization of the algorithms we present to models with different-sized state-spaces is immediate. The switch state itself evolves according to the discrete Markov transition structure specified by the initial state probabilities $P(S_1)$ and the $M \times M$ state transition matrix $P(S_t|S_{t-1})$.

An exact analogy can be made to the mixture-of-experts architecture for modular learning in neural networks (see Figure 3b; Jacobs et al., 1991). Each SSM is a linear expert with gaussian output noise model and linear-gaussian dynamics. The switch state “gates” the outputs of the M SSMs, and therefore plays the role of a gating network with Markovian dynamics.

There are many possible extensions of the model; we shall consider three obvious and straightforward ones:

(Ex1) Differing output covariances, $R^{(m)}$, for each SSM;

(Ex2) Differing output means, $\mu_Y^{(m)}$, for each SSM, such that each model is allowed to capture observations in a different operating range

(Ex3) Conditioning on a sequence of observed input vectors, $\{U_t\}$

4 Learning

An efficient learning algorithm for the parameters of a switching SSM can be derived by generalizing the EM algorithm (Baum et al., 1970; Dempster et al., 1977). EM alternates between optimizing a distribution over the hidden states (the E-step) and optimizing the parameters given the distribution over hidden states (the M-step). Any distribution over the hidden states, $Q(\{S_t, X_t\})$, where $X_t = [X_t^{(1)}, \dots, X_t^{(M)}]$ is the combined state of the SSMs, can be used to define a lower bound, \mathcal{B} , on the log probability of the observed data:

$$\log P(\{Y_t\}|\theta) = \log \sum_{\{S_t\}} \int P(\{S_t, X_t, Y_t\}|\theta) d\{X_t\} \quad (4.1)$$

$$= \log \sum_{\{S_t\}} \int Q(\{S_t, X_t\}) \left[\frac{P(\{S_t, X_t, Y_t\}|\theta)}{Q(\{S_t, X_t\})} \right] d\{X_t\} \quad (4.2)$$

$$\begin{aligned} &\geq \sum_{\{S_t\}} \int Q(\{S_t, X_t\}) \log \left[\frac{P(\{S_t, X_t, Y_t\}|\theta)}{Q(\{S_t, X_t\})} \right] d\{X_t\} \\ &= \mathcal{B}(Q, \theta), \end{aligned} \quad (4.3)$$

where θ denotes the parameters of the model and we have made use of Jensen's inequality (Cover & Thomas, 1991) to establish equation 4.3. Both steps of EM increase the lower bound on the log probability of the observed data. The E-step holds the parameters fixed and sets Q to be the posterior distribution over the hidden states given the parameters,

$$Q(\{S_t, X_t\}) = P(\{S_t, X_t\}|\{Y_t\}, \theta). \quad (4.4)$$

This maximizes \mathcal{B} with respect to the distribution, turning the lower bound into an equality, which can be easily seen by substitution. The M-step holds the distribution fixed and computes the parameters that maximize \mathcal{B} for that distribution. Since $\mathcal{B} = \log P(\{Y_t\}|\theta)$ at the start of the M-step and since the E-step does not affect $\log P$, the two steps combined can never decrease $\log P$. Given the change in the parameters produced by the M-step, the distribution produced by the previous E-step is typically no longer optimal, so the whole procedure must be iterated.

Unfortunately, the exact E-step for switching SSMs is intractable. Like the related hybrid models described in section 2.3, the posterior probability of the real-valued states is a gaussian mixture with M^T terms. This can be seen by using the semantics of directed graphs, in particular the d -separation criterion (Pearl, 1988), which implies that the hidden state variables in Figure 3, while marginally independent, become conditionally dependent given the observation sequence. This induced dependency effectively couples all of the real-valued hidden state variables to the discrete switch variable, as a

consequence of which the exact posteriors become Gaussian mixtures with an exponential number of terms.⁶

In order to derive an efficient learning algorithm for this system, we relax the EM algorithm by approximating the posterior probability of the hidden states. The basic idea is that since expectations with respect to P are intractable, rather than setting $Q(\{S_t, X_t\}) = P(\{S_t, X_t\}|\{Y_t\})$ in the E-step, a tractable distribution Q is used to *approximate* P . This results in an EM learning algorithm that maximizes a lower bound on the log-likelihood. The difference between the bound \mathcal{B} and the log-likelihood is given by the Kullback-Liebler (KL) divergence between Q and P :

$$\text{KL}(Q\|P) = \sum_{\{S_t\}} \int Q(\{S_t, X_t\}) \log \left[\frac{Q(\{S_t, X_t\})}{P(\{S_t, X_t\}|\{Y_t\})} \right] d\{X_t\}. \quad (4.5)$$

Since the complexity of exact inference in the approximation given by Q is determined by its conditional independence relations, not by its parameters, we can choose Q to have a tractable structure—a graphical representation that eliminates some of the dependencies in P . Given this structure, the parameters of Q are varied to obtain the tightest possible bound by minimizing equation 4.5. Therefore, the algorithm alternates between optimizing the parameters of the distribution Q to minimize equation 4.5 (the E-step) and optimizing the parameters of P given the distribution over the hidden states (the M-step). As in exact EM, both steps increase the lower bound \mathcal{B} on the log-likelihood; however, equality is not reached in the E-step.

We will refer to the general strategy of using a parameterized approximating distribution as a *variational approximation* and refer to the free parameters of the distribution as *variational parameters*. A completely factorized approximation is often used in statistical physics, where it provides the basis for simple yet powerful mean-field approximations to statistical mechanical systems (Parisi, 1988). Theoretical arguments motivating approximate E-steps are presented in Neal and Hinton (1998; originally in a technical report in 1993). Saul and Jordan (1996) showed that approximate E-steps could be used to maximize a lower bound on the log-likelihood, and proposed the powerful technique of structured variational approximations to intractable probabilistic networks. The key insight of their work, which this article makes use of, is that by judicious use of an approximation Q , exact inference algorithms can be used on the tractable substructures in an intractable network. A general tutorial on variational approximations can be found in Jordan, Ghahramani, Jaakkola, and Saul (1998).

⁶ The intractability of the E-step or smoothing problem in the simpler single-state switching model has been noted by Ackerson and Fu (1970), Chang and Athans (1978), Bar-Shalom and Li (1993), and others.

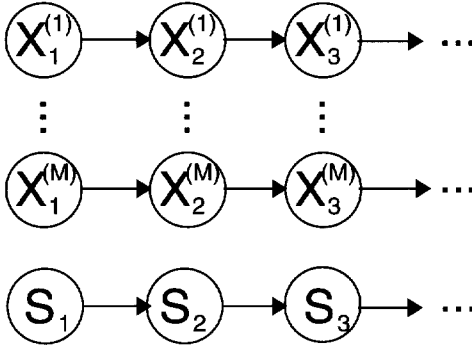


Figure 4: Graphical model representation for the structured variational approximation to the posterior distribution of the hidden states of a switching state-space model.

The parameters of the switching SSM are $\theta = \{A^{(m)}, C^{(m)}, Q^{(m)}, \mu_{X_1}^{(m)}, Q_1^{(m)}, R, \pi, \Phi\}$, where $A^{(m)}$ is the state dynamics matrix for model m , $C^{(m)}$ is its output matrix, $Q^{(m)}$ is its state noise covariance, $\mu_{X_1}^{(m)}$ is the mean of the initial state, $Q_1^{(m)}$ is the covariance of the initial state, R is the (tied) output noise covariance, $\pi = P(S_1)$ is the prior for the discrete Markov process, and $\Phi = P(S_t|S_{t-1})$ is the discrete transition matrix. Extensions (Ex1) through (Ex3) can be readily implemented by substituting $R^{(m)}$ for R , adding means $\mu_Y^{(m)}$ and input matrices $B^{(m)}$.

Although there are many possible approximations to the posterior distribution of the hidden variables that one could use for learning and inference in switching SSMs, we focus on the following:

$$Q(\{S_t, X_t\}) = \frac{1}{Z_Q} \left[\psi(S_1) \prod_{t=2}^T \psi(S_{t-1}, S_t) \right] \prod_{m=1}^M \psi(X_1^{(m)}) \cdot \prod_{t=2}^T \psi(X_{t-1}^{(m)}, X_t^{(m)}), \quad (4.6)$$

where the ψ are unnormalized probabilities, which we will call *potential* functions and define soon, and Z_Q is a normalization constant ensuring that Q integrates to one. Although Q has been written in terms of potential functions rather than conditional probabilities, it corresponds to the simple graphical model shown in Figure 4. The terms involving the switch variables S_t define a discrete Markov chain, and the terms involving the state vectors $X_t^{(m)}$ define M *uncoupled* SSMs. As in mean-field approximations, we have approximated the stochastically coupled system by removing some of the

couplings of the original system. Specifically, we have removed the stochastic coupling between the chains that results from the fact that the observation at time t depends on all the hidden variables at time t . However, we retain the coupling between the hidden variables at successive time steps since these couplings can be handled exactly using the forward-backward and Kalman smoothing recursions. This approximation is therefore structured, in the sense that not all variables are uncoupled.

The discrete switching process is defined by

$$\psi(S_1 = m) = P(S_1 = m) q_1^{(m)} \quad (4.7)$$

$$\psi(S_{t-1}, S_t = m) = P(S_t = m | S_{t-1}) q_t^{(m)}, \quad (4.8)$$

where the $q_t^{(m)}$ are variational parameters of the Q distribution. These parameters scale the probabilities of each of the states of the switch variable at each time step, so that $q_t^{(m)}$ plays exactly the same role that the observation probability $P(Y_t | S_t = m)$ would play in a regular HMM. We will soon see that minimizing $KL(Q||P)$ results in an equation for $q_t^{(m)}$ that supports this intuition.

The uncoupled SSMs in the approximation Q are also defined by potential functions that are related to probabilities in the original system. These potentials are the prior and transition probabilities for $X^{(m)}$ multiplied by a factor that changes these potentials to try to account for the data:

$$\psi(X_1^{(m)}) = P(X_1^{(m)}) \left[P(Y_1 | X_1^{(m)}, S_1 = m) \right]^{h_1^{(m)}} \quad (4.9)$$

$$\psi(X_{t-1}^{(m)}, X_t^{(m)}) = P(X_t^{(m)} | X_{t-1}^{(m)}) \left[P(Y_t | X_t^{(m)}, S_t = m) \right]^{h_t^{(m)}}, \quad (4.10)$$

where the $h_t^{(m)}$ are variational parameters of Q . The vector h_t plays a role very similar to the switch variable S_t . Each component $h_t^{(m)}$ can range between 0 and 1. When $h_t^{(m)} = 0$, the posterior probability of $X_t^{(m)}$ under Q does not depend on the observation at time Y_t . When $h_t^{(m)} = 1$, the posterior probability of $X_t^{(m)}$ under Q includes a term that assumes that SSM m generated Y_t . We call $h_t^{(m)}$ the *responsibility* assigned to SSM m for the observation vector Y_t . The difference between $h_t^{(m)}$ and $S_t^{(m)}$ is that $h_t^{(m)}$ is a deterministic parameter, while $S_t^{(m)}$ is a stochastic random variable.

To maximize the lower bound on the log-likelihood, $KL(Q||P)$ is minimized with respect to the variational parameters $h_t^{(m)}$ and $q_t^{(m)}$ separately for each sequence of observations. Using the definition of P for the switching state-space model (equations 3.1 and 3.2) and the approximating distribution Q , the minimum of KL satisfies the following fixed-point equations for the variational parameters (see appendix B):

$$h_t^{(m)} = Q(S_t = m) \quad (4.11)$$

$$q_t^{(m)} = \exp \left\{ -\frac{1}{2} \left\langle \left(Y_t - C^{(m)} X_t^{(m)} \right)' R^{-1} \left(Y_t - C^{(m)} X_t^{(m)} \right) \right\rangle \right\}, \quad (4.12)$$

where $\langle \cdot \rangle$ denotes expectation over the Q distribution. Intuitively, the responsibility, $h_t^{(m)}$ is equal to the probability under Q that SSM m generated observation vector Y_t , and $q_t^{(m)}$ is an unnormalized gaussian function of the expected squared error if SSM m generated Y_t .

To compute $h_t^{(m)}$ it is necessary to sum Q over all the S_τ variables not including S_t . This can be done efficiently using the forward-backward algorithm on the switch state variables, with $q_t^{(m)}$ playing exactly the same role as an observation probability associated with each setting of the switch variable. Since $q_t^{(m)}$ is related to the prediction error of model m on data Y_t , this has the intuitive interpretation that the switch state associated with models with smaller expected prediction error on a particular observation will be favored at that time step. However, the forward-backward algorithm ensures that the final responsibilities for the models are obtained after considering the entire sequence of observations.

To compute $q_t^{(m)}$, it is necessary to calculate the expectations of $X_t^{(m)}$ and $X_t^{(m)} X_t^{(m)'}$ under Q . We see this by expanding equation 4.12:

$$q_t^{(m)} = \exp \left\{ -\frac{1}{2} Y_t' R^{-1} Y_t + Y_t' R^{-1} C^{(m)} \langle X_t^{(m)} \rangle - \frac{1}{2} \text{tr} \left[C^{(m)'} R^{-1} C^{(m)} \langle X_t^{(m)} X_t^{(m)'} \rangle \right] \right\}, \quad (4.13)$$

where tr is the matrix trace operator, and we have used $\text{tr}(AB) = \text{tr}(BA)$. The expectations of $X_t^{(m)}$ and $X_t^{(m)} X_t^{(m)'}$ can be computed efficiently using the Kalman smoothing algorithm on each SSM, where for model m at time t , the data are weighted by the responsibilities $h_t^{(m)}$.⁷ Since the h parameters depend on the q parameters, and vice versa, the whole process has to be iterated, where each iteration involves calls to the forward-backward and Kalman smoothing algorithms. Once the iterations have converged, the E-step outputs the expected values of the hidden variables under the final Q .

The M-step computes the model parameters that optimize the expectation of the log-likelihood (see equation B.7), which is a function of the expectations of the hidden variables. For switching SSMs, all the parameter reestimates can be computed analytically. For example, taking derivatives of the expectation of equation B.7 with respect to $C^{(m)}$ and setting to zero,

⁷ Weighting the data by $h_t^{(m)}$ is equivalent to running the Kalman smoother on the unweighted data using a time-varying observation noise covariance matrix $R_t^{(m)} = R/h_t^{(m)}$.

Initialize parameters of the model.

Repeat until bound on log likelihood has converged:

E step Repeat until convergence of $KL(Q||P)$:

E.1 Compute $q_t^{(m)}$ from the prediction error of state-space model m on observation Y_t

E.2 Compute $h_t^{(m)}$ using the forward-backward algorithm on the HMM, with observation probabilities $q_t^{(m)}$

E.3 For $m = 1$ to M

Run Kalman smoothing recursions, using the data weighted by $h_t^{(m)}$

M step

M.1 Re-estimate parameters for each state-space model using the data weighted by $h_t^{(m)}$

M.2 Re-estimate parameters for the switching process using Baum-Welch update equations.

Figure 5: Learning algorithm for switching state-space models.

we get

$$\hat{C}^{(m)} = \left(\sum_{t=1}^T \langle S_t^{(m)} \rangle Y_t \langle X_t^{(m)'} \rangle \right) \left(\sum_{t=1}^T \langle S_t^{(m)} \rangle \langle X_t^{(m)} X_t^{(m)'} \rangle \right)^{-1}, \quad (4.14)$$

which is a weighted version of the reestimation equations for SSMs. Similarly, the reestimation equations for the switch process are analogous to the Baum-Welch update rules for HMMs. The learning algorithm for switching state-space models using the above structured variational approximation is summarized in Figure 5.

4.1 Deterministic Annealing. The KL divergence minimized in the E-step of the variational EM algorithm can have multiple minima in general. One way to visualize these minima is to consider the space of all possible segmentations of an observation sequence of length T , where by segmentation we mean a discrete partition of the sequence between the SSMs. If there are M SSMs, then there are M^T possible segmentations of the sequence. Given one such segmentation, inferring the optimal distribution for the real-valued states of the SSMs is a convex optimization problem, since these real-valued states are conditionally gaussian. So the difficulty in the KL minimization lies in trying to find the best (soft) partition of the data.

As in other combinatorial optimization problems, the possibility of getting trapped in local minima can be reduced by gradually annealing the cost function. We can employ a deterministic variant of the annealing idea by making the following simple modifications to the variational fixed-point equations 4.11 and 4.12:

$$h_t^{(m)} = \frac{1}{\mathcal{T}} Q(S_t = m) \quad (4.15)$$

$$q_t^{(m)} = \exp \left\{ -\frac{1}{2\mathcal{T}} \left\langle \left(Y_t - C^{(m)} X_t^{(m)} \right)' R^{-1} \left(Y_t - C^{(m)} X_t^{(m)} \right) \right\rangle \right\}. \quad (4.16)$$

Here \mathcal{T} is a temperature parameter, which is initialized to a large value and gradually reduced to 1. The above equations maximize a modified form of the bound \mathcal{B} in equation 4.3, where the entropy of Q has been multiplied by \mathcal{T} (Ueda & Nakano, 1995).

4.2 Merging Gaussians. Almost all the approximate inference methods that are described in the literature for switching SSMs are based on the idea of merging, at each time step, a mixture of M gaussians into one gaussian. The merged gaussian is obtained by setting its mean and covariance equal to the mean and covariance of the mixture. Here we briefly describe, as an alternative to the variational approximation methods we have derived, how this more traditional gaussian merging procedure can be applied to the model we have defined.

In the switching state-space models described in section 3, there are M different SSMs, with possibly different state-space dimensionalities, so it would be inappropriate to merge their states into one gaussian. However, it is still possible to apply a gaussian merging technique by considering each SSM separately. In each SSM, m , the hidden state density produces at each time step a mixture of two gaussians: one for the case $S_t = m$ and one for $S_t \neq m$. We merge these two gaussians, weighted the current estimates of $P(S_t = m | Y_1, \dots, Y_t)$ and $1 - P(S_t = m | Y_1, \dots, Y_t)$, respectively. This merged gaussian is used to obtain the gaussian prior for $X_{t+1}^{(m)}$ for the next time step. We implemented a forward-pass version of this approximate inference scheme, which is analogous to the IMM procedure described in Bar-Shalom and Li (1993).

This procedure finds at each time step the “best” gaussian fit to the current mixture of gaussians for each SSM. If we denote the approximating gaussian by Q and the mixture being approximated by P , “best” is defined here as minimizing $\text{KL}(P \| Q)$. Furthermore, gaussian merging techniques are greedy in that the “best” gaussian is computed at every time step and used immediately for the next time step. For a gaussian Q , $\text{KL}(P \| Q)$ has no local minima, and it is very easy to find the optimal Q by computing the first two moments of P . Inaccuracies in this greedy procedure arise because the estimates of $P(S_t | Y_1, \dots, Y_t)$ are based on this single merged gaussian, not on the real mixture.

In contrast, variational methods seek to minimize $\text{KL}(Q\|P)$, which can have many local minima. Moreover, these methods are not greedy in the same sense: they iterate forward and backward in time until obtaining a locally optimal Q .

5 Simulations

5.1 Experiment 1: Variational Segmentation and Deterministic Annealing. The goal of this experiment was to assess the quality of solutions found by the variational inference algorithm and the effect of using deterministic annealing on these solutions. We generated 200 sequences of length 200 from a simple model that switched between two SSMs. These SSMs and the switching process were defined by:

$$X_t^{(1)} = 0.99 X_{t-1}^{(1)} + w_t^{(1)} \quad w_t^{(1)} \sim \mathcal{N}(0, 1) \quad (5.1)$$

$$X_t^{(2)} = 0.9 X_{t-1}^{(2)} + w_t^{(2)} \quad w_t^{(2)} \sim \mathcal{N}(0, 10) \quad (5.2)$$

$$Y_t = X_t^{(m)} + v_t \quad v_t \sim \mathcal{N}(0, 0.1), \quad (5.3)$$

where the switch state m was chosen using priors $\pi^{(1)} = \pi^{(2)} = 1/2$ and transition probabilities $\Phi_{11} = \Phi_{22} = 0.95$; $\Phi_{12} = \Phi_{21} = 0.05$. Five sequences from this data set are shown in in Figure 6, along with the true state of the switch variable.

We compared three different inference algorithms: variational inference, variational inference with deterministic annealing (section 4.1), and inference by gaussian merging (section 4.2). For each sequence, we initialized the variational inference algorithms with equal responsibilities for the two SSMs and ran them for 12 iterations. The nonannealed inference algorithm ran at a fixed temperature of $\mathcal{T} = 1$, while the annealed algorithm was initialized to a temperature of $\mathcal{T} = 100$, which decayed down to 1 over the 12 iterations, using the decay function $\mathcal{T}_{i+1} = \frac{1}{2}\mathcal{T}_i + \frac{1}{2}$. To eliminate the effect of model inaccuracies we gave all three inference algorithms the true parameters of the generative model.

The segmentations found by the nonannealed variational inference algorithm showed little similarity to the true segmentations of the data (see Figure 7). Furthermore, the nonannealed algorithm generally underestimated the number of switches, often converging on solutions with no switches at all. Both the annealed variational algorithm and the gaussian merging method found segmentations that were more similar to the true segmentations of the data. Comparing percentage correct segmentations, we see that annealing substantially improves the variational inference method and that the gaussian merging and annealed variational methods perform comparably (see Figure 8). The average performance of the annealed variational method is only about 1.3% better than gaussian merging.

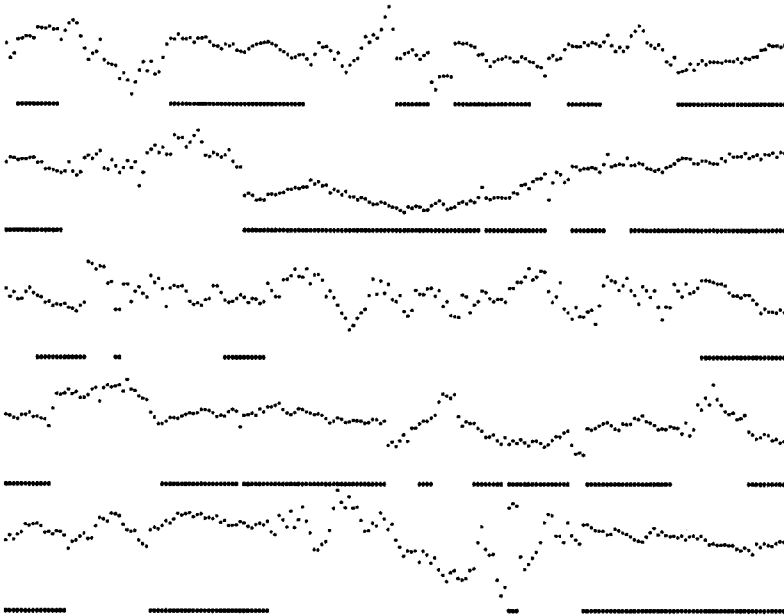


Figure 6: Five data sequences of length 200, with their true segmentations below them. In the segmentations, switch states 1 and 2 are represented with presence and absence of dots, respectively. Notice that it is difficult to segment the sequences correctly based only on knowing the dynamics of the two processes.

5.2 Experiment 2: Modeling Respiration in a Patient with Sleep Apnea. Switching state-space models should prove useful in modeling time series that have dynamics characterized by several different regimes. To illustrate this point, we examined a physiological data set from a patient tentatively diagnosed with sleep apnea, a medical condition in which an individual intermittently stops breathing during sleep. The data were obtained from the repository of time-series data sets associated with the Santa Fe Time Series Analysis and Prediction Competition (Weigend & Gershenfeld, 1993) and is described in detail in Rigney et al. (1993).⁸ The respiration pattern in sleep apnea is characterized by at least two regimes: no breathing and gasping breathing induced by a reflex arousal. Furthermore, in this patient there also seem to be periods of normal rhythmic breathing (see Figure 9).

⁸ The data are available online at <http://www.stern.nyu.edu/~aweigend/Time-Series/SantaFe.html#setB>. We used samples 6201–7200 for training and 5201–6200 for testing.

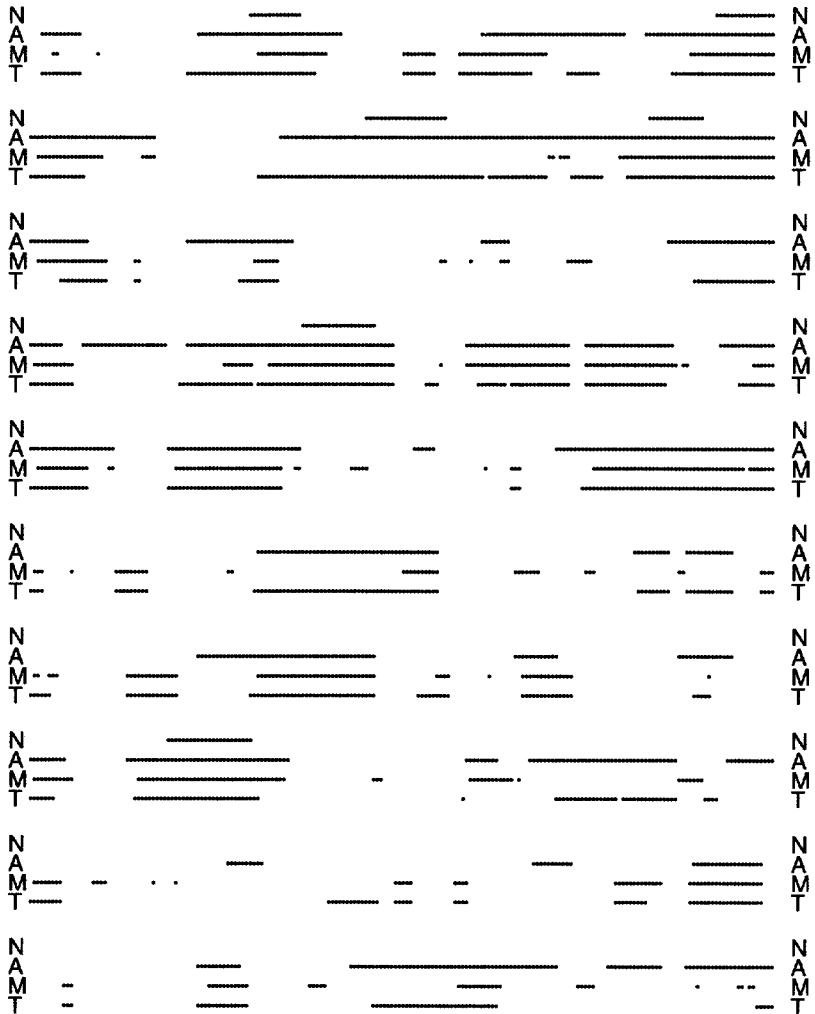


Figure 7: For 10 different sequences of length 200, segmentations are shown with presence and absence of dots corresponding to the two SSMs generating these data. The rows are the segmentations found using the variational method with no annealing (N), the variational method with deterministic annealing (A), the gaussian merging method (M), and the true segmentation (T). All three inference algorithms give real-valued $h_t^{(m)}$; hard segmentations were obtained by thresholding the final $h_t^{(m)}$ values at 0.5. The first five sequences are the ones shown in Figure 6.

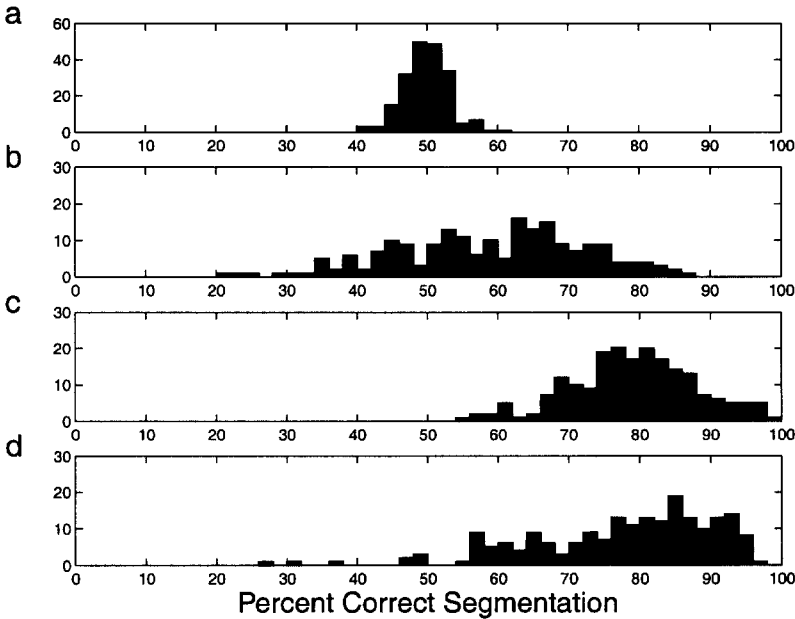


Figure 8: Histograms of percentage correct segmentations: (a) control, using random segmentation; (b) variational inference without annealing; (c) variational inference with annealing; (d) gaussian merging. Percentage correct segmentation was computed by counting the number of time steps for which the true and estimated segmentations agree.

We trained switching SSMs varying the random seed, the number of components in the mixture ($M = 2$ to 5), and the dimensionality of the state-space in each component ($K = 1$ to 10) on a data set consisting of 1000 consecutive measurements of the chest volume. As controls, we also trained simple SSMs (i.e., $M = 1$), varying the dimension of the state-space from $K = 1$ to 10 , and simple HMMs (i.e., $K = 0$), varying the number of discrete hidden states from $M = 2$ to $M = 50$. Simulations were run until convergence or for 200 iterations, whichever came first; convergence was assessed by measuring the change in likelihood (or bound on the likelihood) over consecutive steps of EM.

The likelihood of the simple SSMs and the HMMs was calculated on a test set consisting of 1000 consecutive measurements of the chest volume. For the switching SSMs, the likelihood is intractable, so we calculated the lower bound on the likelihood, \mathcal{B} . The simple SSMs modeled the data very poorly for $K = 1$, and the performance was flat for values of $K = 2$ to 10 (see

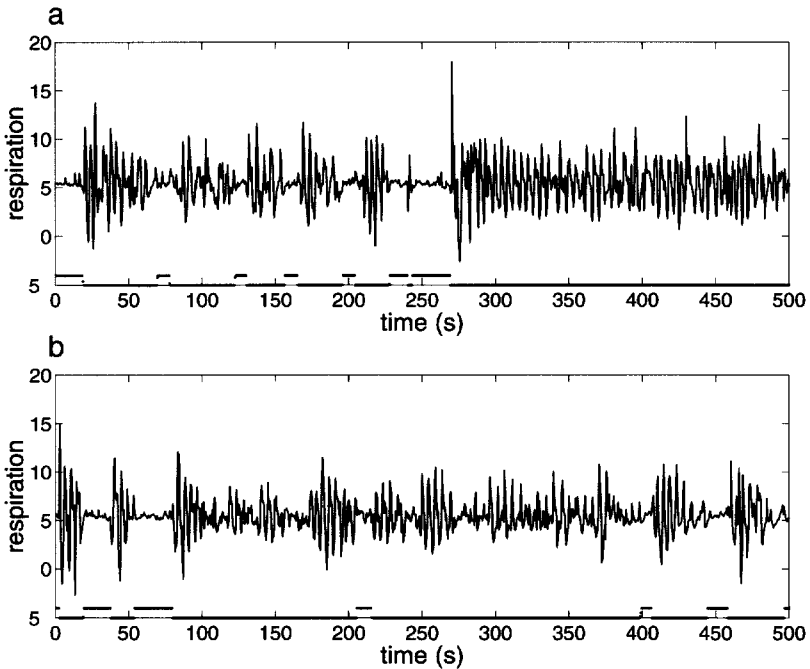


Figure 9: Chest volume (respiration force) of a patient with sleep apnea during two noncontinuous time segments of the same night (measurements sampled at 2 Hz). (a) Training data. Apnea is characterized by extended periods of small variability in chest volume, followed by bursts (gasping). Here we see such behavior around $t = 250$, followed by normal rhythmic breathing. (b) Test data. In this segment we find several instances of apnea and an approximately rhythmic region. (The thick lines at the bottom of each plot are explained in the main text.)

Figure 10a). The large majority of runs of the switching state-space model resulted in models with higher likelihood than those of the simple SMMs (see Figures 10b–10e). One consistent exception should be noted: for values of $M = 2$ and $K = 6$ to 10, the switching SSM performed almost identically to the simple SSM. Exploratory experiments suggest that in these cases, a single component takes responsibility for all the data, so the model has $M = 1$ effectively. This may be a local minimum problem or a result of poor initialization heuristics. Looking at the learning curves for simple and switching SSMs, it is easy to see that there are plateaus at the solutions found by the simple one-component SSMs that the switching SSM can get caught in (see Figure 11).

The likelihoods for HMMs with around $M = 15$ were comparable to those of the best switching SSMs (see Figure 10f). Purely in terms of cod-

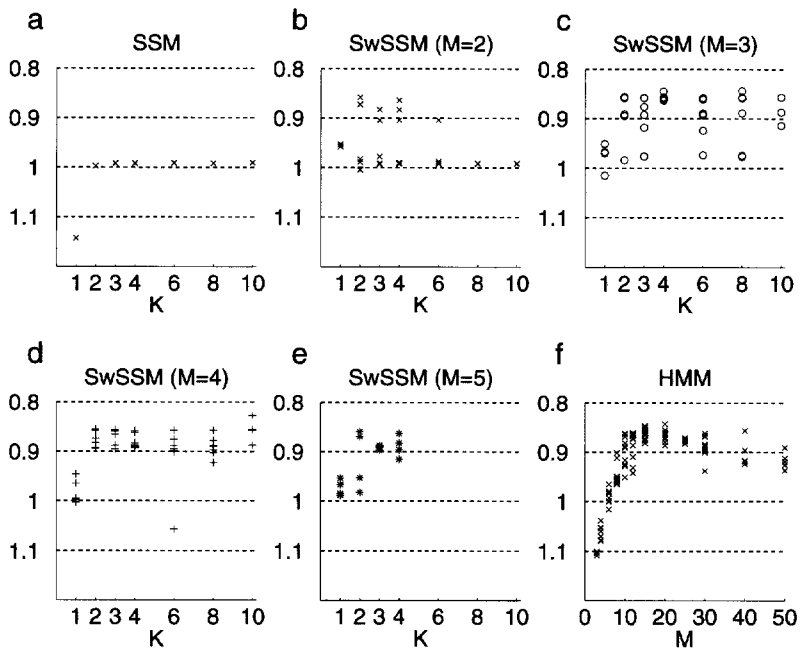


Figure 10: Log likelihood (nats per observation) on the test data from a total of almost 400 runs of simple state-space models (a), switching state-space models with differing numbers of components (b–e), and hidden Markov models (f).

ing efficiency, switching SSMs have little advantage over HMMs on this data.

However, it is useful to contrast the solutions learned by HMMs with the solutions learned by the switching SSMs. The thick dots at the bottom of the Figures 9a and 9b show the responsibility assigned to one of two components in a fairly typical switching SSM with $M = 2$ components of state size $K = 2$. This component has clearly specialized to modeling the data during periods of apnea, while the other component models the gasps and periods of rhythmic breathing. These two switching components provide a much more intuitive model of the data than the 10 to 20 discrete components needed in an HMM with comparable coding efficiency.⁹

⁹ By using further assumptions to constrain the model, such as continuity of the real-valued hidden state at switch times, it should be possible to obtain even better performance on these data.

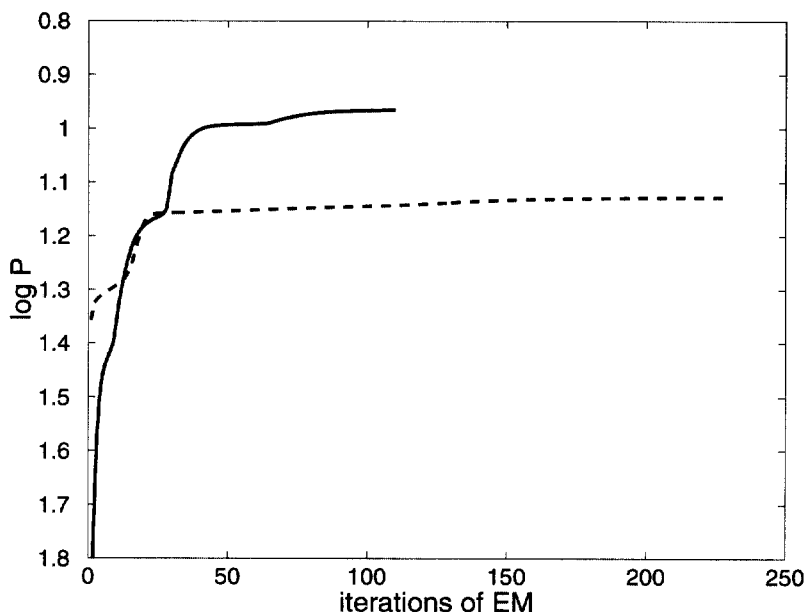


Figure 11: Learning curves for a state space model ($K = 4$) and a switching state-space model ($M = 2, K = 2$).

6 Discussion

The main conclusion we can draw from the first series of experiments is that even when given the correct model parameters, the problem of segmenting a switching time series into its components is difficult. There are combinatorially many alternatives to be considered, and the energy surface suffers from many local minima, so local optimization approaches like the variational method we used are limited by the quality of the initial conditions. Deterministic annealing can be thought of as a sophisticated initialization procedure for the hidden states: the final solution at each temperature provides the initial conditions at the next. We found that annealing substantially improved the quality of the segmentations found.

The first experiment also indicates that the much simpler gaussian merging method performs comparably to annealed variational inference. The gaussian merging methods have the advantage that at each time step, the cost function minimized has no local minima. This may account for how well they perform relative to the nonannealed variational method. On the other hand, the variational methods have the advantage that they iteratively improve their approximation to the posterior, and they define a lower bound

on the likelihood. Our results suggest that it may be very fruitful to use the gaussian merging method to initialize the variational inference procedure. Furthermore, it is possible to derive variational approximations for other switching models described in the literature, and a combination of gaussian merging and variational approximation may provide a fast and robust method for learning and inference in those models.

The second series of experiments suggests that on a real data set believed to have switching dynamics, the switching SSM can indeed uncover multiple regimes. When it captures these regimes, it generalizes to the test set much better than the simple linear dynamical model. Similar coding efficiency can be obtained by using HMMs, which due to the discrete nature of the state-space, can model nonlinear dynamics. However, in doing so, the HMMs had to use 10 to 20 discrete states, which makes their solutions less interpretable.

Variational approximations provide a powerful tool for inference and learning in complex probabilistic models. We have seen that when applied to the switching SSM, they can incorporate within a single framework well-known exact inference methods like Kalman smoothing and the forward-backward algorithm. Variational methods can be applied to many of the other classes of intractable switching models described in section 2.3. However, training more complex models also makes apparent the importance of good methods for model selection and initialization.

To summarize, switching SSMs are a dynamical generalization of mixture-of-experts neural networks, are closely related to well-known models in econometrics and control, and combine the representations underlying HMMs and linear dynamical systems. For domains in which we have some a priori belief that there are multiple, approximately linear dynamical regimes, switching SSMs provide a natural modeling tool. Variational approximations provide a method to overcome the most difficult problem in learning switching SSMs: that the inference step is intractable. Deterministic annealing further improves on the solutions found by the variational method.

Appendix A: Notation

Symbol	Size	Description
Variables		
Y_t	$D \times 1$	observation vector at time t
$\{Y_t\}$	$D \times T$	sequence of observation vectors $[Y_1, Y_2, \dots Y_T]$
$X_t^{(m)}$	$K \times 1$	state vector of state-space model (SSM) m at time t
X_t	$KM \times 1$	entire real-valued hidden state at time t : $X_t = [X_t^{(1)}, \dots, X_t^{(M)}]$

S_t	$M \times 1$	switch state variable (represented either as discrete variable $S_t \in \{1, \dots, M\}$, or as an $M \times 1$ vector $S_t = [S_t^{(1)}, \dots, S_t^{(M)}]'$ where $S_t^{(m)} \in \{0, 1\}$)
-------	--------------	---

Model parameters

$A^{(m)}$	$K \times K$	state dynamics matrix for SSM m
$C^{(m)}$	$D \times K$	output matrix for SSM m
$Q^{(m)}$	$K \times K$	state noise covariance matrix for SSM m
$\mu_{x_1}^{(m)}$	$K \times 1$	initial state mean for SSM m
$Q_1^{(m)}$	$K \times K$	initial state noise covariance matrix for SSM m
R	$D \times D$	output noise covariance matrix
π	$M \times 1$	initial state probabilities for switch state
Φ	$M \times M$	state transition matrix for switch state

Variational parameters

$h_t^{(m)}$	1×1	responsibility of SSM m for Y_t
$q_t^{(m)}$	1×1	related to expected squared error if SSM m generated Y_t

Miscellaneous

X'	matrix transpose of X
$ X $	matrix determinant of X
$\langle X \rangle$	expected value of X under the Q distribution

Dimensions

D	size of observation vector
T	length of a sequence of observation vectors
M	number of state-space models
K	size of state vector in each state-space model

Appendix B: Derivation of the Variational Fixed-Point Equations

In this appendix we derive the variational fixed-point equations used in the learning algorithm for switching SSMs. First, we write out the probability density P defined by a switching SSM. For convenience, we express this probability density in the log domain, through its associated energy function or *hamiltonian*, H . The probability density is related to the hamiltonian through the usual Boltzmann distribution (at a temperature of 1),

$$P(\cdot) = \frac{1}{Z} \exp\{-H(\cdot)\},$$

where Z is a normalization constant required such that $P(\cdot)$ integrates to unity. Expressing the probabilities in the log domain does not affect the resulting algorithm. We then similarly express the approximating distribution Q through its hamiltonian H_Q . Finally, we obtain the variational fixed-point equations by setting to zero the derivatives of the KL divergence between Q and P with respect to the variational parameters of Q .

The joint probability of observations and hidden states in a switching SSM is (see equation 3.1)

$$P(\{S_t, X_t, Y_t\}) = \left[P(S_1) \prod_{t=2}^T P(S_t|S_{t-1}) \right] \prod_{m=1}^M \left[P(X_1^{(m)}) \prod_{t=2}^T P(X_t^{(m)}|X_{t-1}^{(m)}) \right] \cdot \prod_{t=1}^T P(Y_t|X_t, S_t). \quad (\text{B.1})$$

We proceed to dissect this expression into its constituent parts. The initial probability of the switch variable at time $t = 1$ is given by

$$P(S_1) = \prod_{m=1}^M \left(\pi^{(m)} \right)^{S_1^{(m)}}, \quad (\text{B.2})$$

where S_1 is represented by an $M \times 1$ vector $[S_1^{(1)} \dots S_1^{(M)}]$ where $S_1^{(m)} = 1$ if the switch state is in state m , and 0 otherwise. The probability of transitioning from a switch state at time $t - 1$ to a switch state at time t is given by

$$P(S_t|S_{t-1}) = \prod_{m=1}^M \prod_{n=1}^M \left(\Phi^{(m,n)} \right)^{S_t^{(m)} S_{t-1}^{(n)}}. \quad (\text{B.3})$$

The initial distribution for the hidden state variable in SSM m is gaussian with mean $\mu_{X_1}^{(m)}$ and covariance matrix $\mathcal{Q}_1^{(m)}$:

$$P(X_1^{(m)}) = \left| 2\pi \mathcal{Q}_1^{(m)} \right|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \left(X_1^{(m)} - \mu_{X_1}^{(m)} \right)' \left(\mathcal{Q}_1^{(m)} \right)^{-1} \left(X_1^{(m)} - \mu_{X_1}^{(m)} \right) \right\}. \quad (\text{B.4})$$

The probability distribution of the state in SSM m at time t given the state at time $t - 1$ is gaussian with mean $A^{(m)} X_{t-1}^{(m)}$ and covariance matrix $\mathcal{Q}^{(m)}$:

$$P(X_t^{(m)}|X_{t-1}^{(m)}) = \left| 2\pi \mathcal{Q}^{(m)} \right|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \left(X_t^{(m)} - A^{(m)} X_{t-1}^{(m)} \right)' \left(\mathcal{Q}^{(m)} \right)^{-1} \cdot \left(X_t^{(m)} - A^{(m)} X_{t-1}^{(m)} \right) \right\}. \quad (\text{B.5})$$

Finally, using equation 3.2, we can write:

$$P(Y_t|X_t, S_t) = \prod_{m=1}^M \left[|2\pi R|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \left(Y_t - C^{(m)} X_t^{(m)} \right)' R^{-1} \left(Y_t - C^{(m)} X_t^{(m)} \right) \right\} \right]^{S_t^{(m)}} \quad (B.6)$$

since the terms with exponent equal to 0 vanish in the product.

Combining equations B.1 through B.6 and taking the negative of the logarithm, we obtain the hamiltonian of a switching SSM (ignoring constants):

$$\begin{aligned} H = & \frac{1}{2} \sum_{m=1}^M \log |\mathcal{Q}_1^{(m)}| + \frac{1}{2} \sum_{m=1}^M \left(X_1^{(m)} - \mu_{X_1}^{(m)} \right)' \left(\mathcal{Q}_1^{(m)} \right)^{-1} \left(X_1^{(m)} - \mu_{X_1}^{(m)} \right) \\ & + \frac{(T-1)}{2} \sum_{m=1}^M \log |\mathcal{Q}^{(m)}| \\ & + \frac{1}{2} \sum_{m=1}^M \sum_{t=2}^T \left(X_t^{(m)} - A^{(m)} X_{t-1}^{(m)} \right)' \left(\mathcal{Q}^{(m)} \right)^{-1} \left(X_t^{(m)} - A^{(m)} X_{t-1}^{(m)} \right) \\ & + \frac{T}{2} \log |R| + \frac{1}{2} \sum_{m=1}^M \sum_{t=1}^T S_t^{(m)} \left(Y_t - C^{(m)} X_t^{(m)} \right)' R^{-1} \left(Y_t - C^{(m)} X_t^{(m)} \right) \\ & - \sum_{m=1}^M S_1^{(m)} \log \pi^{(m)} - \sum_{t=2}^T \sum_{m=1}^M \sum_{n=1}^M S_t^{(m)} S_{t-1}^{(n)} \log \Phi^{(m,n)}. \end{aligned} \quad (B.7)$$

The hamiltonian for the approximating distribution can be analogously derived from the definition of Q (see equation 4.6):

$$Q(\{S_t, X_t\}) = \frac{1}{Z_Q} \left[\psi(S_1) \prod_{t=2}^T \psi(S_{t-1}, S_t) \right] \prod_{m=1}^M \psi(X_1^{(m)}) \prod_{t=2}^T \psi(X_{t-1}^{(m)}, X_t^{(m)}). \quad (B.8)$$

The potentials for the initial switch state and switch state transitions are

$$\psi(S_1) = \prod_{m=1}^M (\pi^{(m)} q_1^{(m)})^{S_1^{(m)}} \quad (B.9)$$

$$\psi(S_{t-1}, S_t) = \prod_{m=1}^M \prod_{n=1}^M \left(\Phi^{(m,n)} q_t^{(m)} \right)^{S_t^{(m)} S_{t-1}^{(n)}}. \quad (B.10)$$

The potential for the initial state of SSM m is

$$\psi(X_1^{(m)}) = P(X_1^{(m)}) \left[P(Y_1|X_1^{(m)}, S_1 = m) \right]^{h_1^{(m)}}, \quad (B.11)$$

and the potential for the state at time t given the state at time $t - 1$ is

$$\psi(X_{t-1}^{(m)}, X_t^{(m)}) = P(X_t^{(m)} | X_{t-1}^{(m)}) \left[P(Y_t | X_t^{(m)}, S_t = m) \right]^{h_t^{(m)}}. \quad (\text{B.12})$$

The hamiltonian for Q is obtained by combining these terms and taking the negative logarithm:

$$\begin{aligned} H_Q = & \frac{1}{2} \sum_{m=1}^M \log |\mathcal{Q}_1^{(m)}| + \frac{1}{2} \sum_{m=1}^M \left(X_1^{(m)} - \mu_{X_1}^{(m)} \right)' (\mathcal{Q}_1^{(m)})^{-1} \left(X_1^{(m)} - \mu_{X_1}^{(m)} \right) \\ & + \frac{(T-1)}{2} \sum_{m=1}^M \log |\mathcal{Q}^{(m)}| \\ & + \frac{1}{2} \sum_{m=1}^M \sum_{t=2}^T \left(X_t^{(m)} - A^{(m)} X_{t-1}^{(m)} \right)' (\mathcal{Q}^{(m)})^{-1} \left(X_t^{(m)} - A^{(m)} X_{t-1}^{(m)} \right) \\ & + \frac{T}{2} \sum_{m=1}^M \log |R| + \frac{1}{2} \sum_{m=1}^M \sum_{t=1}^T h_t^{(m)} \left(Y_t - C^{(m)} X_t^{(m)} \right)' R^{-1} \left(Y_t - C^{(m)} X_t^{(m)} \right) \\ & - \sum_{m=1}^M S_1^{(m)} \log \pi^{(m)} - \sum_{t=2}^T \sum_{m=1}^M \sum_{n=1}^M S_t^{(m)} S_{t-1}^{(n)} \log \Phi^{(m,n)} \\ & - \sum_{t=1}^T \sum_{m=1}^M S_t^{(m)} \log q_t^{(m)}. \end{aligned} \quad (\text{B.13})$$

Comparing H_Q with H , we see that the interaction between the $S_t^{(m)}$ and the $X_t^{(m)}$ variables has been eliminated, while introducing two sets of variational parameters: the responsibilities $h_t^{(m)}$ and the bias terms on the discrete Markov chain, $q_t^{(m)}$. In order to obtain the approximation Q that maximizes the lower bound on the log-likelihood, we minimize the KL divergence $\text{KL}(Q||P)$ as a function of these variational parameters:

$$\text{KL}(Q||P) = \sum_{\{S_t\}} \int Q(\{S_t, X_t\}) \log \frac{Q(\{S_t, X_t\})}{P(\{S_t, X_t\}|\{Y_t\})} d\{X_t\} \quad (\text{B.14})$$

$$= \langle H - H_Q \rangle - \log Z_Q + \log Z, \quad (\text{B.15})$$

where $\langle \cdot \rangle$ denotes expectation over the approximating distribution Q and Z_Q is the normalization constant for Q . Both Q and P define distributions in the exponential family. As a consequence, the zeros of the derivatives of KL with respect to the variational parameters can be obtained simply by equating derivatives of $\langle H \rangle$ and $\langle H_Q \rangle$ with respect to corresponding sufficient statistics (Ghahramani, 1997):

$$\frac{\partial \langle H_Q - H \rangle}{\partial \langle S_t^{(m)} \rangle} = 0 \quad (\text{B.16})$$

$$\frac{\partial \langle H_Q - H \rangle}{\partial \langle X_t^{(m)} \rangle} = 0 \quad (\text{B.17})$$

$$\frac{\partial \langle H_Q - H \rangle}{\partial \langle P_t^{(m)} \rangle} = 0 \quad (\text{B.18})$$

where $P_t^{(m)} = \langle X_t^{(m)} X_t^{(m)'} \rangle - \langle X_t^{(m)} \rangle \langle X_t^{(m)} \rangle'$ is the covariance of $X_t^{(m)}$ under Q . Many terms cancel when we subtract the two hamiltonians,

$$H_Q - H = \sum_{m=1}^M \sum_{t=1}^T \frac{1}{2} \left(h_t^{(m)} - S_t^{(m)} \right) \left(Y_t - C^{(m)} X_t^{(m)} \right)' R^{-1} \cdot \left(Y_t - C^{(m)} X_t^{(m)} \right) - S_t^{(m)} \log q_t^{(m)}. \quad (\text{B.19})$$

Taking derivatives, we obtain

$$\frac{\partial \langle H_Q - H \rangle}{\partial \langle S_t^{(m)} \rangle} = -\log q_t^{(m)} - \frac{1}{2} \left\langle \left(Y_t - C^{(m)} X_t^{(m)} \right)' R^{-1} \left(Y_t - C^{(m)} X_t^{(m)} \right) \right\rangle \quad (\text{B.20})$$

$$\frac{\partial \langle H_Q - H \rangle}{\partial \langle X_t^{(m)} \rangle} = - \left(h_t^{(m)} - \langle S_t^{(m)} \rangle \right) \left((Y_t - C^{(m)} \langle X_t^{(m)} \rangle)' R^{-1} C^{(m)} \right) \quad (\text{B.21})$$

$$\frac{\partial \langle H_Q - H \rangle}{\partial P_t^{(m)}} = \frac{1}{2} \left(h_t^{(m)} - \langle S_t^{(m)} \rangle \right) \left(C^{(m)'} R^{-1} C^{(m)} \right) \quad (\text{B.22})$$

From equation B.20, we get the fixed-point equation, 4.12, for $q_t^{(m)}$. Both equations B.21 and B.22 are satisfied when $h_t^{(m)} = \langle S_t^{(m)} \rangle$. Using the fact that $\langle S_t^{(m)} \rangle = Q(S_t = m)$ we get equation 4.11.

References

- Ackerson, G. A., & Fu, K. S. (1970). On state estimation in switching environments. *IEEE Transactions on Automatic Control*, AC-15(1):10-17.
- Anderson, B. D. O., & Moore, J. B. (1979). *Optimal filtering*. Englewood Cliffs, NJ: Prentice-Hall.
- Athaide, C. R. (1995). *Likelihood evaluation and state estimation for nonlinear state space models*. Unpublished doctoral dissertation, University of Pennsylvania.
- Baldi, P., Chauvin, Y., Hunkapiller, T., & McClure, M. (1994). Hidden Markov models of biological primary sequence information. *Proc. Nat. Acad. Sci. (USA)*, 91(3), 1059-1063.
- Bar-Shalom, Y., & Li, X.-R. (1993). *Estimation and tracking*. Boston: Artech House.
- Baum, L., Petrie, T., Soules, G., & Weiss, N. (1970). A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains. *Annals of Mathematical Statistics*, 41, 164-171.

- Bengio, Y., & Frasconi, P. (1995). An input-output HMM architecture. In G. Tesauro, D. S. Touretzky, & T. K. Leen (Eds.), *Advances in neural information processing systems*, 7 (pp. 427–434). Cambridge, MA: MIT Press.
- Cacciatore, T. W., & Nowlan, S. J. (1994). Mixtures of controllers for jump linear and non-linear plants. In J. D. Cowan, G. Tesauro, & J. Alspector (Eds.), *Advances in neural information processing systems*, 6 (pp. 719–726). San Mateo: Morgan Kaufmann Publishers.
- Carter, C. K., & Kohn, R. (1994). On Gibbs sampling for state space models. *Biometrika*, 81, 541–553.
- Chaer, W. S., Bishop, R. H., & Ghosh, J. (1997). A mixture-of-experts framework for adaptive Kalman filtering. *IEEE Trans. on Systems, Man and Cybernetics*.
- Chang, C. B., & Athans, M. (1978). State estimation for discrete systems with switching parameters. *IEEE Transactions on Aerospace and Electronic Systems*, AES-14(3), 418–424.
- Cover, T., & Thomas, J. (1991). *Elements of information theory*. New York: Wiley.
- Dean, T., & Kanazawa, K. (1989). A model for reasoning about persistence and causation. *Computational Intelligence*, 5(3), 142–150.
- Dempster, A., Laird, N., & Rubin, D. (1977). Maximum likelihood from incomplete data via the EM algorithm. *J. Royal Statistical Society Series B*, 39, 1–38.
- Deng, L. (1993). A stochastic model of speech incorporating hierarchical non-stationarity. *IEEE Trans. on Speech and Audio Processing*, 1(4), 471–474.
- Digalakis, V., Rohlicek, J. R., & Ostendorf, M. (1993). ML estimation of a stochastic linear system with the EM algorithm and its application to speech recognition. *IEEE Transactions on Speech and Audio Processing*, 1(4), 431–442.
- Elliott, R. J., Aggoun, L., & Moore, J. B. (1995). *Hidden Markov models: Estimation and control*. New York: Springer-Verlag.
- Fraser, A. M., & Dimitriadis, A. (1993). Forecasting probability densities by using hidden Markov models with mixed states. In A. S. Wiegand & N. A. Gershenfeld (Eds.), *Time series prediction: Forecasting the future and understanding the past* (pp. 265–282). Reading, MA: Addison-Wesley.
- Ghahramani, Z. (1997). *On structured variational approximations* (Tech. Rep. No. CRG-TR-97-1). Toronto: Department of Computer Science, University of Toronto. Available online at: <http://www.gatsby.ucl.ac.uk/~zoubin/papers/struct.ps.gz>.
- Ghahramani, Z., & Hinton, G. E. (1996a). *Parameter estimation for linear dynamical systems* (Tech. Rep. No. CRG-TR-96-2). Toronto: Department of Computer Science, University of Toronto. Available online at: <http://www.gatsby.ucl.ac.uk/~zoubin/papers/tr-96-2.ps.gz>.
- Ghahramani, Z., & Hinton, G. E. (1996b). *The EM algorithm for mixtures of factor analyzers* (Tech. Rep. No. CRG-TR-96-1). Toronto: Department of Computer Science, University of Toronto. Available online at: <http://www.gatsby.ucl.ac.uk/~zoubin/papers/tr-96-1.ps.gz>.
- Ghahramani, Z., & Jordan, M. I. (1997). Factorial hidden Markov models. *Machine Learning*, 29, 245–273.
- Goodwin, G., & Sin, K. (1984). *Adaptive filtering prediction and control*. Englewood Cliffs, NJ: Prentice-Hall.

- Hamilton, J. D. (1989). A new approach to the economic analysis of nonstationary time series and the business cycle. *Econometrica*, 57, 357–384.
- Hamilton, J. D. (1994). *Time series analysis*. Princeton, NJ: Princeton University Press.
- Harrison, P. J., & Stevens, C. F. (1976). Bayesian forecasting (with discussion). *J. Royal Statistical Society B*, 38, 205–247.
- Hinton, G. E., Dayan, P., & Revow, M. (1997). Modeling the manifolds of images of handwritten digits. *IEEE Transactions on Neural Networks*, 8, 65–74.
- Jacobs, R. A., Jordan, M. I., Nowlan, S. J., & Hinton, G. E. (1991). Adaptive mixture of local experts. *Neural Computation*, 3, 79–87.
- Jordan, M. I., Ghahramani, Z., Jaakkola, T. S., & Saul, L. K. (1998). An introduction to variational methods in graphical models. In M. I. Jordan (Ed.), *Learning in graphical models*. Norwell, MA: Kluwer.
- Jordan, M. I., & Jacobs, R. (1994). Hierarchical mixtures of experts and the EM algorithm. *Neural Computation*, 6, 181–214.
- Juang, B. H., & Rabiner, L. R. (1991). Hidden Markov models for speech recognition. *Technometrics*, 33, 251–272.
- Kadirkamanathan, V., & Kadirkamanathan, M. (1996). Recursive estimation of dynamic modular RBF networks. In D. Touretzky, M. Mozer, and M. Hasselmo (Eds.), *Advances in neural information processing systems*, 8 (pp. 239–245). Cambridge, MA: MIT Press.
- Kalman, R. E., & Bucy, R. S. (1961). New results in linear filtering and prediction. *Journal of Basic Engineering (ASME)*, 83D, 95–108.
- Kanazawa, K., Koller, D., & Russell, S. J. (1995). Stochastic simulation algorithms for dynamic probabilistic networks. In P. Besnard & S. Hanks (Eds.), *Uncertainty in artificial intelligence. Proceedings of the Eleventh Conference* (pp. 346–351). San Mateo, CA: Morgan Kaufmann.
- Kehagias, A., & Petrides, V. (1997). Time series segmentation using predictive modular neural networks. *Neural Computation*, 9(8), 1691–1710.
- Kim, C.-J. (1994). Dynamic linear models with Markov-switching. *J. Econometrics*, 60, 1–22.
- Lauritzen, S. L., & Spiegelhalter, D. J. (1988). Local computations with probabilities on graphical structures and their application to expert systems. *J. Royal Statistical Society B*, 157–224.
- Ljung, L., & Söderström, T. (1983). *Theory and practice of recursive identification*. Cambridge, MA: MIT Press.
- Meila, M., & Jordan, M. I. (1996). Learning fine motion by Markov mixtures of experts. In D. S. Touretzky, M. C. Mozer, & M. E. Hasselmo (Eds.), *Advances in neural information processing systems*, 8. Cambridge, MA: MIT Press.
- Neal, R. M., & Hinton, G. E. (1998). A new view of the EM algorithm that justifies incremental, sparse, and other variants. In M. I. Jordan (Ed.), *Learning in graphical models*. Norwell, MA: Kluwer.
- Parisi, G. (1988). *Statistical field theory*. Redwood City, CA: Addison-Wesley.
- Pawelzik, K., Kohlmorgen, J., & Müller, K.-R. (1996). Annealed competition of experts for a segmentation and classification of switching dynamics. *Neural Computation*, 8(2), 340–356.

- Pearl, J. (1988). *Probabilistic reasoning in intelligent systems: Networks of plausible inference*. San Mateo, CA: Morgan Kaufmann.
- Rabiner, L. R., & Juang, B. H. (1986). An introduction to hidden Markov models. *IEEE Acoustics, Speech & Signal Processing Magazine*, 3, 4–16.
- Rauch, H. E. (1963). Solutions to the linear smoothing problem. *IEEE Transactions on Automatic Control*, 8, 371–372.
- Rigney, D., Goldberger, A., Ocasio, W., Ichimaru, Y., Moody, G., & Mark, R. (1993). Multi-channel physiological data: Description and analysis. In A. Weigend & N. Gershenfeld (Eds.), *Time series prediction: Forecasting the future and understanding the past* (pp. 105–129). Reading, MA: Addison-Wesley.
- Roweis, S., & Ghahramani, Z. (1999). A unifying review of linear gaussian models. *Neural Computation*, 11(2), 305–345.
- Saul, L. and Jordan, M. I. (1996). Exploiting tractable substructures in intractable networks. In D. Touretzky, M. Mozer, & M. Hasselmo (Eds.), *Advances in neural information processing systems*, 8. Cambridge, MA: MIT Press.
- Shumway, R. H., & Stoffer, D. S. (1982). An approach to time series smoothing and forecasting using the EM algorithm. *J. Time Series Analysis*, 3(4), 253–264.
- Shumway, R. H., & Stoffer, D. S. (1991). Dynamic linear models with switching. *J. Amer. Stat. Assoc.*, 86, 763–769.
- Smyth, P. (1994). Hidden Markov models for fault detection in dynamic systems. *Pattern Recognition*, 27(1), 149–164.
- Smyth, P., Heckerman, D., & Jordan, M. I. (1997). Probabilistic independence networks for hidden Markov probability models. *Neural Computation*, 9, 227–269.
- Ueda, N., & Nakano, R. (1995). Deterministic annealing variant of the EM algorithm. In G. Tesauro, D. Touretzky, & J. Alspector (Eds.), *Advances in neural information processing systems*, 7 (pp. 545–552). San Mateo, CA: Morgan Kaufmann.
- Weigend, A., & Gershenfeld, N. (1993). *Time series prediction: Forecasting the future and understanding the past*. Reading, MA: Addison-Wesley.

Received March 2, 1998; accepted March 16, 1999.