
Variational Bayesian Inference and Learning for Continuous Switching Linear Dynamical Systems

Jack Goffinet

Department of Computer Science
Duke University
Durham, NC 27708
jack.goffinet@duke.edu

David E. Carlson

Department of Computer Science
Duke University
Durham, NC 27708
david.carlson@duke.edu

Abstract

Linear-Gaussian dynamical systems (LDSs) are computationally tractable because all latents and observations are jointly Gaussian. However, these systems are too restrictive to satisfactorily model many dynamical systems of interest. One generalization, the switching linear dynamical system (SLDS), trades analytic tractability for a more expressive model, allowing a discrete set of different linear regimes to model the data. Here we introduce a switching linear dynamical system with a continuum of linear regimes that are traversed continuously in time. We call this model a *continuous switching linear dynamical system* (CSLDS) and derive efficient variational Bayesian methods for inference and model learning.

1 Introduction

Linear dynamical systems (LDSs) are fundamental to probabilistic timeseries modeling due to their analytic tractability and the past several years have seen an increased interest in more expressive extensions of LDSs. The most common among these is the *switching* linear dynamical system (SLDS), which posits a set of discrete modes, each corresponding to a distinct LDS, with Markov transitions between the modes [3, 7]. Several extensions of SLDS exist, including recurrent variants that use both the discrete mode and the continuous state to determine mode dynamics [5, 6, 4].

The choice made by SLDS of maintaining a discrete set of linear regimes is computationally convenient given the Markov model structure of the modes. However, we note that an SLDS with a continuum of linear regimes can also be handled in a computationally convenient way. The resulting model, which we call a *continuous switching linear dynamical systems* (CSLDS), is better suited to modeling systems for which we expect the underlying dynamics vary smoothly in time.

The most similar model to the CSLDS in the literature is the warped autoregressive hidden Markov model (WAR-HMM) [2]. The autoregressive (AR)-HMM can be viewed as a discrete SLDS with a restricted observation model, which the WAR-HMM extends by associating with each discrete mode a univariate latent variable that modulates the mode’s corresponding linear dynamics. The CSLDS, by contrast, modulates its linear dynamics with a continuously varying multivariate latent variable and extends the more flexible observation model of the SLDS.

2 Continuous Switching Linear Dynamical Systems

Setup Let time $t = 1, \dots, T$ have associated observations $\mathbf{y}_t \in \mathbb{R}^N$ for all t . The model has two layers of latent variables: $\boldsymbol{\eta}_t \in \mathbb{R}^K$ for all $t < T$, and $\mathbf{x}_t \in \mathbb{R}^M$ for all t , with $M \geq K$. The graphical model is shown in Figure 1, which has an associated joint distribution of

$$p(\boldsymbol{\eta}_{1:T-1}, \mathbf{x}_{1:T}, \mathbf{y}_{1:T}) = p(\boldsymbol{\eta}_{1:T-1})p(\mathbf{x}_1) \prod_{t=2}^T p(\mathbf{x}_t | \boldsymbol{\eta}_{t-1}, \mathbf{x}_{t-1}) \prod_{t=1}^T p(\mathbf{y}_t | \mathbf{x}_t). \quad (1)$$

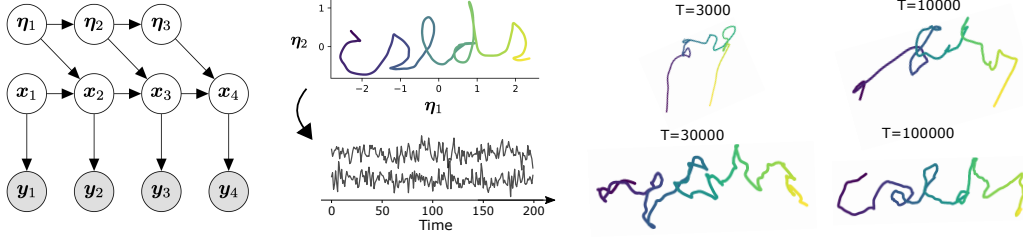


Figure 1: Left: Graphical model shared by the switching linear dynamical system (SLDS) and the proposed continuous switching linear dynamical system (CSLDS). Center: A synthetic dataset was made using continuously varying $\boldsymbol{\eta}$'s which spell out "cslds," where each setting of $\boldsymbol{\eta}$ corresponds to a distinct linear dynamical regime. Right: The proposed model is able to approximately recover the $\boldsymbol{\eta}$'s from the synthetic timeseries, with more accurate results obtained for datasets with more timepoints, T .

We choose an independent Gaussian process prior on each component of $\boldsymbol{\eta}$: $\boldsymbol{\eta}_{1:T-1}^{(k)} \sim \mathcal{GP}(\mathbf{0}, \kappa_k)$, where κ_k is a given kernel. We define the initial condition prior $\mathbf{x}_1 \sim \mathcal{N}(0, Q_{IC}^{-1})$, and the interaction between the $\boldsymbol{\eta}$'s and the \mathbf{x} 's as $\mathbf{x}_t | \boldsymbol{\eta}_{t-1}, \mathbf{x}_{t-1} \sim \mathcal{N}(A(\boldsymbol{\eta}_{t-1})\mathbf{x}_{t-1}, [Q(\boldsymbol{\eta}_{t-1})]^{-1})$. Lastly, we define a linear-Gaussian observation model: $\mathbf{y}_t | \mathbf{x}_t \sim \mathcal{N}(C\mathbf{x}_t + \mathbf{b}, R)$. In this way, the $\boldsymbol{\eta}$'s control the parameters of a linear-Gaussian state space model, producing an LDS with continuously varying parameters. We introduce two concrete choices for $A(\boldsymbol{\eta}_t)$ and $Q(\boldsymbol{\eta}_t)$, a *simplex* definition that is more consistent with the existing SLDS literature, and a *factor* definition that is more readily comparable to factor analysis. The *simplex* definition is

$$\mathbf{s} = \text{Softmax}(B\boldsymbol{\eta} + \mathbf{c}), \quad A(\boldsymbol{\eta}) = \sum_{k=1}^K s^{(k)} A^{(k)}, \quad [Q(\boldsymbol{\eta})]^{-1} = \sum_{k=1}^K s^{(k)} (Q^{(k)})^{-1}, \quad (2)$$

where $A^{(k)}, Q^{(k)} \in \mathbb{R}^{M \times M}$ for $k = 1, \dots, K$. The standard discrete SLDS is the limiting case using only the linear regimes in the corners of the simplex, i.e. \mathbf{s} is constrained to be one-hot. The *factor* definition is given by

$$A(\boldsymbol{\eta}) = A^{(0)} + \sum_{k=1}^K \eta^{(k)} A^{(k)}, \quad \mathbf{s} = \text{Softplus}(B\boldsymbol{\eta} + \mathbf{c}), \quad [Q(\boldsymbol{\eta})]^{-1} = \sum_{k=1}^K s^{(k)} (Q^{(k)})^{-1}. \quad (3)$$

In this work we only consider the *factor* variant and leave the *simplex* for consideration in future work. All told, the model parameters are $\theta = \{C, \mathbf{b}, R, B, \mathbf{c}, \{A^{(k)}\}_{k=0}^K, \{Q^{(k)}\}_{k=0}^K, \{\vartheta_k\}_{k=0}^K, Q_{IC}\}$ where the ϑ 's denote the η GP kernel parameters, discussed below.

Variational Inference and Learning The multiplicative interaction between the $\boldsymbol{\eta}$'s in defining the A matrices and the \mathbf{x} 's produces a joint distribution of $\boldsymbol{\eta}$'s and \mathbf{x} 's that is not jointly Gaussian in general, hindering analytic inference. However, we note that the \mathbf{x} 's and \mathbf{y} 's form a linear-Gaussian state space model *conditioned on* a setting of $\boldsymbol{\eta}_{1:T-1}$. This suggests a model fitting procedure in which the $\boldsymbol{\eta}$'s are sampled from an approximate posterior, the log likelihood of the observations conditioned on the sampled $\boldsymbol{\eta}$ is calculated, and the model parameters are updated by stochastic gradient ascent.

Let $q(\boldsymbol{\eta}_{1:T-1}; \phi)$ denote an arbitrary joint distribution over the $\boldsymbol{\eta}$'s with parameters ϕ that admits sampling and density evaluation. We can lower bound the marginal observation log likelihood using the standard evidence lower bound, or ELBO:

$$\mathcal{L}(\mathbf{y}_{1:T}; \theta, \phi) \triangleq \mathbb{E}_{\boldsymbol{\eta}_{1:T-1} \sim q(\boldsymbol{\eta}_{1:T-1})} \log \left[\frac{p(\mathbf{y}_{1:T}, \boldsymbol{\eta}_{1:T-1})}{q(\boldsymbol{\eta}_{1:T-1})} \right] \leq \log p(\mathbf{y}_{1:T}). \quad (4)$$

$q(\boldsymbol{\eta}_{1:T-1}; \phi)$ is known as an approximate posterior because the inequality becomes equality when $q(\boldsymbol{\eta}_{1:T-1}; \phi) = p(\boldsymbol{\eta}_{1:T-1} | \mathbf{y}_{1:T}; \theta)$, the true posterior. The following sections detail choices that allow for the efficient estimation of an ELBO, which will end up being a lower bound of the standard ELBO in Eq. 4.

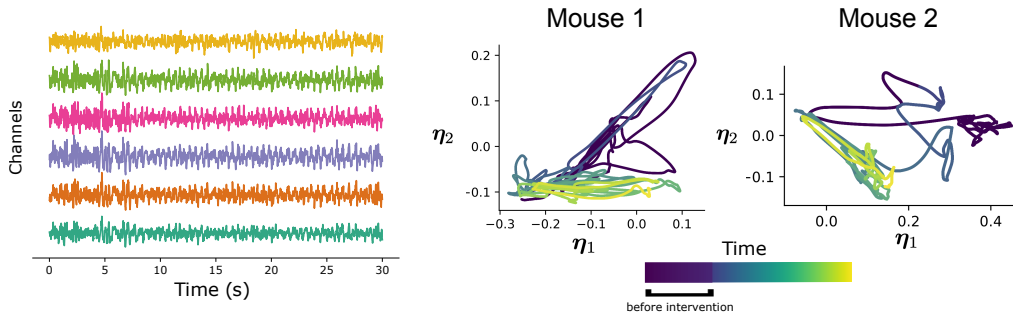


Figure 2: Left: Two CSLDS models were trained on six-channel, 2000 second local field potentials recorded from two mice before and after undergoing a pharmacological intervention. Right: The inferred η traces clearly distinguish between the dynamics before the intervention (dark purple) and after (lighter colors).

Identifiability of η 's We assert under certain conditions, namely that the A matrices have distinct eigenvalues and are sufficiently close to some A_0 , that the true η 's underlying a CSLDS are identifiable up to linear transformation in the slow- η limit. Here we present a brief logical sketch and in Figure 1 we present empirical evidence of this claim, leaving a formal proof for future work. In the slow- η limit, we may window the \mathbf{y} timeseries and fit a separate LDS model for each window corresponding to a specific $\eta \in \mathbb{R}^K$. In general, the dynamics matrix A in an LDS is only identifiable up to matrix similarity [1], which reduces to identifiability up to eigenvalues given the distinct eigenvalue assumption. Under the small perturbation assumption, we can match eigenvalues across windows unambiguously and identify a subspace of eigenvalues. Lastly, in this low perturbation regime, we can use first-order eigenvalue perturbations to confirm that the true η 's are given by a linear transformation of coordinates in this eigenvalue subspace.

3 Experiments

Synthetic Data As a test of the proposed model, we asked whether the CSLDS could recover ground-truth η 's in a synthetic dataset. First, we drew “cslds” in cursive (Figure 1, center) to represent the time-evolution of a 2-dimensional η . Multiple datasets were created that divided the “cslds” path into $T = 3 \cdot 10^3, 10^4, 3 \cdot 10^4,$ and 10^5 timesteps. The ground-truth model was randomly initialized with $K, M, N = 2$ and held fixed between experiments. Given the discussion on the identifiability of the η 's above, we expect that we should be able to infer the correct η 's up to linear transformation with a large amount of data. We fit a CSLDS model to each timeseries and confirmed that it was able to approximately reconstruct the ground-truth η 's (Figure 1, right, shown rotated to best match the ground-truth η 's).

Modeling Multichannel Neural Data We next used CSLDS to model multichannel local field potential data from two freely behaving mice before and after a pharmacological intervention, asking what effect, if any, the intervention has on the inferred η 's. The data for each mouse comprises 6 channels and 2000 seconds sampled at 50Hz for 10^5 total timesteps and the CSLDS model parameters are $K = 2$ and $M = 4$. The inferred η 's for both recordings, shown in Figure 2, clearly distinguish the times before the intervention (dark purple) and after the intervention (lighter colors).

4 Conclusion

We propose the continuous switching linear dynamical systems (CSLDS), which traverses a continuum of linear dynamical systems continuously in time. The model admits computationally feasible methods for inference and learning, allowing us to recover ground-truth latents and model neural recordings in mice. One particularly promising direction of future work is to use the CSLDS to study transitions between nominally discrete states, for example, by modeling the neural dynamics underlying sleep state transitions.

References

- [1] KS Arun and SY Kung. “Balanced approximation of stochastic systems”. In: *SIAM journal on matrix analysis and applications* 11.1 (1990), pp. 42–68.
- [2] Julia C Costacurta et al. “Distinguishing discrete and continuous behavioral variability using warped autoregressive HMMs”. In: *bioRxiv* (2022).
- [3] Emily Fox et al. “Nonparametric Bayesian learning of switching linear dynamical systems”. In: *Advances in neural information processing systems* 21 (2008).
- [4] Joshua Glaser et al. “Recurrent switching dynamical systems models for multiple interacting neural populations”. In: *Advances in neural information processing systems* 33 (2020), pp. 14867–14878.
- [5] Scott Linderman et al. “Bayesian learning and inference in recurrent switching linear dynamical systems”. In: *Artificial Intelligence and Statistics*. PMLR, 2017, pp. 914–922.
- [6] Josue Nassar et al. “Tree-Structured Recurrent Switching Linear Dynamical Systems for Multi-Scale Modeling”. In: *International Conference on Learning Representations*. 2019. URL: <https://openreview.net/forum?id=HkzRQhR9YX>.
- [7] Sang Min Oh et al. “Learning and inferring motion patterns using parametric segmental switching linear dynamic systems”. In: *International Journal of Computer Vision* 77.1 (2008), pp. 103–124.
- [8] Virginia Rutten et al. “Non-reversible Gaussian processes for identifying latent dynamical structure in neural data”. In: *Advances in Neural Information Processing Systems* (2020).

ELBO details

This section addresses details necessary to compute the CSLDS ELBO objective efficiently.

An efficient log determinant bound Here we derive an upper bound for the log determinant of the data covariance, Σ_{yy} . We begin with a derivation that applies generally to the ‘‘GP factor analysis’’ literature, and can be found in [8]. In what follows, we prioritize the common case $M < N \ll T$. Applying the Woodbury matrix inversion lemma to Eq. 10, we get

$$\begin{aligned}\Lambda_{yy} &= (R^{-1} \otimes I_T) - (R^{-1}C \otimes I_T)(L^{-\top}L^{-1} + C^\top R^{-1}C \otimes I_T)^{-1}(C^\top R^{-1} \otimes I_T) \\ &= (R^{-\frac{1}{2}} \otimes I_T)(I_{NT} - A^\top B^{-1}A)(R^{-\frac{1}{2}} \otimes I_T),\end{aligned}\quad (12)$$

where $LL^\top = \Sigma_x$, $A = L^\top(C^\top R^{-\frac{1}{2}} \otimes I_T)$, and $B = I + AA^\top$. We can then simplify the log determinant:

$$\begin{aligned}\log\det(\Lambda_{yy}) &= 2 \log\det(R^{-1/2} \otimes I_T) + \log\det(I - A^\top B^{-1}A) \\ &= -T \log\det(R) + \log\det(B - AA^\top) - \log\det(B) \\ &= -T \log\det(R) + \log\det(I) - \log\det(B) \\ &= -T \log\det(R) - \log\det(I + L^\top(C^\top R^{-1}C \otimes I_T)L) \\ \Rightarrow \log\det(\Sigma_{yy}) &= T \log\det(R) + \log\det(I + L^\top(C^\top R^{-1}C \otimes I_T)L)\end{aligned}\quad (13)$$

Note that the matrix in the second logdet term is MT -by- MT , so a naive logdet calculation would take time $\mathcal{O}(M^3T^3)$, which is an improvement on the $\mathcal{O}(N^3T^3)$ required for the full NT -by- NT matrix, but still prohibitive for large T . Instead of performing an exact log determinant calculation, we deviate from the GP factor analysis approach and opt to upper bound the log determinant using the identity $\log\det(A) \leq \text{tr}(A - I)$:

$$\begin{aligned}\log\det(\Sigma_{yy}) &\leq T \log\det(R) + \text{tr}(L^\top(C^\top R^{-1}C \otimes I_T)L) \\ &= T \log\det(R) + \sum_m [C^\top R^{-1}C]_{mm} \text{tr}(L_m L_m^\top) \\ &= T \log\det(R) + \sum_m [C^\top R^{-1}C]_{mm} \text{tr}(\Sigma_{xx,m}) \\ &= T(\log\det(R) + \sum_m [C^\top R^{-1}C]_{mm}) \\ &= T[\log\det(R) + \text{tr}(C^\top R^{-1}C)]\end{aligned}\quad (14)$$

Note that $\text{tr}(\Sigma_{xx,m}) = T$ relies on our choice of kernel, which guarantees $\kappa(t, t) = 1$. If we restrict R to be diagonal, this log determinant bound can be calculated in time $\mathcal{O}(MN)$, which is notably independent of T , the number of timepoints.

In what situations is this bound tight? It can be easily checked that the bound $\log\det(A) \leq \text{tr}(A - I)$ is tight when $A = I$. The bound can be tightened in a flexible way by using a lower triangular matrix W . We have:

$$\log\det(A) = \log\det(WW^\top A) - 2 \log\det(W) \leq \text{tr}(WW^\top A - I) - 2 \log\det(W). \quad (15)$$

The log determinant of W can be calculated efficiently because W is triangular. The bound is tight when $WW^\top A = I$, and for a positive definite A , which is our use case, we are guaranteed the existence of a W such that $WW^\top A = I$. This suggests an optimization procedure in which W is a free parameter to be optimized that keeps the log determinant bound tight during the course of model training. We leave an implementation and evaluation of this adaptive log determinant bound strategy to future work.

Quadratic form calculation A naive calculation of the quadratic form $\mathbf{y}^\top \Lambda_{yy} \mathbf{y}$ takes time $\mathcal{O}(N^2T^2)$, which is prohibitive for large T . The key step is to perform matrix/vector products with

$$\Lambda_{yy} = (R^{-\frac{1}{2}} \otimes I_T)(I_{NT} - A^\top B^{-1}A)(R^{-\frac{1}{2}} \otimes I_T). \quad (16)$$

Multiplication with the first and third terms are easy. Remember $A = L^\top (C^\top R^{-\frac{1}{2}} \otimes I_T)$. Matrix/vector multiplication with $(C^\top R^{-\frac{1}{2}} \otimes I_T)$ is easy, but we do not have direct access to L^\top . However, we can efficiently find $L^{-\top} L^{-1} = \Lambda_{xx}$ by a Cholesky decomposition because Λ_{xx} is block-tridiagonal. Then we have $L^\top \mathbf{b} = \text{Solve}(L^{-\top}, \mathbf{b})$, with an efficient block-bidiagonal solve. Therefore, $A\mathbf{b}$ can be efficiently calculated as $\text{Solve}(L^{-\top}, (C^\top R^{-\frac{1}{2}} \otimes I_T)\mathbf{b})$.

Now for $B^{-1}\mathbf{b}$. We can write $B = L^\top [C^\top R^{-1}C \otimes I_T + \Lambda_{xx}]L$. Then

$$B^{-1} = L^{-1} [C^\top R^{-1}C \otimes I_T + \Lambda_{xx}]^{-1} L^{-\top}, \quad (17)$$

and as we saw before, matrix/vector multiplication with L^{-1} and $L^{-\top}$ is efficient. If P is the permutation matrix that switched the order of the Kronecker product expansion, then $P(C^\top R^{-1}C \otimes I_T)$ is block-diagonal with T separate $C^\top R^{-1}C$ blocks and $P\Lambda_{xx}$ is symmetric block-tridiagonal with diagonal blocks. It follows that $P[C^\top R^{-1}C \otimes I_T + \Lambda_{xx}]$ is block-tridiagonal with $3T - 2$ separate M -by- M blocks. Solve operations can be performed efficiently with block tridiagonal matrices, $\mathcal{O}(M^3T)$ in this case. Note the permutation P can be applied as a $\mathcal{O}(1)$ reshaping operation as opposed to a $\mathcal{O}(M^2T^2)$ matrix/vector multiplication.

To summarize, the quadratic form calculation is given by:

1. $\mathbf{b} \leftarrow (R^{-\frac{1}{2}} \otimes I_T)\mathbf{y}$
2. $\tilde{\mathbf{b}} \leftarrow \text{Solve}(L^{-\top}, (C^\top R^{-\frac{1}{2}} \otimes I_T)\mathbf{b})$
3. $\tilde{\mathbf{b}} \leftarrow P^\top \text{Solve}(P[C^\top R^{-1}C \otimes I_T + \Lambda_{xx}], \tilde{\mathbf{b}})$
4. $\tilde{\mathbf{b}} \leftarrow (R^{-\frac{1}{2}}C \otimes I_T)\text{Solve}(L^{-1}, \tilde{\mathbf{b}})$
5. $\mathbf{b} \leftarrow \mathbf{b} - \tilde{\mathbf{b}}$
6. $\mathbf{b} \leftarrow (R^{-\frac{1}{2}} \otimes I_T)\mathbf{b}$
7. Return $\mathbf{y}^\top \mathbf{b}$

ELBO estimation The complete CSLDS objective function is

$$\begin{aligned} -\frac{1}{2} [T \log \det(R) + T \text{tr}(C^\top R^{-1}C) + \mathbf{y}^\top \Lambda_{yy} \mathbf{y} + NT \log(2\pi)] - D(q(\boldsymbol{\eta}_{1:T-1}) || p(\boldsymbol{\eta}_{1:T-1})) \\ \leq \mathcal{L}(\mathbf{y}_{1:T}; \theta, \phi) \leq \log p(\mathbf{y}_{1:T}), \end{aligned} \quad (18)$$

where the KL divergence $D(q(\boldsymbol{\eta}_{1:T-1}) || p(\boldsymbol{\eta}_{1:T-1}))$ is given in Eq. 7. This is a lower bound to the standard ELBO $\mathcal{L}(\mathbf{y}_{1:T})$ from Eq. 4 (due to the Σ_{yy} log determinant bound in Eq. 14), which is in turn a lower bound to the observation marginal log likelihood $\log p(\mathbf{y}_{1:T})$.