
Deep Variational Bayes Filters: Unsupervised Learning of State Space Models from Raw Data

Maximilian Karl, Maximilian Soelch, Justin Bayer*, Patrick van der Smagt*
Chair of Robotics and Embedded Systems, Department of Informatics,
Technische Universität München, Germany

Abstract

We introduce Deep Variational Bayes Filters (DVBF), a new method for unsupervised learning of latent Markovian state space models. Leveraging recent advances in Stochastic Gradient Variational Bayes, DVBF can overcome intractable inference distributions by means of variational inference. Thus, it can handle highly nonlinear input data with temporal and spatial dependencies such as image sequences without domain knowledge. Our experiments show that enabling backpropagation through transitions enforces state space assumptions and significantly improves information content of the latent embedding. This also enables realistic long-term prediction.

1 Introduction

Estimating probabilistic models for sequential data is central to many domains, such as audio, natural language or physical plants [5, 12, 3, 4, 9]. The goal is to obtain a model $p(\mathbf{x}_{1:T})$ that best reflects a data set of observed sequences $\mathbf{x}_{1:T}$. Recent advances in deep learning have paved the way to powerful models capable of representing high-dimensional sequences with temporal dependencies, e.g. [5, 12, 3, 1].

A typical model assumption in systems theory is that the observed sequence $\mathbf{x}_{1:T}$ is generated by a corresponding latent sequence $\mathbf{z}_{1:T}$. More specifically, *state space models* assume the latent sequence to be Markovian, i.e., \mathbf{z}_t contains all information on the distribution of \mathbf{z}_{t+1} . Moreover, the *emission distribution* of \mathbf{x}_t is assumed to be determined by the corresponding \mathbf{z}_t . In short, we assume a latent state \mathbf{z}_t that holds all information available at time step t . This results in the following assumptions:

$$p(\mathbf{x}_{1:T} \mid \mathbf{z}_{1:T}, \mathbf{u}_{1:T}) = \prod_{t=1}^T p(\mathbf{x}_t \mid \mathbf{z}_t) \quad (1)$$

$$p(\mathbf{z}_{1:T} \mid \boldsymbol{\beta}_{1:T}, \mathbf{u}_{1:T}) = \prod_{t=0}^{T-1} p(\mathbf{z}_{t+1} \mid \mathbf{z}_t, \mathbf{u}_t, \boldsymbol{\beta}_t) \quad (2)$$

with \mathbf{u}_t as current control input and $\boldsymbol{\beta}_t$ as transition parameters.

We consider modeling a time-discrete, non-linear dynamical system with *observations* in some space $\mathcal{X} \subset \mathbb{R}^{n_x}$, depending on *control inputs* (or *actions*) from the space $\mathcal{U} \subset \mathbb{R}^{n_u}$. Elements of \mathcal{X} can be high-dimensional sensory data such as raw images, or any other state observation. With $\mathbf{x}_t \in \mathcal{X}$, let $\mathbf{x}_{1:T} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T)$ be a sequence of length T of observations. Similarly, with $\mathbf{u}_t \in \mathcal{U}$, let $\mathbf{u}_{1:T} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_T)$ be a corresponding sequence of equal length T of control inputs, which we consider as given. We are interested in deriving² a probabilistic model $p(\mathbf{x}_{1:T} \mid \mathbf{u}_{1:T})$.

*Justin Bayer and Patrick van der Smagt are affiliated with Data Lab, Volkswagen Group.

²The case without control inputs can be recovered by setting $\mathcal{U} = \emptyset$, i.e., not conditioning on control inputs.

Efficient inference of such latent states is only partially solved with state-space models. Under strong assumptions on the system, one can derive optimal Bayesian filters, such as the classical Kalman filter [7] for linear Gaussian models (LGMs). Yet, for less restrictive models, posterior distributions $p(\mathbf{z}_{1:T} | \mathbf{x}_{1:T})$ are often intractable.

Leveraging a recently proposed estimator based on variational inference, stochastic gradient variational Bayes (SGVB, [8, 11]), approximate inference of latent variables becomes tractable.

The principle of SGVB has been transferred to time series [1, 3, 6, 10, 12]. Some³ of these models violate Eq. (2), require inference subroutines, only softly encode the state-space assumptions (1) and (2) in the KL-divergence or fail to be a mathematically correct lower bound to the marginal data likelihood.

The contribution of this work is, to our knowledge, the first model that (i) *enforces* the state-space model assumptions in latent space allowing for reliable and plausible long-term prediction of the observable system, (ii) inherits the merit of neural architectures to be trainable on raw data such as images, audio or other sensory inputs and (iii) scales to large data due to optimization of parameters based on stochastic gradient descent [2].

2 Deep Variational Bayes Filters

2.1 Reparametrizing the Transition

Previous approaches emphasized good reconstruction, so that the space only contains information necessary for reconstruction of one time step. Similar to the reparametrization trick from [8, 11], we establish gradient paths through transitions over time so that the transition becomes the driving factor for shaping the latent space, rather than adjusting the transition to the recognition model’s latent space:

$$\mathbf{z}_{t+1} = f(\mathbf{z}_t, \mathbf{u}_t, \beta_t) \quad (3)$$

Given the stochastic parameters β_t , the state transition is deterministic (which in turn means that by marginalizing β_t , we still have a stochastic transition). The immediate and crucial consequence is that errors in reconstruction of \mathbf{x}_t from \mathbf{z}_t are backpropagated directly through time. This is different to the method used in [10], where the transition is optimized by minimizing a KL divergence. No gradient from the generative model is backpropagated through the transitions.

Fig. 1a shows a generic view on our new computational architecture. Fig. 1b shows an example for Eq. (3), a locally linear transition inspired by [12]. In this case we set

$$\mathbf{z}_{t+1} = \mathbf{A}_t \mathbf{z}_t + \mathbf{B}_t \mathbf{u}_t + \mathbf{C}_t \mathbf{w}_t, \quad t = 1, \dots, T, \quad (4)$$

2.2 The Lower Bound Objective Function

In analogy to VAEs [8, 11], we now derive a lower bound to the marginal likelihood $p(\mathbf{x}_{1:T} | \mathbf{u}_{1:T})$. After reflecting the Markov assumptions (1) and (2) in the factorized likelihood and due to the deterministic transition given β_{t+1} , we have:

$$p(\mathbf{x}_{1:T} | \mathbf{u}_{1:T}) = \int p(\beta_{1:T}) \prod_{t=1}^T p_\theta(\mathbf{x}_t | \mathbf{z}_t) \Big|_{\mathbf{z}_t=f(\mathbf{z}_{t-1}, \mathbf{u}_{t-1}, \beta_{t-1})} d\beta_{1:T}$$

We now derive the objective function, a lower bound to the data likelihood:

$$\begin{aligned} \ln p(\mathbf{x}_{1:T} | \mathbf{u}_{1:T}) &\geq \mathbb{E}_{q_\phi}[\ln p_\theta(\mathbf{x}_{1:T} | \mathbf{z}_{1:T})] - \text{KL}(q_\phi(\beta_{1:T} | \mathbf{x}_{1:T}, \mathbf{u}_{1:T}) || p(\beta_{1:T})) \\ &=: \mathcal{L}_{\text{DVBF}}(\mathbf{x}_{1:T}, \theta, \phi | \mathbf{u}_{1:T}) \end{aligned} \quad (5)$$

3 Dynamic Pendulum Experiments

In order to test our algorithm on truly non-Markovian observations of a dynamical system, we simulated a dynamic torque-controlled pendulum governed by the differential equation

$$ml^2 \ddot{\varphi}(t) = -\mu \dot{\varphi}(t) + mgl \sin \varphi(t) + u(t),$$

³Details about the specific differences can be found in the full version on <http://arxiv.org/abs/1605.06432>

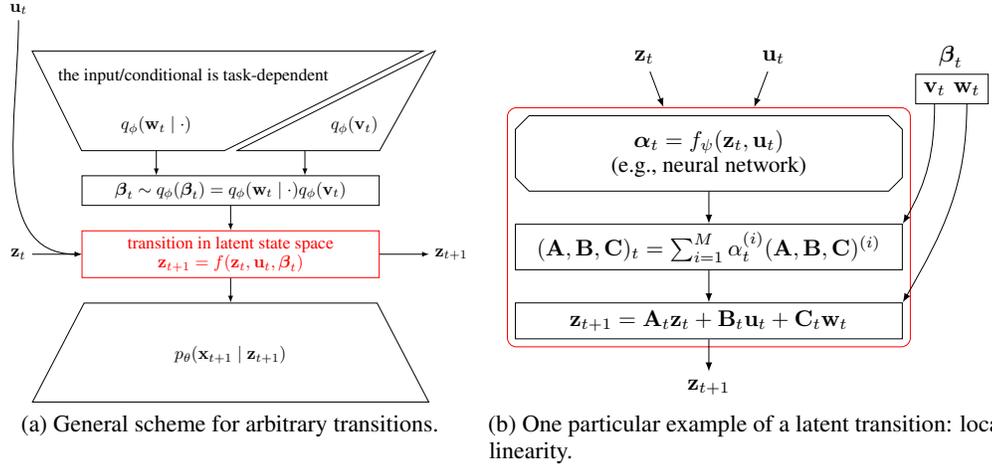


Figure 1: Left: General architecture for DVBF. Stochastic transition parameters β_t are inferred via the recognition model, e.g., a neural network. Based on a sampled β_t , the state transition is computed deterministically. The updated latent state z_{t+1} is used for predicting x_{t+1} . For details, see Section 2.1. Right: Zoom into latent space transition (red box in left figure).

$m = l = 1, \mu = 0.5, g = 9.81$, via numerical integration, and then converted the ground-truth angle φ into an image observation in \mathcal{X} . The one-dimensional control corresponds to angle acceleration (which is proportional to joint torque). Angle and angular velocity fully describe the system.

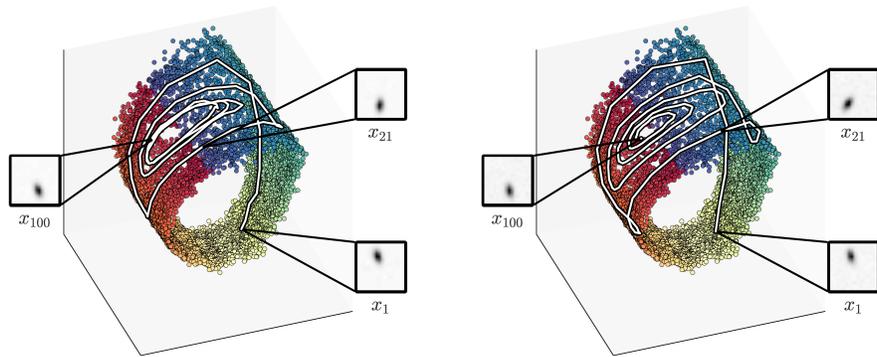
Fig. 2 shows that the strong relation between ground truth and latent state is beneficial for generative sampling. All plots show 100 time steps of a pendulum starting from the exact same latent state and not being actuated. The top row plots show a purely generative walk in the latent space on the left, and a walk in latent space that is corrected by filtering observations on the right. We can see that both follow a similar trajectory to an attractor. The bottom plot shows the first 45 steps of the corresponding observations (top row), reconstructions (middle row), and generative samples (without correcting from observations). Interestingly, DVBF works very well even though the sequence is much longer than all training sequences (indicated by the red line).

4 Conclusion

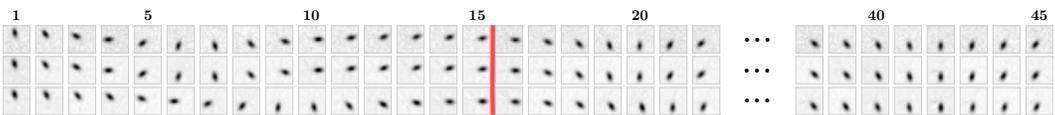
We have proposed Deep Variational Bayes Filters (DVBF), a new method to learn state space models from raw non-Markovian sequence data. DVBFs make use of stochastic gradient variational Bayes to overcome intractable inference and thus naturally scale to large data sets. In a vision-based experiment we demonstrated that latent states can be recovered which identify the underlying physical quantities. The generative model showed stable long-term predictions far beyond the sequence length used during training.

References

- [1] Justin Bayer and Christian Osendorfer. Learning stochastic recurrent networks. *arXiv preprint arXiv:1411.7610*, 2014.
- [2] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010*, pages 177–186. Springer, 2010.
- [3] Junyoung Chung, Kyle Kastner, Laurent Dinh, Kratarth Goel, Aaron C. Courville, and Yoshua Bengio. A recurrent latent variable model for sequential data. *CoRR*, abs/1506.02216, 2015.
- [4] Marc Deisenroth and Carl E Rasmussen. Pilco: A model-based and data-efficient approach to policy search. In *Proceedings of the 28th International Conference on machine learning (ICML-11)*, pages 465–472, 2011.



(a) Generative latent walk. (b) Reconstructive latent walk.



(c) Ground truth (top), reconstructions (middle), generative samples (bottom) from identical initial latent state.

Figure 2: (a) Latent space walk in generative mode. (b) Latent space walk in filtering mode. (c) Ground truth and samples from recognition and generative model. The reconstruction sampling has access to observation sequence and performs filtering. The generative samples only get access to the observations once for creating the initial state while all subsequent samples are predicted from this single initial state. The red bar indicates the length of training sequences. Samples beyond show the generalization capabilities for sequences longer than during training.

[5] Alex Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, 2013.

[6] Matthew J Johnson, David Duvenaud, Alexander B Wiltschko, Sandeep R Datta, and Ryan P Adams. Structured VAEs: Composing probabilistic graphical models and variational autoencoders. *arXiv preprint arXiv:1603.06277*, 2016.

[7] Rudolph E Kalman and Richard S Bucy. New results in linear filtering and prediction theory. *Journal of basic engineering*, 83(1):95–108, 1961.

[8] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.

[9] Jonathan Ko and Dieter Fox. Learning gp-bayesfilters via gaussian process latent variable models. *Autonomous Robots*, 30(1):3–23, 2011.

[10] Rahul G Krishnan, Uri Shalit, and David Sontag. Deep Kalman filters. *arXiv preprint arXiv:1511.05121*, 2015.

[11] Danilo J. Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In Tony Jebara and Eric P. Xing, editors, *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 1278–1286. JMLR Workshop and Conference Proceedings, 2014.

[12] Manuel Watter, Jost Springenberg, Joschka Boedecker, and Martin Riedmiller. Embed to control: A locally linear latent dynamics model for control from raw images. In *Advances in Neural Information Processing Systems*, pages 2728–2736, 2015.