Variational Inference for Switching Dynamics

Variierende Inferenz für wechselnde Dynamiken Master thesis by Thomas Lautenschläger Date of submission: August 31, 2020

 Review: M.Sc. Hany Abdulsamad
 Review: Prof. Jan Peters, Ph.D. Darmstadt



TECHNISCHE UNIVERSITÄT DARMSTADT



Erklärung zur Abschlussarbeit gemäß §22 Abs. 7 und §23 Abs. 7 APB der TU Darmstadt

Hiermit versichere ich, Thomas Lautenschläger, die vorliegende Masterarbeit ohne Hilfe Dritter und nur mit den angegebenen Quellen und Hilfsmitteln angefertigt zu haben. Alle Stellen, die Quellen entnommen wurden, sind als solche kenntlich gemacht worden. Diese Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen.

Mir ist bekannt, dass im Fall eines Plagiats (§38 Abs. 2 APB) ein Täuschungsversuch vorliegt, der dazu führt, dass die Arbeit mit 5,0 bewertet und damit ein Prüfungsversuch verbraucht wird. Abschlussarbeiten dürfen nur einmal wiederholt werden.

Bei der abgegebenen Thesis stimmen die schriftliche und die zur Archivierung eingereichte elektronische Fassung gemäß §23 Abs. 7 APB überein.

Bei einer Thesis des Fachbereichs Architektur entspricht die eingereichte elektronische Fassung dem vorgestellten Modell und den vorgelegten Plänen.

Darmstadt, 31. August 2020

Thomas Lautenschläger

Abstract

Learning and inferring the dynamics of a system ideally requires a model design that expresses the system's behavior, does not overfit to the training data and uses few computational resources. Bayesian probabilistic models have the powerful feature to express models as simple as possible while maximizing the model's information. Moreover, being Bayesian makes the model's uncertainty accessible. This quantity can be advantageous by integrating a learned dynamics model into policy search for planning. In this thesis, we provide a Bayesian extension for recurrent hidden Markov models and analyze their performance capabilities by training them with data sampled from dynamical systems. Since hidden Markov models capture switching regime behavior from data emitted by dynamical systems, we introduce the policy extension Switching actor (SWAC), a policy extension that incorporates the switching regime detection into policy search. This extension supplies additional structure to the policy. We present a performance comparison of SWAC with the non SWAC equipped baselines. In addition, we give an outlook to future work. The complete mathematical derivation is attached as supplementary material.

Zusammenfassung

Bayes'sche probabilistische Modelle haben die mächtige Eigenschaft, die Modelle so simpel wie möglich zu Gestalten und gleichzeitig den Informationsgehalt der Modelle zu maximieren. Außerdem bieten Bayes'sche probabilistische Modelle auf die Ungewissheit des Modells zu quantifizieren und diese auch zugänglich zu machen. In dieser Thesis definieren wir eine Bayes'sche Erweiterung für rekurrente verdeckte Markowmodelle und analysieren deren Leistungsverhalten. Zuvor werden diese Modelle mit Daten trainiert, die aus dynamischen Systemen stammen. Da verdeckte Markowmodelle die Eigenschaft haben, ein wechselndes Schaltverhalten aus Daten, die aus einem dynamischen System emittiert sind, festzustellen, integrieren wir diese Modelle für die Policysuche. Diese Methode führen wir unter der Bezeichnung SWAC ein. Diese Erweiterung stattet eine Policy mit zusätzlicher Struktur aus. Wir präsentieren einen Leistungsvergleich zwischen SWAC und den nicht mit SWAC ausgestatteten Referenzmodellen. Die gesamte mathematische Herleitung ist im Anhang beigefügt.

Contents

1.	Introduction			
2.	Related work2.1. Bayesian HMMs2.2. Model-based reinforcement learning2.3. Switching linear dynamical systems2.4. Pathwise gradient estimator	4 4 5 5		
3.	Background 3.1. Hidden Markov models	7 7 7 9 9 11 12		
	3.2.1. The evidence lower bound 1 3.2.2. Mean-field approximation 1 3.3. Minimizing the gradient variance 1	13 13 14		
4.	Bayesian nidden Markov models 4.1. Variational Bayes HMM	16 16 18 19 19 21 26		

	4.2.	Variational Bayes ARHMM4.2.1. Model priors4.2.2. Target distribution4.2.3. Mean-field approximation4.2.4. Variational E-step4.2.5. Variational M-step4.2.6. The evidence lower boundVariational Bayes rAR-HMM4.3.1. Model priors4.3.2. Target distribution4.3.3. Mean-field approximation4.3.4. Variational E-step4.3.5. Variational M-step	27 29 29 30 32 36 36 36 38 38 38 39 39
5.	Swit	ching actor	41
0.	5.1. 5.2. 5.3.	Implementation Extending A2C with SWAC Extending SAC with SWAC Extending SAC with SWAC	41 42 43
6.	Eval	uation	44
	 6.1. 6.2. 6.3. 6.4. 6.5. 	Bayesian HMMs . 6.1.1. Viterbi state detection . 6.1.2. Model learning performance . Learning control on rAR-HMM . 6.2.1. Policy evaluation on the pendulum task . 6.2.2. Policy evaluation on the cart-pole task . SWAC-A2C . SWAC-SAC . Discussion .	44 45 46 47 48 49 50 51
7.	Outle	ook	52
8.	Cond	clusion	53
Α.	Sum A.1. A.2.	Imary VBEM steps Imary VBEM steps VBHMM EM steps Image: Steps VBAR-HMM EM steps Image: Steps	60 60 61

	A.3. V	rAR-HMM EM steps	53
В.	Deriva B.1. D B.2. D	on VBHMMs 6 rivation VBHMM	55 56 72
C.	Distrib	tions 7	79
	C.1. D	ichlet distribution	79
	C.2. U	iform distibution	30
	C.3. M	ltivariate normal distribution	31
		short distribution	აი

Figures

List of Figures

3.1.	Three different types of Hidden Markov model (HMM)s. The models expressiveness is ordered from left to right.	8
6.1.	Randomly sampled trajectories on the pendulum (left) and cart-pole (right) environment. Each color represents a detected latent state from the Variational Bayes recurrent autoregressive hidden Markov model (VBrAR-HMM) using the Viterbi algorithm. The state transitions match the repeating dynamic patterns	45
6.2.	Out-of-sample predictive performance of the Variational Bayes autoregres- sive hidden Markov model (VBAR-HMM)s and VBrAR-HMMs to state-of- the-art (SOTA) models. Evaluation on the pendulum (left) and cart-pole (right) environment.	46
6.3.	Policy evaluation on the rAR-HMM dynamics (left) and on the true dynamics (right) of the pendulum environment. The policy training observations are sampled from the rAR-HMM dynamics.	47
6.4.	Policy evaluation on the rAR-HMM dynamics (left) and on true dynamics (right) of the pendulum environment. The policy training observations are sampled from the rAR-HMM dynamics	48
6.5.	Policy evaluation performance comparison of Advantage actor-critic (A2C), SWAC-A2C and LAX on the pendulum (left) and on the cart-pole (right) environment.	49

Abbreviations

List of Abbreviations

Notation	Description			
VBAR-HMM	Variational Bayes autoregressive hidden Markov model			
VBHMM	Variational Bayes hidden Markov model			
VBrAR-HMM	Variational Bayes recurrent autoregressive hidden Markov model			
A2C	Advantage actor-critic			
AR-HMM	Autoregressive hidden Markov model			
BNN	Bayesian neural network			
ELBO	Evidence lower bound			
EM	Expectation maximization			
FNN	Feed-forward neural network			
GMM	Gaussian mixture model			
GP	Gaussian process			

Hidden Markov model			
independently and identically distributed			
Long short-term memory neural networks			
Maximum a posteriori probability			
Model-based reinforcement learning			
Markov decision property			
Recurrent autoregressive hidden Markov model			
Recurrent neural network			
Soft actor-critic			
switching linear-Gaussian dynamical systems			
state-of-the-art			
Switching actor			
Variational autoencoders			

1. Introduction

The optimal control for a plane to fly stable depends on specific target scenarios. The plane is starting, being in target height or on landing approach. As the scenarios differ, the control needs to be adapted for each scenario. This example motivates this thesis's goal to find a model that is able to detect the different scenarios and apply the optimal control on each.

Policy optimization for control of complex dynamical systems requires a policy search framework that is capable of capturing the system's dynamics. Since the dynamical systems can be broken into simpler units [1], the policy can be extended with additional structure that makes use of the simpler system units. This is comparable to a car's transmission system where the gear is chosen according to the specific scenario (driving backwards, accelerating, constant speed etc.).

This thesis proposes a Bayesian extension of HMM. These models are able to capture the complex system dynamics into simpler units with switching latent states. Furthermore, we combine these models with policy search to obtain a policy with additional structure.

In Chapter 2, we give an overview of the work related to ours.

In Chapter 3, we explain the basics of HMMs by showing the model structure, how to train and do inference with them. Moreover, we outline the method of variational inference and how this method can be applied by making use of mean-field approximation. Lastly, we describe a recent method that is able to reduce variance for score-function estimator gradient [2].

In Chapter 4, we derive the Bayesian HMMs. Starting the derivation from the classical Variational Bayes hidden Markov model (VBHMM), extending this derivation for the VBAR-HMM and finalizing it with the VBrAR-HMM.

In Chapter 5, we propose a policy search algorithm, called SWAC that requires a switching state space model. SWAC is internally structured with multiple policies. The number of

policies corresponds to the number of the switching state spaces. Besides, we show how to extend A2C and Soft actor-critic (SAC) with SWAC.

Chapter 6 gives a switching state detection and a prediction performance analysis of the introduced Bayesian HMMs from Chapter 4. We compare their prediction performance to SOTA models. In addition, we train Recurrent autoregressive hidden Markov model (rAR-HMM)s from true environment observations. We perform policy search on the rAR-HMM dynamics and evaluate the trained policies on the true dynamics. Lastly, we compare the learning performance of A2C and SAC with the extension of SWAC and without this extension. We use the pendulum and the cart-pole [3] environments for the evaluation.

Chapter 7 gives an outlook of future work and in Chapter 8 we conclude this thesis.

2. Related work

2.1. Bayesian HMMs

[4] derives a solution for Bayesian HMMs with discrete observations. He showed the advantages of being Bayesian compared to non Bayesian HMMs. The HMM is Bayesian in the sense that a prior distribution is put over the parameters and the Evidence lower bound (ELBO) is maximized with respect to (w.r.t.) the assumed distribution parameters. Beal shows the performance of automatically detecting the amount of hidden states that are at least required to describe the data. He assumes a Dirichlet distribution over the state transitions. Makes use of the property of the Dirichlet distribution that cancels out the number of states that are not required to express the data without overfitting. In addition, being Bayesian has the advantage to express the uncertainty of the learned model and is thus is less prone to overfitting. The obtained uncertainty can be used for exploration tasks.

2.2. Model-based reinforcement learning

The goal of reinforcement learning is to compute a policy that returns an optimal control signal for an underlying dynamical model to achieve a defined task such as a stabilization task.

Learning the dynamical model and integrating it into the policy search procedure is called Model-based reinforcement learning (MBRL). Making use of the learned model can result to a more sample efficient policy search.

[5] combines the learned dynamical model and uses a deterministic policy to obtain an analytic gradient for the policy parameter update. They use Gaussian process (GP) to capture the model dynamics. Moreover, they make use of the uncertainty of the GPs and

incorporate this quantity in a planning step. The planning step involves to use the learned the dynamics model to predict the dynamics behavior with a control signal obtained from the policy.

[6] formulate a general formulation for the combination of a deterministic policy and learning the model dynamics. The long-term reward is $J_{\theta} = \sum_{t=1}^{T} \gamma^t r(\mathbf{x}_t)$ where $\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t)$. f(.) is the learned dynamical model, \mathbf{x} are the observations and \mathbf{u} is the control signal from the deterministic policy. Applying the chain rule on J_{θ} w.r.t. to the policy parameters θ gives:

$$\frac{\partial J_{\boldsymbol{\theta}}}{\partial \boldsymbol{\theta}} = \sum_{t=1}^{T} \gamma^{t} \frac{\partial r(\mathbf{x}_{t})}{\partial \mathbf{x}_{t}} \left(\frac{\partial \mathbf{x}_{t}}{\partial \mathbf{x}_{t}} \frac{\partial \mathbf{x}_{t}}{\partial \boldsymbol{\theta}} + \frac{\partial \mathbf{x}_{t}}{\partial \mathbf{u}_{t}} \frac{\partial \mathbf{u}_{t}}{\partial \boldsymbol{\theta}} \right)$$
(2.1)

where we can differentiate through the learned dynamics model and the deterministic policy.

Recent research makes use of Variational autoencoders (VAE) [7, 8, 9] as a general model that captures the model dynamics. VAEs enable the differentiation through the learned dynamics model in a simple fashion that does not require specific model specification [10].

2.3. Switching linear dynamical systems

Learning time series data emitted from a dynamical system, requires a model design that is capable of capturing the dynamical behavior over time. [1] introduce the extension of a switching linear-Gaussian dynamical systems (SLDS)[11, 12, 13, 14, 15] with a recurrent connection where the next latent discrete state depends on the current continuous latent states. The recurrent connection is captured using logistic regression. This extension results in a more interpretable model. The data generated with this model better expresses the true dynamics.

2.4. Pathwise gradient estimator

Gradient estimation for parameter optimization of a distribution function can lead to gradients with a large variance when using estimators such as the score function estimator [2]. A less general solution are pathwise gradient estimators. They make use of the

structure of the problem and can thus lead to minimizing the variance of the estimated gradient.

The idea is to push in the drawn sample from the target distribution into the cost function and differentiate the cost function w.r.t. the distribution function parameters. The drawn sample needs to be reparameterized such that the differentiation is possible. The sampling process becomes:

$$\hat{\mathbf{x}} \sim p(\mathbf{x}; \boldsymbol{\theta}) \equiv \hat{\mathbf{x}} = g(\hat{\boldsymbol{\epsilon}}, \boldsymbol{\theta}), \ \hat{\boldsymbol{\epsilon}} \sim p(\boldsymbol{\epsilon})$$

where g(.) is a differentiable function with parameters θ and $p(\epsilon)$ is independent of the parameters θ .

[16] shows the reparameterization procedure for Gaussian distribution for continuous data. [17] and [18] simultaneously introduced a relaxation of the uniform distribution that made gradient passing possible for discrete data.

[19] and [20] make use of the pathwise gradient estimator to profit from the reduced variance of this gradient estimator. In Section 3.3 we explain their approach in detail.

3. Background

3.1. Hidden Markov models

HMMs are probabilistic models that expresses the switching dynamics of time series data under the Markov decision property (MDP) assumption. HMMs find their application such as in a wide range of fields such as capturing robot dynamics, weather forecasting, speech recognition, financial time series prediction.

This section provides the basics for HMMs. We present the basic parameters and show how to apply HMMs for learning and inference on Gaussian independently and identically distributed (i.i.d.) data.

3.1.1. Overview of HMM types

In this thesis we cover three types of HMMs. The classical HMM depicted in Figure 3.1a where the observation \mathbf{x}_t are emitted from the discrete latent states z_t . A linear dependence from the current observation \mathbf{x}_t to the next observation \mathbf{x}_{t+1} is the Autoregressive hidden Markov model (AR-HMM) as represented in Figure 3.1b. Figure 3.1c shows the rAR-HMM which extends the AR-HMM with a recurrent link of the current observation \mathbf{x}_t to the next discrete latent state z_{t+1} .

3.1.2. Structure of HMMs

A HMM has captures the transitions between the latent variables z. The observations x are emissions from the latent variables. Figure 3.1a shows the dependencies of a HMM.



Figure 3.1.: Three different types of HMMs. The models expressiveness is ordered from left to right.

In this thesis we focus on HMMs with Gaussian i.i.d. observations. The latent variables z absorb the Gaussian parameters. This model is similar to a Gaussian mixture model (GMM). The difference is the dependence to discrete time steps between of the data points.

Joint-distribution

The joint-distribution of a HMM is [21]:

$$p(\mathbf{z}_{1:T}, \mathbf{x}_{1:T}) = p(z_1) \prod_{t=2}^{T} p(z_t | z_{t-1}) p(\mathbf{x}_t | z_t)$$

where $p(z_1)$ is the initial state probability vector, $p(z_t|z_{t-1})$ is the transition probability and $p(\mathbf{x}_t|z_t)$ are the observation probabilities. The initial state density is:

$$p(z_1) = \prod_{k=1}^K \pi_k^{\mathbb{I}(z_1=k)}$$

The density for Gaussian observation becomes:

$$p(\mathbf{x}_t|z_t) = \prod_{k=1}^{K} \mathcal{N}(\mathbf{x}_t|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{\mathbb{I}(z_t=k)}$$

and the transition density becomes:

$$p(z_t|z_{t-1}) = \prod_{k=1}^K \prod_{j=1}^K a_{kj}^{\mathbb{I}(z_t=j,z_{t-1}=k)}, \quad a_{kj} = p(z_t=j|z_{t-1}=k)$$

		-
	-	-
		_

where $z_t \in 1, ..., K$. *K* is the maximum number of discrete latent states in the model and **A** is the transition matrix with $a_{kj} \in \mathbf{A}$. We abbreviate the model parameters with $\boldsymbol{\theta} = (\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{A})$

3.1.3. Baum-Welch algorithm - learning a HMM

The Baum-Welch algorithm first introduced in [22] is the extension of the EM algorithm to HMMs. The Baum-Welch algorithm iteratively learns estimates of the model probabilities shown in the joint-distribution equation above.

3.1.3.1. E-step

Using the above defined model densities, we get the estimated log complete data likelihood from: κ

$$\log p(\mathbf{z}_{1:T}, \mathbf{x}_{1:T} | \boldsymbol{\theta}) = \sum_{k=1}^{K} \mathbb{E} \left[\mathbb{I}(z_1 = k) \right] \log \pi_k$$

+
$$\sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} \mathbb{E} \left[\mathbb{I}(z_{t-1} = k, z_t = j) \right] \log a_{kj} \qquad (3.1)$$

+
$$\sum_{t=1}^{T} \sum_{k=1}^{K} \mathbb{E} \left[\mathbb{I}(z_t = k) \right] \log p(\mathbf{x}_t | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

We make use of dynamic programming to solve the estimates from Equation (3.1). First, we compute the forward probabilities where we define the log posterior of the hidden states as follows:

$$\alpha_t(k) = p(z_t = k | \mathbf{x}_{1:t}) = \frac{p(z_t = k | \mathbf{x}_{1:t-1}) p(\mathbf{x}_t | z_t = k)}{\sum_{j=1}^{K} p(z_t = j | \mathbf{x}_{1:t-1}) p(\mathbf{x}_t | z_t = j)}$$

where $p(z_t = k | \mathbf{x}_{1:t-1})$ is the one-step ahead predictive from the start t = 1 until density given by:

$$p(z_t = k | \mathbf{x}_{1:t-1}) = \sum_{j=1}^{K} p(z_t = k | z_{t-1} = j) p(z_{t-1} = j | \mathbf{x}_{1:t-1})$$
$$= \sum_{j=1}^{K} p(z_t = k | z_{t-1} = j) \alpha_{t-1}(j)$$

We set the prior $\alpha_1(z_1) = 1$ and obtain the final forward pass solution:

$$\alpha_t(k) \propto p(\mathbf{x}_t | z_t = k) \sum_{j=1}^K p(z_t = k | z_{t-1} = j) \alpha_{t-1}(j)$$

Similar to the forward message passing, we compute the backward messages from t = T, ..., 1 which is the conditional likelihood of future evidence:

$$\begin{split} \beta_t(k) &= p(\mathbf{x}_{t+1:T}|z_t) \\ &= \sum_{j=1}^K p(z_t = j, \mathbf{x}_t, \mathbf{x}_{t+1:T}|z_{t-1} = k) \\ &= \sum_{j=1}^K p(\mathbf{x}_{t+1:T}|z_t = j) p(z_t = j, \mathbf{x}_t|z_{t-1} = k) \\ &= \sum_{j=1}^K p(\mathbf{x}_{t+1:T}|z_t = j) p(\mathbf{x}_t|z_t = j) p(z_t = j|z_{t-1} = k) \\ &= \sum_{j=1}^K \beta_t(j) p(\mathbf{x}_t|z_t = j) p(z_t = j|z_{t-1} = k) \end{split}$$

where the end condition is $\beta_T(k) = 1$.

We combine α_t and β_t and obtain the state marginal:

$$p(z_t = k | \mathbf{x}_{1:T}) \propto p(z_t = k | \mathbf{x}_{1:t}) p(\mathbf{x}_{t+1:T} | z_t = k)$$

= $\alpha_t(k) \beta_t(k)$ (3.2)

and the pairwise state marginal:

$$p(z_{t-1} = k, z_t = k' | \mathbf{x}_{1:T}) \propto p(z_{t-1} = k | \mathbf{x}_{1:t-1}) p(\mathbf{x}_t | z_t = k') p(\mathbf{x}_t | z_t = k') p(\mathbf{x}_{t+1:T} | z_t = k')$$

= $\alpha_{t-1}(k) p(z_t = k' | z_t = k) p(\mathbf{x}_t | z_t = k') \beta_t(k')$
(3.3)

Using the results from the forward-backward messages, we can solve the expectations from equation (3.1). We rewrite the expectations:

$$\gamma_t(k) = \mathbb{E} \left[\mathbb{I}(z_t = k) \right]$$
$$\xi_{t-1,t}(k,k') = \mathbb{E} \left[\mathbb{I}(z_{t-1} = k, z_t = k') \right]$$

and solve using the results from Equations (3.2) and (3.3):

$$\gamma_t(k) = \frac{\alpha_t(k)\beta_t(k)}{\sum_{j=1}^{K} \alpha_t(j)\beta_t(j)}$$

which is the posterior density of being in state k at time step t and

$$\xi_{t-1,t}(k,k') = \frac{\alpha_{t-1}(k)a_{kk'}p(\mathbf{x}_t|z_t = k', \boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'})\beta_t(k')}{\sum_{j=1}^{K}\sum_{j'}^{K}\alpha_{t-1}(j)a_{jj'}p(\mathbf{x}_t|z_t = j', \boldsymbol{\mu}_{j'}, \boldsymbol{\Sigma}_{j'})\beta_t(j')}$$

which is the state transition posterior density.

3.1.3.2. M-step

In the M-step we update the model parameters to a maximize the likelihood of the observations. We use the expectation results from the E-step for the model parameter updates. The M-step becomes:

$$\hat{\boldsymbol{\pi}} : \hat{\boldsymbol{\pi}}_{k} = \gamma_{1}(k)$$

$$\hat{\boldsymbol{\mu}} : \hat{\boldsymbol{\mu}}_{k} = \frac{\sum_{t=1}^{T} \gamma_{t}(k) \mathbf{x}_{t}}{\sum_{t=1}^{T} \gamma_{t}(k)}$$

$$\hat{\boldsymbol{\Sigma}} : \hat{\boldsymbol{\Sigma}}_{k} = \frac{\sum_{t=1}^{T} \gamma_{t}(k) \mathbf{x}_{t} \mathbf{x}_{t}^{T} - \sum_{t=1}^{T} \gamma_{t}(k) \hat{\boldsymbol{\mu}}_{k} \hat{\boldsymbol{\mu}}_{k}^{T}}{\sum_{t=1}^{T} \gamma_{t}(k)}$$

$$\hat{\mathbf{A}} : \hat{\boldsymbol{a}}_{kk'} = \frac{\sum_{t=2}^{T} \xi_{t-1,t}(k,k')}{\sum_{t=2}^{T} \sum_{k'=1}^{K} \xi_{t-1,t}(k,k')}$$

This EM routine converges to a local optimum regarding the model likelihood given the observations.

3.1.4. Viterbi algorithm - maximum likely states

The Viterbi algorithm comes in hand if we require the maximum likely latent state sequence given the observations and the trained HMM:

$$\hat{\mathbf{z}} = \operatorname*{arg\,max}_{\mathbf{z}_{1:T}} p(\mathbf{z}_{1:T} | \mathbf{x}_{1:T})$$

where $\hat{\mathbf{z}}$ is the most likely path from the trained model.

3.2. Variational inference

Variational inference [23, 24, 25, 26, 27, 28] tackles the problem to find an approximate density posterior model given an assumed prior model. We assume a probabilistic graphical model with observations \mathbf{x} that depend on parameterized distributions. The parameters are absorbed in the latent variables \mathbf{z} . If we are interested in the posterior $p(\mathbf{z}|\mathbf{x})$, we need to apply Bayes' rule and solve:

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{\int p(\mathbf{x}|\mathbf{z}_i)p(\mathbf{z}_i) \,\mathrm{d}\mathbf{z}_i}$$

$$= \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$$
(3.4)

Marginalizing out z in the denominator can be intractable to compute, because the number of the latent variables z can be infinitely large.

One way to obtain a posterior solution is to utilize Monte Carlo sampling methods. This however can be slow to compute, because of the large amount of samples required to approximate or converge to the exact posterior [29].

Variational inference comes at hand if the goal is to approximate a posterior solution in a decent amount of time. To approximate the posterior, we need to introduce a variational distribution q and seek to minimize the Kullback–Leibler divergence (KL divergence) of p and q. This reformulates the integration problem of finding a posterior solution to an optimization problem as follows:

$$q^{*}(\mathbf{z}) = \underset{q(\mathbf{z})\in\mathcal{D}}{\arg\min} D_{\mathrm{KL}}(q(\mathbf{z})|| p(\mathbf{z}|\mathbf{x}))$$

where \mathcal{D} is a set of approximate distributions.

In Addition, we assume the priors are conjugate to the posteriors which is the case for distributions of the exponential families. Variational inference for non-conjugate models is discussed in [30].

3.2.1. The evidence lower bound

The objective (3.4) is hard to compute since the KL divergence is:

$$D_{\mathrm{KL}}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x})) = \mathbb{E}_q\left[\log q(\mathbf{z})\right] - \mathbb{E}_q\left[\log p(\mathbf{z},\mathbf{x})\right] + \log p(\mathbf{x})$$
(3.5)

where we have the dependence on $\log p(\mathbf{x})$.

An equivalent formulation of the variational optimization problem instead of minimizing (3.5) is to maximize the evidence lower bound (ELBO) defined as:

$$\text{ELBO}(q) = \mathbb{E}_q \left[\log p(\mathbf{z}, \mathbf{x}) \right] - \mathbb{E}_q \left[\log q(\mathbf{z}) \right]$$
(3.6)

A property of the ELBO is that it lower-bounds the log evidence $\log p(\mathbf{x}) \ge \text{ELBO}(q)$ which explains the name.

As [27] shows, combining the equations (3.5) and (3.6) using Jensen's inequality gives:

$$\log p(\mathbf{x}) = D_{\mathrm{KL}}(q(\mathbf{z}) || p(\mathbf{z} | \mathbf{x})) + \mathrm{ELBO}(q)$$

This makes it clear that maximizing the ELBO is equivalent to minimizing the KL divergence of (3.5). We focus on the finding a solution w.r.t. maximizing the ELBO.

3.2.2. Mean-field approximation

We need to restrict the set of approximate distributions $q(\mathbf{z})$ to obtain a tractable solution for this optimization problem. Assuming that $q(\mathbf{z})$ factorizes as follows:

$$q(\mathbf{z}) = \prod_{i=1}^{M} q_i(\mathbf{z}_i)$$
(3.7)

13

where the elements of z are partitioned into disjoint groups.

Using this as the only assumption of q(z). Bishop [31] derives a general solution for the variational optimization problem. The general solution is:

$$\log q_j(\mathbf{z}_j) = \mathbb{E}_{i \neq j}[\log p(\mathbf{x}, \mathbf{z})] + \text{const.}$$
(3.8)

The intuition of the general solution (3.8) is that we fix the the disjoint variational solution for $\log q_j(\mathbf{z}_j)$ and average over all other disjoint groups. The non fixed disjoint groups are called the free energy in this context as the word "free" describes the non fixed property.

3.3. Minimizing the gradient variance

We seek to compute the optimal control policy parameters θ of a policy π_{θ} using gradient based methods. This section gives an overview of the problems that arise applying gradient-based methods on model-free policy search and how to overcome them.

Score-function gradient estimator (REINFORCE)

The score-function estimator [2] is a general applicable solution to obtain policy parameter gradients:

$$\frac{\partial J(\pi_{\boldsymbol{\theta}})}{\partial \boldsymbol{\theta}} = \sum_{t=1}^{T} \frac{\partial \log \pi_{\boldsymbol{\theta}}(a_t|s_t)}{\partial \boldsymbol{\theta}} \sum_{t'=t}^{T} r(a_{t'}, s_{t'})$$

This is an unbiased estimator but it yields high variance because it does not use any information of r(s) (reward function) to compute the gradient.

Variance reduction using control variates

Since r(s) can point in arbitrary directions and thus causes a high variate gradient, we need to limit the arbitrariness of r(s).

Control variates tackle this problem. As [32] shows, the ideal control variate for this variance problem correlates at the maximum with r(s). To obtain this result, we can use the state value function, that we can train using additive path reward samples of r(s).

LAX surrogate gradient

To remove additional variance from the estimated gradient, [20] introduced a control variate in form of a neural network with the objective to minimize the variance of the estimated reinforce gradient. The extension of reinforce with LAX is as follows:

$$\hat{g}_{\boldsymbol{\theta}}^{LAX} = \sum_{t=1}^{T} \frac{\partial \log \pi_{\boldsymbol{\theta}}(a_t|s_t)}{\partial \boldsymbol{\theta}} \left[\sum_{t'=t}^{T} r(a_{t'}, s_{t'}) - b(s_{t'}) - c_{\boldsymbol{\phi}}(a_{t'}, s_{t'}) \right] + \frac{\partial c_{\boldsymbol{\phi}}(a_{t'}, s_{t'})}{\partial \boldsymbol{\theta}}$$

As a_t is sampled from π_{θ} , we need to make reparameterizable to obtain the gradient of of c_{ϕ} w.r.t. the parameters θ . This reparameterization yields the following form of $a_t = a(\epsilon_t, s_t, \theta)$ where $\epsilon_t \sim p(\epsilon_t)$ and ϵ_t does not depend on θ . The last term is forms a pathwise derivative

We compute the gradient for the neural network control variate as follows:

$$\hat{g}_{\phi} = \frac{\partial \left(\hat{g}_{\theta}^{LAX}\right)^2}{\partial \phi}$$

This formulation of the gradient computation is implicitly defined to minimize the variance of the reinforce gradient.

4. Bayesian hidden Markov models

A HMM from Section 3.1 is able to approximate the dynamics of a model. Since we approximate the model dynamics, we do not know how uncertain the model behaves comparing to the true dynamics. Knowing the model uncertainty would be useful when we integrate the trained dynamics model into policy search e.g. for planning or exploration [33, 34, 5, 6, 35].

Making the model Bayesian by putting prior distributions on the model parameters enables to access uncertainty of the model.

In this chapter we derive Bayesian HMMs for all HMMs from Figure 3.1.

4.1. Variational Bayes HMM

In this section, we define a Variational Bayes HMM and derive the Expectation maximization (EM) update steps. The model we derive in this section extends the HMM model from section 3.1. We extend the model with priors on the model parameters. The complete derivation is in the Appendix B.1.

4.1.1. Model priors

We extend the model from section 3.1 by putting priors over the parameters π , A, μ and Σ .

A natural choice for the priors on π and on the rows of **A** are Dirichlet priors. In addition, the Dirichlet prior has the effect to implicitly adjust the number of the latent states that are required for the HMM to be express the data. The latent states that are not used

are virtually canceled out [4]. Details of the Dirichlet distribution can be found in the appendix C.1. For $p(\pi)$ giving:

$$p(\boldsymbol{\pi}) = Dir(\boldsymbol{\pi}|\boldsymbol{\omega_0^{(\pi)}}) = C(\boldsymbol{\omega_0^{(\pi)}}) \prod_{k=1}^{K} \boldsymbol{\pi}_k^{\boldsymbol{\omega}_0^{(\pi)}-1}$$

and for $p(\mathbf{A})$ we get:

$$p(\mathbf{A}) = \prod_{k=1}^{K} Dir(a_{k1}, ..., a_{kK} | \boldsymbol{\omega_0^{(\mathbf{A})}}) = \prod_{k=1}^{K} C(\boldsymbol{\omega_0^{(\mathbf{A})}}) \prod_{j=1}^{K} a_{kj}^{\boldsymbol{\omega_0^{(\mathbf{A})}}-1}$$

where the parameters are $\omega_0^{(\pi)}$ and $\omega_0^{(\mathbf{A})}$. The normalization constant is defined by $C(\omega_0)$. A small value of ω_0 affects the posterior rather by the data than by the prior and this applies vice versa for large value of ω_0 .

Moreover, we assume a Gaussian-Wishart prior on the mean vector and precision matrix for each Gaussian component, given by:

$$p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = p(\boldsymbol{\mu} | \boldsymbol{\Sigma}) p(\boldsymbol{\Sigma})$$
$$= \prod_{k=1}^{K} \mathcal{N} \left(\boldsymbol{\mu}_{k} | \mathbf{m}_{0}, (\beta_{0} \boldsymbol{\Sigma}_{k})^{-1} \right) \mathcal{W} \left(\boldsymbol{\Sigma}_{k} | \mathbf{W}_{0}, \nu_{0} \right)$$

Details of the Wishart Distribution are defined in C.4.

Finally, the graphical model is given by:

$$\begin{aligned} \boldsymbol{\pi} &\sim Dir(\boldsymbol{\omega}_0^{(\boldsymbol{\pi})}) \\ \mathbf{A} &\sim Dir(\boldsymbol{\omega}_0^{(\mathbf{A})}) \\ \boldsymbol{\Sigma} &\sim \mathcal{W}(\mathbf{W}_0, \nu_0) \\ \boldsymbol{\mu} &\sim \mathcal{N}(\mathbf{m}_0, (\beta_0, \boldsymbol{\Sigma})^{-1}) \\ \mathbf{x} &\sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}) \end{aligned}$$

Regarding the prior parameters of the graphical model, the posterior parameters with

dimensions are:

$$\boldsymbol{\omega}_{k}^{(\boldsymbol{\pi})}:1$$
$$\boldsymbol{\omega}_{k}^{(\mathbf{A})}:K$$
$$\mathbf{W}_{k}:D\times D$$
$$\boldsymbol{\nu}_{k}:1$$
$$\mathbf{m}_{k}:D\times 1$$
$$\boldsymbol{\beta}_{k}:1$$

where k denotes the latent state of the VBAR-HMM, D the observation data dimension.

4.1.2. Target distribution

Using the above defined priors, the log target distribution of the VBHMM becomes:

$$\log p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}, \boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) p(\mathbf{z}_{1:T} | \boldsymbol{\pi}, \mathbf{A}) + \log p(\mathbf{A}) + \log p(\boldsymbol{\pi}) + \log p(\boldsymbol{\mu} | \boldsymbol{\Sigma}) + \log p(\boldsymbol{\Sigma}) \propto \sum_{t=1}^{T} \sum_{k=1}^{K} [z_t = k] \log \mathcal{N} \left(\mathbf{x}_t | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k^{-1} \right) + \sum_{k=1}^{K} [z_1 = k] \log \pi_k + \sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} [z_{t-1} = k] [z_t = j] \log a_{kj} + \log Dir(\boldsymbol{\pi} | \boldsymbol{\omega}_0^{(\boldsymbol{\pi})}) + \sum_{k=1}^{K} \log Dir(\mathbf{A} | \boldsymbol{\omega}_{0k}^{(\mathbf{A})}) + \log \sum_{k=1}^{K} \mathcal{N} \left(\boldsymbol{\mu}_k | \mathbf{m}_0, (\beta_0 \boldsymbol{\Sigma}_k)^{-1} \right) + \log \sum_{k=1}^{K} \mathcal{W} \left(\boldsymbol{\Sigma}_k | \mathbf{W}_0, \nu_0 \right)$$

With this assumption of the true model, we can define a variational distribution in the next section.

4.1.3. Mean-field approximation

We consider a variational posterior that factorizes between the latent variables and the model parameters, giving:

$$p(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{z}_{1:T} | \mathbf{x}_{1:T}) \approx q(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) q(\mathbf{z}_{1:T})$$

Using this factorization, we apply the general solution from (3.8) to derive a solution for $q(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ and $q(\mathbf{z}_{1:T})$ in the next section.

4.1.4. Variational E-step

_

To compute the expectations of parameters that are required for the E-step updates, we need to apply the general result from Equation (3.8) to solve $q(\mathbf{z}_{1:T})$ giving:

$$\begin{split} \log q(\mathbf{z}_{1:T}) &\propto \mathbb{E}_{q(\boldsymbol{\pi},\boldsymbol{\mu},\boldsymbol{\Sigma},\mathbf{A})}[\log p(\mathbf{x}_{1:T},\mathbf{z}_{1:T},\boldsymbol{\pi},\mathbf{A},\boldsymbol{\mu},\boldsymbol{\Sigma})] \\ &= \mathbb{E}_{q(\boldsymbol{\pi})} \left[\sum_{k=1}^{K} \mathbb{I}(z_{1}=k) \log \boldsymbol{\pi}_{k} \right] + \mathbb{E}_{q(\mathbf{A})} \left[\sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} \mathbb{I}(z_{t-1}=k,z_{t}=j) \log a_{kj} \right] \\ &+ \mathbb{E}_{q(\boldsymbol{\mu},\boldsymbol{\Sigma})} \left[\sum_{t=1}^{T} \sum_{k=1}^{K} \mathbb{I}(z_{t}=k) \log \mathcal{N} \left(\mathbf{x}_{t} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}^{-1} \right) \right] \\ &= \sum_{k=1}^{K} \mathbb{I}(z_{1}=k) \mathbb{E}_{q(\boldsymbol{\pi})} \left[\log \boldsymbol{\pi}_{k} \right] + \sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} \mathbb{I}(z_{t-1}=k,z_{t}=j) \mathbb{E}_{q(\mathbf{A})} \left[\log a_{kj} \right] \\ &+ \sum_{t=1}^{T} \sum_{k=1}^{K} \mathbb{I}(z_{t}=k) \\ &\left\{ \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\Sigma})} \left[\log |\boldsymbol{\Sigma}_{k}| \right] - \frac{D}{2} \log (2\pi) - \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\mu},\boldsymbol{\Sigma})} \left[(\mathbf{x}_{t} - \boldsymbol{\mu}_{k})^{T} \boldsymbol{\Sigma}_{k} \left(\mathbf{x}_{t} - \boldsymbol{\mu}_{k} \right) \right] \right\} \end{split}$$

It is simple to solve the expectations $\mathbb{E}_{q(\pi)}$, $\mathbb{E}_{q(\mathbf{A})}$ and $\mathbb{E}_{q(\Sigma)}$ because we can look up the expectations for the corresponding distributions that are written in C. This expectation formulations yield:

$$\log \widetilde{\boldsymbol{\pi}}_{k} \equiv \mathbb{E}_{q(\boldsymbol{\pi})} \left[\log \boldsymbol{\pi}_{k}\right] = \int Dir(\boldsymbol{\pi}_{k} | \boldsymbol{\omega}^{(\boldsymbol{\pi})}) \log \boldsymbol{\pi}_{k} \, \mathrm{d}\boldsymbol{\pi}_{k}$$
$$= \psi(\boldsymbol{\omega}_{k}^{(\boldsymbol{\pi})}) - \psi(\sum_{k=1}^{K} \boldsymbol{\omega}_{k}^{(\boldsymbol{\pi})}), \quad \sum_{k=1}^{K} \widetilde{\boldsymbol{\pi}}_{k} \leq 1$$
$$\log \widetilde{\boldsymbol{a}}_{kj} \equiv \mathbb{E}_{q(\mathbf{A})} \left[\log a_{kj}\right] = \int Dir(a_{kj} | \boldsymbol{\omega}^{(\mathbf{A})}) \log a_{kj} \, \mathrm{d}a_{kj}$$
$$= \psi(\boldsymbol{\omega}_{kj}^{(\mathbf{A})}) - \psi(\sum_{j=1}^{K} \boldsymbol{\omega}_{kj}^{(\mathbf{A})}), \quad \sum_{j=1}^{K} \widetilde{\boldsymbol{a}}_{kj} \leq 1$$
$$\log \widetilde{\boldsymbol{\Sigma}}_{k} \equiv \mathbb{E}_{q(\boldsymbol{\Sigma})} \left[\log |\boldsymbol{\Sigma}_{k}|\right] = \int \mathcal{W} \left(\boldsymbol{\Sigma}_{k} | \mathbf{W}_{\mathbf{k}}, \boldsymbol{\nu}_{k}\right) \log |\boldsymbol{\Sigma}_{k}| \, \mathrm{d}\boldsymbol{\Sigma}_{k}$$
$$= \sum_{i=1}^{D} \psi \left(\frac{\boldsymbol{\nu}_{k} + 1 - i}{2}\right) + D \log 2 + \log |\mathbf{W}_{k}|$$

We derive a solution $\mathbb{E}_{q(\pmb{\mu},\pmb{\Sigma})}$ as follows:

$$\mathbb{E}_{q(\boldsymbol{\mu},\boldsymbol{\Sigma})} \left[(\mathbf{x}_{t} - \boldsymbol{\mu}_{k})^{T} \boldsymbol{\Sigma}_{k} (\mathbf{x}_{t} - \boldsymbol{\mu}_{k}) \right] \\ = \int \mathbb{E}_{q(\boldsymbol{\mu})} \left[(\boldsymbol{\mu}_{k} - \mathbf{x}_{t})^{T} \boldsymbol{\Sigma}_{k} (\boldsymbol{\mu}_{k} - \mathbf{x}_{t}) \right] q(\boldsymbol{\Sigma}_{k}) \, \mathrm{d}\boldsymbol{\Sigma}_{k}$$

Using the equation (380) from [36] to solve the inner expectation with respect to μ yields:

$$\mathbb{E}_{q(\boldsymbol{\mu})}\left[(\boldsymbol{\mu}_{k}-\mathbf{x}_{t})^{T}\boldsymbol{\Sigma}_{k}(\boldsymbol{\mu}_{k})-\mathbf{x}_{t}\right] = (\mathbf{m}_{k}-\mathbf{x}_{t})^{T}\boldsymbol{\Sigma}_{k}(\mathbf{m}_{k}-\mathbf{x}_{t}) + \mathbf{Tr}\left(\boldsymbol{\Sigma}_{k}(\boldsymbol{\beta}_{k}^{-1}\boldsymbol{\Sigma}_{k})^{-1}\right) \\ = (\mathbf{m}_{k}-\mathbf{x}_{t})^{T}\boldsymbol{\Sigma}_{k}(\mathbf{m}_{k}-\mathbf{x}_{t}) + D\boldsymbol{\beta}_{k}^{-1}$$

Plugging this term back in the equation

$$\mathbb{E}_{q(\boldsymbol{\mu},\boldsymbol{\Sigma})} \left[(\mathbf{x}_{t} - \boldsymbol{\mu}_{k})^{T} \boldsymbol{\Sigma}_{k} (\mathbf{x}_{t} - \boldsymbol{\mu}_{k}) \right]$$

=
$$\int \left\{ (\mathbf{m}_{k} - \mathbf{x}_{t})^{T} \boldsymbol{\Sigma}_{k} (\mathbf{m}_{k} - \mathbf{x}_{t}) + D\beta_{k}^{-1} \right\} q(\boldsymbol{\Sigma}_{k}) \, \mathrm{d}\boldsymbol{\Sigma}_{k}$$

=
$$D\beta_{k}^{-1} + \nu_{k} (\mathbf{m}_{k} - \mathbf{x}_{t})^{T} \mathbf{W}_{k} (\mathbf{m}_{k} - \mathbf{x}_{t})$$

We used the expectation of the Wishart distribution C.4 for the solution.

4.1.5. Variational M-step

For the M-step updates we need to solve $\log q(\pi, \mathbf{A}, \mu, \Sigma)$. Again make use of the the general solution from equation (3.8) giving:

$$\log q(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\log p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}, \boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma})\right]$$

$$= \log Dir(\boldsymbol{\pi} | \boldsymbol{\omega}_{0}^{(\boldsymbol{\pi})}) + \sum_{k=1}^{K} \log Dir(\mathbf{A} | \boldsymbol{\omega}_{0k}^{(\mathbf{A})})$$

$$+ \sum_{k=1}^{K} \log \mathcal{N} \left(\boldsymbol{\mu}_{k} | \boldsymbol{m}_{0}, (\beta_{0} \boldsymbol{\Sigma}_{k})^{-1}\right) + \sum_{k=1}^{K} \log \mathcal{W} \left(\boldsymbol{\Sigma}_{k} | \mathbf{W}_{0}, \boldsymbol{\nu}_{0}\right)$$

$$+ \sum_{t=1}^{T} \sum_{k=1}^{K} \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\mathbb{I}(z_{t} = k)\right] \log \mathcal{N} \left(\mathbf{x}_{t} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}^{-1}\right)$$

$$+ \sum_{k=1}^{K} \left\{\mathbb{E}_{q(z_{1})} \left[\mathbb{I}(z_{1} = k)\right] \log \pi_{k}\right\}$$

$$+ \sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\mathbb{I}(z_{t-1} = k, z_{t} = j)\right] \log a_{kj}$$

$$(4.1)$$

Since we have a HMM structure, we can solve the expectations from equation (B.1) using message passing. We use solutions from the equations (??) and (??). We rewrite equation (B.1):

$$\log q(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \log Dir(\boldsymbol{\pi} | \boldsymbol{\omega}_{\mathbf{0}}^{(\boldsymbol{\pi})}) + \sum_{k=1}^{K} \log Dir(\mathbf{A} | \boldsymbol{\omega}_{0k}^{(\mathbf{A})}) + \sum_{k=1}^{K} \log \mathcal{N} \left(\boldsymbol{\mu}_{k} | \boldsymbol{m}_{0}, (\beta_{0} \boldsymbol{\Sigma}_{k})^{-1} \right) + \sum_{k=1}^{K} \log \mathcal{W} \left(\boldsymbol{\Sigma}_{k} | \mathbf{W}_{\mathbf{0}}, \boldsymbol{\nu}_{0} \right) + \sum_{t=1}^{T} \sum_{k=1}^{K} \gamma_{t}(k) \log \mathcal{N} \left(\mathbf{x}_{t} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}^{-1} \right) + \sum_{k=1}^{K} \gamma_{1}(k) \log \pi_{k} + \sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} \xi_{t-1,t}(k, j) \log a_{kj}$$
(4.2)

Onwards we directly plug in the densities from the message passing step when expectations are of the above form.

We observe the right-hand side decomposes into a sum of terms only involving π together, A together and μ and Σ together. This observation implies that the variational posterior $q(\pi, \mathbf{A}, \mu, \Sigma)$ factorizes to give $q(\pi)q(\mathbf{A})q(\mu, \Sigma)$

$$q(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = q(\boldsymbol{\pi})q(\mathbf{A})\prod_{k=1}^{K}q(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$
(4.3)

Using this factorization and combine it with the results of the applied mean-field factorization, we obtain the solutions for the VBM step.

Regarding the right-hand side of equation (4.2) that depend on π , we get

$$\log q(\boldsymbol{\pi}) \propto (\omega_0^{(\boldsymbol{\pi})} - 1) \sum_{k=1}^K \log \boldsymbol{\pi}_k + \sum_{k=1}^K \gamma_1(k) \log \boldsymbol{\pi}_k$$
$$= \sum_{k=1}^K \log \boldsymbol{\pi}_k \left\{ \left(\omega_0^{(\boldsymbol{\pi})} - 1 \right) + \gamma_1(k) \right\}$$

Taking the exponential of both sides, we obtain

$$q(\boldsymbol{\pi}) = \text{const.} \cdot \prod_{k=1}^{K} \pi_{k}^{\omega_{0}^{(\boldsymbol{\pi})} - 1 + \gamma_{1}(k)}$$
$$= \text{const.} \cdot \prod_{k=1}^{K} \pi_{k}^{\omega_{k}^{(\boldsymbol{\pi})} - 1}$$
$$= Dir(\boldsymbol{\pi} | \boldsymbol{\omega}^{(\boldsymbol{\pi})})$$

where const. is the normalization term of the Dirichlet distribution and $\omega^{(\pi)}$ has components $\omega_k^{(\pi)}$ given by:

$$\omega_k^{(\boldsymbol{\pi})} = \omega_0^{(\boldsymbol{\pi})} + \gamma_1(k)$$

Applying the same principle on $\log q(\mathbf{A})$ we get

$$\log q(\mathbf{A}) \propto \sum_{k=1}^{K} \sum_{j=1}^{K} (\omega_0^{(\mathbf{A})} - 1) \log a_{kj} + \sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} \xi_{t-1,t}(k,j) \log a_{kj}$$
$$= \sum_{k=1}^{K} \sum_{j=1}^{K} \log a_{kj} \left\{ (\omega_0^{(\mathbf{A})} - 1) + \sum_{t=2}^{T} \xi_{t-1,t}(k,j) \right\}$$

Taking the exponential on both sides gives

$$q(\mathbf{A}) = \prod_{k=1}^{K} \text{const.} \cdot \prod_{j=1}^{K} a_{kj}^{\omega_{0}^{(\mathbf{A})} - 1 + \sum_{t=2}^{T} \xi_{t-1,t}(k,j)}$$
$$= \prod_{k=1}^{K} \text{const.} \cdot \prod_{j=1}^{K} a_{kj}^{\omega_{kj}^{(\mathbf{A})} - 1}$$
$$= \prod_{k=1}^{K} Dir(\mathbf{A}|\boldsymbol{\omega}_{k}^{(\mathbf{A})})$$

where $\pmb{\omega}_{\pmb{k}}^{(\mathbf{A})}$ has components $\omega_{kj}^{(\mathbf{A})}$ given by:

$$\omega_{kj}^{(\mathbf{A})} = \omega_0^{(\mathbf{A})} + \sum_{t=2}^T \xi_{t-1,t}(k,j)$$

To obtain a solution for $q(\mu_k, \Sigma_k)$, we use the result from (4.2) and take only the terms that depend on μ_k and Σ_k .

$$\begin{split} \log q(\boldsymbol{\mu}_{k},\boldsymbol{\Sigma}_{k}) \propto \log \mathcal{N}\left(\boldsymbol{\mu}_{k}|\boldsymbol{m}_{0},(\beta_{0}\boldsymbol{\Sigma}_{k})^{-1}\right) + \log \mathcal{W}\left(\boldsymbol{\Sigma}_{k}|\mathbf{W}_{0},\nu_{0}\right) \\ &+ \sum_{t=1}^{T} \gamma_{t}(k) \log \mathcal{N}\left(\mathbf{x}_{t}|\boldsymbol{\mu}_{k},\boldsymbol{\Sigma}_{k}^{-1}\right) \\ &= -\frac{D}{2} \log(2\pi) + \frac{1}{2} \log |\boldsymbol{\Sigma}_{k}| - \frac{\beta_{0}}{2}(\boldsymbol{\mu}_{k} - \boldsymbol{m}_{0})^{T}\boldsymbol{\Sigma}_{k}(\boldsymbol{\mu}_{k} - \boldsymbol{m}_{0}) \\ &- \frac{1}{2} \mathbf{Tr}\left\{\boldsymbol{\Sigma}_{k}\mathbf{W}_{0}^{-1}\right\} + \frac{(\nu_{0} - D - 1)}{2} \log |\boldsymbol{\Sigma}_{k}| \\ &- \frac{1}{2} \sum_{t=1}^{T} \gamma_{t}(k)(\mathbf{x}_{t} - \boldsymbol{\mu}_{k})^{T}\boldsymbol{\Sigma}_{k}(\mathbf{x}_{t} - \boldsymbol{\mu}_{k}) \\ &+ \frac{1}{2} \left(\sum_{t=1}^{T} \gamma_{t}(k)\right) \log |\boldsymbol{\Sigma}_{k}| \end{split}$$

We make use of the product rule to obtain $\log q(\mu_k, \Sigma_k) = \log q(\mu_k | \Sigma_k) + \log q(\Sigma_k)$. First we write $\log q(\mu_k | \Sigma_k)$ by only making use of the terms that depend on μ_k . We separate the dependencies in quadratic and linear terms.

$$\log q(\boldsymbol{\mu}_{k}|\boldsymbol{\Sigma}_{k}) = \underbrace{-\frac{1}{2} \left(\boldsymbol{\mu}_{k}^{T} \left(\beta_{0}\boldsymbol{\Sigma}_{k}\right)\boldsymbol{\mu}_{k}\right) - \frac{1}{2} \sum_{t=1}^{T} \gamma_{t}(k)\boldsymbol{\mu}_{k}^{T}\boldsymbol{\Sigma}\boldsymbol{\mu}_{k}}_{\text{quadratic term}} \\ + \underbrace{\boldsymbol{\mu}_{k}^{T} \left(\beta_{0}\boldsymbol{\Sigma}_{k}\mathbf{m}_{0}\right) + \sum_{t=1}^{T} \gamma_{t}(k)\boldsymbol{\mu}_{k}^{T} \left(\boldsymbol{\Sigma}_{k}\mathbf{x}_{t}\right)}_{\text{linear term}} \\ = \underbrace{-\frac{1}{2} \left\{ \boldsymbol{\mu}_{k}^{T} \left(\beta_{0} + \sum_{t=1}^{T} \gamma_{t}(k)\right)\boldsymbol{\Sigma}_{k}\boldsymbol{\mu}_{k} \right\}}_{\text{quadratic term}} \\ + \underbrace{\boldsymbol{\mu}_{k}^{T}\boldsymbol{\Sigma}_{k} \left(\beta_{0}\mathbf{m}_{0} + \sum_{t=1}^{T} \gamma_{t}(k)\mathbf{x}_{t}\right)}_{\text{linear term}}$$

The defined form of $q(\boldsymbol{\mu}_k|\boldsymbol{\Sigma}_k)$ is Gaussian, giving

$$q(\boldsymbol{\mu}_k | \boldsymbol{\Sigma}_k) = \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{m}_k, \beta_k \boldsymbol{\Sigma}_k)$$

We separate this Gaussian in quadratic and linear terms and use this result to obtain equations to compute β_k and \mathbf{m}_k . Starting to equate the quadratic terms yields:

$$\beta_k = \beta_0 + \sum_{t=1}^T \gamma_t(k)$$

We repeat this procedure to obtain a solution for m_k by equating the linear terms.

$$\mathbf{m}_{k} = \frac{1}{\beta_{k}} \left(\beta_{0} \mathbf{m}_{0} + \sum_{t=1}^{T} \gamma_{t}(k) \mathbf{x}_{t} \right)$$

To obtain $q(\Sigma_k)$ we subtract $q(\mu_k | \Sigma_k)$ from $q(\mu_k, \Sigma_k)$ focusing only on the terms that depend on Σ_k

$$\begin{split} \log q(\boldsymbol{\Sigma}_{k}) &= \log \mathcal{W}\left(\boldsymbol{\Sigma}_{k} | \mathbf{W}_{k}, \nu_{k}\right) \\ &= \log q(\boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) - \log q(\boldsymbol{\mu}_{k} | \boldsymbol{\Sigma}_{k}) \\ &\propto \frac{1}{2} \log |\boldsymbol{\Sigma}_{k}| - \frac{\beta_{0}}{2} (\boldsymbol{\mu}_{k} - \boldsymbol{m}_{0})^{T} \boldsymbol{\Sigma}_{k} (\boldsymbol{\mu}_{k} - \boldsymbol{m}_{0}) \\ &- \frac{1}{2} \mathbf{Tr} \left\{ \boldsymbol{\Sigma}_{k} \mathbf{W}_{0}^{-1} \right\} + \frac{(\nu_{0} - D - 1)}{2} \log |\boldsymbol{\Sigma}_{k}| \\ &- \frac{1}{2} \sum_{t=1}^{T} \gamma_{t} (k) (\mathbf{x}_{t} - \boldsymbol{\mu}_{k})^{T} \boldsymbol{\Sigma}_{k} (\mathbf{x}_{t} - \boldsymbol{\mu}_{k}) \\ &+ \frac{1}{2} \left(\sum_{t=1}^{T} \gamma_{t} (k) \right) \log |\boldsymbol{\Sigma}_{k}| \\ &+ \frac{\beta_{k}}{2} (\boldsymbol{\mu}_{k} - \mathbf{m}_{k})^{T} \boldsymbol{\Sigma}_{k} (\boldsymbol{\mu}_{k} - \mathbf{m}_{k}) - \frac{1}{2} \log |\boldsymbol{\Sigma}_{k}| \\ &= \frac{(\nu_{k} - D - 1)}{2} \log |\boldsymbol{\Sigma}_{k}| - \frac{1}{2} \mathbf{Tr} \left(\boldsymbol{\Sigma}_{k} \mathbf{W}_{k}^{-1} \right) \end{split}$$

where \mathbf{W}_k^{-1} is defined by:

$$\mathbf{W}_{k}^{-1} = \mathbf{W}_{0}^{-1} N_{k} \mathbf{S}_{k} + \frac{\beta_{0} N_{k}}{\beta_{0} + N_{k}} (\overline{\mathbf{x}}_{k} - \mathbf{m}_{0}) (\overline{\mathbf{x}}_{k} - \mathbf{m}_{0})^{T}$$

where we defined:

$$N_k = \sum_{t=1}^T \gamma_t(k)$$
$$\overline{\mathbf{x}}_k = \frac{1}{N_k} \sum_{t=1}^T \gamma_t(k) \mathbf{x}_t$$
$$\mathbf{S}_k = \frac{1}{N_k} \sum_{t=1}^T \gamma_t(k) (\mathbf{x}_t - \overline{\mathbf{x}}_k) (\mathbf{x}_t - \overline{\mathbf{x}}_k)^T$$

We arrange the terms to obtain the form of a Wishart distribution. This rearrangement gives the solution for ν_k :

$$\nu_k = \nu_0 + \sum_{t=1}^T \gamma_t(k)$$

The EM updates for the VBHMM are summarized in the Appendix A.1.

4.1.6. The lower bound

The ELBO for the VBHMM is:

$$\begin{split} \text{ELBO}_{\text{VBHMM}} &= \int \int \int \int q\left(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}\right) \left[\log \frac{p(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma})}{q(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma})} \right. \\ &+ \sum_{\mathbf{z}_{1:T}} q(\mathbf{z}_{1:T}) \log \frac{p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T} | \boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma})}{q(\mathbf{z}_{1:T})} \right] \, \mathrm{d}\boldsymbol{\pi} \, \mathrm{d}\mathbf{A} \, \mathrm{d}\boldsymbol{\mu} \, \mathrm{d}\boldsymbol{\Sigma} \\ &= \mathbb{E} \left[\log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) \right] + \mathbb{E} \left[\log p(\mathbf{z}_{1:T} | \boldsymbol{\pi}, \mathbf{A}) \right] + \mathbb{E} \left[\log p(\boldsymbol{\pi}) \right] \\ &+ \mathbb{E} \left[\log p(\mathbf{A}) \right] + \mathbb{E} \left[\log p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \right] - \mathbb{E} \left[\log q(\boldsymbol{\pi}) \right] - \mathbb{E} \left[\log q(\mathbf{A}) \right] \\ &- \mathbb{E} \left[\log q(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \right] - \mathbb{E} \left[\log q(\mathbf{z}_{1:T}) \right] \end{split}$$

All expectations are w.r.t. $q(\mathbf{z}_{1:T}, \boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$. The solutions for the expectations are derived in the variational E- and M-steps.
4.2. Variational Bayes ARHMM

The VBAR-HMM is an extension of the VBHMM. We add a linear dependency between each observation and its successive observation as illustrated in Figure 3.1b. In addition, the next observation becomes dependent on a control signal **u**. In this section we define and derive a solution for the model extension. The complete derivation is in the Appendix B.2.

4.2.1. Model priors

Since the VBAR-HMM is an extension of the VBHMM, we can carry over the model priors from Section 4.1.1. In addition, we assume a linear weight matrix U for every component k s.t. $\mathbf{x}_t = \mathbf{U}_k \hat{\mathbf{x}}_{t-1}$ for $z_t = k$ where:

$$\hat{\mathbf{x}}_{t-1} = [\mathbf{x}_{t-1} \ \mathbf{u}_t \ 1]^T$$

We stack the previous observation \mathbf{x}_{t-1} , the successive control signal \mathbf{u}_t with the constant value 1 to obtain a linear transition matrix U that contains all dependency information for \mathbf{x}_{t+1}

We make use of this reformulation to implicitly learn the linear regression constant in **U**. Moreover

A Matrix-Normal prior is a typical choice for linear weight matrix as [37] shows. We choose a Matrix Normal prior with unknown row covariance matrix V giving:

$$p(\mathbf{U}, \mathbf{V}) = p(\mathbf{U} | \mathbf{V}) p(\mathbf{V})$$
$$= \prod_{k=1}^{K} \mathcal{MN} \left(\mathbf{U}_{k} | \mathbf{M}_{0}, \mathbf{V}_{k}^{-1}, \mathbf{K}_{0} \right) \mathcal{W} \left(\mathbf{V}_{k} | \mathbf{P}_{0}, \eta_{0} \right)$$

where we assume a Wishart prior on V. The prior parameter choices of the P_0 and η_0 are equivalent to the choices for the Wishart prior parameters of Equation (4.1.1). For symmetry, we naturally set $M_0 = 0$. Small values for the column covariance matrix K_0 affect the density rather by the data than by the prior.

The updated graphical model for the VBAR-HMM is now given by:

$$\begin{aligned} \boldsymbol{\pi}_{k} &\sim Dir(\boldsymbol{\omega}_{0}^{(\boldsymbol{\pi})}) \\ \mathbf{A}_{k} &\sim Dir(\boldsymbol{\omega}_{0}^{(\mathbf{A})}) \\ \boldsymbol{\Sigma}_{k} &\sim \mathcal{W}(\mathbf{W}_{0}, \nu_{0}) \\ \boldsymbol{\mu}_{k} &\sim \mathcal{N}(\mathbf{m}_{0}, (\beta_{0}, \boldsymbol{\Sigma}_{k})^{-1}) \\ \mathbf{V}_{k} &\sim \mathcal{W}(\mathbf{P}_{0}, \eta_{0}) \\ \mathbf{U}_{k} &\sim \mathcal{M}\mathcal{N}(\mathbf{M}_{0}, \mathbf{V}_{k}^{-1}, \mathbf{K}_{0}) \\ \mathbf{x}_{1} &\sim \mathcal{N}(\boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}^{-1}) \\ \mathbf{x}_{2:\mathbf{T}} &\sim \mathcal{N}(\mathbf{U}_{k} \hat{\mathbf{x}}_{1:T-1}, \mathbf{V}_{k}^{-1}) \end{aligned}$$

The resulting posterior parameters with dimensions for this graphical model are:

$$\begin{aligned}
 \omega_k^{(\boldsymbol{\pi})} &: 1 \\
 \omega_k^{(\mathbf{A})} &: K \\
 \mathbf{W}_k &: D \times D \\
 \nu_k &: 1 \\
 \mathbf{m}_k &: D \times 1 \\
 \beta_k &: 1 \\
 \eta_k &: 1 \\
 \mathbf{M}_k &: D \times (D + C + 1) \\
 \mathbf{K}_k &: (D + C + 1) \times (D + C + 1)
 \end{aligned}$$

where k denotes the latent state of the VBAR-HMM, D the observation data dimension and C the dimension of the control signal **u**.

4.2.2. Target distribution

With the above the defined priors, the log target distribution of the VBAR-HMM becomes:

$$\log p(\mathbf{x}_{1:T}, \hat{\mathbf{x}}_{2:T}, \mathbf{z}_{1:T}, \boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V})$$

$$= \log p(\mathbf{x}_{1}|z_{1}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \log p(\mathbf{x}_{2:T}|\hat{\mathbf{x}}_{2:T}, \mathbf{z}_{2:T}, \mathbf{U}, \mathbf{V})$$

$$+ \log p(\mathbf{z}_{1:T}|\boldsymbol{\pi}, \mathbf{A}) + \log p(\boldsymbol{\pi}) + \log p(\mathbf{A})$$

$$+ \log p(\boldsymbol{\mu}|\boldsymbol{\Sigma}) + \log p(\boldsymbol{\Sigma}) + \log p(\mathbf{U}|\mathbf{V}) + \log p(\mathbf{V})$$

$$\propto \sum_{k=1}^{K} [z_{1} = k] \log \mathcal{N} \left(\mathbf{x}_{1}|\boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}^{-1}\right)$$

$$+ \sum_{t=2}^{T} \sum_{k=1}^{K} [z_{t} = k] \log \mathcal{N} \left(\mathbf{x}_{t}|\mathbf{U}_{k}\hat{\mathbf{x}}_{t-1}, \mathbf{V}_{k}^{-1}\right)$$

$$+ \sum_{k=1}^{K} \{\mathbb{I}(z_{1} = k) \log \pi_{k}\} + \sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} \mathbb{I}(z_{t-1} = k, z_{t} = j) \log a_{kj}$$

$$+ \log Dir(\boldsymbol{\pi}|\boldsymbol{\omega}_{0}^{(\boldsymbol{\pi})}) + \sum_{k=1}^{K} \log Dir(\mathbf{A}|\boldsymbol{\omega}_{0k}^{(\mathbf{A})})$$

$$+ \log \sum_{k=1}^{K} \mathcal{N} \left(\boldsymbol{\mu}_{k}|\mathbf{m}_{0}, (\beta_{0}\boldsymbol{\Sigma}_{k})^{-1}\right) + \log \sum_{k=1}^{K} \mathcal{W} \left(\boldsymbol{\Sigma}_{k}|\mathbf{W}_{0}, \nu_{0}\right)$$

$$+ \sum_{k=1}^{K} \log \mathcal{M} \mathcal{N} \left(\mathbf{U}_{k}|\mathbf{M}_{0}, \mathbf{V}_{k}^{-1}, \mathbf{K}_{0}\right) + \sum_{k=1}^{K} \log \mathcal{W} \left(\mathbf{V}_{k}|\mathbf{P}_{0}, \eta_{0}\right)$$

$$(4.4)$$

4.2.3. Mean-field approximation

We consider the same factorization assumption as in Section 4.1.3 giving:

$$p(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V}\mathbf{z}_{1:T} | \mathbf{x}_{1:T}) \approx q(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V})q(\mathbf{z}_{1:T})$$

Again, we apply the general result from Equation (3.8) to obtain solutions for the variational distributions.

4.2.4. Variational E-step

To compute the expectations of parameters that are required for the M-step updates, we need to apply the general result from Equation (3.8) to solve $q(\mathbf{z}_{1:T})$ giving:

$$\begin{split} \log q(\mathbf{z}_{1:T}) \propto \mathbb{E}_{q(\pi,\mathbf{A},\boldsymbol{\mu},\boldsymbol{\Sigma},\mathbf{U},\mathbf{V})} \left[\log p(\mathbf{x}_{1:T},\mathbf{x}_{2:T}^{*},\mathbf{z}_{1:T},\boldsymbol{\pi},\mathbf{A},\boldsymbol{\mu},\boldsymbol{\Sigma},\mathbf{U},\mathbf{V})\right] \\ &= \mathbb{E}_{q(\pi)} \left[\sum_{k=1}^{K} \mathbb{I}(z_{1}=k)\log\left(\pi_{k}\right)\right] \\ &+ \mathbb{E}_{q(\mathbf{A})} \left[\sum_{t=2}^{T}\sum_{k=1}^{K}\sum_{j=1}^{K} \mathbb{I}(z_{t-1}=k,z_{t}=j)\log a_{kj}\right] \\ &+ \mathbb{E}_{q(\boldsymbol{\mu},\boldsymbol{\Sigma})} \left[\sum_{k=1}^{K}[z_{1}=k]\log\mathcal{N}(\mathbf{x}_{1}|\boldsymbol{\mu}_{k},\boldsymbol{\Sigma}_{k})\right] \\ &+ \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\sum_{t=2}^{T}\sum_{k=1}^{K} \mathbb{I}(z_{t}=k)\log\mathcal{N}\left(\mathbf{x}_{t}|\mathbf{U}_{k}\hat{\mathbf{x}}_{t-1},\mathbf{V}_{k}^{-1}\right)\right] \\ &= \sum_{k=1}^{K} \mathbb{I}(z_{1}=k)\mathbb{E}_{q(\pi)}\left[\log \pi_{k}\right] + \sum_{t=2}^{T}\sum_{k=1}^{K}\sum_{j=1}^{K} \mathbb{I}(z_{t-1}=k,z_{t}=j)\mathbb{E}_{q(\mathbf{A})}\left[\log a_{kj}\right] \\ &+ \sum_{k=1}^{K} \mathbb{I}(z_{1}=k) \\ &\left\{\frac{1}{2}\mathbb{E}_{q(\boldsymbol{\Sigma})}\left[\log|\boldsymbol{\Sigma}_{k}|\right] - \frac{D}{2}\log\left(2\pi\right) - \frac{1}{2}\mathbb{E}_{q(\boldsymbol{\mu},\boldsymbol{\Sigma})}\left[\left(\mathbf{x}_{1}-\boldsymbol{\mu}_{k}\right)^{T}\boldsymbol{\Sigma}_{k}\left(\mathbf{x}_{1}-\boldsymbol{\mu}_{k}\right)\right]\right\} \\ &+ \sum_{t=2}^{T}\sum_{k=1}^{K} \mathbb{I}(z_{t}=k) \\ &\left\{\frac{1}{2}\mathbb{E}_{q(\mathbf{V})}\left[\log|\mathbf{V}_{k}|\right] - \frac{D}{2}\log\left(2\pi\right) \\ &- \frac{1}{2}\mathbb{E}_{q(\mathbf{U},\mathbf{V})}\left[\left(\mathbf{x}_{t}-\mathbf{U}_{k}\hat{\mathbf{x}_{t-1}}\right)^{T}\mathbf{V}_{k}\left(\mathbf{x}_{t}-\mathbf{U}_{k}\hat{\mathbf{x}_{t-1}}\right)\right]\right\} \end{split}$$

We can simply take the expectation results from Section 4.1.4 that are:

$$\log \widetilde{\boldsymbol{\pi}}_k \equiv \mathbb{E}_{q(\boldsymbol{\pi})} \left[\log \boldsymbol{\pi}_k\right] = \psi(\omega_k^{(\boldsymbol{\pi})}) - \psi(\sum_{k=1}^K \omega_k^{(\boldsymbol{\pi})}), \quad \sum_{k=1}^K \widetilde{\boldsymbol{\pi}}_k \le 1$$
(4.5)

$$\log \widetilde{a}_{kj} \equiv \mathbb{E}_{q(\mathbf{A})} \left[\log a_{kj} \right] = \psi(\omega_{kj}^{(\mathbf{A})}) - \psi(\sum_{j=1}^{K} \omega_{kj}^{(\mathbf{A})}), \quad \sum_{j=1}^{K} \widetilde{a}_{kj} \le 1$$
(4.6)

$$\log \widetilde{\boldsymbol{\Sigma}}_k \equiv \mathbb{E}_{q(\boldsymbol{\Sigma})} \left[\log |\boldsymbol{\Sigma}_k| \right] = \sum_{i=1}^{D} \psi \left(\frac{\nu_k + 1 - i}{2} \right) + D \log 2 + \log |\mathbf{W}_k| \quad (4.7)$$

$$\mathbb{E}_{q(\boldsymbol{\mu},\boldsymbol{\Sigma})}\left[(\mathbf{x}_{1}-\boldsymbol{\mu}_{k})^{T}\boldsymbol{\Sigma}_{k}(\mathbf{x}_{1}-\boldsymbol{\mu}_{k})\right] = D\beta_{k}^{-1} + \nu_{k}(\mathbf{m}_{k}-\mathbf{x}_{1})^{T}\mathbf{W}_{k}(\mathbf{m}_{k}-\mathbf{x}_{1})$$
(4.8)

Since we put a Wishart prior on V_k the solution is equivalent to Equation (4.7) giving:

$$\log \widetilde{\mathbf{V}}_k \equiv \mathbb{E}_{q(\mathbf{V})} \left[\log |\mathbf{V}_k| \right] = \sum_{i=1}^{D} \psi \left(\frac{\eta_k + 1 - i}{2} \right) + D \log 2 + \log |\mathbf{P}_k|$$

We are left to solve the expectation for:

$$\mathbb{E}_{q(\mathbf{U},\mathbf{V})}\left[\left(\mathbf{x}_{t}-\mathbf{U}_{k}\hat{\mathbf{x}}_{t-1}\right)^{T}\mathbf{V}_{k}\left(\mathbf{x}_{t}-\mathbf{U}_{k}\hat{\mathbf{x}}_{t-1}\right)\right]$$

We rewrite:

$$\mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[(\mathbf{x}_{t} - \mathbf{U}_{k} \hat{\mathbf{x}}_{t-1})^{T} \mathbf{V}_{k} (\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1}) \right]$$

$$= \int \mathbb{E}_{q(\mathbf{U})} \left[(\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t})^{T} \mathbf{V}_{k} (\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t}) \right] q(\mathbf{V}_{k}) \, \mathrm{d}\mathbf{V}_{k}$$
(4.9)

The inner expectation corresponds to a Matrix-Normal distribution. Thus, we make use of transformations derived in [38] and get two expectations:

$$\begin{split} \mathbb{E}\left[\mathbf{U}_{k}\right] &= \mathbf{M}_{k}\\ \mathbb{E}\left[\mathbf{U}_{k}\mathbf{Q}\mathbf{U}_{k}^{T}\right] &= \mathbf{M}_{k}\mathbf{U}\mathbf{M}_{k}^{T} + \mathbf{Tr}\left\{\mathbf{K}_{k}^{-1}\mathbf{Q}\right\}\mathbf{V}_{k}^{-1} \end{split}$$

Applying this transformations on the inner expectation from Equation (B.2) gives:

$$\begin{split} \mathbb{E}_{q(\mathbf{U})} \left[(\mathbf{U}_k \hat{\mathbf{x}}_{t-1} - \mathbf{x}_t)^T \mathbf{V}_k (\mathbf{U}_k \hat{\mathbf{x}}_{t-1} - \mathbf{x}_t] \\ &= (\mathbf{M}_k \hat{\mathbf{x}}_{t-1} - \mathbf{x}_t)^T \mathbf{V}_k (\mathbf{M}_k \hat{\mathbf{x}}_{t-1} - \mathbf{x}_t) + \mathbf{Tr} \left\{ \mathbf{K}_k^{-1} \hat{\mathbf{x}}_{t-1} \hat{\mathbf{x}}_{t-1}^T \right\} \end{split}$$

31

Substituting back yields the final solution:

$$\mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\left(\mathbf{x}_t - \mathbf{U}_k \hat{\mathbf{x}}_{t-1} \right)^T \mathbf{V}_k \left(\mathbf{U}_k \hat{\mathbf{x}}_{t-1} \right) \right] \\ = \eta_k (\mathbf{M}_k \hat{\mathbf{x}}_{t-1} - \mathbf{x}_t)^T \mathbf{P}_k (\mathbf{M}_k \hat{\mathbf{x}}_{t-1} - \mathbf{x}_t) + \mathbf{Tr} \left\{ \mathbf{K}_k^{-1} \hat{\mathbf{x}}_{t-1} \hat{\mathbf{x}}_{t-1}^T \right\}$$

4.2.5. Variational M-step

For the M-step updates, we need to solve $\log q(\pi, \mathbf{A}, \mu, \Sigma, \mathbf{U}, \mathbf{V})$. Applying the general result from Equation (3.8) gives:

$$\log q(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V}) \propto \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\log p(\mathbf{x}_{1:T}, \mathbf{u}_{2:T}, \mathbf{z}_{1:T}, \boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V})\right]$$

$$= \log Dir(\boldsymbol{\pi} | \boldsymbol{\omega}_{0}^{(\boldsymbol{\pi})}) + \sum_{k=1}^{K} \log Dir(\mathbf{A} | \boldsymbol{\omega}_{0k}^{(\mathbf{A})})$$

$$+ \sum_{k=1}^{K} \gamma_{1}(k) \log \mathcal{N} \left(\mathbf{x}_{1} | \boldsymbol{\mu}_{k}, (\beta_{0} \boldsymbol{\Sigma}_{k})^{-1}\right)$$

$$+ \sum_{t=2}^{T} \sum_{k=1}^{K} \gamma_{t}(k) \log \mathcal{N} \left(\mathbf{x}_{t} | \mathbf{U}_{k} \hat{\mathbf{x}}_{t-1}, \mathbf{V}_{k}^{-1}\right)$$

$$+ \sum_{k=1}^{K} \{\gamma_{1}(k) \log \pi_{k}\}$$

$$+ \sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} \xi_{t-1,t}(k, j) \log a_{ij}$$

$$+ \sum_{k=1}^{K} \log \mathcal{M} \mathcal{N} \left(\mathbf{U}_{k} | \mathbf{M}_{0}, \mathbf{V}_{k}^{-1}, \mathbf{K}_{0}\right) + \sum_{k=1}^{K} \log \mathcal{W} \left(\mathbf{V}_{k} | \mathbf{P}_{0}, \eta_{0}\right)$$

$$(4.10)$$

We observe the right-hands side decomposes into a sum of terms only involving π together, A together, μ and Σ together and U and V together. This observation implies that the variational posterior $q(\pi, \mathbf{A}, \mu, \Sigma, \mathbf{U}, \mathbf{V})$ factorizes to:

$$q(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V}) = q(\boldsymbol{\pi})q(\mathbf{A})\prod_{k=1}^{K}q(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\prod_{k=1}^{K}q(\mathbf{U}_k, \mathbf{V}_k)$$
(4.11)

As we already derived solutions for $q(\pi)$, $q(\mathbf{A})$ and $q(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ in Section 4.1.4, we are left to solve $q(\mathbf{U}_k, \mathbf{V}_k)$.

Regarding the target distribution from Equation (4.3.2), we assume the model for $q(\mathbf{U}, \mathbf{V})$ has the form:

$$q(\mathbf{U}, \mathbf{V}) = \prod_{k=1}^{K} q(\mathbf{U}_{k}, \mathbf{V}_{k}) = \prod_{k=1}^{K} q(\mathbf{U}_{k} | \mathbf{V}_{k}) q(\mathbf{V}_{k})$$
$$= \prod_{k=1}^{K} \mathcal{MN} \left(\mathbf{U}_{k} | \mathbf{M}_{k}, \mathbf{V}_{k}^{-1}, \mathbf{K}_{k} \right) \mathcal{W} \left(\mathbf{V}_{k} | \mathbf{P}_{k}, \eta_{k} \right)$$
$$= \prod_{k=1}^{K} \frac{|\mathbf{K}_{k}|^{d/2} |\mathbf{V}_{k}|^{m/2}}{(2\pi)^{m/2}} \exp \left\{ -\frac{1}{2} \mathbf{Tr} \left\{ (\mathbf{U}_{k} - \mathbf{M}_{k})^{T} \mathbf{V}_{k} (\mathbf{U}_{k} - \mathbf{M}_{k}) \mathbf{K}_{k} \right\} \right\}$$
$$B(\mathbf{P}_{k}, \eta_{k}) |\mathbf{V}_{k}|^{(\eta_{k} - D - 1)/2} \exp \left\{ -\frac{1}{2} \mathbf{Tr} \left\{ \mathbf{P}_{k}^{-1} \mathbf{V}_{k} \right\} \right\}$$

where $B(\mathbf{P}_k, \eta_k)$ is defined in C.4.

Solving $q(\mathbf{U}, \mathbf{V})$ requires to make use of result from Equation (4.11) and take only the terms from Equation (4.10) that depend on \mathbf{U}_k and \mathbf{V}_k giving:

$$\begin{split} \log q(\mathbf{U}_k, \mathbf{V}_k) &\propto \log \mathcal{M} \mathcal{N} \left(\mathbf{U}_k | \mathbf{M}_0, \mathbf{V}_k^{-1}, \mathbf{K}_0 \right) + \log \mathcal{W} \left(\mathbf{V}_k | \mathbf{P}_0, \eta_0 \right) \\ &+ \sum_{t=2}^T \gamma_t(k) \log \mathcal{N} \left(\mathbf{x}_t | \mathbf{U}_k \hat{\mathbf{x}}_{t-1}, \mathbf{V}_k^{-1} \right) \\ &= -\frac{1}{2} \mathbf{Tr} \left\{ (\mathbf{U}_k - \mathbf{M}_0)^T \mathbf{V}_k (\mathbf{U}_k - \mathbf{M}_0) \mathbf{K}_0 \right\} - \frac{m}{2} \log(|\mathbf{V}_k|) \\ &- \frac{1}{2} \mathbf{Tr} \left\{ \mathbf{V}_k \mathbf{W}_0^{-1} \right\} + \frac{(\eta_0 - D - 1)}{2} \log |\mathbf{V}_k| \\ &- \frac{1}{2} \sum_{t=2}^T \gamma_t(k) (\mathbf{x}_t - \mathbf{U}_k \hat{\mathbf{x}}_{t-1})^T \mathbf{V}_k (\mathbf{x}_t - \mathbf{U}_k \hat{\mathbf{x}}_{t-1}) \\ &+ \frac{1}{2} \left(\sum_{t=2}^T \gamma_t(k) \right) \log |\mathbf{V}_k| \end{split}$$

We make use of the product rule to obtain $\log q(\mathbf{U}_k, \mathbf{V}_k) = \log q(\mathbf{U}_k | \mathbf{V}_k) + \log q(\mathbf{V}_k)$. First we write $\log q(\mathbf{U}_k | \mathbf{V}_k)$ by only making use of the terms that depend on \mathbf{U}_k . We gather all terms in (4.2.5) that contain \mathbf{U}_k giving:

$$\begin{split} \log q(\mathbf{U}_k | \mathbf{V}_k) \propto &-\frac{1}{2} \mathbf{Tr} \left\{ (\mathbf{U}_k - \mathbf{M}_0)^T \mathbf{V}_k (\mathbf{U}_k - \mathbf{M}_0) \mathbf{K}_0 \right\} \\ &- \frac{1}{2} \sum_{t=2}^T \gamma_t (k) (\mathbf{x}_t - \mathbf{U}_k \hat{\mathbf{x}}_t)^T \mathbf{V}_k (\mathbf{x}_t - \mathbf{U}_k \hat{\mathbf{x}}_t) \\ &= -\frac{1}{2} \mathbf{Tr} \left\{ \mathbf{V}_k (\mathbf{U}_k \mathbf{K}_0 \mathbf{U}_k^T - \mathbf{U}_k \mathbf{K}_0 \mathbf{M}_0^T - \mathbf{M}_0 \mathbf{K}_0 \mathbf{U}_k^T + \mathbf{M}_0 \mathbf{K}_0 \mathbf{M}_0^T) \right\} \\ &- \frac{1}{2} \sum_{t=2}^T \gamma_t (k) \mathbf{Tr} \left\{ \mathbf{V}_k (\mathbf{x}_t \mathbf{x}_t^T - \mathbf{x}_t \hat{\mathbf{x}}_t^T \mathbf{U}_k^T - \mathbf{U}_k \hat{\mathbf{x}}_t \mathbf{x}_t^T + \mathbf{U}_k \hat{\mathbf{x}}_t \hat{\mathbf{x}}_t^T \mathbf{U}_k^T) \right\} \end{split}$$

We separate the dependencies in quadratic and linear terms which yields:

$$\log q(\mathbf{U}_{k}|\mathbf{V}_{k}) = \underbrace{-\frac{1}{2}\mathbf{Tr}\left\{\mathbf{V}_{k}\mathbf{U}_{k}\left[\sum_{t=2}^{T}\gamma_{t}(k)\hat{\mathbf{x}}_{t}\hat{\mathbf{x}}_{t}^{T} + \mathbf{K}_{0}\right]\mathbf{U}_{k}^{T}\right\}}_{\text{quadratic term}} \\ \underbrace{-\frac{1}{2}\mathbf{Tr}\left\{\mathbf{V}_{k}\mathbf{U}_{k}\left[-\sum_{t=2}^{T}\gamma_{t}(k)\mathbf{x}_{t}\hat{\mathbf{x}}_{t-1}^{T} - \mathbf{M}_{0}\mathbf{K}_{0}\right]\right\}}_{\text{linear term}}$$

The defined form of $q(\mathbf{U}_k | \mathbf{V}_k)$ is Matrix Normal giving:

$$\mathcal{MN}\left(\mathbf{U}_k|\mathbf{M}_k,\mathbf{V}_k^{-1},\mathbf{K}_k
ight)$$

We obtain a solution for the posterior parameter updates using the defined form and bringing Using this form, we obtain a solution for the posterior parameter updates.

$$\mathbf{K}_{k} = \sum_{t=2}^{T} \gamma_{t}(k) \hat{\mathbf{x}}_{t-1} \hat{\mathbf{x}}_{t-1}^{T} + \mathbf{K}_{0}$$

We apply this procedure on the linear terms to obtain a solution for \mathbf{M}_k

$$\mathbf{M}_{k} = \left[\sum_{t=2}^{T} \gamma_{t}(k) \mathbf{x}_{t} \hat{\mathbf{x}}_{t-1}^{T} + \mathbf{M}_{0} \mathbf{K}_{0}\right] \mathbf{K}_{k}^{-1}$$

To obtain the parameters \mathbf{P}_k and η_k of the Wishart distribution $q(\mathbf{V}_k)$, we subtract the conditional distribution from the joint distribution giving:

$$\log q(\mathbf{V}_k) = \log \mathcal{W}(\mathbf{U}_k | \mathbf{P}_k, \eta_k)$$

= log q(U_k, V_k) - log q(U_k | V_k)
 $\propto \frac{(\eta_k - D - 1)}{2} \log |\mathbf{V}_k| - \frac{1}{2} \operatorname{Tr} \left\{ \mathbf{V}_k \mathbf{P}_k^{-1} \right\}$

where \mathbf{P}_k^{-1} is defined as:

$$\mathbf{P}_{k}^{-1} = \mathbf{P}_{0}^{-1} + \mathbf{M}_{0}\mathbf{K}_{0}\mathbf{M}_{0}^{T} + \sum_{t=2}^{T}\gamma_{t}(k)\mathbf{x}_{t}\mathbf{x}_{t}^{T} - \mathbf{M}_{k}\mathbf{K}_{k}\mathbf{M}_{k}^{T}$$

We arange the terms to obtain a solution for η_k and \mathbf{V}_k . For η_k we get:

$$\eta_k = \eta_0 + \sum_{t=2}^T \gamma_t(k)$$

The EM updates for the VBAR-HMM are summarized in the Appendix A.2.

4.2.6. The evidence lower bound

The ELBO for the VBAR-HMM is:

$$\begin{split} \text{ELBO}_{\text{VBAR-HMM}} &= \int \int \int \int \int \int q\left(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V}\right) \left[\log \frac{p(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V})}{q(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}), \mathbf{U}, \mathbf{V}} \\ &+ \sum_{\mathbf{z}_{1:T}} q(\mathbf{z}_{1:T}) \log \frac{p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T} | \boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V})}{q(\mathbf{z}_{1:T})} \right] \\ &\quad d\boldsymbol{\pi} \, d\mathbf{A} \, d\boldsymbol{\mu} \, d\boldsymbol{\Sigma} \, d\mathbf{U} \, d\mathbf{V} \\ &= \mathbb{E} \left[\log p(\mathbf{x}_1 | \boldsymbol{z}_1, \boldsymbol{\mu}, \boldsymbol{\Sigma})\right] + \mathbb{E} \left[\log p(\mathbf{x}_{2:T} | \hat{\mathbf{x}}_{1:T-1} \mathbf{z}_{2:T}, \mathbf{U}, \mathbf{V})\right] \\ &\quad + \mathbb{E} \left[\log p(\mathbf{z}_{1:T} | \boldsymbol{\pi}, \mathbf{A})\right] + \mathbb{E} \left[\log p(\boldsymbol{\pi})\right] + \mathbb{E} \left[\log p(\mathbf{A})\right] \\ &\quad + \mathbb{E} \left[\log p(\boldsymbol{\mu}, \boldsymbol{\Sigma})\right] + \mathbb{E} \left[\log p(\mathbf{U}, \mathbf{V})\right] - \mathbb{E} \left[\log q(\boldsymbol{\pi})\right] - \mathbb{E} \left[\log q(\mathbf{A})\right] \\ &\quad - \mathbb{E} \left[\log q(\boldsymbol{\mu}, \boldsymbol{\Sigma})\right] - \mathbb{E} \left[\log q(\mathbf{U}, \mathbf{V})\right] - \mathbb{E} \left[\log q(\mathbf{z}_{1:T})\right] \end{split}$$

All expectations are w.r.t. $q(\mathbf{z}_{1:T}, \pi, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V})$. The solutions for the expectations are derived in the variational E- and M-steps.

4.3. Variational Bayes rAR-HMM

The VBrAR-HMM similar to the VBAR-HMM. The difference is the added dependency link the observation \mathbf{x}_t to the next state z_{t+1} as Figure 3.1c illustrates. Since z_{t+1} now depends on \mathbf{x}_t and z_t , the latent state transition function has the property of being nonlinear.

4.3.1. Model priors

We introduce the function $f_{\phi}(.)$ for the recurrent state transition which defines a Bayesian neural network with parameters ϕ . We can take the model priors from Section 4.2.1 and need to change the Dirichlet transition prior with a Gaussian prior we put on the neural

network weights ϕ giving the graphical model:

$$\begin{split} \boldsymbol{\pi}_{k} &\sim Dir(\boldsymbol{\omega}_{0}^{(\boldsymbol{\pi})}) \\ \boldsymbol{\phi} &\sim \mathcal{N}(\boldsymbol{\zeta}_{0}, \mathbf{Y}_{0}) \\ \boldsymbol{\Sigma}_{k} &\sim \mathcal{W}(\mathbf{W}_{0}, \nu_{0}) \\ \boldsymbol{\mu}_{k} &\sim \mathcal{N}(\mathbf{m}_{0}, (\beta_{0}, \boldsymbol{\Sigma}_{k})^{-1}) \\ \mathbf{V}_{k} &\sim \mathcal{W}(\mathbf{P}_{0}, \eta_{0}) \\ \mathbf{U}_{k} &\sim \mathcal{M}\mathcal{N}(\mathbf{M}_{0}, \mathbf{V}_{k}^{-1}, \mathbf{K}_{0}) \\ \mathbf{x}_{1} &\sim \mathcal{N}(\boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}^{-1}) \\ \mathbf{x}_{2:\mathbf{T}} &\sim \mathcal{N}(\mathbf{U}_{k} \hat{\mathbf{x}}_{1:T-1}, \mathbf{V}_{k}^{-1}) \end{split}$$

The posterior parameters are equal to the posterior parameters from Subsection 4.2.1. Except the Dirichlet posterior $\omega^{(A)}$ is exchanged with the Bayesian neural network (BNN) posteriors ζ and \mathbf{Y} . The dimensions of the BNN parameter posteriors corresponds to the number of hidden layers and parameters per layer.

4.3.2. Target distribution

With the above the defined priors, the log target distribution of the VBAR-HMM becomes:

$$\begin{split} \log p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}, \pi, \phi, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V}) \\ &= \log p(\mathbf{x}_1 | z_1, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \log p(\mathbf{x}_{2:T} | \hat{\mathbf{x}}_{2:T}, \mathbf{z}_{2:T}, \mathbf{U}, \mathbf{V}) \\ &+ \log p(\mathbf{z}_{1:T} | \pi, \hat{\mathbf{x}}_{2:T}, \phi) + \log p(\pi) + \log p(\phi) \\ &+ \log p(\boldsymbol{\mu} | \boldsymbol{\Sigma}) + \log p(\boldsymbol{\Sigma}) + \log p(\mathbf{U} | \mathbf{V}) + \log p(\mathbf{V}) \\ &\propto \sum_{k=1}^{K} [z_1 = k] \log \mathcal{N} \left(\mathbf{x}_1 | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k^{-1} \right) \\ &+ \sum_{t=2}^{T} \sum_{k=1}^{K} [z_t = k] \log \mathcal{N} \left(\mathbf{x}_t | \mathbf{U}_k \hat{\mathbf{x}}_{t-1}, \mathbf{V}_k^{-1} \right) \\ &+ \sum_{k=1}^{K} \{ \mathbb{I}(z_1 = k) \log \pi_k \} + \sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} \mathbb{I}(z_{t-1} = k, z_t = j) \log \psi_{kj} \\ &+ \log Dir(\pi | \boldsymbol{\omega}_0^{(\pi)}) + \sum_{k=1}^{K} \log \mathcal{N} \left(\boldsymbol{\phi} | \boldsymbol{\zeta}_0, \mathbf{Y}_0 \right) \\ &+ \log \sum_{k=1}^{K} \mathcal{N} \left(\boldsymbol{\mu}_k | \mathbf{m}_0, (\beta_0 \boldsymbol{\Sigma}_k)^{-1} \right) + \log \sum_{k=1}^{K} \mathcal{W} \left(\mathbf{\Sigma}_k | \mathbf{W}_0, \nu_0 \right) \\ &+ \sum_{k=1}^{K} \log \mathcal{M} \mathcal{N} \left(\mathbf{U}_k | \mathbf{M}_0, \mathbf{V}_k^{-1}, \mathbf{K}_0 \right) + \sum_{k=1}^{K} \log \mathcal{W} \left(\mathbf{V}_k | \mathbf{P}_0, \eta_0 \right) \end{split}$$

where $\psi_{k,j}$ is defined as follows:

$$\psi_{k,j} = \frac{\exp(f_{\phi}(\mathbf{x}_{t-1}, \mathbf{u}_t)_{k,j})}{\sum_{k'=1}^{K} \exp(f_{\phi}(\mathbf{x}_{t-1}, \mathbf{u}_t)_{k,k'})}$$
(4.12)

We sum over the columns for each row of the output matrix to obtain normalized transition probabilities.

4.3.3. Mean-field approximation

The model of the rARHMM only changes w.r.t. the transition function. The rARHMM model is equipped with a non-linear function. Thus, the factorization model remains

identical in all parts but the transition function. Applying the mean-field approximation to the rARHMM we get the following result:

4.3.4. Variational E-step

$$\begin{split} \log q(\mathbf{z}_{1:T}) &\propto \mathbb{E}_{q(\boldsymbol{\pi},\boldsymbol{\phi},\boldsymbol{\mu},\boldsymbol{\Sigma},\mathbf{U},\mathbf{V})} \left[\log p(\mathbf{x}_{1:T},\mathbf{z}_{1:T},\boldsymbol{\pi},\boldsymbol{\phi},\boldsymbol{\mu},\boldsymbol{\Sigma},\mathbf{U},\mathbf{V})\right] \\ &= \sum_{k=1}^{K} \mathbb{I}(z_{1}=k) \mathbb{E}_{q(\boldsymbol{\pi})} \left[\log \boldsymbol{\pi}_{k}\right] + \sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} \mathbb{I}(z_{t-1}=k,z_{t}=j) \mathbb{E}_{q(\boldsymbol{\phi})} \left[\log \psi_{kj}\right] \\ &+ \sum_{k=1}^{K} \mathbb{I}(z_{1}=k) \\ &\left\{ \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\Sigma})} \left[\log |\boldsymbol{\Sigma}_{k}|\right] - \frac{D}{2} \log (2\pi) - \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\mu},\boldsymbol{\Sigma})} \left[(\mathbf{x}_{1}-\boldsymbol{\mu}_{k})^{T} \boldsymbol{\Sigma}_{k} (\mathbf{x}_{1}-\boldsymbol{\mu}_{k}) \right] \right\} \\ &+ \sum_{t=2}^{T} \sum_{k=1}^{K} \mathbb{I}(z_{t}=k) \\ &\left\{ \frac{1}{2} \mathbb{E}_{q(\mathbf{V})} \left[\log |\mathbf{V}_{k}|\right] - \frac{D}{2} \log (2\pi) \\ &- \frac{1}{2} \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[(\mathbf{x}_{t}-\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1})^{T} \mathbf{V}_{k} (\mathbf{x}_{t}-\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1}) \right] \right\} \end{split}$$

Since the rAR-HMM changes w.r.t. the transition function, only the estimation of the transition function differs from the AR-HMM model. We can not solve the estimation term in closed form due to its non linearity property as equation (4.12) shows. We make use of a non-linear function approximation framework to approximate the optimal solution. To remain in the Bayesian fashion, we choose a BNN [39, 40, 41] for this task.

We solve the estimation for $\mathbb{E}_{q(\phi)} [\log \psi_{ij}]$ by plugging in the observations \mathbf{x}_{t-1} and \mathbf{u}_t into the bayesian neural network $f_{\phi}(.)$ and normalize the result to get proper probabilities.

4.3.5. Variational M-step

To optimize the parameters ϕ of the BNN, we make use of stochastic gradient function approximations. For this task, we choose to maximize the part of the lower bound that

depends on the network output. We write the objective as follows:

$$\max_{\phi} \log p(\mathbf{z}_{2:T} | \phi) = \max_{\phi} \sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} \xi_{t-1,t}(k,j) \log \psi_{kj}$$
(4.13)

This objective can be plugged in any gradient solver. Thus, a BNN is not mandatory for this task. The EM updates for the VBrAR-HMM are summarized in the Appendix A.3.

4.3.6. The evidence lower bound

The ELBO for the VBAR-HMM is:

$$\begin{split} \text{ELBO}_{\text{VBAR-HMM}} &= \int \int \int \int \int \int q\left(\boldsymbol{\pi}, \boldsymbol{phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V}\right) \left[\log \frac{p(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V})}{q(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}), \mathbf{U}, \mathbf{V}} \\ &+ \sum_{\mathbf{z}_{1:T}} q(\mathbf{z}_{1:T}) \log \frac{p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T} | \boldsymbol{\pi}, \boldsymbol{\phi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V})}{q(\mathbf{z}_{1:T})} \right] \\ &\quad d\boldsymbol{\pi} \, d\boldsymbol{\phi} \, d\boldsymbol{\mu} \, d\boldsymbol{\Sigma} \, d\mathbf{U} \, d\mathbf{V} \\ &= \mathbb{E} \left[\log p(\mathbf{x}_1 | \boldsymbol{z}_1, \boldsymbol{\mu}, \boldsymbol{\Sigma})\right] + \mathbb{E} \left[\log p(\mathbf{x}_{2:T} | \hat{\mathbf{x}}_{2:T} \mathbf{z}_{2:T}, \mathbf{U}, \mathbf{V})\right] \\ &\quad + \mathbb{E} \left[\log p(\mathbf{z}_{1:T} | \boldsymbol{\pi}, \boldsymbol{\phi})\right] + \mathbb{E} \left[\log p(\boldsymbol{\pi})\right] + \mathbb{E} \left[\log p(\boldsymbol{\phi})\right] \\ &\quad + \mathbb{E} \left[\log p(\boldsymbol{\mu}, \boldsymbol{\Sigma})\right] + \mathbb{E} \left[\log p(\mathbf{U}, \mathbf{V})\right] - \mathbb{E} \left[\log q(\boldsymbol{\pi})\right] - \mathbb{E} \left[\log q(\boldsymbol{\phi})\right] \\ &\quad - \mathbb{E} \left[\log q(\boldsymbol{\mu}, \boldsymbol{\Sigma})\right] - \mathbb{E} \left[\log q(\mathbf{U}, \mathbf{V})\right] - \mathbb{E} \left[\log q(\mathbf{z}_{1:T})\right] \end{split}$$

The expectations have the same solutions as VBAR-HMM in Subsection 4.2.6. Only the expectation of the BNN parameters changed. The solution is derived in the EM-steps.

5. Switching actor

The idea of SWAC is to train multiple policies with few parameters rather than training one global policy with a large number of parameters. This procedure is inspired by the option-critic architecture from [42]. Ideally, we can choose linear policies that capture local linearities. Having such policies can be seen as a single global policy with extended internal structure. To achieve this policy transformation, we make use of a trained HMM that mimics the system dynamics. This adaption makes SWAC to a MBRL because of the dynamics model we plug in this algorithm.

In this section, we show the implementation of SWAC and extend one of the recent actor-critic algorithms with SWAC.

5.1. Implementation

The design of SWAC makes SWAC is generally applicable to all actor-critic frameworks where we collect environment sample trajectories with a policy and maximize the reward signal.

Algorithm 1 shows the implementation of SWAC. The training procedure of SWAC in 1 is similar to REINFORCE or any actor-critic framework. The main difference to typical model free policy search frameworks is that we make use of a dynamics model and we initialize multiple policies that are disjointly attached to the latent states of the dynamics model.

To generate samples with the policies, we first compute a belief vector (normalized probabilities for each state). Using the belief vector **b** we sample the policy that samples the action for the current observation.

Input: step size: α , number of switching states: Z, max horizon: H 1 τ_{slds} \leftarrow collect randomly sampled trajectories from environment; $slds \leftarrow$ train dynamics model with Z switching states from τ_{slds} ; 2 $\pi_{\theta}^{(z)} \leftarrow \text{initialize } Z \text{ policies where } z \in \{1, ..., Z\};$ 3 while not converged do 4 for max horizon not reached do 5 $\mathbf{b} \leftarrow$ compute state belief vector with *slds*; 6 $z \leftarrow \text{randomly choose state from } \mathbf{b};$ 7 $\tau \leftarrow \text{collect samples with policy } \pi_{\theta}^{(z)}(\mathbf{a}_t | \mathbf{s}_t);$ 8 9 end foreach latent state z do 10 $\hat{g}_{\theta}^{(z)} \leftarrow \sum_{t=1}^{T} \mathbb{I}(z_t = z) \nabla_{\theta} \log \pi_{\theta}^{(z)}(\mathbf{a}_t | \mathbf{s}_t) \left[\sum_{t'=t}^{T} \mathbb{I}(z_{t'} = z) r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right];$ 11 $\boldsymbol{\theta}^{(z)} \leftarrow \boldsymbol{\theta}^{(z)} + \alpha \hat{g}_{\boldsymbol{\theta}}^{(z)};$ 12 end 13 14 end

Algorithm 1: SWAC: Switching actor

For the policy parameter update step, we only take the actions into account that are sampled from the policies. The indicator functions ensure the mapping to the matching policies.

In addition, we can exchange the reward with a discounted reward, add a baseline, exchange the reward with an advantage function [43] etc. For the evaluation task, we test SWAC on A2C [44] and on SAC [45]. Both extension use a single global value function.

5.2. Extending A2C with SWAC

The implementation of A2C with SWAC is similar to the proposed algorithm 1. We refer to this extension as SWAC-A2C. A2C requires the computation of the advantage function. Thus, we add a global value function network and initialize multiple policies. The gradient computation becomes:

$$\hat{g}_{\boldsymbol{\theta}}^{(z)} \leftarrow \sum_{t=1}^{T} \mathbb{I}(z_t = z) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}^{(z)}(\mathbf{a}_t | \mathbf{s}_t) \left[Q(\mathbf{s}_t, \mathbf{a}_t) - V(s_t) \right]$$

where Q(.) is the state-action value function obtained by computing the discounted reward.

If the HMM detects a latent state less frequently, the policy linked to that state samples only few samples which can lead to gradient distortion. To prevent from this behavior, we attach each policy with a own learn rate. During each policy update step, we measure the policy change by computing the KL divergence of the policy parameters before and after the policy update. The KL divergence needs to remain within a fixed value which is set as a hyper parameter. A too small KL divergence increases the learn rate and a too large KL divergence increases the learn rate for the next iterations.

5.3. Extending SAC with SWAC

The extension of SAC with SWAC is similar to the SWAC extension of A2C. For this implementation we choose a fixed learn rate since the policy gradient has naturally less variance than A2C due to its architecture.

In addition, we need to adapt the computation of the policy log probabilities from the soft value function stated in equation (6) in [45]. Our solution for this problem is to average the log probability over all policies and thus relinquish to use exclusive policies for this computation.

Using the exclusive policy log probabilities affected the accuracy of the soft value function and lead to bad performance.

6. Evaluation

For the experiments, we show the switching state behavior of the VBrAR-HMMs and compare their prediction performance to state of the art neural network architectures.

Moreover, we train policies to control previously trained SLDS and evaluate their performance on the true dynamics.

Finally, we equip actor-critic architectures with SWAC, measure the learning performance and compare it to the non SWAC equipped variants.

We evaluate the experiments on the pendulum and the cart-pole [46] environments. These are classical benchmark systems from the control literature. The swing up task is set as the objective for both environments. For the observation model we represent the angles in the trigonometric space.

6.1. Bayesian HMMs

The experiments of this section exemplarily show the dynamic switching behavior of the VBrAR-HMMs. Additionally, we compare the prediction performance of the VBAR-HMM and VBrAR-HMM to state-of-the-art models.

6.1.1. Viterbi state detection

In this experiment, we evaluate the Viterbi algorithm on randomly sampled trajectories using the likelihoods of previously trained VBrAR-HMMs. The VBrAR-HMMs are initialized with seven latent states to capture the switching dynamics for each environment. From the pendulum environment, the VBrAR-HMM is given 25 training rollouts with a horizon



Figure 6.1.: Randomly sampled trajectories on the pendulum (left) and cart-pole (right) environment. Each color represents a detected latent state from the VBrAR-HMM using the Viterbi algorithm. The state transitions match the repeating dynamic patterns.

of 200 each. The VBrAR-HMM imitating the cart-pole dynamics, is given 45 training rollouts with a horizon of 250 each.

Figure 3.1.4 shows the Viterbi state detection of the VBrAR-HMMs. The horizon for the pendulum trajectory is 200 and 500 for the cart-pole trajectory. Each color is disjointly mapped to a latent state of the model. The VBrAR-HMM only uses 4 of its 7 states to capture the switching dynamics of this trajectory. Similarly, the VBrAR-HMM requires 6 of the total 7 states to express the cart-pole dynamics of the given trajectory.

We observe that the VBrAR-HMM learns a sensible assignment of the states, especially regarding the repeating patterns of the pendulum trajectory.

6.1.2. Model learning performance

In this experiment, we compare the out-of-sample predictive performance of the VBAR-HMM, VBrAR-HMM. Evaluating two types of the VBrAR-HMMs. One uses a BNN and the other uses a classic Feed-forward neural network (FNN) for the latent state transition. We compare the performance of our models with FNN, Recurrent neural network (RNN) and Long short-term memory neural networks (LSTM)[47]. A prediction performance overview of the non Bayesian HMMs can be found in [46].



Figure 6.2.: Out-of-sample predictive performance of the VBAR-HMMs and VBrAR-HMMs to SOTA models. Evaluation on the pendulum (left) and cart-pole (right) environment.

Figure 6.2 shows the evaluation results for the out-of-sample prediction task on the pendulum and cart-pole environment. The s, RNNs and LSTMs use a single hidden layer with 64 neurons each. The VBrAR-HMMs use a single layer with 24 neurons for the transition network. The predictive performance is evaluated on the horizons $h \in \{1, 5, 10, 15, 20, 25\}$. Each model is trained with a 10 trajectories with 250 steps each. The performance is averaged over 20 differently trained models. For the pendulum task, the other models perform better than ours. All three of our presented models perform equally on this task. The logistic link of the VBrAR-HMM does not improve the models predictive accuracy. However, for the cart-pole task, the logistic link improves the prediction accuracy but only for the variant without using the BNN. Our models outperform the other models on this task except the performs better.

6.2. Learning control on rAR-HMM

In this experiment, we perform model-free policy search on trained rAR-HMMs. We use the non Bayesian rAR-HMMs from [46] because they are well tested on the pendulum and cart-pole environment. We train the models with these environments. The goal is to showcase a potential performance gain of SWAC by using a reliable SLDS model.



Figure 6.3.: Policy evaluation on the rAR-HMM dynamics (left) and on the true dynamics (right) of the pendulum environment. The policy training observations are sampled from the rAR-HMM dynamics.

We evaluate the performance of the trained policy on hybrid environments with rAR-HMM dynamics and on the true dynamics models.

The approximation gaps of the trained models can lead to variance increasement of the policy gradients. Thus, we use LAX [20] as the policy search framework for this task.

6.2.1. Policy evaluation on the pendulum task

For the pendulum environment, the rAR-HMM is trained with 25 trajectories each with a horizon of 200. The policy is trained with 800 epochs, each with 5000 samples.

Figure 6.3 illustrates the pendulum policy evaluation on the rAR-HMM dynamics and on the true dynamics. The latent state transitions are detected with the trained rAR-HMM that is used as the hybrid dynamics model on the left. The model dynamics are nearly identical for both systems. Thus, the policy manages to succeed to swing up and stabilize the pendulum on both systems. The policy achieves an expected reward with **-2310.73** \pm **1983.48** on the rAR-HMM dynamics and on the true dynamics the expected reward



Figure 6.4.: Policy evaluation on the rAR-HMM dynamics (left) and on true dynamics (right)of the pendulum environment. The policy training observations are sampled from the rAR-HMM dynamics.

using the same policy is **-2300.90** \pm **2513.69**. We averaged over 20 rollouts, each with a horizon of 1000.

6.2.2. Policy evaluation on the cart-pole task

For the cart-pole environment, the rAR-HMM is trained with 45 trajectories each with a horizon of 250. The policy is trained with 500 epochs, each with 5000 samples.

Figure 6.4 illustrates the cart-pole policy evaluation on the rAR-HMM dynamics and on the true dynamics. The latent state transitions are detected with the trained rAR-HMM that is used as the hybrid dynamics model on the left. Other than in the pendulum task, a difference in the dynamics between the two models is observable. Owing to a more difficult prediction of the cart-pole dynamics. However, the policy still manages to achieve the swing up and stabilization task on both models. The policy achieves an expected



Figure 6.5.: Policy evaluation performance comparison of A2C, SWAC-A2C and LAX on the pendulum (left) and on the cart-pole (right) environment.

reward with **-2993.20** \pm **1918.37** on the rAR-HMM dynamics and on the true dynamics the expected reward using the same policy is **-3389.59** \pm **1785.15**. We averaged over 20 rollouts, each with a horizon of 5000.

6.3. SWAC-A2C

In this experiment, we compare the learning performance of SWAC-A2C, A2C and LAX. The performance is evaluated on 20 different seeds for each algorithm on each environment. The underlying rAR-HMMs of each environment for SWAC-A2C are trained with the same settings as in Section 6.2. For a fair comparison, we have added the learn rate regularization of the KL divergence between the old policy and the updated policy parameters to A2C and (LAX). to A2C.

Figure 6.5 shows the learning performance on both environments with a confidence interval at 95%. Since rAR-HMMs are trained with 7 states, SWAC-A2C is equipped with 7 policies each with a single hidden layer and 14 neurons. A2C and LAX use policies with 2 hidden layers and 64 neurons for each layer. The value function neural network is the same for all three algorithms. 2 hidden layers each with 64 neurons. For each training episode 5000 training samples are generated. The target KL divergence for the policy parameter



Figure 6.6.: Policy evaluation performance comparison of SAC and SWAC-SAC on the pendulum (left) and on the cart-pole (right) environment.

updates is set to 0.02 for all algorithms. SWAC-A2C clearly outperforms the baseline algorithms on the pendulum swing up task. Whereas, on the cart-pole environment SWAC-A2C shows a worse performance at the beginning of the training but performs better in the long run compared to the baseline. LAX shows the same performance as A2C on both environments.

6.4. SWAC-SAC

In this experiment, we compare the learning performance of SWAC-SAC and SAC. The performance is evaluated on 20 different seeds for each algorithm on each environment. The underlying rAR-HMMs of each environment for SWAC-SAC are trained with the same settings as in Section 6.2.

Figure 6.6 shows the learning performance on both environments with a confidence interval at 95%. As in experiment of subsection 5.2 the rAR-HMMs are trained with 7 states, SWAC-SAC is equipped with 7 policies each with a single hidden layer and 14 neurons. SAC uses a policy network with 2 hidden layers and 64 neurons for each layer. The value function neural network is the same for all three algorithms. 2 hidden layers each with 64 neurons. 5,000 training samples are generated for the pendulum and 10,000

training samples are generated for the cart-pole task for each training episode. SWAC-SAC and SAC perform identically on the pendulum task. For the cart-pole task SWAC-SAC outperforms SAC. SWAC-SAC succeeds in approximately 4 training episodes to achieve the task and remains stable for the further training episodes. SAC on the other hand approximately requires 10 training episodes to achieve the cart-pole task.

6.5. Discussion

The results show that the variational Bayes variants of the HMMs are capable of capturing the system dynamics of a given model. The weaker performance of the VBAR-HMM-BNN compared to the VBAR-HMM and VBrAR-HMM-FNN can be due to the difficulties of hyper-parameter tuning for the BNNs [48]. However, using the BNN we can access the uncertainty of the network which could be interesting for further work. Moreover, we have shown that SLDS can be successfully incorporated into policy search problems which is evidenced by our results. The actor-critic frameworks extended with SWAC show a similar or better performance but no worse performance compared to the non SWAC extended actor-critic frameworks. A weakness of the SWAC extension is the computation overhead that occurs during the gradient propagation for each policy. This computation can be parallelized but still requires additional computation resources. Interestingly, LAX shows the same performance as A2C on our selected control environments. In the published paper of LAX [20] the authors used more complex environments to evaluate LAX. Observing this behavior, we conclude that LAX plays its strength for gradient variance reduction on more complex environments.

7. Outlook

The Bayesian HMMs have the property make the access to the models uncertainty possible. For future work it would be interesting to incorporate this quantity into a policy search framework such that exploration can be enhanced. Moreover, further investigation of these models for the planning step in a MBRL setup can improve the learning performance as [5] have done with GPs. A requirement to achieve this task, is to make differentiating through the whole model possible.

For future work of SWAC, it would be interesting to find a way to get SWAC work by only using multiple linear policies. A problem of SWAC is that when a latent state is rarely detected by the switching model, the policy mapped to that latent state produces not enough samples for training. This behavior can be avoided by using as few latent states as possible.

8. Conclusion

In this thesis, we have tackled the problem of equipping a policy with additional internal structure. Initially, we have formulated a Bayesian HMM that makes the model's uncertainty accessible. This uncertainty can be used as an additional quantity for policy search problems. To succeed this thesis's goal, we have successfully integrated a non Bayesian rAR-HMMs to a policy framework that we have defined in this work. This policy framework profits from the switching behavior of the rAR-HMMs and shows superior performance compared to the variants that are not equipped with switching structure of the rAR-HMMs.

In Chapter 4, we have proposed a Bayesian formulation of HMMs, AR-HMMs and rAR-HMMs by making use of variational inference and the mean-field approximation. We have defined the priors of the model parameters and have derived all required steps to make an implementation of these models straightforward using the Baum-Welch algorithm.

In Chapter 5, we have introduced SWAC, an extension for policy search algorithms that brings additional structure to the policy by integrating a switching state space model, in our case rAR-HMMs, into the policy search procedure.

In Chapter 6, we have illustrated the switching behavior of the VBrAR-HMMs.

Moreover, we have compared the predictive performance of the VBAR-HMM and the VBrAR-HMM. On the simpler pendulum environment, the reference models performed better than our proposed models. Whereas, on the more complex cart-pole environment, our models outperformed the reference models except the LSTM.

To show the advantage of minimizing sample complexity by using a learned dynamics model, we have performed policy search on a trained rAR-HMM. The obtained policy is evaluated on the true and on the rAR-HMM dynamics. The policy performed minimally better on the rAR-HMM dynamics.

As a proof of concept for SWAC, we have measured the learning performance of the SWAC to the non SWAC equipped counterparts. For a fair comparison, we have used the same hyper-parameters for both variants. The results of SWAC have shown a better or an equal performance compared to the non SWAC variant.

Finally, we have reviewed a policy search approach that benefits from the Bayesian HMMs by using the model for planning and accessing the model's uncertainty for exploration. Regarding SWAC, we have pointed on the possibility to make it run by using linear policies only.

Bibliography

- S. Linderman, M. Johnson, A. Miller, R. Adams, D. Blei, and L. Paninski, "Bayesian Learning and Inference in Recurrent Switching Linear Dynamical Systems," in *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics* (A. Singh and J. Zhu, eds.), vol. 54 of *Proceedings of Machine Learning Research*, (Fort Lauderdale, FL, USA), pp. 914–922, PMLR, 20–22 Apr 2017.
- [2] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Machine Learning*, vol. 8, pp. 229–256, 1992.
- [3] A. G. Barto, R. S. Sutton, and C. W. Anderson, "Neuronlike adaptive elements that can solve difficult learning control problems," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. SMC-13, no. 5, pp. 834–846, 1983.
- [4] M. J. Beal, "Variational algorithms for approximate bayesian inference," tech. rep., 2003.
- [5] M. Deisenroth and C. Rasmussen, "Pilco: A model-based and data-efficient approach to policy search," in *Proceedings of the 28th International Conference on Machine Learning, ICML 2011*, pp. 465–472, Omnipress, 2011.
- [6] M. P. Deisenroth*, G. Neumann*, and J. Peters, "A survey on policy search for robotics," *Foundations and Trends in Robotics*, pp. 388–403, 2013.
- [7] D. P. Kingma and M. Welling, "An introduction to variational autoencoders," *CoRR*, vol. abs/1906.02691, 2019.
- [8] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," in 2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings (Y. Bengio and Y. LeCun, eds.), 2014.
- [9] C. Doersch, "Tutorial on variational autoencoders," 2016. cite arxiv:1606.05908.

- [10] P. Becker-Ehmck, M. Karl, J. Peters, and P. van der Smagt, "Learning to fly via deep model-based reinforcement learning," 2020.
- [11] G. Ackerson and K. Fu, "On state estimation in switching environments," *IEEE Transactions on Automatic Control*, vol. 15, no. 1, pp. 10–17, 1970.
- [12] C. B. Chang and M. Athans, "State estimation for discrete systems with switching parameters," *IEEE Transactions on Aerospace and Electronic Systems*, vol. AES-14, no. 3, pp. 418–425, 1978.
- [13] J. D. Hamilton, "Analysis of time series subject to changes in regime," 1990.
- [14] E. Fox, E. B. Sudderth, M. I. Jordan, and A. S. Willsky, "Nonparametric bayesian learning of switching linear dynamical systems," in *Advances in Neural Information Processing Systems 21* (D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, eds.), pp. 457–464, Curran Associates, Inc.
- [15] Z. Ghahramani and G. E. Hinton, "Switching state-space models," tech. rep., King's College Road, Toronto M5S 3H5, 1996.
- [16] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," 2013. cite arxiv:1312.6114.
- [17] E. Jang, S. Gu, and B. Poole, "Categorical reparametrization with gumbel-softmax," in *Proceedings International Conference on Learning Representations 2017*, OpenReviews.net, Apr. 2017.
- [18] C. J. Maddison, A. Mnih, and Y. W. Teh, "The concrete distribution: A continuous relaxation of discrete random variables," 2016. cite arxiv:1611.00712.
- [19] G. Tucker, A. Mnih, C. J. Maddison, J. Lawson, and J. Sohl-Dickstein, "Rebar: Lowvariance, unbiased gradient estimates for discrete latent variable models," in *Advances in Neural Information Processing Systems 30* (I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, eds.), pp. 2627–2636, Curran Associates, Inc., 2017.
- [20] W. Grathwohl, D. Choi, Y. Wu, G. Roeder, and D. Duvenaud, "Backpropagation through the void: Optimizing control variates for black-box gradient estimation," 2017. cite arxiv:1711.00123Comment: Published at ICLR 2018.
- [21] K. P. Murphy, *Machine learning : a probabilistic perspective*. Cambridge, Mass. [u.a.]: MIT Press, 2013.

- [22] L. E. Baum and T. Petrie, "Statistical inference for probabilistic functions of finite state markov chains," *Ann. Math. Statist.*, vol. 37, pp. 1554–1563, 12 1966.
- [23] V. Smídl and A. Quinn, *The Variational Bayes Method in Signal Processing (Signals and Communication Technology)*. Berlin, Heidelberg: Springer-Verlag, 2005.
- [24] T. S. Jaakkola and M. I. Jordan, "Bayesian parameter estimation via variational methods," *Statistics and Computing*, vol. 10, no. 1, pp. 25–37, 2000.
- [25] T. S. Jaakkola and M. I. Jordan, "Computing upper and lower bounds on likelihoods in intractable networks," pp. 340–348, Morgan Kaufmann, 1996.
- [26] M. J. Wainwright and M. I. Jordan, "Graphical models, exponential families, and variational inference," *Found. Trends Mach. Learn.*, vol. 1, p. 1–305, Jan. 2008.
- [27] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe, "Variational inference: A review for statisticians," 2018.
- [28] M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul, "An introduction to variational methods for graphical models," *Mach. Learn.*, vol. 37, p. 183–233, Nov. 1999.
- [29] S. Mohamed, M. Rosca, M. Figurnov, and A. Mnih, "Monte carlo gradient estimation in machine learning," 2019. cite arxiv:1906.10652Comment: 59 pages, under review.
- [30] C. Wang and D. M. Blei, "Variational inference in nonconjugate models," *J. Mach. Learn. Res.*, vol. 14, p. 1005–1031, Apr. 2013.
- [31] C. M. Bishop, Pattern Recognition and Machine Learning. Springer, 2006.
- [32] S. M. Ross, Simulation, Fourth Edition. USA: Academic Press, Inc., 2006.
- [33] R. Dearden, N. Friedman, and D. Andre, "Model-based bayesian exploration," CoRR, vol. abs/1301.6690, 2013.
- [34] P. Abbeel, M. Quigley, and A. Y. Ng, "Using inaccurate models in reinforcement learning," in *Proceedings of the 23rd International Conference on Machine Learning*, ICML '06, (New York, NY, USA), p. 1–8, Association for Computing Machinery, 2006.
- [35] J. C. Spall and J. A. Cristion, "Model-free control of nonlinear stochastic systems with discrete-time measurements," *IEEE Transactions on Automatic Control*, vol. 43, no. 9, pp. 1198–1210, 1998.

- [36] K. B. Petersen and M. S. Pedersen, "The matrix cookbook," Oct. 2008. Version 20081110.
- [37] T. P. Minka, "Bayesian linear regression," tech. rep., 3594 Security Ticket Control, 1999.
- [38] D. V. rosen, "Moments for matrix normal variables," *Statistics*, vol. 19, no. 4, pp. 575–583, 1988.
- [39] D. J. MacKay, "Bayesian neural networks and density networks," Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment, vol. 354, no. 1, pp. 73 – 80, 1995. Proceedings of the Third Workshop on Neutron Scattering Data Analysis.
- [40] C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra, "Weight uncertainty in neural network," in *Proceedings of the 32nd International Conference on Machine Learning* (F. Bach and D. Blei, eds.), vol. 37 of *Proceedings of Machine Learning Research*, (Lille, France), pp. 1613–1622, PMLR, 07–09 Jul 2015.
- [41] D. P. Kingma, T. Salimans, and M. Welling, "Variational dropout and the local reparameterization trick," in *Advances in Neural Information Processing Systems* 28 (C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, eds.), pp. 2575–2583, Curran Associates, Inc.
- [42] P.-L. Bacon, J. Harb, and D. Precup, "The option-critic architecture," in *Proceedings* of the Thirty-First AAAI Conference on Artificial Intelligence, AAAI'17, p. 1726–1734, AAAI Press, 2017.
- [43] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel, "High-dimensional continuous control using generalized advantage estimation," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2016.
- [44] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," *CoRR*, vol. abs/1602.01783, 2016.
- [45] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *Proceedings of the* 35th International Conference on Machine Learning (J. Dy and A. Krause, eds.), vol. 80 of *Proceedings of Machine Learning Research*, (Stockholmsmässan, Stockholm Sweden), pp. 1861–1870, PMLR, 10–15 Jul 2018.

- [46] H. Abdulsamad and J. Peters, "Hierarchical decomposition of nonlinear dynamics and control for system identification and policy distillation," in *Proceedings of the 2nd Conference on Learning for Dynamics and Control* (A. M. Bayen, A. Jadbabaie, G. Pappas, P. A. Parrilo, B. Recht, C. Tomlin, and M. Zeilinger, eds.), vol. 120 of *Proceedings of Machine Learning Research*, (The Cloud), pp. 904–914, PMLR, 10–11 Jun 2020.
- [47] F. A. Gers, J. A. Schmidhuber, and F. A. Cummins, "Learning to forget: Continual prediction with lstm," *Neural Comput.*, vol. 12, p. 2451–2471, Oct. 2000.
- [48] B. Lakshminarayanan, A. Pritzel, and C. Blundell, "Simple and scalable predictive uncertainty estimation using deep ensembles," in *Advances in Neural Information Processing Systems 30* (I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, eds.), pp. 6402–6413, Curran Associates, Inc., 2017.

A. Summary VBEM steps

A.1. VBHMM EM steps

E-step updates

$$\log \tilde{\boldsymbol{\pi}}_{k} \equiv \mathbb{E}_{q(\boldsymbol{\pi})} \left[\log \boldsymbol{\pi}_{k}\right] = \psi(\boldsymbol{\omega}_{k}^{(\boldsymbol{\pi})}) - \psi(\sum_{k=1}^{K} \boldsymbol{\omega}_{k}^{(\boldsymbol{\pi})}), \quad \sum_{k=1}^{K} \tilde{\boldsymbol{\pi}}_{k} \leq 1$$
$$\log \tilde{\boldsymbol{a}}_{kj} \equiv \mathbb{E}_{q(\mathbf{A})} \left[\log a_{kj}\right] = \psi(\boldsymbol{\omega}_{kj}^{(\mathbf{A})}) - \psi(\sum_{j=1}^{K} \boldsymbol{\omega}_{kj}^{(\mathbf{A})}), \quad \sum_{j=1}^{K} \tilde{\boldsymbol{a}}_{kj} \leq 1$$
$$\log \tilde{\boldsymbol{\Sigma}}_{k} \equiv \mathbb{E}_{q(\boldsymbol{\Sigma})} \left[\log |\boldsymbol{\Sigma}_{k}|\right] = \sum_{i=1}^{D} \psi\left(\frac{\nu_{k}+1-i}{2}\right) + D\log 2 + \log |\mathbf{W}_{k}|$$
$$\mathbb{E}_{q(\boldsymbol{\mu},\boldsymbol{\Sigma})} \left[(\mathbf{x}_{t} - \boldsymbol{\mu}_{k})^{T} \boldsymbol{\Sigma}_{k} (\mathbf{x}_{t} - \boldsymbol{\mu}_{k})\right] = D\beta_{k}^{-1} + \nu_{k} (\mathbf{m}_{k} - \mathbf{x}_{t})^{T} \mathbf{W}_{k} (\mathbf{m}_{k} - \mathbf{x}_{t})$$

M-step updates

$$\omega_k^{(\pi)} = \omega_0^{(\pi)} + \gamma_1(k)$$

$$\omega_{kj}^{(\mathbf{A})} = \omega_0^{(\mathbf{A})} + \sum_{t=2}^T \xi_{t-1,t}(k,j)$$

$$\beta_k = \beta_0 + N_k$$

$$\mathbf{m}_k = \frac{1}{\beta_k} \left(\beta_0 \mathbf{m}_0 + N_k \overline{\mathbf{x}}_t\right)$$

$$\mathbf{W}_k^{-1} = \mathbf{W}_0^{-1} N_k \mathbf{S}_k + \frac{\beta_0 N_k}{\beta_0 + N_k} (\overline{\mathbf{x}}_k - \mathbf{m}_0) (\overline{\mathbf{x}}_k - \mathbf{m}_0)^T$$

$$\nu_k = \nu_0 + N_k$$

A.2. VBAR-HMM EM steps

E-step updates

$$\log \tilde{\boldsymbol{\pi}}_{k} \equiv \mathbb{E}_{q(\boldsymbol{\pi})} \left[\log \boldsymbol{\pi}_{k}\right] = \psi(\boldsymbol{\omega}_{k}^{(\boldsymbol{\pi})}) - \psi(\sum_{k=1}^{K} \boldsymbol{\omega}_{k}^{(\boldsymbol{\pi})}), \quad \sum_{k=1}^{K} \tilde{\boldsymbol{\pi}}_{k} \leq 1$$
$$\log \tilde{\boldsymbol{\omega}}_{kj} \equiv \mathbb{E}_{q(\mathbf{A})} \left[\log a_{kj}\right] = \psi(\boldsymbol{\omega}_{kj}^{(\mathbf{A})}) - \psi(\sum_{j=1}^{K} \boldsymbol{\omega}_{kj}^{(\mathbf{A})}), \quad \sum_{j=1}^{K} \tilde{\boldsymbol{\alpha}}_{kj} \leq 1$$
$$\log \tilde{\boldsymbol{\Sigma}}_{k} \equiv \mathbb{E}_{q(\boldsymbol{\Sigma})} \left[\log |\boldsymbol{\Sigma}_{k}|\right] = \sum_{i=1}^{D} \psi\left(\frac{\nu_{k}+1-i}{2}\right) + D\log 2 + \log |\mathbf{W}_{k}|$$
$$\mathbb{E}_{q(\boldsymbol{\mu},\boldsymbol{\Sigma})} \left[(\mathbf{x}_{1}-\boldsymbol{\mu}_{k})^{T} \boldsymbol{\Sigma}_{k}(\mathbf{x}_{1}-\boldsymbol{\mu}_{k})\right] = D\beta_{k}^{-1} + \nu_{k}(\mathbf{m}_{k}-\mathbf{x}_{1})^{T} \mathbf{W}_{k}(\mathbf{m}_{k}-\mathbf{x}_{1})$$
$$\log \tilde{\mathbf{V}}_{k} \equiv \mathbb{E}_{q(\mathbf{V})} \left[\log |\mathbf{V}_{k}|\right] = \sum_{i=1}^{D} \psi\left(\frac{\eta_{k}+1-i}{2}\right) + D\log 2 + \log |\mathbf{P}_{k}|$$
$$\mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[(\mathbf{x}_{t}-\mathbf{U}_{k}\hat{\mathbf{x}}_{t-1})^{T} \mathbf{V}_{k}(\mathbf{x}_{t}-\mathbf{U}_{k}\hat{\mathbf{x}}_{t-1})\right] = \eta_{k}(\mathbf{M}_{k}\hat{\mathbf{x}}_{t-1}-\mathbf{x}_{t})^{T} \mathbf{P}_{k}(\mathbf{M}_{k}\hat{\mathbf{x}}_{t-1}-\mathbf{x}_{t}) + \mathbf{Tr}\left\{\mathbf{K}_{k}^{-1}\hat{\mathbf{x}}_{t-1}\hat{\mathbf{x}}_{t-1}^{T}\right\}$$

M-step updates

$$\begin{split} \omega_k^{(\boldsymbol{\pi})} &= \omega_0^{(\boldsymbol{\pi})} + \gamma_1(k) \\ \omega_{kj}^{(\boldsymbol{\Lambda})} &= \omega_0^{(\boldsymbol{\Lambda})} + \sum_{t=2}^T \xi_{t-1,t}(k,j) \\ \beta_k &= \beta_0 + N_k \\ \boldsymbol{\mathbf{m}}_k &= \frac{1}{\beta_k} \left(\beta_0 \boldsymbol{\mathbf{m}}_0 + N_k \boldsymbol{\overline{\mathbf{x}}}_t \right) \\ \boldsymbol{\mathbf{W}}_k^{-1} &= \boldsymbol{\mathbf{W}}_0^{-1} N_k \boldsymbol{\mathbf{S}}_k + \frac{\beta_0 N_k}{\beta_0 + N_k} (\boldsymbol{\overline{\mathbf{x}}}_k - \boldsymbol{\mathbf{m}}_0) (\boldsymbol{\overline{\mathbf{x}}}_k - \boldsymbol{\mathbf{m}}_0)^T \\ \nu_k &= \nu_0 + N_k \\ \boldsymbol{\mathbf{K}}_k &= \sum_{t=2}^T \gamma_t(k) \hat{\boldsymbol{\mathbf{x}}}_{t-1} \hat{\boldsymbol{\mathbf{x}}}_{t-1}^T + \boldsymbol{\mathbf{K}}_0 \\ \boldsymbol{\mathbf{M}}_k &= \left[\sum_{t=2}^T \gamma_t(k) \boldsymbol{\mathbf{x}}_t \hat{\boldsymbol{\mathbf{x}}}_{t-1}^T + \boldsymbol{\mathbf{M}}_0 \boldsymbol{\mathbf{K}}_0 \right] \boldsymbol{\mathbf{K}}_k^{-1} \\ \boldsymbol{\mathbf{P}}_k^{-1} &= \boldsymbol{\mathbf{P}}_0^{-1} + \boldsymbol{\mathbf{M}}_0 \boldsymbol{\mathbf{K}}_0 \boldsymbol{\mathbf{M}}_0^T + \sum_{t=2}^T \gamma_t(k) \boldsymbol{\mathbf{x}}_t \boldsymbol{\mathbf{x}}_t^T - \boldsymbol{\mathbf{M}}_k \boldsymbol{\mathbf{K}}_k \boldsymbol{\mathbf{M}}_k^T \\ \eta_k &= \eta_0 + \sum_{t=2}^T \gamma_t(k) \end{split}$$
A.3. VBrAR-HMM EM steps

E-step updates

$$\log \tilde{\boldsymbol{\pi}}_{k} \equiv \mathbb{E}_{q(\boldsymbol{\pi})} \left[\log \boldsymbol{\pi}_{k}\right] = \psi(\boldsymbol{\omega}_{k}^{(\boldsymbol{\pi})}) - \psi(\sum_{k=1}^{K} \boldsymbol{\omega}_{k}^{(\boldsymbol{\pi})}), \quad \sum_{k=1}^{K} \tilde{\boldsymbol{\pi}}_{k} \leq 1$$
$$\mathbb{E}_{q(\boldsymbol{\phi})} \left[\log \psi_{ij}\right] = f_{\boldsymbol{\phi}}(\mathbf{x}, \mathbf{u})$$
$$\log \tilde{\boldsymbol{\Sigma}}_{k} \equiv \mathbb{E}_{q(\boldsymbol{\Sigma})} \left[\log |\boldsymbol{\Sigma}_{k}|\right] = \sum_{i=1}^{D} \psi\left(\frac{\nu_{k}+1-i}{2}\right) + D\log 2 + \log |\mathbf{W}_{k}|$$
$$\mathbb{E}_{q(\boldsymbol{\mu},\boldsymbol{\Sigma})} \left[(\mathbf{x}_{1}-\boldsymbol{\mu}_{k})^{T} \boldsymbol{\Sigma}_{k}(\mathbf{x}_{1}-\boldsymbol{\mu}_{k})\right] = D\beta_{k}^{-1} + \nu_{k}(\mathbf{m}_{k}-\mathbf{x}_{1})^{T} \mathbf{W}_{k}(\mathbf{m}_{k}-\mathbf{x}_{1})$$
$$\log \tilde{\mathbf{V}}_{k} \equiv \mathbb{E}_{q(\mathbf{V})} \left[\log |\mathbf{V}_{k}|\right] = \sum_{i=1}^{D} \psi\left(\frac{\eta_{k}+1-i}{2}\right) + D\log 2 + \log |\mathbf{P}_{k}|$$
$$\mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[(\mathbf{x}_{t}-\mathbf{U}_{k}\hat{\mathbf{x}}_{t-1})^{T} \mathbf{V}_{k}\left(\mathbf{x}_{t}-\mathbf{U}_{k}\hat{\mathbf{x}}_{t-1}\right)\right] = \eta_{k}(\mathbf{M}_{k}\hat{\mathbf{x}}_{t-1}-\mathbf{x}_{t})^{T} \mathbf{P}_{k}(\mathbf{M}_{k}\hat{\mathbf{x}}_{t-1}-\mathbf{x}_{t})$$
$$+ \mathbf{Tr}\left\{\mathbf{K}_{k}^{-1}\hat{\mathbf{x}}_{t-1}\hat{\mathbf{x}}_{t-1}^{T}\right\}$$

M-step updates

$$\begin{split} \omega_k^{(\pi)} &= \omega_0^{(\pi)} + \gamma_1(k) \\ \max_{\boldsymbol{\phi}} \log p(\mathbf{z}_{2:T} | \boldsymbol{\phi}) &= \max_{\boldsymbol{\phi}} \sum_{t=2}^T \sum_{k=1}^K \sum_{j=1}^K \xi_{t-1,t}(k,j) \log \psi_{kj} \\ \beta_k &= \beta_0 + N_k \\ \mathbf{m}_k &= \frac{1}{\beta_k} \left(\beta_0 \mathbf{m}_0 + N_k \overline{\mathbf{x}}_t \right) \\ \mathbf{W}_k^{-1} &= \mathbf{W}_0^{-1} N_k \mathbf{S}_k + \frac{\beta_0 N_k}{\beta_0 + N_k} (\overline{\mathbf{x}}_k - \mathbf{m}_0) (\overline{\mathbf{x}}_k - \mathbf{m}_0)^T \\ \nu_k &= \nu_0 + N_k \\ \mathbf{K}_k &= \sum_{t=2}^T \gamma_t(k) \hat{\mathbf{x}}_{t-1} \hat{\mathbf{x}}_{t-1}^T + \mathbf{K}_0 \\ \mathbf{M}_k &= \left[\sum_{t=2}^T \gamma_t(k) \mathbf{x}_t \hat{\mathbf{x}}_{t-1}^T + \mathbf{M}_0 \mathbf{K}_0 \right] \mathbf{K}_k^{-1} \\ \mathbf{P}_k^{-1} &= \mathbf{P}_0^{-1} + \mathbf{M}_0 \mathbf{K}_0 \mathbf{M}_0^T + \sum_{t=2}^T \gamma_t(k) \mathbf{x}_t \mathbf{x}_t^T - \mathbf{M}_k \mathbf{K}_k \mathbf{M}_k^T \\ \eta_k &= \eta_0 + \sum_{t=2}^T \gamma_t(k) \end{split}$$

B. Derivation VBHMMs

B.1. Derivation VBHMM

Variational E-step full derivation Solving $q(\mathbf{z}_{1:T})$ with mean-field approximation:

$$\begin{split} \log q(\mathbf{z}_{1:T}) &\propto \mathbb{E}_{q(\pi,\mu,\boldsymbol{\Sigma},\mathbf{A})}[\log p(\mathbf{x}_{1:T},\mathbf{z}_{1:T},\pi,\mathbf{A},\mu,\boldsymbol{\Sigma})] \\ &= \mathbb{E}_{q(\pi,\mu,\boldsymbol{\Sigma},\mathbf{A})}[p(\mathbf{x}_{1:T}|\mathbf{z}_{1:T},\mu,\boldsymbol{\Sigma})p(\mathbf{z}_{1:T}|\pi,\mathbf{A})p(\mathbf{A})p(\pi)p(\mu|\boldsymbol{\Sigma})p(\boldsymbol{\Sigma})] \\ &= \mathbb{E}_{q(\pi,\mu,\boldsymbol{\Sigma},\mathbf{A})}[\log p(\mathbf{z}_{1:T}|\pi,\mathbf{A})p(\mathbf{x}_{1:T}|\mathbf{z}_{1:T},\mu,\boldsymbol{\Sigma})] \\ &= \mathbb{E}_{q(\pi,\mathbf{A})}[\log p(\mathbf{z}_{1:T}|\pi,\mathbf{A})] + \mathbb{E}_{q(\mu,\boldsymbol{\Sigma})}[\log p(\mathbf{x}_{1:T}|\mathbf{z}_{1:T},\mu,\boldsymbol{\Sigma})] \\ &= \mathbb{E}_{q(\pi,\mathbf{A})}[\log p(\mathbf{z}_{1:T}|\pi,\mathbf{A})] + \mathbb{E}_{q(\mu,\boldsymbol{\Sigma})}[\log p(\mathbf{x}_{1:T}|\mathbf{z}_{1:T},\mu,\boldsymbol{\Sigma})] \\ &= \mathbb{E}_{q(\pi,\mathbf{A})}\left[\log p(\mathbf{z}_{1}|\pi)p(\mathbf{z}_{2:T}|\mathbf{z}_{1},\mathbf{A})\right] + \mathbb{E}_{q(\mu,\boldsymbol{\Sigma})}[\log p(\mathbf{x}_{1:T}|\mathbf{z}_{1:T},\mu,\boldsymbol{\Sigma})] \\ &= \mathbb{E}_{q(\pi,\mathbf{A})}\left[\log \prod_{k=1}^{K} \left(\pi_{k}^{\mathbb{I}(z_{1}=k)}\right)\right] \prod_{t=2}^{T} \prod_{k=1}^{K} \prod_{j=1}^{K} a_{kj}^{\mathbb{I}(z_{t-1}=k,z_{t}=j)}\right] \\ &+ \mathbb{E}_{q(\mu,\boldsymbol{\Sigma})}\left[\log \prod_{t=1}^{T} \prod_{k=1}^{K} \mathcal{N}\left(\mathbf{x}_{t}|\mu_{k},\boldsymbol{\Sigma}_{k}^{-1}\right)^{\mathbb{I}(z_{t}=k)}\right] \\ &= \mathbb{E}_{q(\pi)}\left[\log \prod_{k=1}^{K} \left(\pi_{k}^{\mathbb{I}(z_{1}=k)}\right)\right] + \mathbb{E}_{q(\mathbf{A})}\left[\log \prod_{t=2}^{T} \prod_{k=1}^{K} a_{kj}^{\mathbb{I}(z_{t-1}=k,z_{t}=j)}\right] \\ &+ \mathbb{E}_{q(\mu,\boldsymbol{\Sigma})}\left[\log \prod_{t=1}^{T} \prod_{k=1}^{K} \mathcal{N}\left(\mathbf{x}_{t}|\mu_{k},\boldsymbol{\Sigma}_{k}^{-1}\right)^{\mathbb{I}(z_{t}=k)}\right] \\ &= \mathbb{E}_{q(\pi)}\left[\sum_{k=1}^{K} \mathbb{I}(z_{1}=k)\log \pi_{k}\right] \\ &+ \mathbb{E}_{q(\mathbf{A})}\left[\sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} \mathbb{I}(z_{t-1}=k,z_{t}=j)\log a_{kj}\right] \\ &+ \mathbb{E}_{q(\mu,\boldsymbol{\Sigma})}\left[\sum_{t=2}^{T} \sum_{k=1}^{K} \mathbb{I}(z_{t}=k)\log \mathcal{N}\left(\mathbf{x}_{t}|\mu_{k},\boldsymbol{\Sigma}_{k}^{-1}\right)\right] \end{split}$$

Continuation of the derivation for $\log q(\mathbf{z}_{1:T})$:

$$\log q(\mathbf{z}_{1:T}) \propto \sum_{k=1}^{K} \mathbb{I}(z_1 = k)$$

$$\mathbb{E}_{q(\boldsymbol{\pi})} \left[\log \boldsymbol{\pi}_k\right] + \sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} \mathbb{I}(z_{t-1} = k, z_t = j) \mathbb{E}_{q(\mathbf{A})} \left[\log a_{kj}\right]$$

$$+ \sum_{t=1}^{T} \sum_{k=1}^{K} \mathbb{I}(z_t = k)$$

$$\mathbb{E}_{q(\boldsymbol{\mu}, \boldsymbol{\Sigma})} \left[\frac{1}{2} \log |\boldsymbol{\Sigma}_k| - \frac{D}{2} \log 2\pi - \frac{1}{2} (\mathbf{x}_t - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k (\mathbf{x}_t - \boldsymbol{\mu}_k)\right]$$

$$= \sum_{k=1}^{K} \mathbb{I}(z_1 = k)$$

$$\mathbb{E}_{q(\boldsymbol{\pi})} \left[\log \boldsymbol{\pi}_k\right] + \sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} \mathbb{I}(z_{t-1} = k, z_t = j) \mathbb{E}_{q(\mathbf{A})} \left[\log a_{kj}\right]$$

$$+ \sum_{t=1}^{T} \sum_{k=1}^{K} \mathbb{I}(z_t = k)$$

$$\left\{\frac{1}{2} \mathbb{E}_{q(\boldsymbol{\Sigma})} \left[\log |\boldsymbol{\Sigma}_k|\right] - \frac{D}{2} \log (2\pi) - \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\mu}, \boldsymbol{\Sigma})} \left[(\mathbf{x}_t - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k (\mathbf{x}_t - \boldsymbol{\mu}_k) \right] \right\}$$

Derving $\mathbb{E}_{q(\boldsymbol{\mu}, \boldsymbol{\Sigma})}$:

$$\begin{split} \mathbb{E}_{q(\boldsymbol{\mu},\boldsymbol{\Sigma})} \left[(\mathbf{x}_{t} - \boldsymbol{\mu}_{k})^{T} \boldsymbol{\Sigma}_{k} (\mathbf{x}_{t} - \boldsymbol{\mu}_{k}) \right] \\ &= \int \int \left((\mathbf{x}_{t} - \boldsymbol{\mu}_{k})^{T} \boldsymbol{\Sigma}_{k} (\mathbf{x}_{t} - \boldsymbol{\mu}_{k}) \right) q(\boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) \, \mathrm{d}\boldsymbol{\mu}_{k} \, \mathrm{d}\boldsymbol{\Sigma}_{k} \\ &= \int \left\{ \int (\mathbf{x}_{t} - \boldsymbol{\mu}_{k})^{T} \boldsymbol{\Sigma}_{k} (\mathbf{x}_{t} - \boldsymbol{\mu}_{k}) q(\boldsymbol{\mu}_{k} | \boldsymbol{\Sigma}_{k}) \, \mathrm{d}\boldsymbol{\mu}_{k} \right\} q(\boldsymbol{\Sigma}_{k}) \, \mathrm{d}\boldsymbol{\Sigma}_{k} \\ &= \int \mathbb{E}_{q(\boldsymbol{\mu})} \left[(\boldsymbol{\mu}_{k} - \mathbf{x}_{t})^{T} \boldsymbol{\Sigma}_{k} (\boldsymbol{\mu}_{k} - \mathbf{x}_{t}) \right] q(\boldsymbol{\Sigma}_{k}) \, \mathrm{d}\boldsymbol{\Sigma}_{k} \end{split}$$

Using the equation (380) from [36] to solve the inner expectation with respect to μ yields:

$$\mathbb{E}_{q(\boldsymbol{\mu})}\left[(\boldsymbol{\mu}_{k}-\mathbf{x}_{t})^{T}\boldsymbol{\Sigma}_{k}(\boldsymbol{\mu}_{k})-\mathbf{x}_{t}\right] = (\mathbf{m}_{k}-\mathbf{x}_{t})^{T}\boldsymbol{\Sigma}_{k}(\mathbf{m}_{k}-\mathbf{x}_{t}) + \mathbf{Tr}\left(\boldsymbol{\Sigma}_{k}(\boldsymbol{\beta}_{k}^{-1}\boldsymbol{\Sigma}_{k})^{-1}\right) \\ = (\mathbf{m}_{k}-\mathbf{x}_{t})^{T}\boldsymbol{\Sigma}_{k}(\mathbf{m}_{k}-\mathbf{x}_{t}) + D\boldsymbol{\beta}_{k}^{-1}$$

Plugging this term back in the equation

$$\begin{split} \mathbb{E}_{q(\boldsymbol{\mu},\boldsymbol{\Sigma})} \left[(\mathbf{x}_{t} - \boldsymbol{\mu}_{k})^{T} \boldsymbol{\Sigma}_{k} (\mathbf{x}_{t} - \boldsymbol{\mu}_{k}) \right] \\ &= \int \left\{ (\mathbf{m}_{k} - \mathbf{x}_{t})^{T} \boldsymbol{\Sigma}_{k} (\mathbf{m}_{k} - \mathbf{x}_{t}) + D\beta_{k}^{-1} \right\} q(\boldsymbol{\Sigma}_{k}) \, \mathrm{d}\boldsymbol{\Sigma}_{k} \\ &= D\beta_{k}^{-1} + \mathbb{E}_{q(\boldsymbol{\Sigma})} \left[(\mathbf{m}_{k} - \mathbf{x}_{t})^{T} \boldsymbol{\Sigma}_{k} (\mathbf{m}_{k} - \mathbf{x}_{t}) \right] \\ &= D\beta_{k}^{-1} + \mathbb{E}_{q(\boldsymbol{\Sigma})} \left[\mathbf{Tr} \left\{ \boldsymbol{\Sigma}_{k} (\mathbf{m}_{k} - \mathbf{x}_{t}) (\mathbf{m}_{k} - \mathbf{x}_{t})^{T} \right\} \right] \\ &= D\beta_{k}^{-1} + \mathbf{Tr} \left\{ \mathbb{E}_{q(\boldsymbol{\Sigma})} \left[\boldsymbol{\Sigma}_{k} \right] (\mathbf{m}_{k} - \mathbf{x}_{t}) (\mathbf{m}_{k} - \mathbf{x}_{t})^{T} \right\} \\ &= D\beta_{k}^{-1} + \mathbf{Tr} \left\{ \nu_{k} \mathbf{W}_{k} (\mathbf{m}_{k} - \mathbf{x}_{t}) (\mathbf{m}_{k} - \mathbf{x}_{t})^{T} \right\} \\ &= D\beta_{k}^{-1} + \nu_{k} (\mathbf{m}_{k} - \mathbf{x}_{t})^{T} \mathbf{W}_{k} (\mathbf{m}_{k} - \mathbf{x}_{t}) \end{split}$$

Variational M-step full derivation Solving $\log q(\pi, \mathbf{A}, \mu, \Sigma)$ with mean-field approximation:

$$\begin{split} \log q(\pi, \mathbf{A}, \mu, \Sigma) &\propto \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}, \pi, \mathbf{A}, \mu, \Sigma)] \\ &= \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma) p(\mathbf{z}_{1:T} | \pi, \mathbf{A}) p(\pi) p(\mathbf{A}) p(\mu | \Sigma) p(\Sigma)] \\ &= \log p(\pi) + \log p(\mathbf{A}) + \sum_{k=1}^{K} \log p(\mu_{k} | \Sigma_{k}) p(\Sigma_{k}) \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] + \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{z}_{1:T} | \pi, \mathbf{A})] \\ &= \log p(\pi) + \log p(\mathbf{A}) + \sum_{k=1}^{K} \log p(\mu_{k} | \Sigma_{k}) p(\Sigma_{k}) \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{z}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{z}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{z}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{x}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{z}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{z}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\log p(\mathbf{z}_{1:T} | \mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} [\mathbb{E}_{q(\mathbf{z}_{1:T})} [\mathbb{E}_{q(\mathbf{z}_{1:T})} [\mathbb{E}_{q(\mathbf{z}_{1:T}) [\mathbb{E}_{q(\mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} \sum_{\mathbf{z}_{1:T}} \mathbb{E}_{q(\mathbf{z}_{1:T})} [\mathbb{E}_{q(\mathbf{z}_{1:T})} [\mathbb{E}_{q(\mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} \sum_{\mathbf{z}_{1:T}} \mathbb{E}_{q(\mathbf{z}_{1:T})} [\mathbb{E}_{q(\mathbf{z}_{1:T}) [\mathbb{E}_{q(\mathbf{z}_{1:T}, \mu, \Sigma)] \\ &+ \mathbb{E}_{q(\mathbf{z}_{1:T})} \sum_{\mathbf{z}_{1:T}} \mathbb{E}_{q(\mathbf{z}_{$$

$$\log q(\mathbf{A}) \propto \sum_{k=1}^{K} \sum_{j=1}^{K} (\omega_{0}^{(\mathbf{A})} - 1) \log a_{kj}$$

$$+ \sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\mathbb{I}(z_{t-1} = k, z_{t} = j) \right] \log a_{kj}$$

$$= \sum_{k=1}^{K} \sum_{j=1}^{K} \log a_{kj} \left\{ (\omega_{0}^{(\mathbf{A})} - 1) + \sum_{t=2}^{T} \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\mathbb{I}(z_{t-1} = k, z_{t} = j) \right] \right\}$$

$$- \frac{1}{2} \left(\boldsymbol{\mu}_{k}^{T} (\beta_{k} \boldsymbol{\Sigma}_{k}) \boldsymbol{\mu}_{k} \right) = - \frac{1}{2} \left\{ \boldsymbol{\mu}_{k}^{T} \left(\beta_{0} + \sum_{t=1}^{T} \gamma_{t}(k) \right) \boldsymbol{\Sigma}_{k} \boldsymbol{\mu}_{k} \right\}$$

$$\beta_{k} = \beta_{0} + \sum_{t=1}^{T} \gamma_{t}(k)$$

$$\boldsymbol{\mu}_{k}^{T} \boldsymbol{\Sigma}_{k} \beta_{k} \mathbf{m}_{k} = \boldsymbol{\mu}_{k}^{T} \boldsymbol{\Sigma}_{k} \left(\beta_{0} \mathbf{m}_{0} + \sum_{t=1}^{T} \gamma_{t}(k) \mathbf{x}_{t} \right)$$

$$\mathbf{m}_{k} = \frac{1}{\beta_{k}} \left(\beta_{0} \mathbf{m}_{0} + \sum_{t=1}^{T} \gamma_{t}(k) \mathbf{x}_{t} \right)$$

Deriving ν_k by rearanging the term to obtain the form of the Wishart distribution:

$$\frac{(\nu_k - D - 1)}{2} \log |\mathbf{\Sigma}_k| = \frac{(\nu_0 - D - 1)}{2} \log |\mathbf{\Sigma}_k| + \frac{1}{2} (\sum_{t=1}^T \gamma_t(k)) \log |\mathbf{\Sigma}_k|$$
$$\nu_k = \nu_0 + \sum_{t=1}^T \gamma_t(k)$$

70

B.2. Derivation VBAR-HMM

Variational E-step full derivation

$$\begin{split} \log q(\mathbf{z}_{1:T}) &\propto \mathbb{E}_{q(\pi,\mathbf{A},\mu,\Sigma,\mathbf{U},\mathbf{V})} \left[\log p(\mathbf{x}_{1:T}, \hat{\mathbf{x}}_{2:T}, \mathbf{z}_{1:T}, \boldsymbol{\pi}, \mathbf{A}, \mu, \Sigma, \mathbf{U}, \mathbf{V}) \right] \\ &= \mathbb{E}_{q(\pi,\mathbf{A},\mu,\Sigma,\mathbf{U},\mathbf{V})} \left[\log p(\mathbf{x}_{1:T} | \hat{\mathbf{x}}_{2:T}, \mathbf{z}_{1:T}, \mu, \Sigma \mathbf{U}, \mathbf{V}) p(\mathbf{z}_{1:T} | \boldsymbol{\pi}, \mathbf{A}) p(\mathbf{A}) p(\pi) \\ &\quad p(\mu|\Sigma) p(\Sigma) p(\mathbf{U}|V) p(\mathbf{V}) \right] \\ &= \mathbb{E}_{q(\pi,\mathbf{A})} \left[\log p(z_1 | \pi) p(z_{2:T} | \mathbf{z}_{1:T}, \mathbf{U}, \mathbf{V}) \right] \\ &\quad + \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\log p(\mathbf{x}_{1:T} | \hat{\mathbf{x}}_{2:T}, \mathbf{z}_{1:T}, \mathbf{U}, \mathbf{V}) \right] \\ &= \mathbb{E}_{q(\pi,\mathbf{A})} \left[\log \prod_{k=1}^{K} \left(\pi_{k}^{\mathbb{I}(z_{1}=k)} \right) \prod_{t=2}^{T} \prod_{i=1}^{K} \prod_{j=1}^{K} a_{ij}^{\mathbb{I}(z_{i}=j,z_{t-1}=i)} \right] \\ &\quad + \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\log \prod_{k=1}^{K} \mathcal{N} \left(\mathbf{x}_{1} | \mu_{k}, \Sigma_{k}^{-1} \right)^{\mathbb{I}(z_{1}=k)} \right] \\ &\quad + \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\log \prod_{k=1}^{T} \prod_{k=1}^{K} \mathcal{N} \left(\mathbf{x}_{1} | u_{k} \hat{\mathbf{x}}_{t-1}, \mathbf{V}_{k}^{-1} \right)^{\mathbb{I}(z_{i}=k)} \right] \\ &\quad + \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\log \prod_{i=2,k=1}^{K} \mathcal{N} \left(\mathbf{x}_{1} | \mu_{k}, \Sigma_{k}^{-1} \right)^{\mathbb{I}(z_{1}=k)} \right] \\ &\quad + \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\log \prod_{k=1}^{K} \mathcal{N} \left(\mathbf{x}_{1} | \mu_{k}, \Sigma_{k}^{-1} \right)^{\mathbb{I}(z_{i}=k)} \right] \\ &\quad + \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\log \prod_{k=1}^{T} \prod_{k=1}^{K} \mathcal{N} \left(\mathbf{x}_{1} | \mu_{k}, \Sigma_{k}^{-1} \right)^{\mathbb{I}(z_{i}=k)} \right] \\ &\quad + \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\log \prod_{i=1,k=1}^{T} \mathcal{N} \left(\mathbf{x}_{i} | u_{k} \hat{\mathbf{x}}_{i-1}, \mathbf{V}_{k}^{-1} \right)^{\mathbb{I}(z_{i}=k)} \right] \\ &\quad + \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\log \prod_{i=1,k=1}^{T} \mathcal{N} \left(\mathbf{x}_{i} | u_{k} \hat{\mathbf{x}}_{i-1}, \mathbf{V}_{k}^{-1} \right)^{\mathbb{I}(z_{i}=k)} \right] \\ &\quad + \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\log \prod_{i=1,k=1}^{T} \mathcal{N} \left(\mathbf{x}_{i} | u_{k} \hat{\mathbf{x}}_{i-1}, \mathbf{V}_{k}^{-1} \right)^{\mathbb{I}(z_{i}=k)} \right] \\ &\quad + \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\sum_{k=1}^{K} \mathbb{I} \left(z_{1} = k \right) \log \mathcal{N} \left(\mathbf{x}_{i} | \mu_{k}, \Sigma_{k}^{-1} \right) \right] \\ &\quad + \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\sum_{k=2}^{K} \sum_{i=1}^{K} \sum_{j=1}^{K} \mathbb{I} \left(z_{i} = j, z_{i-1} = i \right) \log a_{ij} \right] \\ &\quad + \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\sum_{k=2}^{K} \sum_{i=1}^{K} \mathbb{I} \left(z_{i} = k \right) \log \mathcal{N} \left(\mathbf{x}_{i} | u_{k} \hat{\mathbf{x}}_{i-1}, \mathbf{V}_{k}^{-1} \right) \right] \\ &\quad + \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\sum_{k=2}^{K} \sum_{i=1}^{K} \sum_{j=1}^{K} \mathbb{I} \left(z_{i} = k \right) \log \mathcal{N} \left(\mathbf{x}_{i} | u_{k} \hat{\mathbf{x}}_{i-1} \right) \right] \\ &\quad +$$

Continuation of the derivation for $\log q(\mathbf{z}_{1:T})$

$$\begin{split} \log q(\mathbf{z}_{1:T}) \propto \mathbb{E}_{q(\boldsymbol{\pi})} \left[\sum_{k=1}^{K} \mathbb{I}(z_{1}=k) \log(\boldsymbol{\pi}_{k}) \right] \\ &+ \mathbb{E}_{q(\mathbf{A})} \left[\sum_{t=2}^{T} \sum_{i=1}^{K} \sum_{j=1}^{K} \mathbb{I}(z_{t}=j, z_{t-1}=i) \log a_{ij} \right] \\ &+ \sum_{k=1}^{K} \mathbb{I}(z_{1}=k) \\ &\mathbb{E}_{q(\boldsymbol{\mu},\boldsymbol{\Sigma})} \left[\frac{1}{2} \log |\boldsymbol{\Sigma}_{k}| - \frac{D}{2} \log 2\boldsymbol{\pi} - \frac{1}{2} (\mathbf{x}_{1}-\boldsymbol{\mu}_{k})^{T} \boldsymbol{\Sigma}_{k} (\mathbf{x}_{1}-\boldsymbol{\mu}_{k}) \right] \\ &+ \sum_{t=2}^{T} \sum_{k=1}^{K} \mathbb{I}(z_{t}=k) \\ &\mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\frac{1}{2} \log |\mathbf{V}_{k}| - \frac{D}{2} \log 2\boldsymbol{\pi} - \frac{1}{2} (\mathbf{x}_{t}-\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1})^{T} \mathbf{V}_{k} (\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1}) \right] \\ &= \sum_{k=1}^{K} \mathbb{I}(z_{1}=k) \mathbb{E}_{q(\boldsymbol{\pi})} \left[\log \boldsymbol{\pi}_{k} \right] + \sum_{t=2}^{T} \sum_{k=1}^{K} \sum_{j=1}^{K} \mathbb{I}(z_{t-1}=k, z_{t}=j) \mathbb{E}_{q(\mathbf{A})} \left[\log a_{kj} \right] \\ &+ \sum_{k=1}^{K} \mathbb{I}(z_{1}=k) \\ &\left\{ \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\Sigma})} \left[\log |\boldsymbol{\Sigma}_{k}| \right] - \frac{D}{2} \log (2\boldsymbol{\pi}) - \frac{1}{2} \mathbb{E}_{q(\boldsymbol{\mu},\boldsymbol{\Sigma})} \left[(\mathbf{x}_{1}-\boldsymbol{\mu}_{k})^{T} \boldsymbol{\Sigma}_{k} (\mathbf{x}_{1}-\boldsymbol{\mu}_{k}) \right] \right\} \\ &+ \sum_{t=2}^{T} \sum_{k=1}^{K} \mathbb{I}(z_{t}=k) \\ &\left\{ \frac{1}{2} \mathbb{E}_{q(\mathbf{V})} \left[\log |\mathbf{V}_{k}| \right] - \frac{D}{2} \log (2\boldsymbol{\pi}) \\ &- \frac{1}{2} \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[(\mathbf{x}_{t}-\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1})^{T} \mathbf{V}_{k} (\mathbf{x}_{t}-\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1}) \right] \right\} \end{split}$$

$$\begin{split} \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\left(\mathbf{x}_{t} - \mathbf{U}_{k} \hat{\mathbf{x}}_{t-1} \right)^{T} \mathbf{V}_{k} \left(\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1} \right) \right] \\ &= \int \int \left(\left(\mathbf{x}_{t} - \mathbf{U}_{k} \hat{\mathbf{x}}_{t-1} \right)^{T} \mathbf{V}_{k} \left(\mathbf{x}_{t} - \mathbf{U}_{k} \hat{\mathbf{x}}_{t-1} \right) \right) q(\mathbf{U}_{k}, \mathbf{V}_{k}) \, \mathrm{d}\mathbf{U}_{k} \, \mathrm{d}\mathbf{V}_{k} \\ &= \int \left\{ \int (\mathbf{x}_{t} - \mathbf{U}_{k} \hat{\mathbf{x}}_{t-1})^{T} \mathbf{V}_{k} (\mathbf{x}_{t} - \mathbf{U}_{k} \hat{\mathbf{x}}_{t-1}) q(\mathbf{U}_{k} | \mathbf{V}_{k}) \, \mathrm{d}\mathbf{U}_{k} \right\} q(\mathbf{V}_{k}) \, \mathrm{d}\mathbf{V}_{k} \\ &= \int \mathbb{E}_{q(\mathbf{U})} \left[(\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t})^{T} \mathbf{V}_{k} (\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t}) \right] q(\mathbf{V}_{k}) \, \mathrm{d}\mathbf{V}_{k} \end{split}$$

The inner expectation corresponds to a Matrix-Normal distribution. Thus, we make use of transformations derived in [38] and get two expectations:

$$\begin{split} \mathbb{E}\left[\mathbf{U}_{k}\right] &= \mathbf{M}_{k} \\ \mathbb{E}\left[\mathbf{U}_{k}\mathbf{Q}\mathbf{U}_{k}^{T}\right] &= \mathbf{M}_{k}\mathbf{U}\mathbf{M}_{k}^{T} + \mathbf{Tr}\left\{\mathbf{K}_{k}^{-1}\mathbf{Q}\right\}\mathbf{V}_{k}^{-1} \end{split}$$

Applying this transformations on the inner expectation from Equation (B.2) gives:

$$\begin{split} \mathbb{E}_{q(\mathbf{U})} \left[(\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t})^{T} \mathbf{V}_{k} (\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t}] \\ &= \mathbb{E}_{q(\mathbf{U})} \left[\mathbf{Tr} \left\{ \mathbf{V}_{k} (\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t}) (\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t})^{T} \right\} \right] \\ &= \mathbf{Tr} \left\{ \mathbf{V}_{k} (\mathbb{E}_{q(\mathbf{U})} \left[\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1} \hat{\mathbf{x}}_{t-1}^{T} \mathbf{U}_{k}^{T} \right] \\ &- \mathbb{E}_{q(\mathbf{U})} \left[\mathbf{U}_{k} \right] \hat{\mathbf{x}}_{t-1} \mathbf{x}_{t}^{T} - \mathbf{x}_{t} \hat{\mathbf{x}}_{t-1}^{T} \mathbb{E}_{q(\mathbf{U})} \left[\mathbf{U}_{k}^{T} \right] + \mathbf{x}_{t} \mathbf{x}_{t}^{T} \right] \right\} \\ &= \mathbf{Tr} \left\{ \mathbf{V}_{k} (\mathbf{M}_{k} \hat{\mathbf{x}}_{t-1} \hat{\mathbf{x}}_{t-1}^{T} \mathbf{M}_{k}^{T} + (\mathbf{Tr} \left\{ \mathbf{K}_{k}^{-1} \hat{\mathbf{x}}_{t-1} \hat{\mathbf{x}}_{t-1}^{T} \right\} \mathbf{V}_{k}^{-1}) \\ &- \mathbf{M}_{k} \hat{\mathbf{x}}_{t-1} \mathbf{x}_{t}^{T} - \mathbf{x}_{t} \hat{\mathbf{x}}_{t-1}^{T} \mathbf{M}_{k}^{T} + \mathbf{x}_{t} \mathbf{x}_{t}^{T} \right) \right\} \\ &= \mathbf{Tr} \left\{ \mathbf{V}_{k} \left((\mathbf{M}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t}) (\mathbf{M}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t})^{T} + (\mathbf{Tr} \left\{ \mathbf{K}_{k}^{-1} \hat{\mathbf{x}}_{t-1} \hat{\mathbf{x}}_{t-1}^{T} \right\} \mathbf{V}_{k}^{-1}) \right\} \\ &= (\mathbf{M}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t})^{T} \mathbf{V}_{k} (\mathbf{M}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t}) + \mathbf{Tr} \left\{ \mathbf{K}_{k}^{-1} \hat{\mathbf{x}}_{t-1} \hat{\mathbf{x}}_{t-1}^{T} \right\} \end{split}$$

$$\begin{split} \mathbb{E}_{q(\mathbf{U},\mathbf{V})} \left[\left(\mathbf{x}_{t} - \mathbf{U}_{k} \hat{\mathbf{x}}_{t-1} \right)^{T} \mathbf{V}_{k} \left(\mathbf{U}_{k} \hat{\mathbf{x}}_{t-1} \right) \right] \\ &= \int \left[\left(\mathbf{M}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t} \right)^{T} \mathbf{V}_{k} \left(\mathbf{M}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t} \right) + \mathbf{Tr} \left\{ \mathbf{K}_{k}^{-1} \hat{\mathbf{x}}_{t-1} \hat{\mathbf{x}}_{t-1}^{T} \right\} \right] q(\mathbf{V}_{k}) \, \mathrm{d} \mathbf{V}_{k} \\ &= \mathbb{E}_{q(\mathbf{V})} \left[\left(\mathbf{M}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t} \right)^{T} \mathbf{V}_{k} \left(\mathbf{M}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t} \right) + \mathbf{Tr} \left\{ \mathbf{K}_{k}^{-1} \hat{\mathbf{x}}_{t-1} \hat{\mathbf{x}}_{t-1}^{T} \right\} \right] \\ &= \mathbf{Tr} \left\{ \mathbb{E}_{q(\mathbf{V})} \left[\mathbf{V}_{k} \right] \left(\mathbf{M}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t} \right) \left(\mathbf{M}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t} \right)^{T} \right\} + \mathbf{Tr} \left\{ \mathbf{K}_{k}^{-1} \hat{\mathbf{x}}_{t-1} \hat{\mathbf{x}}_{t-1}^{T} \right\} \\ &= \mathbf{Tr} \left\{ \eta_{k} \mathbf{P}_{k} \left(\mathbf{M}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t} \right) \left(\mathbf{M}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t} \right)^{T} \right\} + \mathbf{Tr} \left\{ \mathbf{K}_{k}^{-1} \hat{\mathbf{x}}_{t-1} \hat{\mathbf{x}}_{t-1}^{T} \right\} \\ &= \eta_{k} \left(\mathbf{M}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t} \right)^{T} \mathbf{P}_{k} \left(\mathbf{M}_{k} \hat{\mathbf{x}}_{t-1} - \mathbf{x}_{t} \right) + \mathbf{Tr} \left\{ \mathbf{K}_{k}^{-1} \hat{\mathbf{x}}_{t-1} \hat{\mathbf{x}}_{t-1}^{T} \right\} \end{split}$$

Variational M-step full derivation

$$\begin{split} \log q(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V}) &\propto \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\log p(\mathbf{x}_{1:T}, \hat{\mathbf{x}}_{2:T}, \mathbf{z}_{1:T}, \boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V}) \right] \\ &= \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\log (p(\mathbf{x}_{1:T} | \hat{\mathbf{x}}_{2:T}, \mathbf{z}_{1:T}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V}) \right. \\ &\quad p(\mathbf{z}_{1:T} | \boldsymbol{\pi}, \mathbf{A}) p(\boldsymbol{\pi}) p(\mathbf{A}) p(\mathbf{U} | \boldsymbol{\mu}) p(\boldsymbol{\Sigma}) p(\mathbf{U} | \mathbf{V}) p(\mathbf{V})) \right] \\ &= \log p(\boldsymbol{\pi}) + \log p(\mathbf{A}) \\ &\quad + \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\log p(\mathbf{x}_{1:T} | \hat{\mathbf{x}}_{1:T}, \mathbf{z}_{1:T}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V}) \right] \\ &\quad + \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\log p(\mathbf{z}_{1:T} | \boldsymbol{\pi}, \mathbf{A}) \right] + \sum_{k=1}^{K} \log p(\boldsymbol{\mu}_{k} | \boldsymbol{\Sigma}_{k}) p(\boldsymbol{\Sigma}_{k}) \\ &\quad + \sum_{k=1}^{K} \log p(\mathbf{U}_{k} | \mathbf{V}_{k}) p(\mathbf{V}_{k}) \\ &= \log p(\boldsymbol{\pi}) + \log p(\mathbf{A}) \\ &\quad + \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\log p(\mathbf{x}_{1} | \mathbf{z}_{1:T}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) \right] \\ &\quad + \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\log p(\mathbf{x}_{2:T} | \hat{\mathbf{x}}_{2:T}, \mathbf{z}_{2:T}, \mathbf{U}, \mathbf{V}) \right] \\ &\quad + \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\log p(\mathbf{z}_{1} | \boldsymbol{\pi}) p(\mathbf{z}_{2:T} | \mathbf{z}_{1}, \mathbf{A}) \right] \\ &\quad + \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\log p(\boldsymbol{\mu}_{k} | \boldsymbol{\Sigma}_{k}) p(\boldsymbol{\Sigma}_{k}) + \sum_{k=1}^{K} \log p(\mathbf{U}_{k} | \mathbf{V}_{k}) p(\mathbf{V}_{k}) \right] \end{split}$$

Continuation of deriving $\log q(\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathbf{U}, \mathbf{V})$

$$\log q(\boldsymbol{\pi}, \mathbf{A}, \mathbf{U}, \mathbf{V}) \propto \log Dir(\boldsymbol{\pi} | \boldsymbol{\omega}_{\mathbf{0}}^{(\boldsymbol{\pi})}) + \sum_{k=1}^{K} \log Dir(\mathbf{A} | \boldsymbol{\omega}_{0k}^{(\mathbf{A})}) \\ + \sum_{k=1}^{K} \mathbb{E}_{q(z_{1})} \left[\mathbb{I}(z_{1} = k) \right] \log \mathcal{N} \left(\mathbf{x}_{1} | \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}^{-1} \right) \\ + \sum_{t=2}^{T} \sum_{k=1}^{K} \mathbb{E}_{q(\mathbf{z}_{2:T})} \left[\mathbb{I}(z_{t} = k) \right] \log \mathcal{N} \left(\mathbf{x}_{t} | \mathbf{U}_{k} \hat{\mathbf{x}}_{t-1}, \mathbf{V}_{k}^{-1} \right) \\ + \sum_{k=1}^{K} \left\{ \mathbb{E}_{q(z_{1})} \left[\mathbb{I}(z_{1} = k) \right] \log \pi_{k} \right\} \\ + \sum_{t=2}^{T} \sum_{i=1}^{K} \sum_{j=1}^{K} \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\mathbb{I}(z_{t} = j, z_{t-1} = i) \right] \log a_{ij} \\ + \sum_{k=1}^{K} \log \mathcal{N} \left(\boldsymbol{\mu}_{k} | \boldsymbol{m}_{0}, (\beta_{0} \boldsymbol{\Sigma}_{k})^{-1} \right) + \sum_{k=1}^{K} \log \mathcal{W} \left(\boldsymbol{\Sigma}_{k} | \mathbf{W}_{0}, \boldsymbol{\nu}_{0} \right) \\ + \sum_{k=1}^{K} \log \mathcal{M} \mathcal{N} \left(\mathbf{U}_{k} | \mathbf{M}_{0}, \mathbf{V}_{k}^{-1}, \mathbf{K}_{0} \right) + \sum_{k=1}^{K} \log \mathcal{W} \left(\mathbf{V}_{k} | \mathbf{P}_{0}, \eta_{0} \right)$$

$$\begin{split} \log q(\mathbf{U}_{k}|\mathbf{V}_{k}) \propto &-\frac{1}{2}\mathbf{Tr}\left\{(\mathbf{U}_{k}-\mathbf{M}_{0})^{T}\mathbf{V}_{k}(\mathbf{U}_{k}-\mathbf{M}_{0})\mathbf{K}_{0}\right\} \\ &-\frac{1}{2}\sum_{t=2}^{T}\gamma_{t}(k)(\mathbf{x}_{t}-\mathbf{U}_{k}\hat{\mathbf{x}}_{t-1})^{T}\mathbf{V}_{k}(\mathbf{x}_{t}-\mathbf{U}_{k}\hat{\mathbf{x}}_{t-1}) \\ &=-\frac{1}{2}\mathbf{Tr}\left\{\mathbf{V}_{k}(\mathbf{U}_{k}-\mathbf{M}_{0})\mathbf{K}_{0}(\mathbf{U}_{k}-\mathbf{M}_{0})^{T}\right\} \\ &-\frac{1}{2}\sum_{t=2}^{T}\gamma_{t}(k)\mathbf{Tr}\left\{\mathbf{V}_{k}(\mathbf{x}_{t}-\mathbf{U}_{k}\hat{\mathbf{x}}_{t-1})(\mathbf{x}_{t}-\mathbf{U}_{k}\hat{\mathbf{x}}_{t-1})^{T}\right\} \\ &=-\frac{1}{2}\mathbf{Tr}\left\{\mathbf{V}_{k}(\mathbf{U}_{k}\mathbf{K}_{0}\mathbf{U}_{k}^{T}-\mathbf{U}_{k}\mathbf{K}_{0}\mathbf{M}_{0}^{T}-\mathbf{M}_{0}\mathbf{K}_{0}\mathbf{U}_{k}^{T}+\mathbf{M}_{0}\mathbf{K}_{0}\mathbf{M}_{0}^{T})\right\} \\ &-\frac{1}{2}\sum_{t=2}^{T}\gamma_{t}(k)\mathbf{Tr}\left\{\mathbf{V}_{k}(\mathbf{x}_{t}\mathbf{x}_{t}^{T}-\mathbf{x}_{t}\hat{\mathbf{x}}_{t-1}^{T}\mathbf{U}_{k}^{T}-\mathbf{U}_{k}\hat{\mathbf{x}}_{t-1}\mathbf{x}_{t}^{T}+\mathbf{U}_{k}\hat{\mathbf{x}}_{t-1}\hat{\mathbf{x}}_{t-1}^{T}\mathbf{U}_{k}^{T})\right\} \end{split}$$

$$-\frac{1}{2}\mathbf{Tr}\left\{\mathbf{V}_{k}\mathbf{U}_{k}\mathbf{K}_{k}\mathbf{U}_{k}^{T}\right\} = -\frac{1}{2}\mathbf{Tr}\left\{\mathbf{V}_{k}\mathbf{U}_{k}\left[\sum_{t=2}^{T}\mathbb{E}_{q}\left[\mathbb{I}(z_{t}=k)\right]\hat{\mathbf{x}}_{t-1}\hat{\mathbf{x}}_{t-1}^{T} + \mathbf{K}_{0}\right]\mathbf{U}_{k}^{T}\right\}$$
$$\mathbf{V}_{k}\mathbf{U}_{k}\mathbf{K}_{k}\mathbf{U}_{k}^{T} = \mathbf{V}_{k}\mathbf{U}_{k}\left[\sum_{t=2}^{T}\mathbb{E}_{q}\left[\mathbb{I}(z_{t}=k)\right]\hat{\mathbf{x}}_{t-1}\hat{\mathbf{x}}_{t-1}^{T} + \mathbf{K}_{0}\right]\mathbf{U}_{k}^{T}$$
$$\mathbf{K}_{k} = \sum_{t=2}^{T}\mathbb{E}_{q}\left[\mathbb{I}(z_{t}=k)\right]\hat{\mathbf{x}}_{t-1}\hat{\mathbf{x}}_{t-1}^{T} + \mathbf{K}_{0}$$

$$-\frac{1}{2} \operatorname{Tr} \left\{ \mathbf{V}_{k} \mathbf{U}_{k} \mathbf{M}_{k} \mathbf{K}_{k} \right\} = -\frac{1}{2} \operatorname{Tr} \left\{ \mathbf{V}_{k} \mathbf{U}_{k} \left[-\sum_{t=2}^{T} \mathbb{E}_{q} \left[\mathbb{I}(z_{t}=k) \right] \mathbf{x}_{t} \hat{\mathbf{x}}_{t-1}^{T} - \mathbf{M}_{0} \mathbf{K}_{0} \right] \right\}$$
$$\mathbf{V}_{k} \mathbf{U}_{k} \mathbf{M}_{k} \mathbf{K}_{k} = \mathbf{V}_{k} \mathbf{U}_{k} \left[\sum_{t=2}^{T} \mathbb{E}_{q} \left[\mathbb{I}(z_{t}=k) \right] \mathbf{x}_{t} \hat{\mathbf{x}}_{t-1}^{T} + \mathbf{M}_{0} \mathbf{K}_{0} \right]$$
$$\mathbf{M}_{k} \mathbf{K}_{k} = \left[\sum_{t=2}^{T} \mathbb{E}_{q} \left[\mathbb{I}(z_{t}=k) \right] \mathbf{x}_{t} \hat{\mathbf{x}}_{t-1}^{T} + \mathbf{M}_{0} \mathbf{K}_{0} \right]$$
$$\mathbf{M}_{k} = \left[\sum_{t=2}^{T} \mathbb{E}_{q} \left[\mathbb{I}(z_{t}=k) \right] \mathbf{x}_{t} \hat{\mathbf{x}}_{t-1}^{T} + \mathbf{M}_{0} \mathbf{K}_{0} \right]$$

$$\begin{split} \log q(\mathbf{V}_k) &= \log \mathcal{W}\left(\mathbf{U}_k | \mathbf{P}_k, \eta_k\right) \\ &= \log q(\mathbf{U}_k, \mathbf{V}_k) - \log q(\mathbf{U}_k | \mathbf{V}_k) \\ &= -\frac{1}{2} \mathbf{Tr} \left\{ (\mathbf{U}_k - \mathbf{M}_0)^T \mathbf{V}_k (\mathbf{U}_k - \mathbf{M}_0) \mathbf{K}_0 \right\} - \frac{m}{2} \log(|\mathbf{V}_k|) \\ &- \frac{1}{2} \mathbf{Tr} \left\{ \mathbf{V}_k \mathbf{W}_0^{-1} \right\} + \frac{(\eta_0 - D - 1)}{2} \log |\mathbf{V}_k| \\ &- \frac{1}{2} \sum_{t=2}^T \mathbb{E}_{q(\mathbf{z}_{2:T})} \left[\mathbb{I}(z_t = k) \right] (\mathbf{x}_t - \mathbf{U}_k \hat{\mathbf{x}}_{t-1})^T \mathbf{V}_k (\mathbf{x}_t - \mathbf{U}_k \hat{\mathbf{x}}_{t-1}) \\ &+ \frac{1}{2} \left(\sum_{t=2}^T \mathbb{E}_{q(\mathbf{z}_{2:T})} \left[\mathbb{I}(z_t = k) \right] \right) \log |\mathbf{V}_k| \\ &+ \frac{m}{2} \log |\mathbf{V}_k| + \frac{1}{2} \mathbf{Tr} \left\{ (\mathbf{U}_k - \mathbf{M}_k)^T \mathbf{V}_k (\mathbf{U}_k - \mathbf{M}_k) \mathbf{K}_k \right\} \\ &\propto \frac{(\eta_k - D - 1)}{2} \log |\mathbf{V}_k| - \frac{1}{2} \mathbf{Tr} \left\{ \mathbf{V}_k \mathbf{P}_k^{-1} \right\} \end{split}$$

$$\mathbf{P}_{k}^{-1} = \sum_{t=2}^{T} \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\mathbb{I}(z_{t}=k) \right] (\mathbf{x}_{t} - \mathbf{U}_{k} \hat{\mathbf{x}}_{t-1}) (\mathbf{x}_{t} - \mathbf{U}_{k} \hat{\mathbf{x}}_{t-1})^{T} + (\mathbf{U}_{k} - \mathbf{M}_{0}) \mathbf{K}_{0} (\mathbf{U}_{k} - \mathbf{M}_{0})^{T} + \mathbf{P}_{0}^{-1} - (\mathbf{U}_{k} - \mathbf{M}_{k}) \mathbf{K}_{k} (\mathbf{U}_{k} - \mathbf{M}_{k})^{T} = \mathbf{P}_{0}^{-1} + \mathbf{M}_{0} \mathbf{K}_{0} \mathbf{M}_{0}^{T} + \sum_{t=2}^{T} \gamma_{t}(k) \mathbf{x}_{t} \mathbf{x}_{t}^{T} - \mathbf{M}_{k} \mathbf{K}_{k} \mathbf{M}_{k}^{T}$$

$$\frac{(\eta_k - D - 1)}{2} \log |\mathbf{V}_k| = \frac{(\eta_0 - D - 1)}{2} \log |\mathbf{V}_k| + \frac{1}{2} \left(\sum_{t=1}^T \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\mathbb{I}(z_t = k) \right] \right) \log |\mathbf{V}_k|$$
$$\eta_k = \eta_0 + \sum_{t=1}^T \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\mathbb{I}(z_t = k) \right]$$

C. Distributions

C.1. Dirichlet distribution

Notation and parameters

$$\boldsymbol{\pi} \sim Dir(\boldsymbol{\alpha})$$
$$\boldsymbol{\alpha} = \{\alpha_1, ..., \alpha_k\}$$
$$\alpha_j > 0; \ \alpha_0 = \sum_{j=1}^K \alpha_j$$

Density function

$$p(\boldsymbol{\pi}|\boldsymbol{\alpha}) = C(\boldsymbol{\alpha})\pi_1^{\alpha_1-1}...\pi_k^{\alpha_k-1}$$
$$\pi_1, ..., \pi_k \ge 0; \sum_{j=1}^K \pi_j = 1$$
$$C(\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1)...\Gamma(\alpha_k)}$$

Expectations

$$\mathbb{E}[\boldsymbol{\pi}] = \boldsymbol{\alpha}/\alpha_0$$
$$\log \mathbb{E}[\pi_j] = \psi(\alpha_j) - \psi(\alpha_0)$$

79

KL-divergence

$$D_{\mathrm{KL}}\left(\tilde{\boldsymbol{\alpha}}||\boldsymbol{\alpha}\right) = \ln\frac{\Gamma(\tilde{\alpha}_{0})}{\Gamma(\alpha_{0})} - \sum_{j=1}^{K} \left[\ln\frac{\Gamma(\tilde{\alpha}_{j})}{\Gamma(\alpha_{j})} - (\tilde{\alpha}_{j} - \alpha_{j})(\psi(\tilde{\alpha}_{j}) - \psi(\tilde{\alpha}_{0}))\right]$$

Entropy

$$H[\boldsymbol{\pi}] = -\sum_{k=1}^{K} (\alpha_k - 1)(\psi(\alpha_k) - \psi(\alpha_0)) - \ln C(\boldsymbol{\alpha})$$

C.2. Uniform distibution

Notation and parameters

$$x \sim U(a, b)$$

boundaries a, b
with $b > a$

Density function

$$p(x|a,b) = \frac{1}{b-a}, \ x \in [a,b]$$

Expectations

$$\mathbb{E}\left[x\right] = \frac{a+b}{2}$$

Entropy

$$H\left[x\right] = \ln(b-a)$$

C.3. Multivariate normal distribution

Notation and parameters

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

 $\boldsymbol{\mu}$ mean vector
 $\boldsymbol{\Sigma}$ covariance matrix
 $\boldsymbol{\Sigma}^{-1}$ precision matrix

Density function

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-D/2} |\boldsymbol{\Sigma}|^{-1/2} \exp\left\{-\frac{1}{2}\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right)\right\}$$

Expectations

$$\mathbb{E}\left[\mathbf{x}
ight] = oldsymbol{\mu}$$

 $\mathbb{E}\left[\mathbf{x}\mathbf{x}^T
ight] = oldsymbol{\Sigma}$

KL-divergence

$$D_{\mathrm{KL}}\left(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}} || \boldsymbol{\mu}, \boldsymbol{\Sigma}\right) = -\frac{1}{2} \left(\ln |\tilde{\boldsymbol{\Sigma}} \boldsymbol{\Sigma}^{-1}| + \mathrm{Tr} \left\{ I - \left[\tilde{\boldsymbol{\Sigma}} + (\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu})(\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu})^T \right] \boldsymbol{\Sigma}^{-1} \right\} \ln e \right)$$

Entropy

$$H\left[\mathbf{x}\right] = \frac{1}{2}\ln|\mathbf{\Sigma}| + \frac{D}{2}(1+\ln(2\pi))$$

C.4. Wishart distribution

Notation and parameters

$$\Sigma \sim W(\mathbf{W}, \nu)$$

W precision matrix
 ν degree of freedom

Density function

$$p(\boldsymbol{\Sigma}|\mathbf{W},\nu) = B(\mathbf{W},\nu)|\boldsymbol{\Sigma}|^{(\nu-D-1)/2} \exp\left(-\frac{1}{2}\mathbf{Tr}\left\{\mathbf{W}^{-1}\boldsymbol{\Sigma}\right\}\right)$$
$$B(\mathbf{W},\nu) = |\mathbf{W}|^{-\nu/2} \left(2^{\nu D/2}\pi^{D(D-1)/4}\prod_{i=1}^{D}\Gamma\left(\frac{\nu+1-i}{2}\right)\right)^{-1}$$

Expectations

$$\mathbb{E}\left[\mathbf{\Sigma}\right] = \mathbf{W}\nu$$
$$\mathbb{E}\left[\ln|\mathbf{\Sigma}|\right] = \sum_{i=1}^{D} \psi\left(\frac{\nu+1-i}{2}\right) + D\ln 2 + \ln|\mathbf{W}|$$

KL-divergence

$$D_{\mathrm{KL}}\left(\tilde{\mathbf{W}}, \tilde{\nu} || \mathbf{W}, \nu\right) = \ln \frac{B(\mathbf{W}, \nu)}{\tilde{\mathbf{W}}, \tilde{\nu}} + \frac{\tilde{\nu} - \nu}{2} \mathbb{E}\left[\ln |\mathbf{\Sigma}|\right] + \frac{1}{2} \tilde{\nu} \mathbf{Tr} \left\{ \mathbf{W}^{-1} \tilde{\mathbf{W}} - I \right\}$$

Entropy

$$H[\mathbf{\Sigma}] = -\ln B(\mathbf{W}, \nu) - \frac{\nu - D - 1}{2} \mathbb{E}[|\mathbf{\Sigma}|] + \frac{\nu D}{2}$$