

A Supplementary Material

This text is the Supplementary Material of the paper “*Approximate inference in latent Gaussian-Markov models from continuous time observations*” by B. Cseke, M. Oppen and G. Sanguinetti (*Neural Information Processing Systems* 2013).

A.1 Variational formulation using expectation constraints

In this section we formulate an expectation constraints based approximate inference scheme [Heskes et al., 2005] for our model. It turns out that, when only discrete time observation are present, the inference results in an expectation propagation type algorithm whereas when only continuous type observations are present it collapses to the variational approach. These two approaches can be combined into a joint inference scheme.

A.1.1 Discrete time observations

In the case of only discrete time observations no time discretisation is needed for the formal manipulation of the distributions. The propagation algorithm we arrive to can be viewed as an EP on a latent Gaussian model where the (partial) matrix inversion for computing marginal means and variances is replaced by solving the forward-backward differential equations.

In the following we present an expectation constrained free energy optimisation that leads to the EP algorithm. Here we use the concept of the variational formulation by free energies [e.g. Yedidia et al., 2000]. In a similar spirit as in Heskes et al. [2005], instead of approximating

$$p(\{\mathbf{x}_t\}|\{\mathbf{y}_{t_i}^d\}_i) \propto p_0(\{\mathbf{x}_t\}) \times \prod_i p(\mathbf{y}_{t_i}^d|\mathbf{x}_{t_i})$$

with an OU process we define the free energy as function of a family of approximate marginals $\mathcal{Q} = \{q_0(\{\mathbf{x}_t\}), \{q_{ds}^i(\mathbf{x}_{t_i})\}_i, \{q_d^i(\mathbf{x}_{t_i})\}_i\}$ constrained by expectation constraints. The densities q_d^i will be assigned to the factors corresponding to the likelihood terms, q_0 will assigned to the factor p_0 . With some abuse in notation the family \mathcal{Q} can be viewed as representing an approximation q having the form

$$q(\{\mathbf{x}_t\}) \propto \frac{q_0(\{\mathbf{x}_t\}) \prod_i q_d^i(\mathbf{x}_{t_i})}{\prod_i q_{ds}^i(\mathbf{x}_{t_i})}.$$

Note that in the graphical model formalism the densities q_{ds}^i correspond to the densities defined over the variables of the separator sets in graphical models [Lauritzen, 1996]. The expectation constraints are defined over the function $\mathbf{f}(\mathbf{z}) = (\mathbf{z}, -\mathbf{z}\mathbf{z}^T/2)$ and ensure that the corresponding marginals of the members of \mathcal{Q} are consistent up to second moments, i.e., their marginal means and covariances are equal. Given the above assumptions, one can define an approximation of the $D[q||p]$ divergence called free energy that reads as

$$\begin{aligned} F(\mathcal{Q}) = & -\langle \log p_0(\{\mathbf{x}_t\}) \rangle_{q_0} - \sum_i \left\langle \log p(\mathbf{y}_{t_i}^d|\mathbf{x}_{t_i}) \right\rangle_{q_d^i} \\ & + \langle \log q_0(\{\mathbf{x}_t\}) \rangle_{q_0} + \sum_i \left[\left\langle \log q_d^i(\mathbf{x}_{t_i}) \right\rangle_{q_d^i} - \left\langle \log q_{ds}^i(\mathbf{x}_{t_i}) \right\rangle_{q_{ds}^i} \right] \end{aligned} \quad (16)$$

and specify the expectation constraints

$$\langle \mathbf{f}(\mathbf{x}_{t_i}) \rangle_{q_0} = \langle \mathbf{f}(\mathbf{x}_{t_i}) \rangle_{q_{ds}^i} \quad \text{and} \quad \langle \mathbf{f}(\mathbf{x}_{t_i}) \rangle_{q_d^i} = \langle \mathbf{f}(\mathbf{x}_{t_i}) \rangle_{q_{ds}^i} \quad \text{for all } t_i \in T_d.$$

The stationary equations of the corresponding Lagrangian

$$L(\mathcal{Q}, \Lambda) = F(\mathcal{Q}) + \sum_i \left[\lambda_{t_i}^0 \cdot [\langle \mathbf{f}(\mathbf{x}_{t_i}) \rangle_{q_{ds}^i} - \langle \mathbf{f}(\mathbf{x}_{t_i}) \rangle_{q_0}] + \lambda_{t_i}^d \cdot [\langle \mathbf{f}(\mathbf{x}_{t_i}) \rangle_{q_d^i} - \langle \mathbf{f}(\mathbf{x}_{t_i}) \rangle_{q_{ds}^i}] \right] \quad (17)$$

result in $q_0(\{\mathbf{x}_t\}) \propto p_0(\{\mathbf{x}_t\}) \times \exp(\sum_i \lambda_{t_i}^0 \cdot \mathbf{f}(\mathbf{x}_{t_i}))$, $q_d^i(\mathbf{x}_{t_i}) \propto p(\mathbf{y}_{t_i}^d|\mathbf{x}_{t_i}) \times \exp(\lambda_{t_i}^d \cdot \mathbf{f}(\mathbf{x}_{t_i}))$ and $q_{ds}^i(\mathbf{x}_{t_i}) \propto \exp([\lambda_{t_i}^0 + \lambda_{t_i}^d] \cdot \mathbf{f}(\mathbf{x}_{t_i}))$. The differential w.r.t. the Lagrange multipliers $\lambda_{t_i}^0$ and $\lambda_{t_i}^d$ lead to the above mentioned expectation constraints. Since the expectation constraints are defined over the sufficient statistics $\mathbf{f}(\mathbf{z})$ and the optimal q_{ds}^i belongs to the exponential (Gaussian) family defined by \mathbf{f} , we can rewrite these constraints into canonical forms. These read as

$$\lambda_{t_i}^d + \lambda_{t_i}^0 = \text{Collapse}(q_0(\mathbf{x}_{t_i}); \mathbf{f}) \quad \text{and} \quad \lambda_{t_i}^d + \lambda_{t_i}^0 = \text{Collapse}(q_d^i(\mathbf{x}_{t_i}); \mathbf{f}).$$

Here $\text{Collapse}(q(\mathbf{z}); \mathbf{f})$ denotes the (unique) moment matching canonical parameters, in other words the Kullback-Leibler projection $\text{Collapse}(q(\mathbf{z}); \mathbf{f}) = \text{argmin}_{\boldsymbol{\theta}} D[q(\mathbf{z}) \parallel \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{z}) - \log Z(\boldsymbol{\theta}))]$. From the moment matching constraints one can introduce the fixed point iteration

$$[\boldsymbol{\lambda}_{t_i}^0]^{new} = \text{Collapse}(q_d^i(\mathbf{x}_{t_i}); \mathbf{f}) - \boldsymbol{\lambda}_{t_i}^d, \quad (18)$$

$$[\boldsymbol{\lambda}_{t_i}^d]^{new} = \text{Collapse}(q_0(\mathbf{x}_{t_i}); \mathbf{f}) - \boldsymbol{\lambda}_{t_i}^0, \quad (19)$$

which corresponds to a (parallel) EP algorithm in a latent Gaussian model where the latent Gaussian is given by an OU process. It can be shown that the free energy in (16) is finite, details are given in Section A.2. The collapse operation $\text{Collapse}(q_d^i(\mathbf{x}_{t_i}); \mathbf{f})$ is computed by moment matching—computing the corresponding canonical parameters—while $\text{Collapse}(q_0(\mathbf{x}_{t_i}); \mathbf{f})$ is computed by using the moment parameters of the marginals resulting from the forward-backward equations in Section B.1, with $\boldsymbol{\lambda}_{t_i}^0 = (\mathbf{h}_{t_i}^d, \mathbf{Q}_{t_i}^d)$. Readers familiar with the EP presented in Opper and Winther [2000] or [Minka, 2001] can identify the multipliers $\boldsymbol{\lambda}_{t_i}^0$ as the canonical parameters of the so called term approximations whereas the $\boldsymbol{\lambda}_{t_i}^d$ s correspond to the canonical parameters of the so called cavity distributions. Equations (18) and (19) correspond to the updates of the term approximation and the cavity distribution through moment matching. The free energy approach presented above starts from the variational formulation $D[q \parallel p]$ where instead of single specially chosen Gaussian q , a family of approximate marginals \mathcal{Q} is introduced. The EP style iterative moment matching minimizations in the $D[p \parallel \cdot]$ sense corresponding to $\text{Collapse}(\cdot; \mathbf{f})$, arise from the satisfaction of moment matching constraints.

Clearly, any method that computes the marginal means and covariances of q_0 at the time-point in T_d suffices to keep the iteration running, and thus, when possible, one should solve the differential equations between observation points analytically. We can also opt for the alternative generic approach of computing the covariance matrix corresponding to the variables $\{\mathbf{x}_{t_i}\}_{t_i \in T_d}$ and opt for the equivalent Gaussian process (OU covariance function) expectation propagation in Opper and Winther [2000] or Minka [2001]. In the latter case the marginal means and variances for $t \notin T_d$ can be computed by using the conditional independencies in the model and computing the predictive distributions.

A.1.2 Continuous time observations

In this section we extend the approach to the case when only continuous time observations are present. The task is to approximate a posterior distribution having the form

$$p(\{\mathbf{x}_t\} | \{\mathbf{y}_t^c\}) \propto p(\{\mathbf{x}_t\}) \times \exp \left\{ - \int_0^1 dt V(t, \mathbf{y}_t^c, \mathbf{x}_t) \right\}.$$

In order to simplify notation, we will omit the dependence of V on \mathbf{y}_t^c . In an similar fashion as in the previous section we introduce a family of marginals $\mathcal{Q} = \{q_0(\{\mathbf{x}_t\}), q_{cs}(\{\mathbf{x}_t\}), q_c(\{\mathbf{x}_t\})\}$ and define the free energy as

$$F(\mathcal{Q}) = - \langle \log p_0(\{\mathbf{x}_t\}) \rangle_{q_0} + \left\langle \int_0^1 dt V(t, \mathbf{x}_t) \right\rangle_{q_c} + \langle \log q_0(\{\mathbf{x}_t\}) \rangle_{q_0} + \langle \log q_c(\{\mathbf{x}_t\}) \rangle_{q_c} - \langle \log q_{cs}(\{\mathbf{x}_t\}) \rangle_{q_{cs}} \quad (20)$$

The moment matching constraints will be defined as

$$\langle \mathbf{f}(\mathbf{x}_t) \rangle_{q_0} = \langle \mathbf{f}(\mathbf{x}_t) \rangle_{q_{cs}} \quad \text{and} \quad \langle \mathbf{f}(\mathbf{x}_t) \rangle_{q_c} = \langle \mathbf{f}(\mathbf{x}_t) \rangle_{q_{cs}} \quad \text{for all } t \in [0, 1].$$

which imply using Lagrange multiplier terms of the form

$$C(\mathcal{Q}, \mathcal{M}) = \int_0^1 dt \boldsymbol{\mu}_t^0 \cdot [\langle \mathbf{f}(\mathbf{x}_t) \rangle_{q_{cs}} - \langle \mathbf{f}(\mathbf{x}_t) \rangle_{q_0}] + \int_0^1 dt \boldsymbol{\mu}_t^c \cdot [\langle \mathbf{f}(\mathbf{x}_t) \rangle_{q_{cs}} - \langle \mathbf{f}(\mathbf{x}_t) \rangle_{q_c}]. \quad (21)$$

In order to carry out the computations and show that the above quantities exist, we discretise the time domain by using the time-points $T = \{t_0 = 0, t_1, \dots, t_{K-1}, t_K = 1\}$ with the lags $\Delta t_k = t_{k+1} - t_k$ and represent the process $\{\mathbf{x}_t\}_t$ by the matrix $\mathbf{x} = [\mathbf{x}_{t_0}, \dots, \mathbf{x}_{t_K}]$. By using this discretisation, we approximate all integrals using the corresponding Euler discretisation and, view $p_0(\mathbf{x})$ as a multivariate Gaussian. We use the indexing T to highlight the discretisation. We define the Lagrangian as

$$L(\mathcal{Q}_T, \mathcal{M}_T) = F(\mathcal{Q}_T) + C(\mathcal{Q}_T, \mathcal{M}_T). \quad (22)$$

The stationary conditions (22) corresponding to the differentiation w.r.t q_0 , q_{cs} and q_c , result in

$$q_0(\mathbf{x}) \propto p_0(\mathbf{x}) \times \exp \left\{ \sum_k \Delta t_k \boldsymbol{\mu}_{t_k}^0 \cdot \mathbf{f}(\mathbf{x}_{t_k}) \right\}, \quad (23)$$

$$q_c(\mathbf{x}) \propto \exp \left\{ \sum_k \Delta t_k [-V(t_k, \mathbf{x}_{t_k}) + \boldsymbol{\mu}_{t_k}^c \cdot \mathbf{f}(\mathbf{x}_{t_k})] \right\}, \quad (24)$$

$$q_{cs}(\mathbf{x}) \propto \exp \left\{ \sum_k \Delta t_k [\boldsymbol{\mu}_{t_k}^0 + \boldsymbol{\mu}_{t_k}^c] \cdot \mathbf{f}(\mathbf{x}_{t_k}) \right\}. \quad (25)$$

Due to the factorisation of q_{cs} and q_c and the Gaussian nature of $q_{cs}(\mathbf{x}_{t_k})$, the stationary conditions corresponding to the moment constraints can be rewritten as

$$\Delta t_k [\boldsymbol{\mu}_{t_k}^0 + \boldsymbol{\mu}_{t_k}^c] = \text{Collapse}(q_0(\mathbf{x}_{t_k}); \mathbf{f}) \quad \text{and} \quad \Delta t_k [\boldsymbol{\mu}_{t_k}^0 + \boldsymbol{\mu}_{t_k}^c] = \text{Collapse}(q_c(\mathbf{x}_{t_k}); \mathbf{f}) \quad \text{for all, } t_k \in T. \quad (26)$$

Taking the limit $\Delta t_k \rightarrow 0$ is not feasible at this point because the marginals of both q_{cs} and q_c collapse into delta distributions. However, we can observe that $\text{Collapse}(q_0(\mathbf{x}_{t_k}); \mathbf{f})$ should always be finite and well defined. We use the alias $\boldsymbol{\mu}_{t_k} = \text{Collapse}(q_0(\mathbf{x}_{t_k}); \mathbf{f})$ and we eliminate $\boldsymbol{\mu}_{t_k}^c$ from the formulae above. In an similar spirit as in Section A.1.1, we use the moment matching constraint to define the fixed point iteration

$$[\boldsymbol{\mu}_{t_k}^0]^{new} = \boldsymbol{\mu}_{t_k}^0 + \frac{1}{\Delta t_k} [\text{Collapse}(q_c(\mathbf{x}_{t_k}); \mathbf{f}) - \boldsymbol{\mu}_{t_k}], \quad (27)$$

where, due to eliminating $\boldsymbol{\mu}_{t_k}^c$, we have $q_c(\mathbf{x}_{t_k}) \propto \exp \{-\Delta t_k [V(t_k, \mathbf{x}_{t_k}) + \boldsymbol{\mu}_{t_k}^0 \cdot \mathbf{f}(\mathbf{x}_{t_k})] + \boldsymbol{\mu}_{t_k} \cdot \mathbf{f}(\mathbf{x}_{t_k})\}$. The $\Delta t_k \rightarrow 0$ limit is presented in Section 2.2.3 of the paper and by using $\boldsymbol{\mu}_t^0 = (\mathbf{h}_t^c, \mathbf{Q}_t^c)$, it results in the updates

$$[\mathbf{h}_t^c]^{new} = -\partial_{\mathbf{m}_t} \langle V(t, \mathbf{x}_t) \rangle_{q_0(\mathbf{x}_t)} + 2\partial_{\mathbf{V}_t} \langle V(t, \mathbf{x}_t) \rangle_{q_0(\mathbf{x}_t)} \mathbf{m}_t \quad \text{and} \quad [\mathbf{Q}_t^c]^{new} = \partial_{\mathbf{V}_t} \langle V(t, \mathbf{x}_t) \rangle_{q_0(\mathbf{x}_t)} \quad (28)$$

$$[q_0(\{\mathbf{x}_t\})]^{new} \propto p_0(\{\mathbf{x}_t\}) \times \exp \left\{ \int_0^1 dt [\mathbf{x}_t^T \mathbf{h}_t^c - \frac{1}{2} \mathbf{x}_t^T \mathbf{Q}_t^c \mathbf{x}_t] \right\}$$

where the marginal moments of $q_0(\{\mathbf{x}_t\})$ are computed by using the Kalman-Bucy algorithm (Section B.1).

A.1.3 The joint approximation scheme

Now that we have derived the approximation scheme for both discrete and continuous time observations, we can show that they can be easily combined to obtain a joint approximation. Without loss of generality, we can assume that $T_d \subset T$. By defining the joint family as $\mathcal{Q} = \{q_0(\mathbf{x}), \{q_d^i(\mathbf{x}_{t_i})\}_i, \{q_{ds}^i(\mathbf{x}_{t_i})\}_i, q_c(\mathbf{x})q_{cs}(\mathbf{x}), q_c(\mathbf{x})\}$ and the free energy as

$$F(\mathcal{Q}) = -\langle \log p_0(\mathbf{x}) \rangle_{q_0} - \sum_i \left\langle \log p(\mathbf{y}_{t_i}^d | \mathbf{x}_{t_i}) \right\rangle_{q_d^i} + \left\langle \int_0^1 dt V(t, \mathbf{x}_t) \right\rangle_{q_c} \quad (29)$$

$$+ \langle \log q_0(\mathbf{x}) \rangle_{q_0} + \sum_i [\langle \log q_d^i(\mathbf{x}_{t_i}) \rangle_{q_d^i} - \langle \log q_{ds}^i(\mathbf{x}_{t_i}) \rangle_{q_{ds}^i}] + \langle \log q_c(\mathbf{x}) \rangle_{q_c} - \langle \log q_{cs}(\mathbf{x}) \rangle_{q_{cs}}$$

we can construct the Lagrangian by using the multiplier terms from (17) and the Euler discretisation of (21). The fixed point iteration follows (18), (19) and (28) where q_0 is defined by

$$[q_0(\{\mathbf{x}_t\})]^{new} \propto p_0(\{\mathbf{x}_t\}) \times \exp \left\{ \sum_i [\mathbf{x}_{t_i}^T \mathbf{h}_{t_i}^d - \frac{1}{2} \mathbf{x}_{t_i}^T \mathbf{Q}_{t_i}^d \mathbf{x}_{t_i}] + \int_0^1 dt [\mathbf{x}_t^T \mathbf{h}_t^c - \frac{1}{2} \mathbf{x}_t^T \mathbf{Q}_t^c \mathbf{x}_t] \right\}. \quad (30)$$

We use the forward-backward equations (Section B) to compute the marginals of q_0 . As a result we have an algorithm behaves like a EP/variational hybrid: the parameters corresponding to the discrete time observations follow an EP style update (18) and (19), while the ones corresponding to the continuous observations are updated in a variational fashion according to (28).

A.2 The computation of the free energy

In the following we show that the free energy exist when $\Delta t_k \rightarrow 0$.

A.2.1 Continuous time observations

The expression of the free energy in (20) after discretisation is

$$F(\mathcal{Q}_T) = -\langle \log p_0(\mathbf{x}) \rangle_{q_0} + \sum_k \Delta t_k \langle V(t_k, \mathbf{x}_{t_k}) \rangle_{q_c} + \langle \log q_0(\mathbf{x}) \rangle_{q_0} + \langle \log q_c(\mathbf{x}) \rangle_{q_c} - \langle \log q_{cs}(\mathbf{x}) \rangle_{q_{cs}}$$

by substituting (23), (24) and (25) into $F(\mathcal{Q}_T)$ we find that

$$F(\mathcal{Q}_T) = -\log Z_0(\{\boldsymbol{\mu}_{t_k}^0\}) - \log Z_c(\{\boldsymbol{\mu}_{t_k}^c\}) + \log Z_{cs}(\{\boldsymbol{\mu}_{t_k}^0 + \boldsymbol{\mu}_{t_k}^c\})$$

where Z_0 , Z_c and Z_{cs} stand for the corresponding normalisation constants. By using the Legendre duality we can write

$$-\log Z_0(\{\boldsymbol{\mu}_{t_k}^0\}) = D[q_0(\mathbf{x})||p_0(\mathbf{x})] - \sum_k \Delta t_k \boldsymbol{\mu}_{t_k}^0 \cdot \mathbf{f}(\mathbf{x}_{t_k})$$

and by using the expansion in (11) in Section 2.2.3 of the paper, we find that

$$-\log Z_c(\{\boldsymbol{\mu}_{t_k}^c\}) + \log Z_{cs}(\{\boldsymbol{\mu}_{t_k}^0 + \boldsymbol{\mu}_{t_k}^c\}) \simeq \sum_k \Delta t_k \langle V(t_k, \mathbf{x}_{t_k}) + \boldsymbol{\mu}_{t_k}^0 \cdot \mathbf{f}(\mathbf{x}_{t_k}) \rangle_{q_0}.$$

Since q_0 corresponds to an OU process (see Section B.3 for its parametric form), we can take the limit $\Delta t_k \rightarrow 0$ and obtain

$$\begin{aligned} F(\mathcal{Q}) &= D[q_0(\{\mathbf{x}_t\})||p_0(\{\mathbf{x}_t\})] + \int_0^1 dt \langle V(t, \mathbf{x}_t) \rangle_{q_0} \\ &= \frac{1}{2} \int_0^1 dt \left\langle [(\mathbf{A}_t - \mathbf{A}_t^q)\mathbf{x}_t + (\mathbf{c}_t - \mathbf{c}_t^q)]^T \mathbf{B}_t^{-1} [(\mathbf{A}_t - \mathbf{A}_t^q)\mathbf{x}_t + (\mathbf{c}_t - \mathbf{c}_t^q)] \right\rangle_{q_0} \\ &\quad + D[p_0(\mathbf{x}_0)||q_0(\mathbf{x}_0)] + \int_0^1 dt \langle V(t, \mathbf{x}_t) \rangle_{q_0}, \end{aligned}$$

where \mathbf{A}_t^q and \mathbf{c}_t^q represent the parameters corresponding to q_0 . Computing $D[q_0(\{\mathbf{x}_t\})||p_0(\{\mathbf{x}_t\})]$ when q_0 and p_0 are parameterised OU processes can be done as in [e.g. Archambeau et al., 2007].

A.2.2 Discrete time observations

We use the notation

$$\begin{aligned} q_d^i(\mathbf{x}_{t_i}) &= \frac{1}{Z_d^i} p(\mathbf{y}_i^d|\mathbf{x}_{t_i}) \times \exp(\boldsymbol{\lambda}_i^d \cdot \mathbf{f}(\mathbf{x}_{t_i})), \\ q_{ds}^i(\mathbf{x}_{t_i}) &= \frac{1}{Z_{ds}^i} \exp([\boldsymbol{\lambda}_i^0 + \boldsymbol{\lambda}_i^d] \cdot \mathbf{f}(\mathbf{x}_{t_i})) \end{aligned}$$

and by using the Legendre duality as above, we can rewrite (16) as

$$F(\mathcal{Q}) = D[q_0(\{\mathbf{x}_t\})||p_0(\{\mathbf{x}_t\})] - \sum_i \boldsymbol{\lambda}_i^0 \cdot \langle \mathbf{f}(\mathbf{x}_{t_i}) \rangle_{q_{ds}^i} - \sum_i [\log Z_d^i - \log Z_{ds}^i],$$

which is a finite, computable quantity. The joint free energy follows from combining the discrete and continuous free energies according to (29), that is,

$$\begin{aligned} F(\mathcal{Q}) &= D[p_0(\mathbf{x}_0)||q_0(\mathbf{x}_0)] + \frac{1}{2} \int_0^1 dt \left\langle [(\mathbf{A}_t - \mathbf{A}_t^q)\mathbf{x}_t + (\mathbf{c}_t - \mathbf{c}_t^q)]^T \mathbf{B}_t^{-1} [(\mathbf{A}_t - \mathbf{A}_t^q)\mathbf{x}_t + (\mathbf{c}_t - \mathbf{c}_t^q)] \right\rangle_{q_0} \\ &\quad + \int_0^1 dt \langle V(t, \mathbf{x}_t) \rangle_{q_0} - \sum_i \boldsymbol{\lambda}_i^0 \cdot \langle \mathbf{f}(\mathbf{x}_{t_i}) \rangle_{q_{ds}^i} - \sum_i [\log Z_d^i - \log Z_{ds}^i]. \end{aligned}$$

Note that after convergence we have $q_{ds}^i(\mathbf{x}_t) = q_0(\mathbf{x}_t)$.

B Computations related to the Kalman-Bucy forward-backward algorithm

This section contains the computations that complement the material presented in Section 2.1. For reasons of simplicity, in Sections B.2 and B.3 we focus on the continuous time case, the computations related to the additional discrete time terms follow naturally.

B.1 The Kalman-Bucy forward-backward equations

By using the Euler discretisation and first order expansions as in [e.g. Särkkä, 2006] one can show that the forward and backward filtering equations satisfy

$$\begin{aligned}\frac{d}{dt}\mathbf{V}_t^{fw} &= \mathbf{A}_t\mathbf{V}_t^{fw} + \mathbf{V}_t^{fw}\mathbf{A}_t^T + \mathbf{B}_t - \mathbf{V}_t^{fw}\mathbf{Q}_t^c\mathbf{V}_t^{fw}, & \mathbf{m}_{t_i+}^{fw} &= (\mathbf{I} + \mathbf{V}_{t_i}^{fw}\mathbf{Q}_{t_i}^d)^{-1}(\mathbf{m}_{t_i}^{fw} + \mathbf{V}_{t_i}^{fw}\mathbf{h}_{t_i}^d), \\ \frac{d}{dt}\mathbf{m}_t^{fw} &= \mathbf{A}_t\mathbf{m}_t^{fw} + \mathbf{c}_t + \mathbf{V}_t^{fw}[\mathbf{h}_t^c - \mathbf{Q}_t^c\mathbf{m}_t^{fw}], & \mathbf{V}_{t_i+}^{fw} &= (\mathbf{I} + \mathbf{V}_{t_i}^{fw}\mathbf{Q}_{t_i}^d)^{-1}\mathbf{V}_{t_i}^{fw}, \\ \frac{d}{dt}\mathbf{V}_t^{bw} &= \mathbf{A}_t\mathbf{V}_t^{bw} + \mathbf{V}_t^{bw}\mathbf{A}_t^T - \mathbf{B}_t + \mathbf{V}_t^{bw}\mathbf{Q}_t^c\mathbf{V}_t^{bw}, & \mathbf{m}_{t_i-}^{bw} &= (\mathbf{I} + \mathbf{V}_{t_i}^{bw}\mathbf{Q}_{t_i}^d)^{-1}(\mathbf{m}_{t_i}^{bw} + \mathbf{V}_{t_i}^{bw}\mathbf{h}_{t_i}^d), \\ \frac{d}{dt}\mathbf{m}_t^{bw} &= \mathbf{A}_t\mathbf{m}_t^{bw} + \mathbf{c}_t - \mathbf{V}_t^{bw}[\mathbf{h}_t^c - \mathbf{Q}_t^c\mathbf{m}_t^{bw}], & \mathbf{V}_{t_i-}^{bw} &= (\mathbf{I} + \mathbf{V}_{t_i}^{bw}\mathbf{Q}_{t_i}^d)^{-1}\mathbf{V}_{t_i}^{bw}.\end{aligned}$$

We solve the equations for \mathbf{m}^{fw} and \mathbf{V}^{fw} in a forward fashion using the initial conditions $N(\mathbf{x}_0; \mathbf{m}_0, \mathbf{V}_0)$, whereas the equations for \mathbf{m}^{bw} and \mathbf{V}^{bw} are solved in a backwards with the initial, or more specifically, the end conditions given by the a non-informative Gaussian. By combining the forward and backward solutions we obtain the posterior marginal density $p(\mathbf{x}_t | \{\mathbf{y}_i^d\}_i, \{\mathbf{y}_i^c\}) \propto N(\mathbf{x}_t; \mathbf{m}_t^{fw}, \mathbf{V}_t^{fw}) \times N(\mathbf{x}_t; \mathbf{m}_t^{bw}, \mathbf{V}_t^{bw})$. Note that the backward equations are often replaced by the so called smoothing equations [e.g. Särkkä, 2006]. These combine the backward equations and the latter computation of the marginals into a pair of differential equations for the mean and the covariance respectively. In some cases one is better off with computing directly the inverse of \mathbf{V}_t^{fw} or \mathbf{V}_t^{bw} , these also follow similar quadratic or linear differential equations as the ones above.

B.2 The variational approach to the Kalman-Bucy problem

In this section we present the computations of $D[q(\{\mathbf{x}_t\})||p(\{\mathbf{x}_t\}|\{\mathbf{y}_i^c\})]$ for the probabilistic model corresponding to the Kalman-Bucy problem defined by the equations

$$d\mathbf{x}_t = (\mathbf{A}_t\mathbf{x}_t + \mathbf{c}_t)dt + \mathbf{B}_t^{1/2}d\mathbf{W}_t, \quad \text{and} \quad d\mathbf{y}_t = \mathbf{H}_t\mathbf{x}_t dt + \mathbf{R}_t^{1/2}d\mathbf{W}_t.$$

We relate its optimum's marginals to the marginals computed by the Kalman-Bucy algorithm. Let us discretise (again) by using the Euler scheme and write

$$\begin{aligned}p(\mathbf{x}|\mathbf{y}^c) &\propto N(\mathbf{x}_0; \mathbf{m}_0, \mathbf{V}_0) \prod_k N(\mathbf{x}_{t_{k+1}}; \mathbf{x}_{t_k} + (\mathbf{A}_{t_k}\mathbf{x}_{t_k} + \mathbf{c}_{t_k})\Delta t_k, \Delta t_k \mathbf{B}_{t_k}) \\ &\quad \times \prod_k N(\mathbf{y}_{t_{k+1}}^c; \mathbf{y}_{t_k}^c + \mathbf{H}_{t_k}\mathbf{x}_{t_k}\Delta t_k, \Delta t_k \mathbf{R}_{t_k})\end{aligned}$$

and assume that we approximate this density by a $q(\mathbf{x})$ which is the discretisation

$$q(\mathbf{x}) \propto N(\mathbf{x}_0; \mathbf{m}_0, \mathbf{V}_0) \prod_k N(\mathbf{x}_{t_{k+1}}; \mathbf{x}_{t_k} + (\mathbf{A}_{t_k}^q\mathbf{x}_{t_k} + \mathbf{c}_{t_k}^q)\Delta t_k, \Delta t_k \mathbf{B}_{t_k})$$

of an approximating $d\mathbf{x}_t = (\mathbf{A}_t^q\mathbf{x}_t + \mathbf{c}_t^q)dt + \mathbf{B}_t^{1/2}d\mathbf{W}_t$, say, with same initial conditions. After some algebra, one can show that

$$\begin{aligned}D[q(\mathbf{x})||p(\mathbf{x}|\mathbf{y})] &= \frac{1}{2} \sum_k \Delta t_k \left\langle [(\mathbf{A}_{t_k} - \mathbf{A}_{t_k}^q)\mathbf{x}_{t_k} + (\mathbf{c}_t - \mathbf{c}_{t_k}^q)]^T \mathbf{B}_{t_k}^{-1} [(\mathbf{A}_{t_k} - \mathbf{A}_{t_k}^q)\mathbf{x}_{t_k} + (\mathbf{c}_{t_k} - \mathbf{c}_{t_k}^q)] \right\rangle_{q(\mathbf{x}_{t_k})} \\ &\quad + \frac{1}{2} \sum_k \Delta t_k \left\langle \left[\frac{\Delta \mathbf{y}_{t_k}^c}{\Delta t_k} - \mathbf{H}_{t_k}\mathbf{x}_{t_k} \right]^T \mathbf{R}_{t_k}^{-1} \left[\frac{\Delta \mathbf{y}_{t_k}^c}{\Delta t_k} - \mathbf{H}_{t_k}\mathbf{x}_{t_k} \right] \right\rangle_{q(\mathbf{x}_{t_k})} + \frac{d}{2} \sum_k \log(2\pi \Delta t_k \mathbf{R}_{t_k}).\end{aligned}$$

Clearly, due to the $\sum_k \log(2\pi \Delta t_k \mathbf{R}_{t_k})$ terms, the limit $\Delta t_k \rightarrow 0$ does not exist but since these are not dependent on the variational parameters, one can still define the free energy

$$\begin{aligned}F(q) &= \frac{1}{2} \int_0^1 dt \left\langle [(\mathbf{A}_t - \mathbf{A}_t^q)\mathbf{x}_t + (\mathbf{c}_t - \mathbf{c}_t^q)]^T \mathbf{B}_t^{-1} [(\mathbf{A}_t - \mathbf{A}_t^q)\mathbf{x}_t + (\mathbf{c}_t - \mathbf{c}_t^q)] \right\rangle_{q(\mathbf{x}_t)} \\ &\quad + \frac{1}{2} \int_0^1 dt \left\langle \left[\frac{d\mathbf{y}_t^c}{dt} - \mathbf{H}_t\mathbf{x}_t \right]^T \mathbf{R}_t^{-1} \left[\frac{d\mathbf{y}_t^c}{dt} - \mathbf{H}_t\mathbf{x}_t \right] \right\rangle_{q(\mathbf{x}_t)}.\end{aligned}$$

Due to the Gaussian nature of the problem, the Kalman-Bucy algorithm provides the marginal means and variances of the optimal q corresponding to a quadratic loss function. As we show in Section B.3 below, the Kalman-Bucy algorithm can also be used to compute the \mathbf{A}_t^q and \mathbf{c}_t^q parameters of the optimal q . The case with additional discrete observations follows naturally.

B.3 The moment matching OU to a Kalman-Bucy solution

Suppose now that we have computed, by Kalman-Bucy smoothing, the marginals of a process q . Since both the prior and observation processes are linear, the posterior process will also be of OU type; however, Kalman-Bucy smoothing only computes marginals of the process, and it may be expedient to compute the SDE formulation of the process. To do that, one needs the drift coefficients A_t , c_t and the diffusion matrix B_t . In order to compute these parameters, we resort to a variational computation.

We consider the (discretised) KL divergence between a posterior process $q(x)$ arising from a Kalman-Bucy problem and an OU process $p(x)$ with parameters A_t , c_t and B_t

$$D[q(x)||p(x)] = \text{const.} + \langle \log q(x) \rangle_q + \frac{1}{2} \sum_t \left\langle [\Delta x_t - (A_t x_t + c_t) \Delta t] [\Delta t B_t]^{-1} [\Delta x_t - (A_t x_t + c_t) \Delta t]^T + \log \det(\Delta t B_t) \right\rangle_q$$

To find the optimal A_t^* , c_t^* and B_t^* we need to compute the expectations needed in the above expression. Let us denote the forward filtering distribution at time t by $N(x_t; m_t^{fw}, V_t^{fw})$ while the backward filtering at $t+\Delta t$ is represented by $N(x_{t+\Delta t}; m_{t+\Delta t}^{bw}, V_{t+\Delta t}^{bw})$; these distributions are known as they are the outcome of the Kalman-Bucy forward-backward filtering. The joint posterior of $(x_t, x_{t+\Delta t})$ can then be written as

$$q(x_t, x_{t+\Delta t}) \propto N(x_t; m_t^{fw}, V_t^{fw}) N(y_{t+\Delta t}^c + H_t x_t \Delta t, \Delta t R_t) \times N(x_{t+\Delta t}; (I + A_t \Delta t) x_t + c_t \Delta t, \Delta t B_t) N(x_{t+\Delta t}; m_{t+\Delta t}^{bw}, V_{t+\Delta t}^{bw}).$$

This density can be rewritten in a conditional form

$$q(x_t, x_{t+\Delta t}) \propto N(x_{t+\Delta t}; U_t x_t + v_t \Delta t, \Delta t Z_t) N(x_t; \hat{m}_t, \hat{V}_t).$$

where the first order approximations of the quantities above are given by

$$\begin{aligned} \hat{m}_t &= m_t^{fw} + \Delta t V_t^{fw} H_t^T (\Delta t H_t V_t^{fw} H_t^T + R_t)^{-1} \left[\frac{\Delta y_t^c}{\Delta t} - H_t m_t^{fw} \right] \\ &\simeq m_t^{fw} + \Delta t V_t^{fw} H_t^T R_t^{-1} \left[\frac{\Delta y_t^c}{\Delta t} - H_t m_t^{fw} \right] \\ \hat{V}_t &= V_t^{fw} - \Delta t V_t^{fw} H_t^T (\Delta t H_t V_t^{fw} H_t^T + R_t)^{-1} H_t V_t^{fw} \\ &\simeq V_t^{fw} - \Delta t V_t^{fw} H_t^T R_t^{-1} H_t V_t^{fw} \\ U_t &= (I + \Delta t B_t [V_{t+\Delta t}^{bw}]^{-1})^{-1} (I + \Delta t A_t) \\ &\simeq I + \Delta t (A_t - B_t [V_{t+\Delta t}^{bw}]^{-1}) \\ v_t &= (I + \Delta t B_t [V_{t+\Delta t}^{bw}]^{-1})^{-1} (c_t + B_t [V_{t+\Delta t}^{bw}]^{-1} m_{t+\Delta t}^{bw}) \\ &\simeq c_t + B_t [V_{t+\Delta t}^{bw}]^{-1} m_{t+\Delta t}^{bw} - \Delta t B_t [V_{t+\Delta t}^{bw}]^{-1} (c_t + B_t [V_{t+\Delta t}^{bw}]^{-1} m_{t+\Delta t}^{bw}) \\ Z_t &\simeq B_t - \Delta t B_t [V_{t+\Delta t}^{bw}]^{-1} B_t. \end{aligned}$$

Using the parameterisation from above, the minimisers of the KL divergence are given by

$$\begin{aligned} A_t^* &= \frac{1}{\Delta t} \left[\langle \Delta x_t x_t^T \rangle - \langle \Delta x_t \rangle \langle x_t \rangle^T \right] \left[\langle x_t x_t^T \rangle - \langle x_t \rangle \langle x_t \rangle^T \right]^{-1} \\ &\simeq \frac{1}{\Delta t} (U - I) \\ &= A_t - B_t [V_{t+\Delta t}^{bw}]^{-1} \end{aligned} \tag{31}$$

$$\begin{aligned} c_t^* &= \frac{1}{\Delta t} \langle x_{t+\Delta t} - x_t \rangle - A_t^* \langle x_t \rangle \\ &\simeq \frac{1}{\Delta t} (U - I) \langle x_t \rangle - A_t^* \langle x_t \rangle + v_t \\ &\simeq c_t + B_t [V_{t+\Delta t}^{bw}]^{-1} m_{t+\Delta t}^{bw} \end{aligned} \tag{32}$$

$$\begin{aligned} B_t^* &= \frac{1}{\Delta t} \left\langle [\Delta x_t - (A_t^* x_t + c_t^*) \Delta t] [\Delta x_t - (A_t^* x_t + c_t^*) \Delta t]^T \right\rangle \\ &\simeq B_t. \end{aligned} \tag{33}$$

We remark that when adding the discrete time observations, the form of (31), (32) and (33) does not change, we only have to make sure that the backward filtering for computing m_t^{bw} and V_t^{bw} does include these terms.

C The inference algorithm

Until convergence do

- (1) Update $\{(\mathbf{h}_{t_i}^d, \mathbf{Q}_{t_i}^d)\}_i$ and $\{(\mathbf{h}_t^c, \mathbf{Q}_t^c)\}_t$ according to

- (1.1) Update $(\mathbf{h}_{t_i}^d, \mathbf{Q}_{t_i}^d)$ by

- (1.1.1) compute the cavity means and variances

$$\mathbf{m}_{t_i}^{\setminus t_i} = (\mathbf{I} - \mathbf{Q}_{t_i}^d \mathbf{V}_{t_1})^{-1}(\mathbf{m}_{t_i} - \mathbf{V}_{t_i} \mathbf{h}_{t_i}^d) \quad \text{and} \quad \mathbf{V}_{t_i}^{\setminus t_i} = \mathbf{V}_{t_i} (\mathbf{I} - \mathbf{Q}_{t_i}^d \mathbf{V}_{t_i})^{-1}$$

- (1.1.2) compute mean $\hat{\mathbf{m}}_{t_i}$ and covariance $\hat{\mathbf{V}}_{t_i}$ of the tilted distribution $q_d^i(\mathbf{x}_{t_i}) \propto p(\mathbf{y}_{t_i}^d | \mathbf{x}_{t_i}) N(\mathbf{x}_{t_i}; \mathbf{m}_{t_i}^{\setminus t_i}, \mathbf{V}_{t_i}^{\setminus t_i})$ either by exact or numerical methods (see EP related references for details)

- (1.1.3) compute $[\mathbf{h}_{t_i}^d]^{new}$ and $[\mathbf{Q}_{t_i}^d]^{new}$ from

$$[\mathbf{h}_{t_i}^d]^{new} = [\hat{\mathbf{V}}_{t_i}]^{-1} \hat{\mathbf{m}}_{t_i} - [\mathbf{V}_{t_i}^{\setminus t_i}]^{-1} \mathbf{m}_{t_i}^{\setminus t_i}$$

$$[\mathbf{Q}_{t_i}^d]^{new} = [\hat{\mathbf{V}}_{t_i}]^{-1} - [\mathbf{V}_{t_i}^{\setminus t_i}]^{-1}$$

- (1.2) Update \mathbf{h}_t^c and \mathbf{Q}_t^c by

$$[\mathbf{h}_t^c]^{new} = -\partial_{\mathbf{m}_t} \langle V(t, \mathbf{x}_t) \rangle_{\mathcal{N}(\mathbf{m}_t, \mathbf{V}_t)} + 2\partial_{\mathbf{V}_t} \langle V(t, \mathbf{x}_t) \rangle_{\mathcal{N}(\mathbf{m}_t, \mathbf{V}_t)} \mathbf{m}_t$$

$$[\mathbf{Q}_t^c]^{new} = \partial_{\mathbf{V}_t} \langle V(t, \mathbf{x}_t) \rangle_{\mathcal{N}(\mathbf{m}_t, \mathbf{V}_t)}$$

- (2) Update \mathbf{m}_t and \mathbf{V}_t by

- (2.1) solve forward starting at $(\mathbf{m}_0, \mathbf{V}_0)$

$$\frac{d}{dt} \mathbf{V}_t^{fw} = \mathbf{A}_t \mathbf{V}_t^{fw} + \mathbf{V}_t^{fw} \mathbf{A}_t^T + \mathbf{B}_t - \mathbf{V}_t^{fw} \mathbf{Q}_t^c \mathbf{V}_t^{fw}, \quad \mathbf{m}_{t_i+}^{fw} = (\mathbf{I} + \mathbf{V}_{t_i}^{fw} \mathbf{Q}_{t_i}^d)^{-1}(\mathbf{m}_{t_i}^{fw} + \mathbf{V}_{t_i}^{fw} \mathbf{h}_{t_i}^c),$$

$$\frac{d}{dt} \mathbf{m}_t^{fw} = \mathbf{A}_t \mathbf{m}_t^{fw} + \mathbf{c}_t + \mathbf{V}_t^{fw} [\mathbf{h}_t^c - \mathbf{Q}_t^c \mathbf{m}_t^{fw}], \quad \mathbf{V}_{t_i+}^{fw} = (\mathbf{I} + \mathbf{V}_t^{fw} \mathbf{Q}_{t_i}^d)^{-1} \mathbf{V}_t^{fw}$$

- (2.2) solve backwards starting, say, at $(\mathbf{0}, 100\mathbf{V}_1^{fw})$

$$\frac{d}{dt} \mathbf{V}_t^{bw} = \mathbf{A}_t \mathbf{V}_t^{bw} + \mathbf{V}_t^{bw} \mathbf{A}_t^T - \mathbf{B}_t + \mathbf{V}_t^{bw} \mathbf{Q}_t^c \mathbf{V}_t^{bw}, \quad \mathbf{m}_{t_i-}^{bw} = (\mathbf{I} + \mathbf{V}_{t_i}^{bw} \mathbf{Q}_{t_i}^d)^{-1}(\mathbf{m}_{t_i}^{bw} + \mathbf{V}_{t_i}^{bw} \mathbf{h}_{t_i}^d),$$

$$\frac{d}{dt} \mathbf{m}_t^{bw} = \mathbf{A}_t \mathbf{m}_t^{bw} + \mathbf{c}_t - \mathbf{V}_t^{bw} [\mathbf{h}_t^c - \mathbf{Q}_t^c \mathbf{m}_t^{bw}], \quad \mathbf{V}_{t_i-}^{bw} = (\mathbf{I} + \mathbf{V}_t^{bw} \mathbf{Q}_{t_i}^d)^{-1} \mathbf{V}_t^{bw}$$

- (2.3) compute \mathbf{m}_t and \mathbf{V}_t from

$$[\mathbf{V}_t]^{-1} = [\mathbf{V}_t^{fw}]^{-1} + [\mathbf{V}_t^{bw}]^{-1} \quad \text{and} \quad \mathbf{m}_t = \mathbf{V}_t [[\mathbf{V}_t^{fw}]^{-1} \mathbf{m}_t^{fw} + [\mathbf{V}_t^{bw}]^{-1} \mathbf{m}_t^{bw}]$$

The above equations are given for illustrative purposes only, one should always avoid inversion, try to stabilise computations by reorganising them and, when necessary, performing matrix inversions through matrix factorisations.

The sequential forward-backward scheduling follows by iteratively solving the forward and backward equations with \mathbf{Q}_t^c and \mathbf{h}_t^c computed according to (1.2) and the update steps (1.1) for \mathbf{Q}_t^d and \mathbf{h}_t^d performed at the jump times t_i .