

# Supplementary Material: The Variational Coupled Gaussian Process Dynamical Model

Dmytro Velychko, Benjamin Knopp, Dominik Endres\*

March 31, 2017

## 1 Partial optimization of variational distribution for simplified ELBO

While optimizing the full variational posterior in augmented Gaussian Processes models the following type of term appears often in the ELBO equation:

$$\begin{aligned}\mathcal{R}(q(\mathbf{u}), q(\mathbf{v})) &= \int q(\mathbf{u}) \left( f(q(\mathbf{v}), \mathbf{u}) + \log \frac{p(\mathbf{u})}{q(\mathbf{u})} \right) d\mathbf{u} \\ &= \int q(\mathbf{u}) f(q(\mathbf{v}), \mathbf{u}) d\mathbf{u} + \int q(\mathbf{u}) \log p(\mathbf{u}) d\mathbf{u} - \int q(\mathbf{u}) \log q(\mathbf{u}) d\mathbf{u}\end{aligned}\tag{S1}$$

To simplify the optimization of such terms, we would like to carry out the optimization with respect to the density  $q(\mathbf{u})$  analytically, so as to remove the dependency on  $q(\mathbf{u})$ . To this end, we calculate for the optimal variational  $q^*(\mathbf{u})$  in the above equation. This approach was suggested in [1], however, it is not well described there. Here we give an extended derivation. A necessary condition for maximality is a vanishing functional derivative under the constraint that the density  $q(\mathbf{u})$  is normalized to one:

$$\int q(\mathbf{u}) d\mathbf{u} - 1 = 0\tag{S2}$$

which is fulfilled at the stationary points of the Lagrangian

$$\mathcal{X}(q(\mathbf{u}), q(\mathbf{v})) = \mathcal{R}(q(\mathbf{u}), q(\mathbf{v})) + \lambda \left( \int q(\mathbf{u}) d\mathbf{u} - 1 \right)\tag{S3}$$

---

\*Theoretical Neuroscience Group, Dept. Psychology, University of Marburg, Gutenbergstr. 18, 35032 Marburg, Germany  
{dmytro.velychko,benjamin.knopp,dominik.endres}@uni-marburg.de

where  $\lambda$  is chosen so that (S2) holds. Taking the derivative of  $\mathcal{X}(q(\mathbf{u}), q(\mathbf{v}))$  and setting it to zero yields

$$\frac{\delta \mathcal{X}(q(\mathbf{u}), q(\mathbf{v}))}{\delta q(\mathbf{u})} = f(q(\mathbf{v}), \mathbf{u}) + \log p(\mathbf{u}) - \log q(\mathbf{u}) + 1 + \lambda = 0 \quad (\text{S4})$$

and therefore, denoting  $Z = \exp(-\lambda - 1)$

$$q^*(\mathbf{u}) = \exp(f(q(\mathbf{v}), \mathbf{u}) + \log p(\mathbf{u}) + 1 + \lambda) \quad (\text{S5})$$

$$q^*(\mathbf{u}) = \frac{1}{Z} p(\mathbf{u}) \exp(f(q(\mathbf{v}), \mathbf{u})) \quad (\text{S6})$$

$$Z = \exp(-\lambda - 1) = \int p(\mathbf{u}) \exp(f(q(\mathbf{v}), \mathbf{u})) d\mathbf{u} \quad (\text{S7})$$

Substituting the optimal  $q^*(\mathbf{u})$  into the original term, we get:

$$\begin{aligned} \mathcal{R}(q(\mathbf{v})) &= \int \frac{1}{Z} p(\mathbf{u}) \exp(f(q(\mathbf{v}), \mathbf{u})) (f(q(\mathbf{v}), \mathbf{u}) + \log \frac{p(\mathbf{u})}{\frac{1}{Z} p(\mathbf{u}) \exp(f(q(\mathbf{v}), \mathbf{u}))}) d\mathbf{u} \\ &= \int \frac{1}{Z} p(\mathbf{u}) \exp(f(q(\mathbf{v}), \mathbf{u})) \left( \log \frac{p(\mathbf{u}) \exp(f(q(\mathbf{v}), \mathbf{u}))}{\frac{1}{Z} p(\mathbf{u}) \exp(f(q(\mathbf{v}), \mathbf{u}))} \right) d\mathbf{u} \\ &= \log(Z) \frac{1}{Z} \int p(\mathbf{u}) \exp(f(q(\mathbf{v}), \mathbf{u})) d\mathbf{u} \\ &= \log \int p(\mathbf{u}) \exp(f(q(\mathbf{v}), \mathbf{u})) d\mathbf{u} \end{aligned} \quad (\text{S8})$$

This is the simplified version of (S1), which depends only on  $q(\mathbf{v})$  because  $q(\mathbf{u}) = q^*(\mathbf{u})$  has been determined by optimization.

## 2 vCGPDM ELBO derivation

Here we give the extended derivation of the ELBO from the main paper. Let's assume we deal with  $M$  parts. The model reads:

$$\mathbf{X} = \{\mathbf{X}^0 \dots \mathbf{X}^M\} \quad (\text{S9})$$

$$\mathbf{X}^i = \{\mathbf{x}_0^i \dots \mathbf{x}_T^i\}; \mathbf{x}_t^i \in \mathbb{R}^{Q^i} \quad (\text{S10})$$

$$\mathbf{Y}^i = \{\mathbf{y}_0^i \dots \mathbf{y}_T^i\}; \mathbf{y}_t^i \in \mathbb{R}^{D^i} \quad (\text{S11})$$

$$f^{j,i}(\mathbf{X}_{-t}^j) \sim \mathcal{GP}(0, k_f^j(\mathbf{x}_{-t}^j, \mathbf{x}_{-t}^{j'})) \quad (\text{S12})$$

$$p(\mathbf{x}_t^i) \sim \prod_{j=1}^M \mathcal{N}(\mathbf{x}_t^i | f^{j,i}(\mathbf{x}_{-t}^j), \mathbf{I} \alpha^{j,i}); \alpha > 0 \quad (\text{S13})$$

$$g_d^i(\mathbf{X}^i) \sim \mathcal{GP}(0, k_g(\mathbf{x}_t^i, \mathbf{x}_t^{i'})) \quad (\text{S14})$$

$$\mathbf{g}_d^i = g_d^i(\mathbf{X}^i) \quad (\text{S15})$$

$$\mathbf{Y}_{:,d}^i \sim \mathcal{N}(\mathbf{g}_d^i, \mathbf{I} \beta^i); \beta > 0 \quad (\text{S16})$$

$$p(\mathbf{x}_0^i) = \mathcal{N}(\mathbf{x}_0^i | 0, \mathbf{I}) \quad (\text{S17})$$

Here and following upper indexes are part-related, lower indexes indicate dimensions.

We have  $M \times M$  latent dynamics mappings, which are combined into  $M$  mappings with the Product of Experts - multiplying and renormalizing the distribution from separate partial predictions. Each of the  $M \times M$  mappings is augmented with some inducing inputs and outputs:  $f^{j,i} : \mathbf{z}^{j,i} \rightarrow \mathbf{u}^{j,i}$ . The full augmented joint distribution of the model reads:

$$\begin{aligned}
p(\mathbf{x}, \mathbf{u}, \mathbf{f}, \mathbf{v}, \mathbf{g}, \mathbf{y}) &= p(\mathbf{y}|\mathbf{g})p(\mathbf{g}|\mathbf{x}, \mathbf{v})p(\mathbf{v})p(\mathbf{x}, \mathbf{f}|\mathbf{u})p(\mathbf{u}) \\
&= \left[ \prod_{i=1}^M p(\mathbf{y}^i|\mathbf{g}^i)p(\mathbf{g}^i|\mathbf{x}^i, \mathbf{v}^i)p(\mathbf{v}^i) \right] p(\mathbf{x}, \mathbf{f}|\mathbf{u})p(\mathbf{u}) \\
&= \left[ \prod_{i=1}^M \left[ \prod_{d=1}^{D_i} p(\mathbf{y}_d^i|\mathbf{g}_d^i)p(\mathbf{g}_d^i|\mathbf{x}^i, \mathbf{v}^i) \right] p(\mathbf{v}^i) \right] \\
&\times \left[ \prod_{t=1}^T \left[ \prod_{i=1}^M p(\mathbf{x}_t^i|\{\mathbf{f}_t^{:,i}, \alpha^{:,i}\}) \right] \left[ \prod_{i=1}^M \prod_{j=1}^M p(\mathbf{f}_t^{j,i}|\mathbf{f}_{1:t-1}^{j,i}, \mathbf{x}_{0:t-1}^j, \mathbf{u}^{j,i}) \right] \right] \\
&\times \left[ \prod_{i=1}^M \prod_{j=1}^M p(\mathbf{u}^{j,i}) \right] \left[ \prod_{i=1}^M p(\mathbf{x}_0^i) \right] \tag{S18}
\end{aligned}$$

The full proposal variational posterior is:

$$\begin{aligned}
q(\mathbf{x}, \mathbf{u}, \mathbf{f}, \mathbf{v}, \mathbf{g}) &= p(\mathbf{g}|\mathbf{x}, \mathbf{v})q(\mathbf{v})p(\mathbf{f}|\mathbf{x}, \mathbf{u})q(\mathbf{x})q(\mathbf{u}) \\
&= p(\mathbf{g}|\mathbf{x}, \mathbf{v})q(\mathbf{v}) \left[ \prod_{t=1}^T \prod_{i=1}^M \prod_{j=1}^M p(\mathbf{f}_t^{j,i}|\mathbf{f}_{1:t-1}^{j,i}, \mathbf{x}_{0:t-1}^j, \mathbf{u}^{j,i}) \right] q(\mathbf{x})q(\mathbf{u}) \tag{S19}
\end{aligned}$$

The ELBO is:

$$\mathcal{L}(\theta) = \int_{\mathbf{x}, \mathbf{v}, \mathbf{g}} q(\mathbf{x}, \mathbf{u}, \mathbf{f}, \mathbf{v}, \mathbf{g}) \log \left( \frac{p(\mathbf{x}, \mathbf{u}, \mathbf{f}, \mathbf{v}, \mathbf{g}, \mathbf{y})}{q(\mathbf{x}, \mathbf{u}, \mathbf{f}, \mathbf{v}, \mathbf{g})} \right) \quad (\text{S20})$$

$$= \sum_{i=1}^M \sum_{d=1}^D \int_{\mathbf{x}^i, \mathbf{v}^i, \mathbf{g}_d^i} p(\mathbf{g}_d^i | \mathbf{x}^i, \mathbf{v}^i) q(\mathbf{x}^i) q(\mathbf{v}^i) \log \frac{p(\mathbf{y}_d^i | \mathbf{g}_d^i)}{q(\mathbf{v}^i)} \quad (\text{S21})$$

$$\begin{aligned} &+ \int q(\mathbf{u}) \left[ \sum_{t=1}^T \int q(\mathbf{x}_t) q(\mathbf{x}_{-t}) \left( \int \left[ \prod_{i=1}^M \prod_{j=1}^M p(\mathbf{f}_t^{j,i} | \mathbf{f}_{1:t-1}^{j,i}, \mathbf{x}_{0:t-1}^j, \mathbf{u}^{j,i}) \right] \right. \right. \\ &\times \left. \left. \log \prod_{i=1}^M p(\mathbf{x}_t^i | \{\mathbf{f}_t^{:,i}, \alpha^{:,i}\}) d\mathbf{f}_t \right) d\mathbf{x}_t d\mathbf{x}_{-t} \right] \\ &+ q(\mathbf{u}) \log \frac{p(\mathbf{u})}{q(\mathbf{u})} d\mathbf{u} \\ &+ \int q(\mathbf{x}_0) \log p(\mathbf{x}_0) d\mathbf{x}_0 + H(q(\mathbf{x})) \quad (\text{S22}) \end{aligned}$$

The (S21) part is the GPLVM ELBO and is given in [2]. Next, we consider only the ELBO component which is relevant for the dynamics and apply the sufficient statistics assumption: knowing  $\mathbf{x}_{-t}^j$  and  $\mathbf{u}^{j,i}$  is sufficient for the  $\mathbf{f}_t^{j,i}$  distribution. The innermost integral is:

$$\begin{aligned} \mathcal{A} &= \int \left[ \prod_{i=1}^M \prod_{j=1}^M p(\mathbf{f}_t^{j,i} | \mathbf{x}_{-t}^j, \mathbf{u}^{j,i}) \right] \log \prod_{i=1}^M p(\mathbf{x}_t^i | \{\mathbf{f}_t^{:,i}, \alpha^{:,i}\}) d\mathbf{f}_t \\ &= \sum_{i=1}^M \int \left[ \prod_{j=1}^M p(\mathbf{f}_t^{j,i} | \mathbf{x}_{-t}^j, \mathbf{u}^{j,i}) \right] \log p(\mathbf{x}_t^i | \{\mathbf{f}_t^{:,i}, \alpha^{:,i}\}) d\mathbf{f}_t^i \\ &= \sum_{i=1}^M \int \left[ \prod_{j=1}^M \mathcal{N}(\mathbf{f}_t^{j,i} | \boldsymbol{\mu}_{\mathbf{f}_t^{j,i}}, \mathbf{S}_{\mathbf{f}_t^{j,i}}) \right] \log \mathcal{N}(\mathbf{x}_t^i | \alpha_i \sum_{j=1}^M \alpha_{j,i}^{-1} \mathbf{f}_t^{j,i}, \mathbf{I} \alpha_i) d\mathbf{f}_t^i \\ &= \sum_{i=1}^M \left[ -\frac{1}{2} \text{tr} \left[ \alpha_i \sum_{j=1}^M (\alpha_{j,i}^{-1})^2 \mathbf{S}_{\mathbf{f}_t^{j,i}} \right] + \log \mathcal{N}(\mathbf{x}_t^i | \alpha_i \sum_{j=1}^M \alpha_{j,i}^{-1} \boldsymbol{\mu}_{\mathbf{f}_t^{j,i}}, \mathbf{I} \alpha_i) \right] \quad (\text{S23}) \end{aligned}$$

$$\boldsymbol{\mu}_{\mathbf{f}_t^{j,i}} = \mathbf{K}_{\mathbf{x}_{-t}^j, \mathbf{z}^{j,i}}^{j,i} \left( \mathbf{K}_{\mathbf{z}^{j,i}, \mathbf{z}^{j,i}}^{j,i} \right)^{-1} \mathbf{u}^{j,i} \quad (\text{S24})$$

$$\mathbf{S}_{\mathbf{f}_t^{j,i}} = \mathbf{K}_{\mathbf{x}_{-t}^j, \mathbf{x}_{-t}^j}^{j,i} - \mathbf{K}_{\mathbf{x}_{-t}^j, \mathbf{z}^{j,i}}^{j,i} \left( \mathbf{K}_{\mathbf{z}^{j,i}, \mathbf{z}^{j,i}}^{j,i} \right)^{-1} \mathbf{K}_{\mathbf{x}_{-t}^j, \mathbf{z}^{j,i}}^{j,i} \quad (\text{S25})$$

$$\alpha_i = \left( \sum_{j=1}^M \alpha_{j,i}^{-1} \right)^{-1} \quad (\text{S26})$$

Now let's take the integral over  $\mathbf{x}$ :

$$\begin{aligned}
\mathcal{B} &= \int q(\mathbf{x}_t)q(\mathbf{x}_{-t})\mathcal{A}d\mathbf{x}_td\mathbf{x}_{-t} \\
&= \sum_{i=1}^M \left( -\frac{1}{2}\alpha_i \sum_{j=1}^M (\alpha_{j,i}^{-1})^2 \text{tr} \left[ \Psi_0^{j,i}(\mathbf{x}_{-t}^j) - \left( \mathbf{K}_{\mathbf{z}^{j,i},\mathbf{z}^{j,i}}^{j,i} \right)^{-1} \Psi_2^{j,i}(\mathbf{x}_{-t}^j) \right] \right. \\
&\quad - \log Z(\mathbf{I}\alpha_i) - \frac{1}{2}\alpha_i^{-1} \left[ \text{tr}(\mathbf{S}_{\mathbf{x}_t^i}) + \boldsymbol{\mu}_{\mathbf{x}_t^i}^T \boldsymbol{\mu}_{\mathbf{x}_t^i} \right] \\
&\quad + \boldsymbol{\mu}_{\mathbf{x}_t^i}^T \left( \sum_{j=1}^M \alpha_{j,i}^{-1} \Psi_1^{j,i}(\mathbf{x}_{-t}^j) \left( \mathbf{K}_{\mathbf{z}^{j,i},\mathbf{z}^{j,i}}^{j,i} \right)^{-1} \mathbf{u}^{j,i} \right) \\
&\quad \left. - \frac{1}{2}\alpha_i \sum_{j=1}^M \sum_{k=1}^M \alpha_{j,i}^{-1} \alpha_{k,i}^{-1} \mathbf{u}^{j,iT} \left( \mathbf{K}_{\mathbf{z}^{j,i},\mathbf{z}^{j,i}}^{j,i} \right)^{-1} \Psi_2^{j,k,i}(\mathbf{x}_{-t}^j, \mathbf{x}_{-t}^k) \left( \mathbf{K}_{\mathbf{z}^{k,i},\mathbf{z}^{k,i}}^{k,i} \right)^{-1} \mathbf{u}^{k,i} \right)
\end{aligned} \tag{S27}$$

The sum over time points:

$$\begin{aligned}
\mathcal{C} &= \sum_{t=1}^T \mathcal{B} \\
&= \sum_{i=1}^M \left( \sum_{t=1}^T \left[ -\frac{1}{2}\alpha_i \sum_{j=1}^M (\alpha_{j,i}^{-1})^2 \text{tr} \left[ \Psi_0^{j,i}(\mathbf{x}_{-t}^j) - \left( \mathbf{K}_{\mathbf{z}^{j,i},\mathbf{z}^{j,i}}^{j,i} \right)^{-1} \Psi_2^{j,i}(\mathbf{x}_{-t}^j) \right] \right. \right. \\
&\quad \left. - \log Z(\mathbf{I}\alpha_i) - \frac{1}{2}\alpha_i^{-1} \left[ \text{tr}(\mathbf{S}_{\mathbf{x}_t^i}) + \boldsymbol{\mu}_{\mathbf{x}_t^i}^T \boldsymbol{\mu}_{\mathbf{x}_t^i} \right] \right] \\
&\quad + \sum_{j=1}^M \alpha_{j,i}^{-1} \left[ \sum_{t=1}^T \boldsymbol{\mu}_{\mathbf{x}_t^i}^T \Psi_1^{j,i}(\mathbf{x}_{-t}^j) \right] \left( \mathbf{K}_{\mathbf{z}^{j,i},\mathbf{z}^{j,i}}^{j,i} \right)^{-1} \mathbf{u}^{j,i} \\
&\quad \left. - \frac{1}{2}\alpha_i \sum_{j=1}^M \sum_{k=1}^M \alpha_{j,i}^{-1} \alpha_{k,i}^{-1} \mathbf{u}^{j,iT} \left( \mathbf{K}_{\mathbf{z}^{j,i},\mathbf{z}^{j,i}}^{j,i} \right)^{-1} \left[ \sum_{t=1}^T \Psi_2^{j,k,i}(\mathbf{x}_{-t}^j, \mathbf{x}_{-t}^k) \right] \left( \mathbf{K}_{\mathbf{z}^{k,i},\mathbf{z}^{k,i}}^{k,i} \right)^{-1} \mathbf{u}^{k,i} \right)
\end{aligned} \tag{S28}$$

For every  $i \in 1 \dots M$  we may stack up the  $\mathbf{u}^{:,i}$  into  $\mathbf{u}^i$  and construct a large block matrices  $\mathcal{F}^i$  and stacked matrices  $\mathcal{G}^i$  with elements

$$\mathcal{F}_{j,k}^i = \alpha_i \alpha_{j,i}^{-1} \alpha_{k,i}^{-1} \left( \mathbf{K}_{\mathbf{z}^{j,i},\mathbf{z}^{j,i}}^{j,i} \right)^{-1} \left[ \sum_{t=1}^T \Psi_2^{j,k,i}(\mathbf{x}_{-t}^j, \mathbf{x}_{-t}^k) \right] \left( \mathbf{K}_{\mathbf{z}^{k,i},\mathbf{z}^{k,i}}^{k,i} \right)^{-1} \tag{S29}$$

$$\mathcal{G}_j^i = \left( \alpha_{j,i}^{-1} \left[ \sum_{t=1}^T \boldsymbol{\mu}_{\mathbf{x}_t^i}^T \Psi_1^{j,i}(\mathbf{x}_{-t}^j) \right] \left( \mathbf{K}_{\mathbf{z}^{j,i},\mathbf{z}^{j,i}}^{j,i} \right)^{-1} \right)^T \tag{S30}$$

For  $j \neq k$ :  $\Psi_2^{j,k,i}(\mathbf{x}_{-t}^j, \mathbf{x}_{-t}^k) = \Psi_1^{j,i}(\mathbf{x}_{-t}^j)\Psi_1^{k,i}(\mathbf{x}_{-t}^k)$ . Otherwise  $\Psi_2^{j,j,i}(\mathbf{x}_{-t}^j, \mathbf{x}_{-t}^j) = \Psi_2^{j,i}(\mathbf{x}_{-t}^j)$ .

The sum over time points must be expressed as a quadratic form w.r.t. the augmenting outputs  $\mathbf{u}$ :

$$\begin{aligned} \mathcal{C} &= \sum_{i=1}^M \left[ -\frac{1}{2} \mathbf{u}^{iT} \mathcal{F}^i \mathbf{u}^i + \mathbf{u}^{iT} \mathcal{G}^i + \mathcal{H}^i \right] \\ &= \sum_{i=1}^M \left[ -\frac{1}{2} (\mathbf{u}^i - \mathcal{F}^{i-1} \mathcal{G}^i)^T \mathcal{F}^i (\mathbf{u}^i - \mathcal{F}^{i-1} \mathcal{G}^i) + \frac{1}{2} \mathcal{G}^{iT} \mathcal{F}^{i-1} \mathcal{G}^i + \mathcal{H}^i \right] \end{aligned} \quad (\text{S31})$$

$$\mathcal{C} = \sum_{i=1}^M \mathcal{C}^i \quad (\text{S32})$$

$$\mathcal{C}^i = -\frac{1}{2} (\mathbf{u}^i - \mathcal{F}^{i-1} \mathcal{G}^i)^T \mathcal{F}^i (\mathbf{u}^i - \mathcal{F}^{i-1} \mathcal{G}^i) + \frac{1}{2} \mathcal{G}^{iT} \mathcal{F}^{i-1} \mathcal{G}^i + \mathcal{H}^i \quad (\text{S33})$$

$$\begin{aligned} \mathcal{H}^i &= \sum_{t=1}^T \left[ -\frac{1}{2} \alpha_i \sum_{j=1}^M (\alpha_{j,i}^{-1})^2 \text{tr} \left[ \Psi_0^{j,i}(\mathbf{x}_{-t}^j) - \left( \mathbf{K}_{\mathbf{z}^{j,i}, \mathbf{z}^{j,i}}^{j,i} \right)^{-1} \Psi_2^{j,i}(\mathbf{x}_{-t}^j) \right] \right. \\ &\quad \left. - \log Z(\mathbf{I} \alpha_i) - \frac{1}{2} \alpha_i^{-1} \left[ \text{tr}(\mathbf{S}_{\mathbf{x}_t^i}) + \boldsymbol{\mu}_{\mathbf{x}_t^i}^T \boldsymbol{\mu}_{\mathbf{x}_t^i} \right] \right] \end{aligned} \quad (\text{S34})$$

Finally, we may write the dynamics ELBO, also accounting for the optimal variational  $q(\mathbf{u})$  (see eq. S8) and using  $p(\mathbf{u}^i) = \prod_{j=1}^M p(\mathbf{u}^{ji}) = \prod_{j=1}^M \mathcal{N}(\mathbf{u}^{ji} | 0, \mathbf{K}_{\mathbf{z}^{j,i}, \mathbf{z}^{j,i}}) = \mathcal{N}(\mathbf{u}^i | 0, \mathbf{K}_{\mathbf{z}^{:,i}, \mathbf{z}^{:,i}})$  where  $\mathbf{K}_{\mathbf{z}^{:,i}, \mathbf{z}^{:,i}}$  is a block-diagonal covariance matrix:

$$\begin{aligned} \mathcal{L}_{dyn}(\boldsymbol{\theta}) &\geq \log \int p(\mathbf{u}) \exp(\mathcal{C}) d\mathbf{u} + H(q(\mathbf{x})) \\ &= \log \prod_{i=1}^M \int p(\mathbf{u}^i) \exp(\mathcal{C}^i) d\mathbf{u}^i + H(q(\mathbf{x})) \\ &= \sum_{i=1}^M \left[ \log \int p(\mathbf{u}^i) \exp\left(-\frac{1}{2} (\mathbf{u}^i - \mathcal{F}^{i-1} \mathcal{G}^i)^T \mathcal{F}^i (\mathbf{u}^i - \mathcal{F}^{i-1} \mathcal{G}^i)\right) d\mathbf{u}^i + \frac{1}{2} \mathcal{G}^{iT} \mathcal{F}^{i-1} \mathcal{G}^i + \mathcal{H}^i \right] + H(q(\mathbf{x})) \\ &= \sum_{i=1}^M \left[ -\log Z(\mathcal{F}^{i-1} + \mathbf{K}_{\mathbf{z}^{:,i}, \mathbf{z}^{:,i}}) - \frac{1}{2} \mathcal{G}^{iT} \mathcal{F}^{i-1} (\mathcal{F}^{i-1} + \mathbf{K}_{\mathbf{z}^{:,i}, \mathbf{z}^{:,i}})^{-1} \mathcal{F}^{i-1} \mathcal{G}^i + \log Z(\mathcal{F}^{i-1}) \right] \\ &\quad + \sum_{i=1}^M \left[ \frac{1}{2} \mathcal{G}^{iT} \mathcal{F}^{i-1} \mathcal{G}^i + \mathcal{H}^i \right] + H(q(\mathbf{x})) \end{aligned} \quad (\text{S35})$$

This is even lower bound on the ELBO due to the sufficient statistics assumption for the  $\mathbf{f}_t^{j,i}$  distribution. It is easy to notice that the optimal proposal  $q(\mathbf{u}) = \prod_{i=1}^M q(\mathbf{u}^i)$  - factorized.

We optimize the full ELBO w.r.t. the parameters of  $q(\mathbf{x})$ , augmenting inputs  $\mathbf{z}$ , kernel parameters, and couplings  $\boldsymbol{\alpha}$ . After the optimization the optimal  $q(\mathbf{u})$

is computed as:

$$\begin{aligned}
q(\mathbf{u}^i) &= \frac{1}{Z} p(\mathbf{u}^i) \exp(\mathcal{C}^i) \\
&= \frac{1}{Z} \mathcal{N}(\mathbf{u}^i | 0, \mathbf{K}_{\mathbf{z}^{:i}, \mathbf{z}^{:i}}} ) \exp\left(-\frac{1}{2} (\mathbf{u}^i - \mathcal{F}^{i-1} \mathcal{G}^i)^T \mathcal{F}^i (\mathbf{u}^i - \mathcal{F}^{i-1} \mathcal{G}^i)\right) \\
&= \mathcal{N}(\mathbf{u}^i | (\mathbf{K}_{\mathbf{z}^{:i}, \mathbf{z}^{:i}}}^{-1} + \mathcal{F}^i)^{-1} \mathcal{G}^i, (\mathbf{K}_{\mathbf{z}^{:i}, \mathbf{z}^{:i}}}^{-1} + \mathcal{F}^i)^{-1}) \tag{S36}
\end{aligned}$$

### 3 ARD RBF kernel $\Psi$ statistics. Full covariance variational parameters case.

Consider the following form of the approximate variational posterior distribution of  $\mathbf{X}$ :

$$q(X) = \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, S_n) \tag{S37}$$

Here we derive the  $\Psi$  statistics for the variational lower bound for the case when the  $\{S_n\}_{n=1 \dots N}$  are *full covariance matrices* and the ARD SE kernel is defined as:

$$k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp\left(-\frac{1}{2} \sum_{q=1}^Q \frac{(\mathbf{x}_q - \mathbf{x}'_q)^2}{\lambda_q}\right) \tag{S38}$$

In matrix notation:

$$\lambda = \text{diag}(\lambda_1 \dots \lambda_Q) \tag{S39}$$

$$k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp\left(-\frac{1}{2} (\mathbf{x} - \mathbf{x}')^T \lambda^{-1} (\mathbf{x} - \mathbf{x}')\right) \tag{S40}$$

where  $\lambda_q$  are the ARD factors.

The  $\Psi_0$  statistic is easy to calculate and it does not depend on the covariance matrix:

$$\begin{aligned}
\Psi_0^n &= \int k(\mathbf{x}_n, \mathbf{x}_n) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, S_n) d\mathbf{x}_n \\
&= \int \sigma_f^2 \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, S_n) d\mathbf{x}_n \\
&= \sigma_f^2 \int \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, S_n) d\mathbf{x}_n \\
&= \sigma_f^2 \tag{S41}
\end{aligned}$$

$$\Psi_0 = \sum_{n=1}^N \Psi_0^n = N \sigma_f^2 \tag{S42}$$

The  $\Psi_1$  statistic:

$$\begin{aligned} (\Psi_1)_{nm} &= \int k(\mathbf{x}_n, \mathbf{z}_m) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, S_n) d\mathbf{x}_n \\ &= \int \sigma_f^2 \exp\left(-\frac{1}{2}(\mathbf{x}_n - \mathbf{z}_m)^T \lambda^{-1}(\mathbf{x}_n - \mathbf{z}_m)\right) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, S_n) d\mathbf{x}_n \end{aligned} \quad (\text{S43})$$

Next, complete the ARD SE kernel to a scaled Gaussian distribution:

$$\begin{aligned} (\Psi_1)_{nm} &= \sigma_f^2 \int \frac{Z(\lambda)}{Z(\lambda)} \sigma_f^2 \exp\left(-\frac{1}{2}(\mathbf{x}_n - \mathbf{z}_m)^T \lambda^{-1}(\mathbf{x}_n - \mathbf{z}_m)\right) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, S_n) d\mathbf{x}_n \\ &= \sigma_f^2 Z(\lambda) \int \mathcal{N}(\mathbf{x}_n | \mathbf{z}_m, \lambda) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, S_n) d\mathbf{x}_n \end{aligned} \quad (\text{S44})$$

Noticing the product of two Gaussians, which is an unnormalized Gaussian, the integral boils down to the scaling coefficient, which is a Gaussian value:

$$\begin{aligned} (\Psi_1)_{nm} &= \sigma_f^2 Z(\lambda) \mathcal{N}(\mathbf{z}_m | \boldsymbol{\mu}_n, \lambda + S_n) \\ &= \sigma_f^2 (2\pi)^{Q/2} \sqrt{|\lambda|} \frac{1}{(2\pi)^{Q/2} \sqrt{|\lambda + S_n|}} \exp\left(-\frac{1}{2}(\mathbf{z}_m - \boldsymbol{\mu}_n)^T (\lambda + S_n)^{-1} (\mathbf{z}_m - \boldsymbol{\mu}_n)\right) \\ &= \sigma_f^2 \frac{\sqrt{|\lambda|}}{\sqrt{|\lambda + S_n|}} \exp\left(-\frac{1}{2}(\mathbf{z}_m - \boldsymbol{\mu}_n)^T (\lambda + S_n)^{-1} (\mathbf{z}_m - \boldsymbol{\mu}_n)\right) \\ &= \sigma_f^2 \sqrt{\frac{\prod_{q=1}^Q \lambda_q}{|\lambda + S_n|}} \exp\left(-\frac{1}{2}(\mathbf{z}_m - \boldsymbol{\mu}_n)^T (\lambda + S_n)^{-1} (\mathbf{z}_m - \boldsymbol{\mu}_n)\right) \end{aligned} \quad (\text{S45})$$

$$(\text{S46})$$

The  $\Psi_2$  statistic integral can be solved in the same manner:

$$\begin{aligned} (\Psi_2)_{mm'} &= \int k(\mathbf{x}_n, \mathbf{z}_m) k(\mathbf{z}_{m'}, \mathbf{x}_n) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, S_n) d\mathbf{x}_n \\ &= (\sigma_f^2 Z(\lambda))^2 \int \mathcal{N}(\mathbf{x}_n | \mathbf{z}_m, \lambda) \mathcal{N}(\mathbf{x}_n | \mathbf{z}_{m'}, \lambda) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, S_n) d\mathbf{x}_n \end{aligned} \quad (\text{S47})$$



Here we have to multiply the Gaussians twice:

$$\begin{aligned}
(\Psi_2)_{mm'} &= (\sigma_f^2 Z(\lambda))^2 \int \mathcal{N}(\mathbf{z}_{m'} | \mathbf{z}_m, 2\lambda) \\
&\quad \times \mathcal{N}(\mathbf{x}_n | (\lambda^{-1} + \lambda^{-1})^{-1}(\lambda^{-1} \mathbf{z}_m + \lambda^{-1} \mathbf{z}_{m'}), (\lambda^{-1} + \lambda^{-1})^{-1}) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, S_n) d\mathbf{x}_n \\
&= (\sigma_f^2 Z(\lambda))^2 \int \mathcal{N}(\mathbf{z}_{m'} | \mathbf{z}_m, 2\lambda) \mathcal{N}(\mathbf{x}_n | \frac{\mathbf{z}_m + \mathbf{z}_{m'}}{2}, \frac{\lambda}{2}) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, S_n) d\mathbf{x}_n \\
&= (\sigma_f^2 Z(\lambda))^2 \mathcal{N}(\mathbf{z}_{m'} | \mathbf{z}_m, 2\lambda) \int \mathcal{N}(\mathbf{x}_n | \frac{\mathbf{z}_m + \mathbf{z}_{m'}}{2}, \frac{\lambda}{2}) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_n, S_n) d\mathbf{x}_n \\
&= (\sigma_f^2 Z(\lambda))^2 \mathcal{N}(\mathbf{z}_{m'} | \mathbf{z}_m, 2\lambda) \int \mathcal{N}(\boldsymbol{\mu}_n | \frac{\mathbf{z}_m + \mathbf{z}_{m'}}{2}, \frac{\lambda}{2} + S_n) \mathcal{N}(\mathbf{x}_n | \text{mean}, \text{precision}) d\mathbf{x}_n \\
&= (\sigma_f^2 Z(\lambda))^2 \mathcal{N}(\mathbf{z}_{m'} | \mathbf{z}_m, 2\lambda) \mathcal{N}(\boldsymbol{\mu}_n | \frac{\mathbf{z}_m + \mathbf{z}_{m'}}{2}, \frac{\lambda}{2} + S_n) \int \mathcal{N}(\mathbf{x}_n | \text{mean}, \text{precision}) d\mathbf{x}_n \\
&= (\sigma_f^2 Z(\lambda))^2 \mathcal{N}(\mathbf{z}_{m'} | \mathbf{z}_m, 2\lambda) \mathcal{N}(\boldsymbol{\mu}_n | \frac{\mathbf{z}_m + \mathbf{z}_{m'}}{2}, \frac{\lambda}{2} + S_n) \\
&= \sigma_f^4 (2\pi)^Q \left( \prod_{q=1}^Q \lambda_q \right) \mathcal{N}(\mathbf{z}_{m'} | \mathbf{z}_m, 2\lambda) \mathcal{N}(\boldsymbol{\mu}_n | \frac{\mathbf{z}_m + \mathbf{z}_{m'}}{2}, \frac{\lambda}{2} + S_n) \quad (\text{S48})
\end{aligned}$$

For the case of the diagonal covariance  $S_n$  the  $\Psi_3$  statistic looks simpler [2] as it does not require inversion of the  $\lambda + S_n$  and  $\frac{\lambda}{2} + S_n$  matrices.

## 4 Relationship between ELBO and MSE

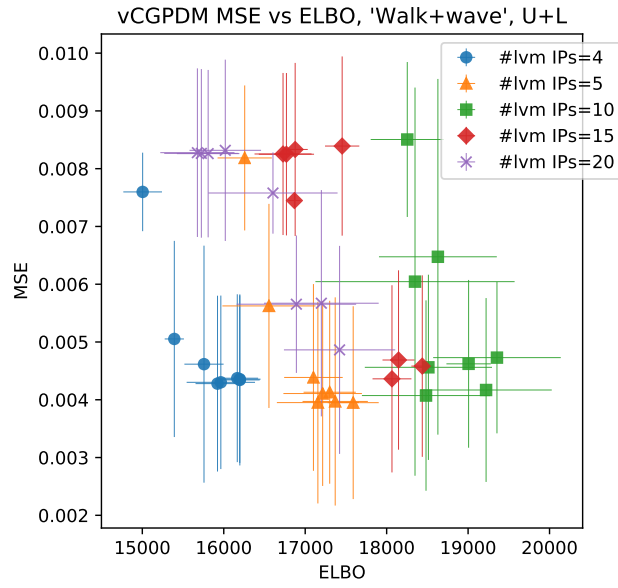


Figure 1: Mean-squared kinematics error, MSE, vs. evidence lower bound, ELBO, for different number of LVM IPs, indicated by symbols. The negative correlation between ELBO and MSE for a given number of LVM IPs is clearly visible. Furthermore, note that the highest ELBO corresponds to an MSE that is very close to the optimal (lowest) one, i.e. ELBO can be used for model selection.

## 5 Learning time of vCGPDM and MAP CGPDM

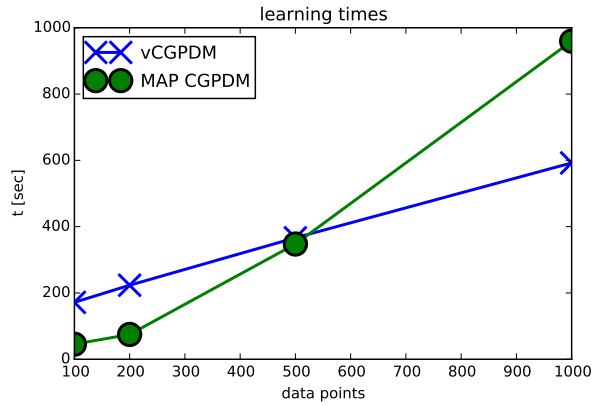


Figure 2: Learning times, including Theano compile times, for a three-part vCGPDM with 10 IPs for all parts and corresponding MAP CGPDM. The linear learning time scaling of the vCGPDM is evident, whereas the MAP CGPDM shows cubic scaling. Thus, the vCGPDM can be used on large data sets which is infeasible for the MAP CGPDM.

## References

- [1] Michalis K. Titsias. Variational learning of inducing variables in sparse gaussian processes. In David A. Van Dyk and Max Welling, editors, *AISTATS*, volume 5 of *JMLR Proceedings*, pages 567–574. JMLR.org, 2009.
- [2] Michalis K. Titsias and Neil D. Lawrence. Bayesian gaussian process latent variable model. In *Proc. 13th AISTATS*, pages 844–851, 2010.