
Learning GPLVM with arbitrary kernels using the unscented transformation:

Supplementary Materials

A THE BAYESIAN GPLVM EVIDENCE LOWER BOUND

The unscented transformation and the competing approximations are applied to terms in the evidence lower bound of the Bayesian GPLVM model derived in Titsias and Lawrence (2010). This section will present a definition of the Bayesian GPLVM, the variational distributions, the evidence lower bound, and where each approximation is applied.

A.1 Prior and variational distributions

The model is defined as:

$$\begin{aligned} \mathbf{Z} &\in \mathbb{R}^{m \times D_x} \\ p(\mathbf{X}) &= \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i \mid \mathbf{0}, \mathbf{I}), \\ p(\mathbf{U} \mid \mathbf{Z}) &= \prod_{d=1}^{D_y} \mathcal{N}(\mathbf{u}_{:d} \mid \mathbf{0}, \mathbf{K}_u), \\ p(\mathbf{F} \mid \mathbf{X}, \mathbf{U}, \mathbf{Z}) &= \prod_{d=1}^{D_y} \mathcal{N}(\mathbf{f}_{:d} \mid \mathbf{K}_{fu} \mathbf{K}_u^{-1} \mathbf{u}_{:d}, \mathbf{K}_f - \mathbf{K}_{fu} \mathbf{K}_u^{-1} \mathbf{K}_{fu}^\top), \\ p(\mathbf{Y} \mid \mathbf{F}) &= \prod_{d=1}^{D_y} \mathcal{N}(\mathbf{y}_{:d} \mid \mathbf{f}_{:d}, \sigma^2 \mathbf{I}), \end{aligned}$$

where:

$$\begin{aligned} [\mathbf{K}_f]_{ij} &= k(\mathbf{x}_i, \mathbf{x}_j), \\ [\mathbf{K}_u]_{ij} &= k(\mathbf{z}_i, \mathbf{z}_j), \\ [\mathbf{K}_{fu}]_{ij} &= k(\mathbf{x}_i, \mathbf{z}_j). \end{aligned}$$

The approximate posterior is defined as:

$$q(\mathbf{F}, \mathbf{U}, \mathbf{X}) = p(\mathbf{F} \mid \mathbf{U}, \mathbf{X}) q(\mathbf{U}) q(\mathbf{X}),$$

where:

$$q(\mathbf{X}) = \prod_{i=1}^N \mathcal{N}(\mathbf{x}_i \mid \check{\boldsymbol{\mu}}_i, \text{Diag}(\check{\boldsymbol{\sigma}}_i^2)),$$

$\text{Diag}(\mathbf{x})$ maps a vector into a corresponding diagonal matrix,

$$q(\mathbf{U}) = \prod_{d=1}^{D_y} \mathcal{N}(\mathbf{u}_{:d} \mid \mathbf{K}_u (\sigma^2 \mathbf{K}_u + \boldsymbol{\Psi}_2)^{-1} \boldsymbol{\Psi}_1^\top \mathbf{y}_{:d}, \sigma^2 \mathbf{K}_u (\sigma^2 \mathbf{K}_u + \boldsymbol{\Psi}_2)^{-1} \mathbf{K}_u).$$

Therefore, the hyperparameters and variational parameters of this model are the output noise variance σ^2 , the kernel hyperparameters, the pseudo-inputs \mathbf{Z} , and the variational parameters $\check{\boldsymbol{\mu}}_i$ and $\check{\boldsymbol{\sigma}}_i^2$ for each i from 1 to n .

A.2 Evidence lower bound

The parameters are learned through joint-optimization of the evidence lower bound (ELBO). The ELBO obtained in Titsias and Lawrence (2010) is:

$$\ln p(\mathbf{Y}) \geq \frac{1}{2} \sum_{d=1}^{D_y} \left[\ln |\mathbf{K}_u| - n \ln(2\pi\sigma^2) - \ln |\mathbf{W}| - \frac{\mathbf{y}_{:d}^\top \mathbf{y}_{:d}}{\sigma^2} + \frac{\mathbf{y}_{:d}^\top \boldsymbol{\Psi}_1 \mathbf{W}^{-1} \boldsymbol{\Psi}_1^\top \mathbf{y}_{:d}}{\sigma^2} - \frac{\psi_0}{\sigma^2} + \frac{\text{tr}(\mathbf{K}_u^{-1} \boldsymbol{\Psi}_2)}{\sigma^2} \right], \quad (1)$$

where:

$$\begin{aligned} \mathbf{W} &= \sigma^2 \mathbf{K}_u + \boldsymbol{\Psi}_2, \\ \psi_0 &= \sum_{i=1}^N \langle k(\mathbf{x}_i, \mathbf{x}_i) \rangle_{q(\mathbf{x}_i)}, \\ [\boldsymbol{\Psi}_1]_{ij} &= \langle k(\mathbf{x}_i, \mathbf{z}_j) \rangle_{q(\mathbf{x}_i)}, \\ [\boldsymbol{\Psi}_2]_{jm} &= \sum_{i=1}^N \langle k(\mathbf{x}_i, \mathbf{z}_j) k(\mathbf{x}_i, \mathbf{z}_m) \rangle_{q(\mathbf{x}_i)}. \end{aligned}$$

The Kullback–Leibler divergence of the approximate posterior with the true posterior is minimized by maximizing the ELBO. Under certain circumstances, when an RBF kernel is used in this model, the Ψ -statistics can be computed analytically. However, for arbitrary kernels, we must apply approximations to compute the expected values.

B UNSCENTED TRANSFORMATION, GAUSS-HERMITE QUADRATURE, AND MONTE CARLO INTEGRATION IN CONTEXT

The only terms that can not be computed analytically on the ELBO are the Ψ -statistics. This section presents how each of the discussed approximation methods can compute an estimate of these terms.

B.1 Unscented transformation

Starting with the unscented transformation, each term is approximated as:

$$\begin{aligned} \psi_0 &\approx \frac{1}{2D_x} \sum_{i=1}^N \sum_{t=1}^{2D_x} k(\mathbf{s}_i^{(t)}, \mathbf{s}_i^{(t)}), \\ [\boldsymbol{\Psi}_1]_{ij} &\approx \frac{1}{2D_x} \sum_{t=1}^{2D_x} k(\mathbf{s}_i^{(t)}, \mathbf{z}_j), \\ [\boldsymbol{\Psi}_2]_{jm} &\approx \frac{1}{2D_x} \sum_{i=1}^N \sum_{t=1}^{2D_x} k(\mathbf{s}_i^{(t)}, \mathbf{z}_j) k(\mathbf{s}_i^{(t)}, \mathbf{z}_m), \end{aligned}$$

where:

$$\mathbf{s}_i^{(t)} = \begin{cases} \check{\boldsymbol{\mu}}_i + \text{Diag}(\sqrt{D_x} \check{\boldsymbol{\sigma}}_i)_{:t} & \text{if } t \in [1, D_x], \\ \check{\boldsymbol{\mu}}_i - \text{Diag}(\sqrt{D_x} \check{\boldsymbol{\sigma}}_i)_{:t} & \text{otherwise.} \end{cases}.$$

These expressions are exactly the empirical mean of the distribution induced by the transformed sigma points.

B.2 Gauss-Hermite quadrature

A simple extension of Gauss-Hermite quadrature can be used to estimate expected values of mappings of multidimensional Gaussian distributed variables. First the number H of evaluation points on each dimension is chosen, then the base H evaluation points and weights are computed from the roots of the Hermite polynomial

$\mathcal{H}_H(x)$. These base evaluation points and weights are organized into a H^{D_x} dimensional grid \mathbf{r}_t and w_t . Finally, the approximation is computed as:

$$\begin{aligned}\psi_0 &\approx \sum_{i=1}^N \sum_{t=1}^{H^{D_x}} w_t \cdot k(\mathbf{s}_i^{(t)}, \mathbf{s}_i^{(t)}), \\ [\Psi_1]_{ij} &\approx \sum_{t=1}^{H^{D_x}} w_t \cdot k(\mathbf{s}_i^{(t)}, \mathbf{z}_j), \\ [\Psi_2]_{jm} &\approx \sum_{i=1}^N \sum_{t=1}^{H^{D_x}} w_t \cdot k(\mathbf{s}_i^{(t)}, \mathbf{z}_j) k(\mathbf{s}_i^{(t)}, \mathbf{z}_m),\end{aligned}$$

where:

$$\mathbf{s}_i^{(t)} = \check{\boldsymbol{\mu}}_i + \text{Diag}(\sqrt{2}\check{\boldsymbol{\sigma}}_i) \mathbf{r}_t.$$

As an concrete example, suppose that $H = 2$ and $D_x = 2$, then, the Hermite polynomial is $\mathcal{H}_2(x) = 4x^2 - 2$, its roots are $\pm \frac{\sqrt{2}}{2}$ with weights $\frac{\sqrt{\pi}}{2}$. Thus, \mathbf{r}_t and w_t are:

$$\begin{aligned}\mathbf{r}_1 &= \left[-\frac{\sqrt{2}}{2}, -\frac{\sqrt{2}}{2} \right], & w_1 &= \frac{\sqrt{\pi}}{2} \cdot \frac{\sqrt{\pi}}{2}, \\ \mathbf{r}_2 &= \left[-\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2} \right], & w_2 &= \frac{\sqrt{\pi}}{2} \cdot \frac{\sqrt{\pi}}{2}, \\ \mathbf{r}_3 &= \left[\frac{\sqrt{2}}{2}, -\frac{\sqrt{2}}{2} \right], & w_3 &= \frac{\sqrt{\pi}}{2} \cdot \frac{\sqrt{\pi}}{2}, \\ \mathbf{r}_4 &= \left[\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2} \right], & w_4 &= \frac{\sqrt{\pi}}{2} \cdot \frac{\sqrt{\pi}}{2}.\end{aligned}$$

B.3 Monte Carlo integration

Finally, we have Monte Carlo integration. First a number of samples T is chosen. Then, the expected values are computed by the empirical mean of each trasformed sample:

$$\begin{aligned}\psi_0 &\approx \frac{1}{T} \sum_{i=1}^N \sum_{t=1}^T k(\mathbf{s}_i^{(t)}, \mathbf{s}_i^{(t)}), \\ [\Psi_1]_{ij} &\approx \frac{1}{T} \sum_{t=1}^T k(\mathbf{s}_i^{(t)}, \mathbf{z}_j), \\ [\Psi_2]_{jm} &\approx \frac{1}{T} \sum_{i=1}^N \sum_{t=1}^T k(\mathbf{s}_i^{(t)}, \mathbf{z}_j) k(\mathbf{s}_i^{(t)}, \mathbf{z}_m),\end{aligned}$$

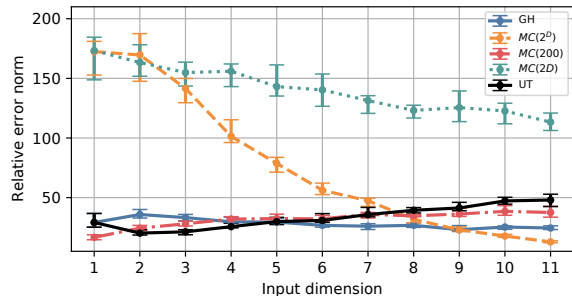
where:

$$\mathbf{s}_i^{(t)} = \check{\boldsymbol{\mu}}_i + \text{Diag}(\check{\boldsymbol{\sigma}}_i) \mathbf{r}_i^{(t)},$$

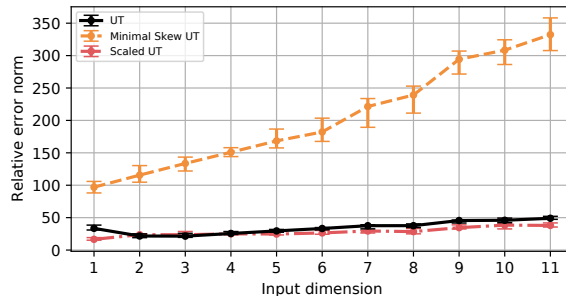
$\mathbf{r}_i^{(t)}$ was sampled from a standard normal distribution.

C PERFORMANCE COMPARISON OF APPROXIMATION METHODS

Alongside the timing experiments from Section 3.3, we also conducted an approximation performance experiment. This experiment's goal is to evaluate how well each approximation evaluates a single Ψ -statistic and how alternatives of the simple symmetric UT might compare to it.



(a) L_1 -norm of the differences between approximations and an MC(10000) estimate. UT maintains a steady accuracy as dimensions increase. Other methods with comparable performance require many more samples.



(b) Uniform UT versus other selections of sigma points and weights. Best parameters for scaled and minimal skew UT were chosen through grid-search and displayed here. Note that our hyper-parameterless UT has a comparable accuracy to best selected scaled UT.

Figure 1: Results of the approximation performance experiment

The tested approximations are the unscented transformation presented in the main paper, Gauss-Hermite with 2 points on each grid for a total of 2^D points, Monte Carlo integration with $2D$, 2^D and 200 samples, and finally, the minimal skew UT and scaled UT as defined by Menegaz et al. (2015).

First, we computed Ψ_1 for a Matérn 3/2 kernel using each approximation on random data ($N = 40, M = 20$) of varying dimensions. Then, we compare the results with an MC approximation that uses 10,000 samples as a proxy for the real value.

In Figure 1a, despite the increasing dimensionality, the UT performs on par with other methods that use many more samples, in contrast with the only method that uses the same amount of samples as it, MC($2D$). As for the comparison between UTs in Figure 1b, our choice of sigma-points outperforms minimal skew UT and does as well as scaled UT. It is worth highlighting that, contrary to its counterparts, this method is hyper-parameter-free.