

Stochastic variational Gaussian process regression

Constructing the ELBO

We are interested in maximizing

$$\log p_{\theta}(\mathbf{y}) = \log \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f})d\mathbf{f} \quad (1)$$

with respect to some set of generative parameters θ (note that we will omit the \cdot_{θ} subscript for notational clarity).

We start by introducing a set of inducing points $\mathbf{z} = \{z_m\}_{m=1}^M$ with corresponding function values (inducing variables) $\mathbf{u} = \{u_m\}_{m=1}^M$ which have a prior distribution $p(\mathbf{u}) = \mathcal{N}(0, \mathbf{K}_{uu})$ and a posterior $p(\mathbf{u}|\mathbf{y})$. We then introduce a variational distribution $q_{\phi}(\mathbf{f}, \mathbf{u}) = p(\mathbf{f}|\mathbf{u})q_{\phi}(\mathbf{u})$, where

$$q_{\phi}(\mathbf{u}) = \mathcal{N}(\mathbf{u}; \mathbf{m}, \mathbf{S}). \quad (2)$$

Here, we have introduced the variational parameters $\phi = \{\mathbf{m}, \mathbf{S}\}$, and we will drop the \cdot_{ϕ} subscript for notational clarity in what follows.

We start by constructing a variational lower bound (ELBO):

$$\log p(\mathbf{y}) \geq \mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\log p(\mathbf{y}|\mathbf{u}, \mathbf{f})] - KL [q(\mathbf{u}, \mathbf{f})||p(\mathbf{u}, \mathbf{f})]. \quad (3)$$

Importantly, maximizing this ELBO is equivalent to minimizing

$$KL [q(\mathbf{f}, \mathbf{u})||p(\mathbf{f}, \mathbf{u}|\mathbf{y})] \quad (4)$$

with respect to the variational parameters ϕ , and this will be important when we compute our predictive distribution in [Equation 23](#).

We then note that $p(\mathbf{y}|\mathbf{u}, \mathbf{f}) = p(\mathbf{y}|\mathbf{f})$. Defining a distribution

$$q(\mathbf{f}) = \int p(\mathbf{f}|\mathbf{u})q(\mathbf{u})d\mathbf{u}, \quad (5)$$

this allows us to integrate out \mathbf{u} in the expectation:

$$\mathbb{E}_{q(\mathbf{u}, \mathbf{f})} [\log p(\mathbf{y}|\mathbf{u}, \mathbf{f})] = \mathbb{E}_{q(\mathbf{f})} [\log p(\mathbf{y}|\mathbf{f})]. \quad (6)$$

We also note that

$$KL [q(\mathbf{u}, \mathbf{f})||p(\mathbf{u}, \mathbf{f})] = \int q(\mathbf{u}, \mathbf{f}) \log \frac{q(\mathbf{u}, \mathbf{f})}{p(\mathbf{u}, \mathbf{f})} d\mathbf{f}d\mathbf{u} \quad (7)$$

$$= \int q(\mathbf{u})p(\mathbf{f}|\mathbf{u}) \log \frac{q(\mathbf{u})p(\mathbf{f}|\mathbf{u})}{p(\mathbf{u})p(\mathbf{f}|\mathbf{u})} d\mathbf{f}d\mathbf{u} \quad (8)$$

$$= \int q(\mathbf{u}) \log \frac{q(\mathbf{u})}{p(\mathbf{u})} d\mathbf{u} \quad (9)$$

$$= KL [q(\mathbf{u})||p(\mathbf{u})]. \quad (10)$$

Putting [Equation 6](#) and [Equation 7](#) together, our ELBO simplifies to

$$\log p(\mathbf{y}) \geq \mathbb{E}_{q(\mathbf{f})} [\log p(\mathbf{y}|\mathbf{f})] - KL [q(\mathbf{u})||p(\mathbf{u})]. \quad (11)$$

Computing the likelihood term

We now turn our attention to the likelihood term $\mathbb{E}_{q(\mathbf{f})} [\log p(\mathbf{y}|\mathbf{f})]$. To evaluate this, we need $q(\mathbf{f})$. We note that everything is jointly Gaussian with a prior given by our kernel

$$p \left(\begin{bmatrix} \mathbf{f} \\ \mathbf{u} \end{bmatrix} \right) = \mathcal{N} \left(0, \begin{bmatrix} \mathbf{K}_{ff} & \mathbf{K}_{fu} \\ \mathbf{K}_{uf} & \mathbf{K}_{uu} \end{bmatrix} \right). \quad (12)$$

From standard Gaussian identities, this gives rise to a conditional distribution

$$p(\mathbf{f}|\mathbf{u}) = \mathcal{N}(\mathbf{f}; \mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{u}, \mathbf{K}_{ff} - \mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{K}_{uf}). \quad (13)$$

We can now compute the marginal distribution $q(\mathbf{f})$ analytically:

$$q(\mathbf{f}) = \int \mathcal{N}(\mathbf{f}; \mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{u}, \mathbf{K}_{ff} - \mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{K}_{uf})\mathcal{N}(\mathbf{u}; \mathbf{m}, \mathbf{S})d\mathbf{u} \quad (14)$$

$$= \mathcal{N}(\mathbf{f}; \mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{m}, \mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{S}\mathbf{K}_{uu}^{-1}\mathbf{K}_{uf} + \mathbf{K}_{ff} - \mathbf{K}_{fu}\mathbf{K}_{uu}^{-1}\mathbf{K}_{uf}). \quad (15)$$

Assuming that the likelihood factorizes (i.e. $p(\mathbf{y}|\mathbf{f}) = \prod_i p(y_i|f_i)$), we can further simplify the likelihood term by writing

$$\mathbb{E}_{q(\mathbf{f})} [\log p(\mathbf{y}|\mathbf{f})] = \sum_i \mathbb{E}_{q(f_i)} [\log p(y_i|f_i)]. \quad (16)$$

We thus only need the marginal distributions $q(f_i)$, which can be computed efficiently (Equation 28). We can then evaluate the 1-dimensional expectations over f_i using either Monte Carlo samples from $q(f_i)$ or with Gauss-Hermite quadrature which has been described extensively elsewhere.

Computing the KL term

We now address the KL term $KL[q(\mathbf{u})||p(\mathbf{u})]$. This has a closed-form expression given by

$$KL[\mathcal{N}(\mathbf{u}; \mathbf{m}, \mathbf{S})||\mathcal{N}(\mathbf{u}; 0, \mathbf{K}_{uu})] = 0.5 (tr(\mathbf{K}_{uu}^{-1}\mathbf{S}) + \mathbf{m}^T \mathbf{K}_{uu}^{-1}\mathbf{m} - k + \log |\mathbf{K}_{uu}| - \log |\mathbf{S}|). \quad (17)$$

Together with Equation 16, this allows us to evaluate the ELBO in a differentiable manner which in turn provides a framework for optimizing the generative parameters θ of our model ($\ell, \sigma_{signal}^2, \sigma_{noise}^2, \dots$).

Inference

To perform inference, we are also interested in the distribution over function values $p(\mathbf{f}^*)$ at some set of test locations \mathbf{x}^* . Recalling that our variational inference procedure gives rise to $q(\mathbf{u}, \mathbf{f}) \approx p(\mathbf{u}, \mathbf{f}|\mathbf{y})$, we can write

$$p(\mathbf{f}^*|\mathbf{y}) = \int p(\mathbf{f}^*, \mathbf{f}, \mathbf{u}|\mathbf{y})d\mathbf{f}d\mathbf{u} \quad (18)$$

$$= \int p(\mathbf{f}^*|\mathbf{f}\mathbf{u})p(\mathbf{f}, \mathbf{u}|\mathbf{y})d\mathbf{f}d\mathbf{u} \quad (19)$$

$$\approx \int p(\mathbf{f}^*|\mathbf{f}\mathbf{u})q(\mathbf{f}, \mathbf{u})d\mathbf{f}d\mathbf{u} \quad (20)$$

$$= \int p(\mathbf{f}^*|\mathbf{f}\mathbf{u})p(\mathbf{f}|\mathbf{u})q(\mathbf{u})d\mathbf{f}d\mathbf{u} \quad (21)$$

$$= \int p(\mathbf{f}^*, \mathbf{f}|\mathbf{u})q(\mathbf{u})d\mathbf{f}d\mathbf{u} \quad (22)$$

$$= \int p(\mathbf{f}^*|\mathbf{u})q(\mathbf{u})d\mathbf{u}. \quad (23)$$

Conveniently, we already know the form of this marginal distribution from Equation 15. It is interesting to note that this predictive distribution does not depend on the data \mathbf{y} but is entirely summarized by the parameters ϕ of our variational distribution $q(\mathbf{u})$.

The ‘whitened’ parameterization

We can simplify the evaluation of our ELBO by parameterizing $q(\mathbf{u}) = \mathcal{N}(\mathbf{u}; \mathbf{K}_{uu}^{\frac{1}{2}}\mathbf{m}, \mathbf{K}_{uu}^{\frac{1}{2}}\mathbf{S}(\mathbf{K}_{uu}^T)^{\frac{1}{2}})$. Here, $\mathbf{K}_{uu}^{\frac{1}{2}}$ can be any matrix factorization with the usual choice being the lower triangular cholesky factor \mathbf{L}_{uu} s.t. $\mathbf{L}_{uu}\mathbf{L}_{uu}^T = \mathbf{K}_{uu}$. Since the KL divergence is invariant to the affine transformation $\mathbf{x} \rightarrow \mathbf{L}^{-1}\mathbf{x}$, we can rewrite KL term from

$$KL[\mathcal{N}(\mathbf{u}; \mathbf{L}\mathbf{m}, \mathbf{L}\mathbf{S}\mathbf{L}^T)||\mathcal{N}(\mathbf{u}; 0, \mathbf{K}_{uu})] \quad (24)$$

to

$$KL[\mathcal{N}(\mathbf{u}; \mathbf{m}, \mathbf{S}) || \mathcal{N}(\mathbf{u}; 0, \mathbf{I})] = 0.5 (tr(\mathbf{S}) + \mathbf{m}^T \mathbf{m} - k - \log |\mathbf{S}|), \quad (25)$$

which gets rid of all terms involving \mathbf{L}_{uu} .

For this whitened parameterization, the distribution $q(\mathbf{f})$ is given by (Equation 15):

$$q(\mathbf{f}) = \mathcal{N}(\mathbf{f}; \mathbf{K}_{fu}(\mathbf{L}^{-1})^T \mathbf{m}, \mathbf{K}_{fu}(\mathbf{L}^{-1})^T \mathbf{S} \mathbf{L}^{-1} \mathbf{K}_{uf} + \mathbf{K}_{ff} - \mathbf{K}_{fu} \mathbf{K}_{uu}^{-1} \mathbf{K}_{uf}). \quad (26)$$

We then write $\mathbf{S} = \mathbf{S}_L \mathbf{S}_L^T$, denote $\mathbf{\Psi} = \mathbf{K}_{fu}(\mathbf{L}^{-1})^T$, and note that for any matrix \mathbf{A}

$$(\mathbf{A}\mathbf{A})_{ii} = \sum_j \mathbf{A}_{ij} \mathbf{A}_{ji}^T = \sum_j \mathbf{A}_{ij}^2. \quad (27)$$

This allows us to compute the marginal variance $\mathbf{\Sigma}_{ii}$ of $q(\mathbf{f})$ as

$$\mathbf{\Sigma}_{ii} = \mathbf{K}_{ff,ii} + \sum_j (\mathbf{\Psi} \mathbf{S}_L)_{ij}^2 - \sum_j \mathbf{\Psi}_{ij}^2. \quad (28)$$

Additionally, the mean parameters $\boldsymbol{\mu}_i$ are easily computed as $(\mathbf{\Psi} \mathbf{m})_i$.

Together, $\boldsymbol{\mu}_i$ and $\mathbf{\Sigma}_{ii}$ fully specify the marginal distributions $q(f_i)$. We also recall from Equation 16 that this is all we need to compute the likelihood term in our ELBO, while the KL term is easily computed according to Equation 25. Finally, these computations are all differentiable, which allows us to perform gradient-based optimization of both the generative parameters θ and the variational parameters ϕ . After convergence, inference is performed according to Equation 23.