

On the paper: Variational Learning of Inducing Variables in Sparse Gaussian Processes (Titsias, 2009)

Thang Bui and Richard Turner

May 25, 2014

Abstract

This summary was prepared for our internal reading club and serves as notes on the sparse GP regression using the variational method (Titsias, 2009). We also discuss why this approximation can be viewed as the corrected version of the Projected Process or Deterministic Training Conditional (DTC) approximation (Seeger, 2003).

Consider the following regression problem: $y_i = f(x_i) + \epsilon_i$, where $f \sim \mathcal{GP}(0, k(x_i, x_j))$ and $\epsilon \sim_{\text{iid}} \mathcal{N}(0, \sigma^2)$. Augment the model with inducing variables \mathbf{u} where $\dim(\mathbf{u}) = m \ll \dim(f) = n$, we have the posterior of \mathbf{f} and \mathbf{u} :

$$p(\mathbf{f}, \mathbf{u} | \mathbf{y}) = \frac{p(\mathbf{f} | \mathbf{u}) p(\mathbf{u}) p(\mathbf{y} | \mathbf{f})}{p(\mathbf{y})}. \quad (1)$$

The reason for this augmentation is that it allows us to produce tractable approximations (in the same spirit with the FITC family). As seen in the above expression, \mathbf{u} can be analytically integrated out to obtain the joint distribution of \mathbf{f} and \mathbf{y} : $p(\mathbf{f}, \mathbf{y}) = \int d\mathbf{u} p(\mathbf{f}, \mathbf{u}, \mathbf{y})$.

The variational approach

- Minimising the KL divergence between $q(\mathbf{f}, \mathbf{u})$ and $p(\mathbf{f}, \mathbf{u} | \mathbf{y})$:

$$\text{KL}(q(\mathbf{f}, \mathbf{u}) || p(\mathbf{f}, \mathbf{u} | \mathbf{y})) = \int d\mathbf{u} d\mathbf{f} q(\mathbf{f}, \mathbf{u}) \log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{f}, \mathbf{u} | \mathbf{y})} \quad (2)$$

$$= \log p(\mathbf{y}) + \int d\mathbf{u} d\mathbf{f} q(\mathbf{f}, \mathbf{u}) \log \frac{q(\mathbf{f}, \mathbf{u})}{p(\mathbf{f}, \mathbf{u}, \mathbf{y})} \quad (3)$$

$$\text{hence, } \log p(\mathbf{y}) = \int d\mathbf{u} d\mathbf{f} q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{f}, \mathbf{u}, \mathbf{y})}{q(\mathbf{f}, \mathbf{u})} + \text{KL}(q(\mathbf{f}, \mathbf{u}) || p(\mathbf{f}, \mathbf{u} | \mathbf{y})), \quad (4)$$

or, $\mathcal{F}(q(\mathbf{u})) = \int d\mathbf{u} d\mathbf{f} q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{f}, \mathbf{u}, \mathbf{y})}{q(\mathbf{f}, \mathbf{u})}$ is the evidence lower bound (ELBO).

- An alternative, equivalent way to obtain the ELBO is by using the Jensen's inequality:

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

$$= \log \int d\mathbf{u} d\mathbf{f} q(\mathbf{f}, \mathbf{u}) \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \quad (6)$$

$$\geq \int d\mathbf{u} d\mathbf{f} q(\mathbf{f}, \mathbf{u}) \log \frac{p(\mathbf{y}, \mathbf{f}, \mathbf{u})}{q(\mathbf{f}, \mathbf{u})} \quad (7)$$

Titsias, 2009 chose a variational distribution $q(\mathbf{f}, \mathbf{u})$ such that,

$$q(\mathbf{f}, \mathbf{u}) = p(\mathbf{f}|\mathbf{u})q(\mathbf{u}). \quad (8)$$

In general $q(\mathbf{f}, \mathbf{u}) = q(\mathbf{f}|\mathbf{u})q(\mathbf{u})$, here the term $q(\mathbf{f}|\mathbf{u})$ is replaced by $p(\mathbf{f}|\mathbf{u})$ which is the prior conditional distribution. This particular choice means that now the only way \mathbf{y} to affect \mathbf{f} is through \mathbf{u} , as opposed to \mathbf{f} separates \mathbf{u} and \mathbf{y} as in the original model. This also helps us mathematically as the optimal distribution $q(\mathbf{u})$ can be obtained analytically. This however limits the applicability of this tricks to derive an extension or a slightly different approximation. Substitute $q(\mathbf{f}, \mathbf{u})$ into the ELBO,

$$\mathcal{F}(q(\mathbf{u})) = \int d\mathbf{u} d\mathbf{f} p(\mathbf{f}|\mathbf{u})q(\mathbf{u}) \log \frac{p(\mathbf{f}|\mathbf{u})p(\mathbf{u})p(\mathbf{y}|\mathbf{f})}{p(\mathbf{f}|\mathbf{u})q(\mathbf{u})} \quad (9)$$

$$= \int d\mathbf{u} d\mathbf{f} p(\mathbf{f}|\mathbf{u})q(\mathbf{u}) \log \frac{p(\mathbf{u})p(\mathbf{y}|\mathbf{f})}{q(\mathbf{u})} \quad (10)$$

$$= \int d\mathbf{u} d\mathbf{f} p(\mathbf{f}|\mathbf{u})q(\mathbf{u}) \log \frac{p(\mathbf{u})}{q(\mathbf{u})} + \int d\mathbf{u} d\mathbf{f} p(\mathbf{f}|\mathbf{u})q(\mathbf{u}) \log p(\mathbf{y}|\mathbf{f}). \quad (11)$$

To find the optimal form of $q(\mathbf{u})$ that maximises the ELBO, consider the derivative of the ELBO w.r.t $q(\mathbf{u})$ with the addition of the Lagrange multiplier,

$$\frac{d}{dq(\mathbf{u})} \mathcal{F} + \lambda = \int d\mathbf{f} p(\mathbf{f}|\mathbf{u}) [\log p(\mathbf{u}) - \log q(\mathbf{u}) - 1] + \int d\mathbf{f} p(\mathbf{f}|\mathbf{u}) \log p(\mathbf{y}|\mathbf{f}) + \lambda \quad (12)$$

Letting 12 equal 0 gives

$$q(\mathbf{u}) = \frac{p(\mathbf{u})}{Z} \exp \left(\int d\mathbf{f} p(\mathbf{f}|\mathbf{u}) \log p(\mathbf{y}|\mathbf{f}) \right), \quad (13)$$

or $q(\mathbf{u}) = p(\mathbf{u})H(\mathbf{y}, \mathbf{u})/Z$ where $H(\mathbf{y}, \mathbf{u}) = \exp \left(\int d\mathbf{f} p(\mathbf{f}|\mathbf{u}) \log p(\mathbf{y}|\mathbf{f}) \right)$. Substitute this into the ELBO in 10,

$$\mathcal{F}(q(\mathbf{u})) = \int d\mathbf{u} d\mathbf{f} p(\mathbf{f}|\mathbf{u})q(\mathbf{u}) \log \frac{p(\mathbf{u})p(\mathbf{y}|\mathbf{f})}{p(\mathbf{u})H(\mathbf{y}, \mathbf{u})/Z} \quad (14)$$

$$= \int d\mathbf{u} d\mathbf{f} p(\mathbf{f}|\mathbf{u})q(\mathbf{u}) [\log p(\mathbf{y}|\mathbf{f}) - \log H(\mathbf{y}, \mathbf{u}) + \log Z] \quad (15)$$

$$= \int d\mathbf{u} q(\mathbf{u}) \left[\log Z - \log H(\mathbf{y}, \mathbf{u}) + \underbrace{\int d\mathbf{f} p(\mathbf{f}|\mathbf{u}) \log p(\mathbf{y}|\mathbf{f})}_{=0} \right] \quad (16)$$

$$= \log Z \quad (17)$$

Consider the integral inside the exponential in the optimal $q(\mathbf{u})$:

$$M = \int d\mathbf{f} p(\mathbf{f}|\mathbf{u}) \log p(\mathbf{y}|\mathbf{f}) \quad (18)$$

$$= \int d\mathbf{f} \mathcal{N}(\mathbf{f}; K_{\mathbf{f}\mathbf{u}} K_{\mathbf{u}\mathbf{u}}^{-1} \mathbf{u}, K_{\mathbf{f}\mathbf{f}} - K_{\mathbf{f}\mathbf{u}} K_{\mathbf{u}\mathbf{u}}^{-1} K_{\mathbf{u}\mathbf{f}}) \log[\mathcal{N}(\mathbf{y}; \mathbf{f}, \sigma^2 I)] \quad (19)$$

$$= \int d\mathbf{f} \mathcal{N}(\mathbf{f}; \mathbf{A}, \mathbf{B}) \left[-\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{f})^\top I (\mathbf{y} - \mathbf{f}) \right] \quad (20)$$

$$= \int d\mathbf{f} \mathcal{N}(\mathbf{f}; \mathbf{A}, \mathbf{B}) \left[-\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \text{Tr}(\mathbf{y}\mathbf{y}^\top - 2\mathbf{y}\mathbf{f}^\top + \mathbf{f}\mathbf{f}^\top) \right] \quad (21)$$

$$= -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \text{Tr}(\mathbf{y}\mathbf{y}^\top - 2\mathbf{y}\mathbf{A}^\top + \mathbf{A}\mathbf{A}^\top + \mathbf{B}) \quad (22)$$

$$= -\frac{1}{2\sigma^2} \text{Tr}(\mathbf{B}) + \log[\mathcal{N}(\mathbf{y}; \mathbf{A}, \sigma^2 I)], \quad (23)$$

where $\mathbf{A} = K_{\mathbf{f}\mathbf{u}}K_{\mathbf{u}\mathbf{u}}^{-1}\mathbf{u}$ and $\mathbf{B} = K_{\mathbf{f}\mathbf{f}} - K_{\mathbf{f}\mathbf{u}}K_{\mathbf{u}\mathbf{u}}^{-1}K_{\mathbf{u}\mathbf{f}}$. Hence the optimal form of $q(\mathbf{u})$ can be found analytically as follows,

$$q(\mathbf{u}) = \mathcal{N}(K_{mn}(K_{nm}K_{mm}^{-1}K_{mn} + \sigma^2 I)\mathbf{y}, K_{mm} - K_{mn}(K_{nm}K_{mm}^{-1}K_{mn} + \sigma^2 I)K_{nm}) \quad (24)$$

$$= \mathcal{N}(\sigma^{-2}K_{mm}\Sigma K_{mn}\mathbf{y}, K_{mm}\Sigma K_{mm}), \quad (25)$$

where $\Sigma = (K_{mm} + \sigma^{-2}K_{mn}K_{nm})^{-1}$. The lower bound on the marginal likelihood is:

$$\mathcal{F}(q(\mathbf{u})) = \underbrace{\log \mathcal{N}(\mathbf{y}|0, \sigma^2 I + K_{nm}K_{mm}^{-1}K_{mn})}_{\text{DTC log marginal likelihood}} - \underbrace{\frac{1}{2\sigma^2} \text{Tr}(K_{nn} - K_{nm}K_{mm}^{-1}K_{mn})}_{\text{regulariser - avoid overfitting}}. \quad (26)$$

Relationship with the DTC approximation

The likelihood approximation presented in Csató and Opper, 2002; Seeger, 2003 can be justified by choosing a likelihood function $q(\mathbf{y}|\mathbf{u})$ to minimise the KL divergence,

$$q(\mathbf{y}|\mathbf{u}) \leftarrow \arg \min_{q(\mathbf{y}|\mathbf{u})} \text{KL}(q(\mathbf{f}, \mathbf{u}|\mathbf{y})||p(\mathbf{f}, \mathbf{u}|\mathbf{y})), \quad (27)$$

where,

$$q(\mathbf{f}, \mathbf{u}|\mathbf{y}) = \frac{q(\mathbf{y}|\mathbf{u})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{y})}, \quad (28)$$

$$p(\mathbf{f}, \mathbf{u}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{p(\mathbf{y})}, \quad (29)$$

$$\text{and, } q(\mathbf{y}) = \int d\mathbf{u}p(\mathbf{u})p(\mathbf{y}|\mathbf{u}). \quad (30)$$

Consider the likelihood $q(\mathbf{y}|\mathbf{u})$ that is a normalised Gaussian or $\int d\mathbf{y}q(\mathbf{y}|\mathbf{u}) = 1$ Seeger, 2003, combining the reversed KL divergence and the normalisation assumption above gives us the Lagrangian:

$$\mathcal{L} = \text{KL}(q(\mathbf{f}, \mathbf{u}|\mathbf{y})||p(\mathbf{f}, \mathbf{u}|\mathbf{y})) + \lambda \left(\int d\mathbf{y}q(\mathbf{y}|\mathbf{u}) - 1 \right) \quad (31)$$

$$= \log p(\mathbf{y}) - \log q(\mathbf{y}) + \int d\mathbf{f} d\mathbf{u} \frac{q(\mathbf{y}|\mathbf{u})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{y})} \log \frac{q(\mathbf{y}|\mathbf{u})}{p(\mathbf{y}|\mathbf{f})} + \lambda \left(\int d\mathbf{y}q(\mathbf{y}|\mathbf{u}) - 1 \right). \quad (32)$$

The derivative \mathcal{L} w.r.t $q(\mathbf{y}|\mathbf{u})$ is,

$$\frac{\partial}{\partial q(\mathbf{y}|\mathbf{u})} = \frac{1}{q(\mathbf{y})}p(\mathbf{u}) + \int d\mathbf{f} \frac{p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{y})} \log \frac{q(\mathbf{y}|\mathbf{u})}{p(\mathbf{y}|\mathbf{f})} + \int d\mathbf{f} \frac{q(\mathbf{y}|\mathbf{u})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})}{q(\mathbf{y})} \frac{1}{q(\mathbf{y}|\mathbf{u})} + \lambda \quad (33)$$

$$= \frac{p(\mathbf{u})}{q(\mathbf{y})} \left[\log q(\mathbf{y}|\mathbf{u}) - \int d\mathbf{f} p(\mathbf{f}|\mathbf{u}) \log p(\mathbf{y}|\mathbf{f}) \right] + \lambda. \quad (34)$$

Setting 34 to zero gives,

$$q(\mathbf{y}|\mathbf{u}) = \exp \left(-\frac{\lambda q(\mathbf{y})}{p(\mathbf{u})} \right) \exp \left(\int d\mathbf{f} p(\mathbf{f}|\mathbf{u}) \log p(\mathbf{y}|\mathbf{f}) \right). \quad (35)$$

Similar as in the variational approach but now we have to use the normalisation constraint, the optimal form for the likelihood approximation is:

$$\boxed{q(\mathbf{y}|\mathbf{u}) = \mathcal{N}(\mathbf{y}; K_{\mathbf{f}\mathbf{u}}K_{\mathbf{u}\mathbf{u}}^{-1}\mathbf{u}, \sigma^2 \mathbf{I})} \quad (36)$$

As noted by Snelson, 2007, the same result can be obtained by optimising the KL divergence between the joint models of y , \mathbf{f} and \mathbf{u} : $\text{KL}(q(\mathbf{y}|\mathbf{u})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})||p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{u})p(\mathbf{u}))$.

Let's remove the normalisation constraint of the likelihood term, this equivalently means that the Lagrange multiplier in the above expression is zero, or the optimal likelihood is:

$$q(\mathbf{y}|\mathbf{u}) = \exp\left(-\frac{1}{2\sigma^2} \text{Tr}(K_{\mathbf{ff}} - K_{\mathbf{fu}}K_{\mathbf{uu}}^{-1}K_{\mathbf{uf}})\right) \mathcal{N}(y; K_{\mathbf{fu}}K_{\mathbf{uu}}^{-1}\mathbf{u}, \sigma^2 \mathbf{I}) \quad (37)$$

Here it becomes clear that why the expression of the posterior of \mathbf{u} in DTC is exactly the same as in Titsias, 2009 and only the approximate marginal likelihood are different. Both are optimising the same KL divergence under the approximate likelihood regime, but Titsias, 2009 allows a free form for the likelihood (which turns out to be easily computed analytically) as opposed to a Gaussian likelihood in Seeger, 2003.

Open questions: 1. tighter bound, how biased is the variational approach is in learning? 2. are \mathbf{u} truly variational parameters?

References

- Csató, Lehel and Opper, Manfred (2002). "Sparse on-line Gaussian processes". In: *Neural computation* 14.3, pp. 641–668.
- Seeger, Matthias (2003). "Bayesian Gaussian process models: PAC-Bayesian generalisation error bounds and sparse approximations". PhD thesis. School of Informatics, College of Science and Engineering, University of Edinburgh.
- Snelson, Edward (2007). "Flexible and efficient Gaussian process models for machine learning". PhD thesis. Gatsby Computational Neuroscience Unit, University College London.
- Titsias, Michalis K. (2009). "Variational Learning of Inducing Variables in Sparse Gaussian Processes". In: *International Conference on Artificial Intelligence and Statistics*, pp. 567–574.