

Sparse Gaussian Processes for Bayesian Optimization

Supplementary Material

March 1, 2016

Here we derive the result (14), the weighted KL divergence. Consider two sparse Gaussian processes, $GP \sim GP(\boldsymbol{\alpha}, C)$, $\widehat{GP} \sim GP(\hat{\boldsymbol{\alpha}}, \hat{C})$, which share a covariance function K . We can therefore write

$$GP \sim \mathcal{N}(\Phi\boldsymbol{\alpha}, I_{\mathcal{F}} + \Phi C \Phi^{\top}) \equiv \mathcal{N}(\boldsymbol{\mu}, \Sigma) , \quad (1)$$

and likewise for \widehat{GP} , where Φ is again the feature space representation of the inducing variables, which are shared. This assumption is made without loss of generality, since the inducing variables of each GP can simply be concatenated into a combined representation. In such a representation, if an inducing variable is only used by one of the GPs, the other will have zeros in the corresponding entries of $\boldsymbol{\alpha}$ and C . This is seen for example in online selection of inducing variables, in which one GP has $m + 1$ inducing variables and the other uses a subset of m of these.

We restate Equation 10 for convenience:

$$D_{KL}^f(P||Q) = \int_{\mathcal{F}} P(\mathbf{x}) \log \left(\frac{P(\mathbf{x})}{Q(\mathbf{x})} \right)^{f(\mathbf{x})} d\mathbf{x} = \int_{\mathcal{F}} f(\mathbf{x}) P(\mathbf{x}) \log \frac{P(\mathbf{x})}{Q(\mathbf{x})} d\mathbf{x} . \quad (2)$$

As before, we define $f^*(\mathbf{x})$ to be the prediction of GP at a point \mathbf{x} in feature space:

$$f^*(\mathbf{x}) = \frac{\mathbf{x}^{\top} \Phi \boldsymbol{\alpha}}{|y^*|} = \frac{\boldsymbol{\mu}^{\top} \mathbf{x}}{|y^*|} .$$

Note that the weighting function f^* uses the full GP prediction rather than the reduced GP, and is therefore constant through the optimization of $D_{KL}^{f^*}(GP||\widehat{GP})$ with respect to \widehat{GP} .

Abusing notation slightly, we write $GP(\mathbf{x})$ to denote the evaluation of the normal distribution corresponding to GP at a point \mathbf{x} in the feature space, and likewise for \widehat{GP} . Our goal is then to find

$$D_{KL}^{f^*}(GP||\widehat{GP}) = \int_{\mathcal{F}} f(\mathbf{x}) GP(\mathbf{x}) \log \frac{GP(\mathbf{x})}{\widehat{GP}(\mathbf{x})} d\mathbf{x} . \quad (3)$$

In the following, we will use $D_{KL}^{f^*}$ as a shorthand for $D_{KL}^{f^*}(GP||\widehat{GP})$.

Evaluating the log term based on the normal distributions of GP and \widehat{GP} , this becomes

$$D_{KL}^{f^*} = \int_{\mathcal{F}} f(\mathbf{x}) GP(\mathbf{x}) \left[\frac{1}{2} (\mathbf{x} - \hat{\boldsymbol{\mu}})^\top \hat{\Sigma}^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) - \frac{1}{2} \log(|\Sigma \hat{\Sigma}^{-1}|) \right] d\mathbf{x} . \quad (4)$$

Taking expectations over $GP \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$, we have

$$\begin{aligned} 2D_{KL}^{f^*} &= E[f^*(\mathbf{x})(\mathbf{x} - \hat{\boldsymbol{\mu}})^\top \hat{\Sigma}^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}})] - \\ &\quad E[f^*(\mathbf{x})(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})] - \log(|\Sigma \hat{\Sigma}^{-1}|) E[f^*(\mathbf{x})] \\ &\Rightarrow 2|y^*| D_{KL}^{f^*} = E[\boldsymbol{\mu}^\top \mathbf{x} (\mathbf{x} - \hat{\boldsymbol{\mu}})^\top \hat{\Sigma}^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}})] - \\ &\quad E[\boldsymbol{\mu}^\top \mathbf{x} (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})] - \log(|\Sigma \hat{\Sigma}^{-1}|) \|\boldsymbol{\mu}\|^2 . \end{aligned} \quad (5)$$

We use the following formula for Gaussian expectation:

$$\begin{aligned} E[\mathbf{x}(\mathbf{x} - \mathbf{m})^\top M(\mathbf{x} - \mathbf{m})] &= \Sigma M(\boldsymbol{\mu} - \mathbf{m}) + \Sigma M^\top (\boldsymbol{\mu} - \mathbf{m}) \\ &\quad + \text{Tr}[\Sigma M^\top] \boldsymbol{\mu} + \boldsymbol{\mu}(\boldsymbol{\mu} - \mathbf{m})^\top M(\boldsymbol{\mu} - \mathbf{m}) . \end{aligned} \quad (6)$$

Since $\hat{\Sigma}$ is symmetric, this allows us to reduce (5) to

$$\begin{aligned} 2|y^*| D_{KL}^{f^*} &= 2\boldsymbol{\mu}^\top \Sigma \hat{\Sigma}^{-1} (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}) + \text{Tr}[\Sigma \hat{\Sigma}^{-1} - I_{\mathcal{F}}] \|\boldsymbol{\mu}\|^2 \\ &\quad + (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}})^\top \hat{\Sigma}^{-1} (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}) \|\boldsymbol{\mu}\|^2 - \log(|\Sigma \hat{\Sigma}^{-1}|) \|\boldsymbol{\mu}\|^2 . \end{aligned} \quad (7)$$

We now use the matrix inversion lemma to write

$$\hat{\Sigma}^{-1} = (I_{\mathcal{F}} + \Phi \hat{C} \Phi^\top)^{-1} = I_{\mathcal{F}} - \Phi (\hat{C}^{-1} + K)^{-1} \Phi^\top , \quad (8)$$

recalling that $K = \Phi^\top \Phi$. Then we can compute

$$\begin{aligned} \Phi^\top \Sigma \hat{\Sigma}^{-1} \Phi &= \\ &\quad \Phi^\top (I_{\mathcal{F}} + \Phi \hat{C} \Phi^\top) (I_{\mathcal{F}} - \Phi (\hat{C}^{-1} + K)^{-1} \Phi^\top) \Phi \\ &= K - K(\hat{C}^{-1} + K)^{-1} K + K C K - K C K (\hat{C}^{-1} + K)^{-1} K \\ &= (I + K C) (\hat{C} + K^{-1})^{-1} , \end{aligned} \quad (9)$$

where in the last line we have used the matrix inversion lemma again in the reverse direction. Likewise, we have

$$\Phi^\top \hat{\Sigma}^{-1} \Phi = K - K(\hat{C}^{-1} + K)^{-1} K = (\hat{C} + K^{-1})^{-1} \equiv \hat{V} , \quad (10)$$

In the same way, we can then compute

$$\text{Tr}[\Sigma \hat{\Sigma}^{-1} - I_{\mathcal{F}}] = \text{Tr}[(C - \hat{C}) \hat{V}] , \quad \log |\Sigma \hat{\Sigma}^{-1}| = \log |(C + K^{-1}) \hat{V}| , \quad (11)$$

as shown in Csato (2002). Note that $\hat{C} = \hat{V}^{-1} - Q$; substituting then gives us w as defined previously:

$$w \equiv \text{Tr}[\Sigma\hat{\Sigma}^{-1} - I_{\mathcal{F}}] - \log(|\Sigma\hat{\Sigma}^{-1}|) = \text{Tr}[(C+Q)\hat{V} - I] - \log(|(C+Q)\hat{V}|) . \quad (12)$$

We can simplify (7) by substituting for $\boldsymbol{\mu}$, Σ , $\hat{\boldsymbol{\mu}}$, and $\hat{\Sigma}$:

$$\begin{aligned} 2|y^*|D_{KL}^{f*} &= 2\boldsymbol{\alpha}^\top \Phi^\top \Sigma \hat{\Sigma} \Phi (\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}}) + \\ &\quad (\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}})^\top \Phi^\top (I_{\mathcal{F}} - \Phi(\hat{C}^{-1} + K)^{-1}\Phi^\top) \Phi (\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}}) \boldsymbol{\alpha}^\top \Phi^\top \Phi \boldsymbol{\alpha} \\ &\quad + w \boldsymbol{\alpha}^\top \Phi^\top \Phi \boldsymbol{\alpha} \\ &= 2\boldsymbol{\alpha}^\top (I + KC)\hat{V}(\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}}) + (\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}})^\top \hat{V}(\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}}) \boldsymbol{\alpha}^\top K \boldsymbol{\alpha} + w \boldsymbol{\alpha}^\top K \boldsymbol{\alpha} . \end{aligned} \quad (13)$$

We define

$$\Gamma = I + \frac{(I + KC)^\top}{\boldsymbol{\alpha}^\top K \boldsymbol{\alpha}} ,$$

and can then write

$$\begin{aligned} \frac{2|y^*|D_{KL}^{f*}}{\boldsymbol{\alpha}^\top K \boldsymbol{\alpha}} &= 2\boldsymbol{\alpha}^\top (\Gamma^\top - I) \hat{V}(\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}}) + \\ &\quad (\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}})^\top \hat{V}(\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}}) + w \\ &= [2\boldsymbol{\alpha}^\top (\Gamma^\top - I) + (\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}})^\top] \hat{V}(\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}}) + w \\ &= (2\Gamma \boldsymbol{\alpha} - (\boldsymbol{\alpha} + \hat{\boldsymbol{\alpha}}))^\top \hat{V}(\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}}) + w . \end{aligned} \quad (14)$$

This is precisely Equation (14), and so we are done.

Similarly, we can compute the weighted KL divergence with arguments reversed:

$$\frac{2|y^*|D_{KL}^{f*}(\widehat{GP} \| GP)}{\hat{\boldsymbol{\alpha}}^\top K \hat{\boldsymbol{\alpha}}} = (2\hat{\Gamma} \hat{\boldsymbol{\alpha}} + (\hat{\boldsymbol{\alpha}} - 3\boldsymbol{\alpha}))^\top V(\hat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}) + \hat{w} , \quad (15)$$

where

$$\begin{aligned} V &= (C + K^{-1})^{-1}, \hat{\Gamma} = I + \frac{(I + K\hat{C})^\top}{\hat{\boldsymbol{\alpha}}^\top K \hat{\boldsymbol{\alpha}}} , \\ \hat{w} &= \text{Tr}[(\hat{C} + K^{-1})V - I] - \log |(\hat{C} + K^{-1})V| . \end{aligned} \quad (16)$$