

KL Corrected Variational Inference for Gaussian Processes

Nathaniel J. King and Neil D. Lawrence

10th June 2005



Overview

- Variational Inference in Gaussian Processes
- Modified Variational Inference
 - ↳ Probabilistic Point Assimilation (PPA) ‘more tractable’
- KL Correction of the Variational Bound
- Results
- Speculation



Notation

- Labels $\mathbf{y} = [y_1 \dots y_N]^T$.
- Input vector $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_N]^T$
- Gaussian distribution over \mathbf{y} is $N(\mathbf{y}|\boldsymbol{\mu}, \Sigma)$ with mean $\boldsymbol{\mu}$ and covariance Σ .
- Process variable (the function) $\mathbf{f} = [f_1 \dots f_N]^T$ and $\bar{\mathbf{f}} = [\bar{f}_1 \dots \bar{f}_N]^T$.
- The notation $\mathbf{f}_{\setminus n}$ represents the vector without the n th element.



Gaussian Process Graph

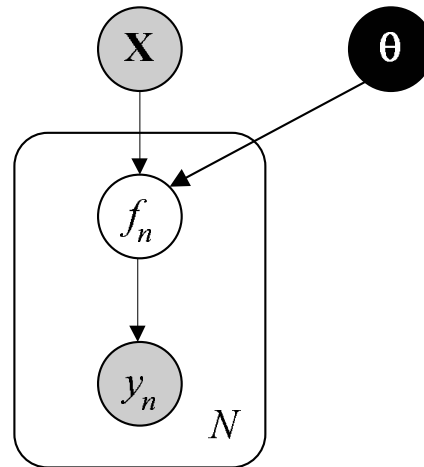


Figure 1: Graphical model of Gaussian process.

$$\log p(\mathbf{y}) = \log \int \prod_{n=1}^N p(y_n | f_n) p(\mathbf{f} | \mathbf{X}, \theta) d\mathbf{f}$$

$p(y_n | f_n)$ is a noise model

$$p(\mathbf{f} | \mathbf{X}, \theta) = N(\mathbf{f} | \mathbf{0}, \mathbf{K})$$

\mathbf{K} is a covariance function parameterised by θ



Variational Inference (Vanilla)

$$\log p(\mathbf{y}) \geq \left\langle \sum_{n=1}^N \log p(y_n | f_n) p(\mathbf{f} | \mathbf{X}, \theta) \right\rangle_{q(\mathbf{f})} - \langle \log q(\mathbf{f}) \rangle_{q(\mathbf{f})}$$
$$q(\mathbf{f}) \propto \prod_{n=1}^N p(y_n | f_n) p(\mathbf{f} | \mathbf{X}, \theta)$$

- Constrain $q(\mathbf{f})$ to be Gaussian — Seeger [2000].
- Constrain covariance of $q(\mathbf{f})$ to have a FA style structure.
- Method is slow and not easily adjusted to new noise models.



Augmented Model — PPA

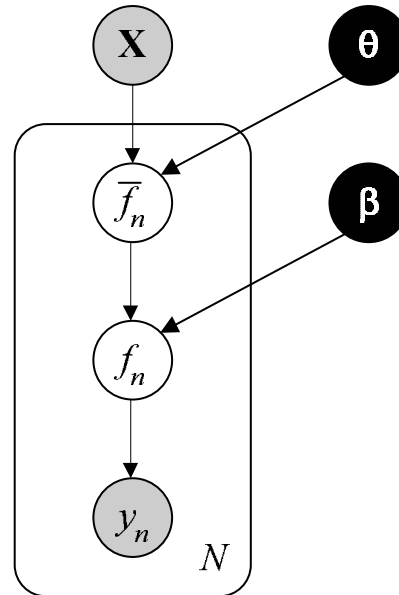


Figure 2: Graph of Augmented Model

$$\log p(\mathbf{y}) = \log \int \prod_{n=1}^N \int p(y_n | f_n) p(f_n | \bar{f}_n, \beta) p(\bar{\mathbf{f}} | \mathbf{X}, \theta) d\mathbf{f}$$

$$p(f_n | \bar{f}_n, \beta) = N(f_n | \bar{f}_n, \beta^{-1})$$



Variational Inference (PPA)

$$\log p(\mathbf{y}) \geq \sum_{n=1}^N \langle \log p(y_n | f_n) p(f_n | \bar{f}_n, \beta) p(\bar{\mathbf{f}} | \mathbf{X}, \theta) \rangle_{\prod_{n=1}^N q(\bar{f}_n) q(\bar{\mathbf{f}})}$$

$$- \sum_{n=1}^N \langle \log q(f_n) \rangle_{q(f_n)} - \langle \log q(\bar{\mathbf{f}}) \rangle_{q(\bar{\mathbf{f}})}$$

Maximised by

$$q(\bar{\mathbf{f}}) \propto \exp \left\langle \sum_{n=1}^N \log p(f_n | \bar{f}_n, \beta) \right\rangle_{\prod_{n=1}^N q(f_n)} p(\bar{\mathbf{f}} | \mathbf{X}, \theta)$$

and

$$q(f_n) \propto \exp \langle \log p(f_n | \bar{f}_n, \beta) \rangle_{q(\bar{\mathbf{f}})} p(y_n | f_n)$$



Expectations of $\log p(f_n | \bar{f}_n, \beta)$

- Since $\exp \langle \log p(f_n | \bar{f}_n, \beta) \rangle_{q(f_n)} \propto N(\langle f_n \rangle | \bar{f}_n, \beta^{-1})$ we have

$$q(\bar{\mathbf{f}}) \propto \prod_{n=1}^N N(\langle f_n \rangle | \bar{f}_n, \beta^{-1}) p(\bar{\mathbf{f}} | \mathbf{X}, \theta)$$

- Since $\exp \langle \log p(f_n | \bar{f}_n, \beta) \rangle_{q(\bar{\mathbf{f}})} \propto N(f_n | \langle \bar{f}_n \rangle, \beta^{-1})$ we have

$$q(f_n) \propto N(f_n | \langle \bar{f}_n \rangle, \beta^{-1}) p(y_n | f_n)$$

- So

- ➔ $q(\bar{\mathbf{f}})$ is a Gaussian process regardless of form of $p(y_n | f_n)$.
- ➔ Moments of $q(f_n)$ are straightforward to compute for any $p(y_n | f_n)$ see *e.g.* *Csató [2002]*



Speed up Variational Method

- Variational Methods can be tediously slow. (yawn!)

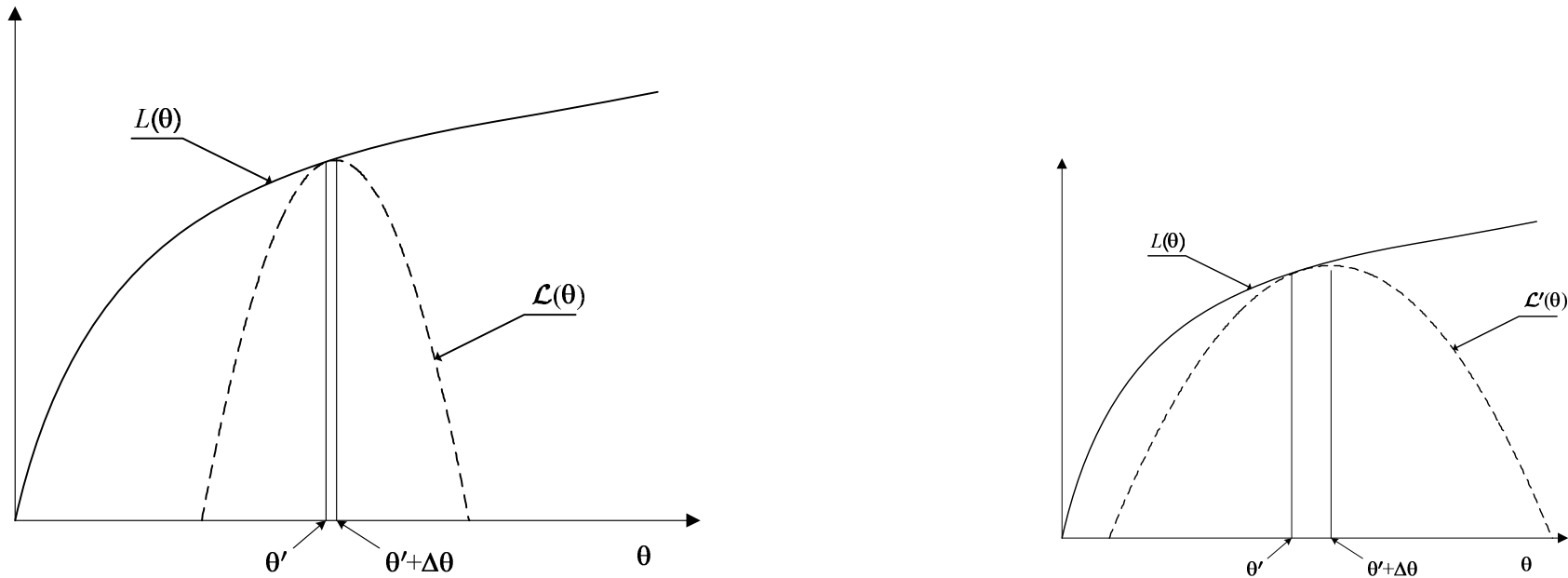


Figure 3: When variational methods are slow.

- Problem occurs when bound's quality degrades rapidly with parameter changes.



KL Corrected Variational Inference

- Updating Parameters: variational lower bound

$$L(\beta, \theta) = \sum_{n=1}^N \langle \log p(f_n | \bar{f}_n, \beta) \rangle_{q(\bar{f}_n)q(f_n)} + \langle \log p(\bar{\mathbf{f}} | \mathbf{X}, \theta) \rangle_{q(\bar{\mathbf{f}})}. \quad (1)$$

- Solution: make the quality of the bound responsive to changes in the parameters.



KL Corrected Variational Inference

- Ideally we would like to optimise the marginal likelihood,

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \beta) = L(\boldsymbol{\theta}) = \log \int \prod_{n=1}^N p(y_n|\bar{f}_n, \beta) p(\bar{\mathbf{f}}|\mathbf{X}, \boldsymbol{\theta}) d\bar{\mathbf{f}}, \quad (2)$$

- Substitute for noise model

$$\begin{aligned} \log p(y_n|\bar{f}_n, \beta) &\geq \langle \log p(y_n|f_n) \rangle_{q(f_n)} + \langle \log p(f_n|\bar{f}_n, \beta) \rangle_{q(f_n)} \\ &\quad - \sum_{n=1}^N \langle \log q(f_n) \rangle_{q(f_n)}, \end{aligned}$$



KL Corrected Lower Bound

- A new lower bound is

$$\log p(\mathbf{y}|\mathbf{X}, \theta, \beta) \geq \log \int \prod_{n=1}^N \exp \langle \log p(f_n | \bar{f}_n, \beta) \rangle_{q(f_n)} p(\bar{\mathbf{f}}|\mathbf{X}, \theta) d\bar{\mathbf{f}} + \text{const}$$

- Which leads to

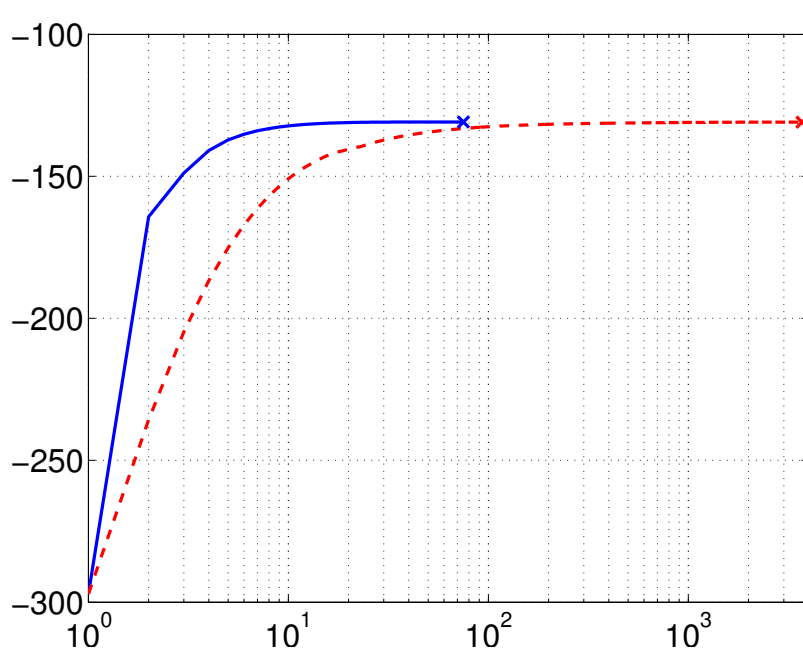
$$\mathcal{L}'(\theta) = \log \int \prod_{n=1}^N N(\langle f_n \rangle | \bar{f}_n, \beta^{-1}) p(\bar{\mathbf{f}}|\mathbf{X}) d\bar{\mathbf{f}} + \text{const},$$

which does not depend on $q(\bar{\mathbf{f}})$.

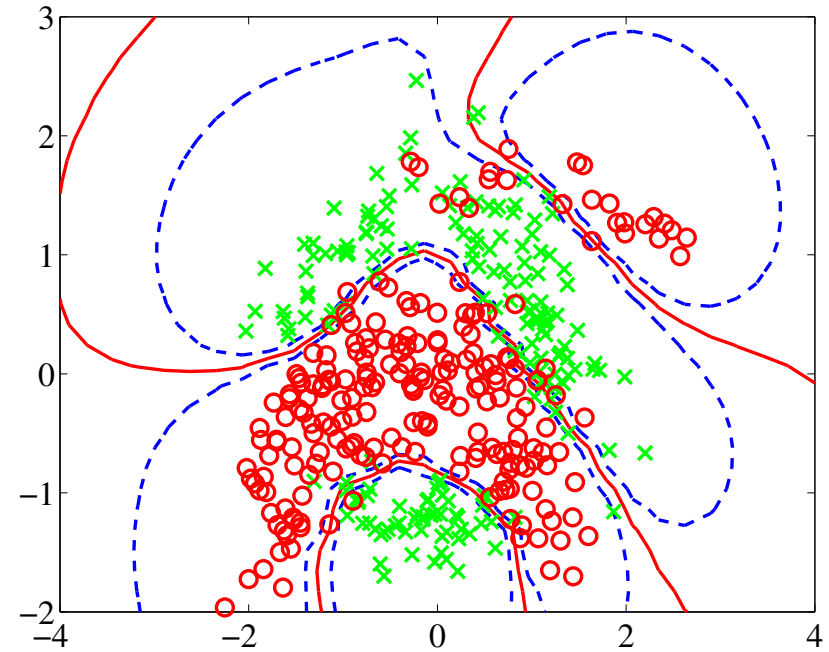
- This is the KL corrected bound.



Convergence Speed



(a)



(b)

Figure 4: (a) Plot of log-likelihood vs iteration number (log-scale) for the KL-corrected objective function (solid line) and the standard variational bound (dashed line). (b) The resulting classification of the **banana** data set.

- KL-corrected requires 74 iterations for convergence, standard variational inference (via PPA) requires 3697 iterations.



Alternative View Point

The marginal likelihood for the augmented model is

$$\log p(\mathbf{y}) = \int \int p(y_n | f_n) p(f_n | \bar{f}_n, \beta) df_n p(\bar{\mathbf{f}} | \mathbf{X}, \theta) d\bar{\mathbf{f}}$$

to make progress we insert a variational lower bound on the inner integral,

$$\begin{aligned} \log \int p(y_n | f_n) p(f_n | \bar{f}_n, \beta) df_n &\geq \langle \log p(y_n | f_n) \rangle_{q(f_n)} \\ &+ \langle \log p(f_n | \bar{f}_n, \beta) \rangle_{q(f_n)} \\ &+ \langle \log q(f_n) \rangle_{q(f_n)} \end{aligned}$$



New Bound

Substituting in this lower bound we have,

$$\begin{aligned}
 \log p(\mathbf{y}) &\geq \log \int \prod_{n=1}^N \exp \langle \log p(f_n | \bar{f}_n, \beta) \rangle p(\bar{\mathbf{f}} | \mathbf{X}, \theta) d\bar{\mathbf{f}} \\
 &\quad + \sum_{n=1}^N \langle \log p(y_n | f_n) \rangle_{q(f_n)} \\
 &\quad + \sum_{n=1}^N \langle \log q(f_n) \rangle_{q(f_n)} \\
 &\doteq \mathcal{L}'(\theta)
 \end{aligned} \tag{3}$$



Minimise directly wrt $q(f_n)$

Bound's dependence on $q(f_n)$ is summarised as

$$\begin{aligned} \mathcal{L}'_n(\theta) &= \frac{1}{2} \left(\beta - \frac{1}{\sigma_n^2} \right) \left(\langle f_n^2 \rangle - \langle f_n \rangle^2 \right) - \langle \log N(f_n | \mu_n, \sigma_n^2) \rangle \\ &\quad + \langle \log p(y_n | f_n) \rangle + \langle \log q(f_n) \rangle_{q(f_n)} + \text{const} \end{aligned}$$

where

$$\mu_n = \mathbf{k}_n^T \left(\mathbf{K}_{\setminus n} + \beta^{-1} \mathbf{I} \right)^{-1} \langle \mathbf{f}_{\setminus n} \rangle$$

and

$$\sigma_n^2 = \beta^{-1} + k_{nn} - \mathbf{k}_n^T \left(\mathbf{K}_{\setminus n} + \beta^{-1} \mathbf{I} \right)^{-1} \mathbf{k}_n$$

where if $\frac{1}{2} \left(\beta - \frac{1}{\sigma_n^2} \right) \left(\langle f_n^2 \rangle - \langle f_n \rangle^2 \right)$ is small then this implies

$$q(f_n) \propto p(y_n | f_n) N(f_n | \mu_n, \sigma_n^2)$$

Which is very similar to the approximating distribution that arises in ... EP



Conclusions

- Variational inference in GPs is practical.
 - ➔ Various noise models can be accommodated.
 - ➔ Slow convergence can be solved.
- Recent (Monday & Tuesday!) analysis suggests connections with EP.



References

Lehel Csató. *Gaussian Processes — Iterative Sparse Approximations*. PhD thesis, Aston University, 2002.

Matthias Seeger. Bayesian model selection for support vector machines, Gaussian processes and other kernel classifiers. In Sara A. Solla, Todd K. Leen, and Klaus-Robert Müller, editors, *Advances in Neural Information Processing Systems*, volume 12, pages 603–609, Cambridge, MA, 2000. MIT Press.

