

# Sparse Gaussian Process Hyperparameters: Optimize or Integrate?

Vidhi Lalchand

Wessel P. Bruinsma

David R. Burt

Carl E. Rasmussen

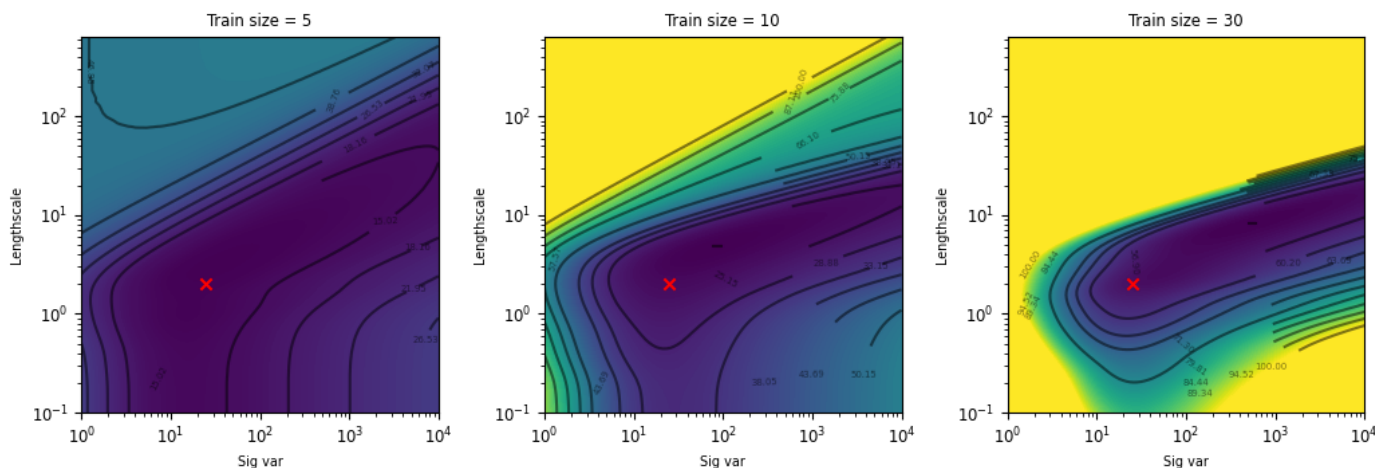
# Motivation

- This work is about Bayesian hyperparameter inference in sparse Gaussian process regression.
- Traditional gradient based optimisation (ML-II) can be extremely sensitive to starting values.
- ML-II hyperparameter estimates are subject to high variability and underestimate prediction uncertainty.
- We propose a novel and computationally efficient scheme Fully Bayesian inference in sparse GPs.

$$\mathbf{y}_n = f(\mathbf{x}_n) + \epsilon_n, \quad \epsilon_n \sim \mathcal{N}(0, \sigma^2), \quad f \sim \mathcal{GP}(0, k_\theta)$$

$$\log p(\mathbf{y}|\boldsymbol{\theta}) = \log \int p(\mathbf{y}|f)p(f|\boldsymbol{\theta})df = \underbrace{c - \frac{1}{2}\mathbf{y}^T (K_\theta + \sigma^2 I)^{-1} \mathbf{y}}_{\text{data fit term}} - \underbrace{\frac{1}{2}|K_\theta + \sigma^2 I|}_{\text{complexity penalty}}$$

Neg. log marginal likelihood (Sig var vs. Lengthscale)



# Mathematical set-up

Inputs  $X = (\mathbf{x}_n)_{n=1}^N \subseteq \mathbb{R}^D$

Latent function prior  $f \sim \mathcal{GP}(0, k_\theta)$

Outputs  $\mathbf{y} = (y_n)_{n=1}^N \subseteq \mathbb{R}$

Factorised Gaussian likelihood  $p(\mathbf{y}|f) = \prod_{n=1}^N \mathcal{N}(y_n|f_n, \sigma^2)$

Inducing variables  $\mathbf{u} = \{f(\mathbf{z}_m)\}_{m=1}^M \subseteq \mathbb{R}$

Inducing locations  $Z = \{\mathbf{z}_m\}_{m=1}^M, \mathbf{z}_m \in \mathbb{R}^D$

---

## Canonical Inference for $\theta$ in sparse GPs:

1. Specify a variational approximation to the posterior over  $(f, \mathbf{u})$
2. Lower bound the GP log-marginal likelihood  $\log p(\mathbf{y}|\theta) \geq \mathcal{L}_{\theta, Z}$
3. Use the closed-form ELBO to learn hyperparameters ( $\theta$ ) and inducing locations ( $Z$ )

Variational approximation to the posterior  $p(f, \mathbf{u}|\mathbf{y}, \theta) \approx q(f, \mathbf{u}|\theta) = p(f|\mathbf{u}, \theta)q(\mathbf{u})$

Hyperparameter inference  $\theta^* \in \arg \max_{\theta, Z} \mathcal{L}_{\theta, Z}$ . ← "Collapsed ELBO"

Titsias [2009] showed that in the case of a Gaussian likelihood the optimal variational distribution  $q^*(\mathbf{u})$  is Gaussian and can be derived in closed-form.

---

## [Ours] Doubly collapsed Inference for $\theta$ in sparse GPs:

1. Specify a variational approximation to the posterior over  $(f, \mathbf{u}, \theta)$  ← "Collapsed ELBO"
2. Lower bound the GP log-marginal likelihood  $\log p(\mathbf{y}) \geq \int q(\theta) \mathcal{L}_{\theta, Z} d\theta - \text{KL}(q(\theta)||p(\theta))$
3. Crucially, we can write down the optimal  $q^*(\theta)$  upto a normalising constant.

Variational approximation to the posterior  $p(f, \mathbf{u}, \theta|\mathbf{y}) \approx q(f, \mathbf{u}, \theta) = p(f|\mathbf{u}, \theta)q(\mathbf{u}|\theta)q(\theta)$

Hyperparameter inference Sample  $\theta^* \sim q^*(\theta)$

# Algorithm

$$\log p(\mathbf{y}) \geq \int q(\boldsymbol{\theta}) \mathcal{L}_{\boldsymbol{\theta}, Z} d\boldsymbol{\theta} - \text{KL}(q(\boldsymbol{\theta}) || p(\boldsymbol{\theta}))$$

Overall, the core training algorithm alternates between two steps:

1. Sampling step for  $\boldsymbol{\theta}$ :  $\theta_j \sim q^*(\boldsymbol{\theta}) \propto \mathcal{L}(\boldsymbol{\theta}, Z_{opt}) + \log p(\boldsymbol{\theta})$ , [Keep  $Z_{opt}$  fixed]

2. Optimisation step for  $Z$ :  $Z_{opt} \leftarrow \text{optim}(\hat{\mathcal{L}})$ , where

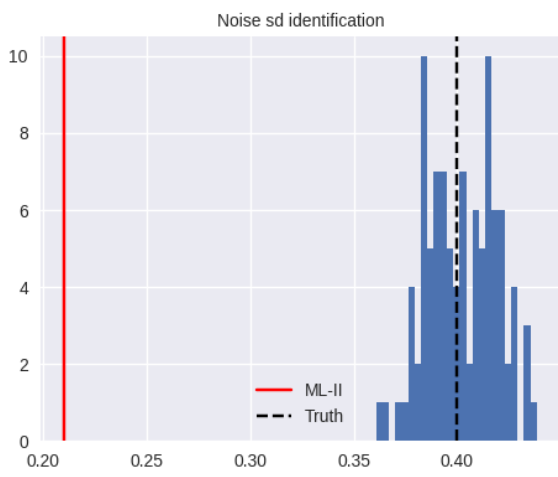
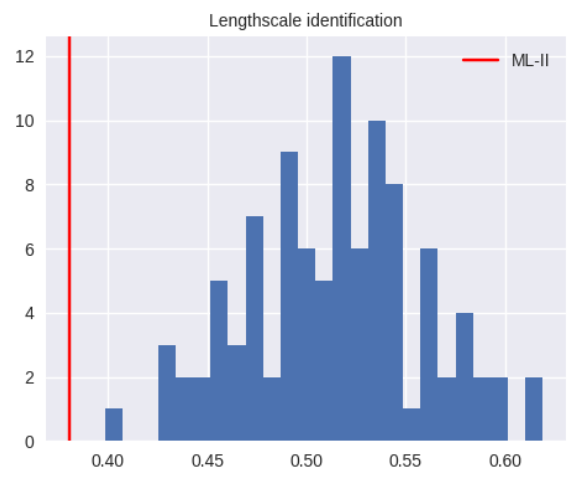
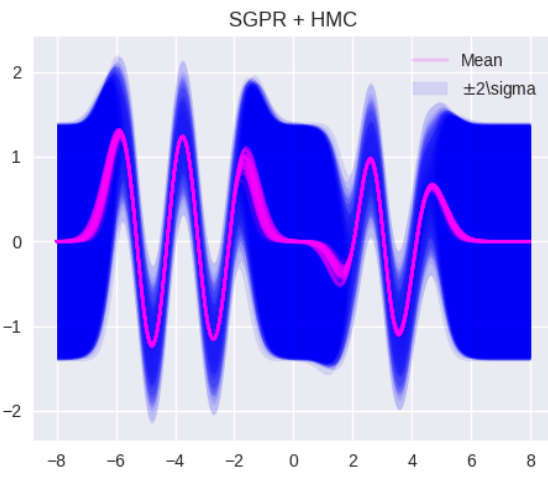
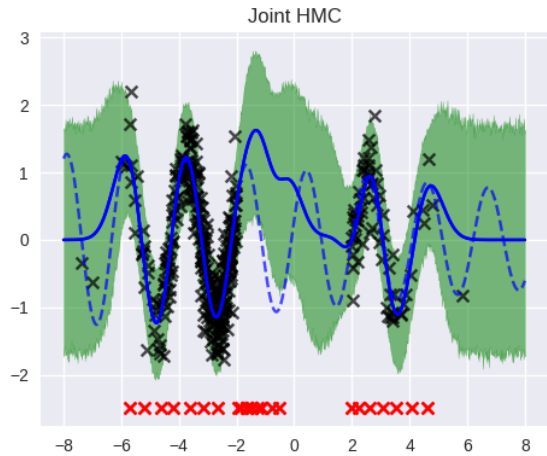
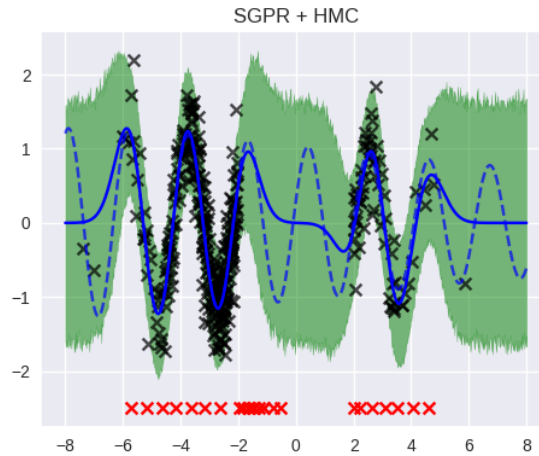
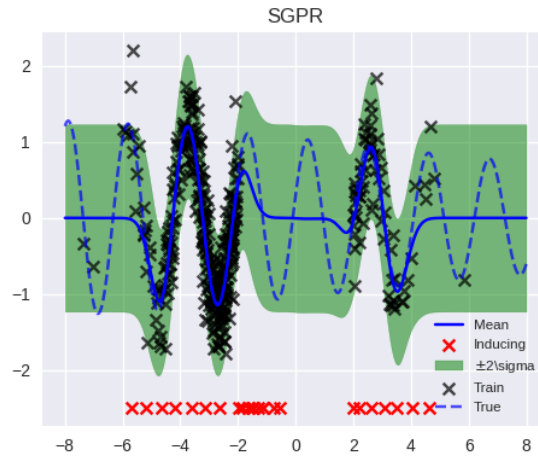
$$\hat{\mathcal{L}} = \mathbb{E}_{q^*(\boldsymbol{\theta})} [\mathcal{L}(\boldsymbol{\theta}, Z)] \approx \frac{1}{J} \sum_{j=1}^J \mathcal{L}(\theta_j, Z_{opt}), \quad [\text{Keep } \boldsymbol{\theta} \text{ fixed}]$$

By sampling from  $q^*(\boldsymbol{\theta})$ , we side-step the need to sample from the joint  $(\mathbf{u}, \boldsymbol{\theta})$ -space yielding a significantly more efficient algorithm in the case of regression with a Gaussian likelihood.

	Approach	Time/it.	Mem./it.	Param / Vars
	Non-collapsed [Hensman et al, 2015]	$m^3$	$m^2$	$n_\theta + m$
	Collapsed (ours)	$nm^2$	$m^2$	$n_\theta$

$n_\theta$  is the number of hyperparameters and  $m$  is the number of inducing variables

# 1d Synthetic Experiment



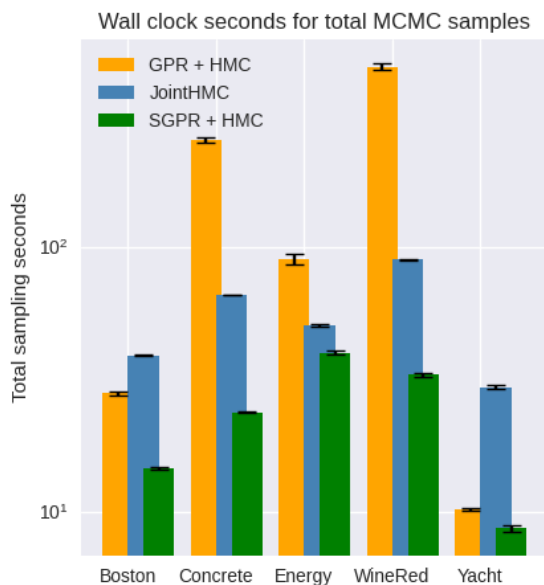
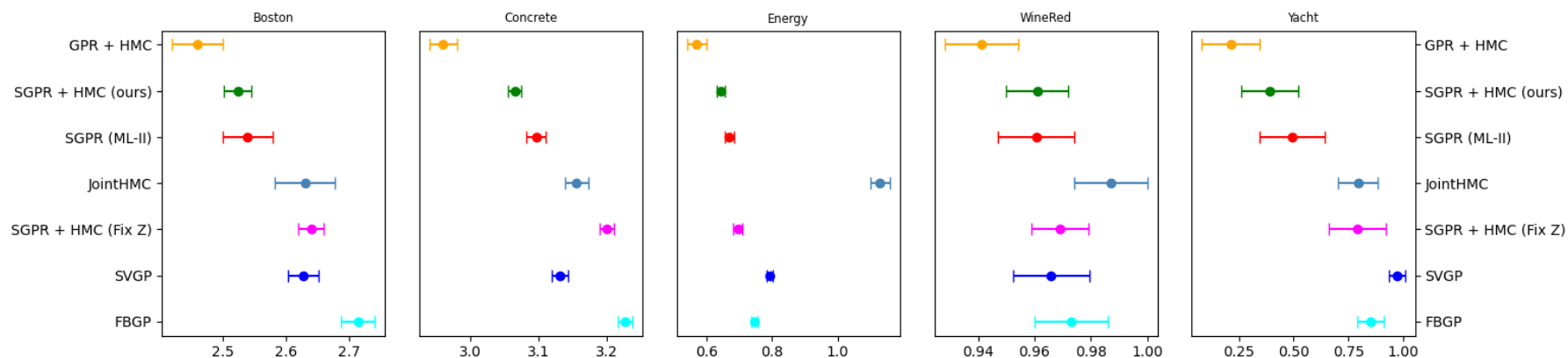
$$f(x) = \sin(3x) + 0.3 \cos(\pi x)$$

with the constraint  $(x < -2)$  and  $(x > 2)$ .

Method	SGPR	SGPR + HMC	JointHMC
RMSE	0.580	<b>0.537</b>	0.682
NLPD	0.214	<b>0.065</b>	0.74

# Sparse GP Benchmarks

Neg. log predictive density (mean  $\pm$  se) on test data, 10 splits.



- Our method, SGPR + HMC (--) outperforms other fully Bayesian benchmarks like jointHMC (--) and FBGP (--) in terms of negative log predictive density on unseen data.
- It is significantly faster relative to jointHMC and Exact GP inference with HMC (--).

Thank you!