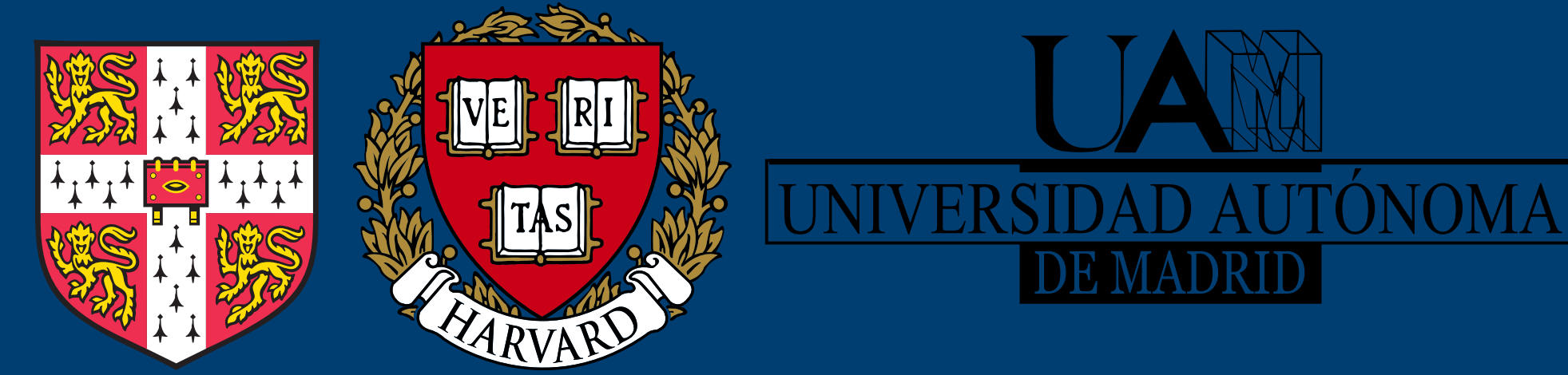


Deep Gaussian Processes for Regression using Approximate Expectation Propagation



Thang D. Bui¹, José Miguel Hernández-Lobato², Daniel Hernández-Lobato³, Yingzhen Li¹, Richard E. Turner¹

University of Cambridge¹, Harvard University², Universidad Autónoma de Madrid³

1 - Contributions

- **Deep GPs for large scale regression:** **scalable** + **flexible** + **calibrated**, largest dataset: **500k** datapoints, **90**-dimensional inputs
- **Novel approx. inference technique:** direct optimisation of the EP energy
- **Rigorous experiments:** an extensive comparison to GPs and Bayesian neural networks, demonstrating state-of-the-art performance of DGPs on many datasets
- **Code:** http://github.com/thangbui/deepGP_approxEP

3 - Deep Gaussian processes for regression

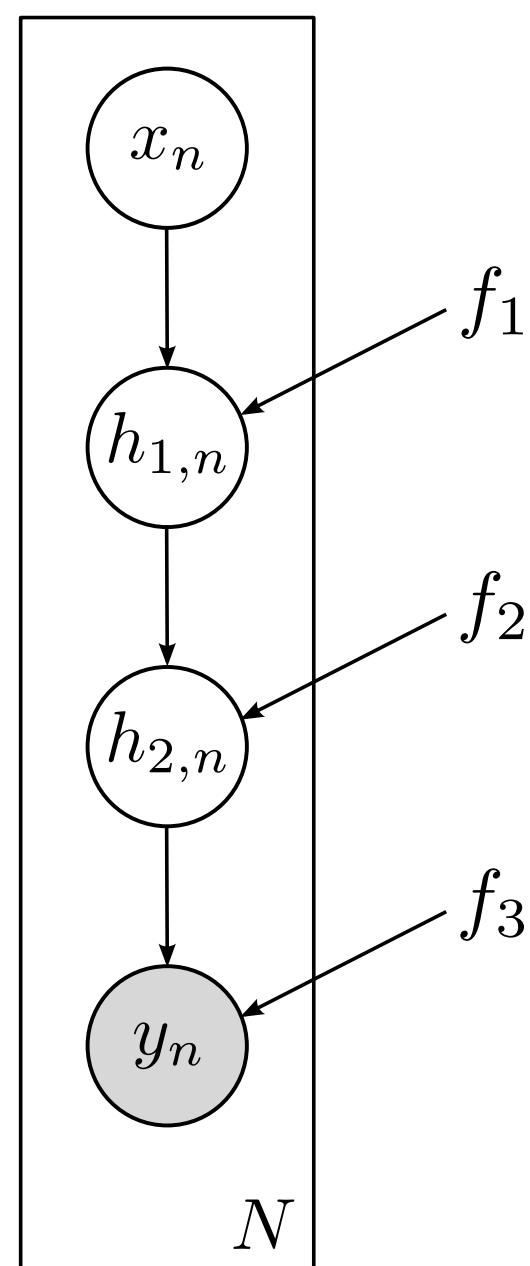
Deep GPs are

- multi-layer generalisation of Gaussian processes, hence,
- equivalent to deep neural networks with infinitely wide hidden layers

$$\begin{aligned} f_l &\sim \mathcal{GP}(0, k(\cdot, \cdot)) \\ h_{l,n} &:= f_l(f_{l-1}(\cdots f_1(\mathbf{x}_n))) \\ y_n &= g(\mathbf{x}_n) = f_L(f_{L-1}(\cdots f_2(f_1(\mathbf{x}_n)))) + \epsilon_n \end{aligned}$$

Advantages of Deep GPs:

- discover useful input warping or dimensionality compression/expansion, i.e. automatic, nonparametric Bayesian kernel design,
- give a non-Gaussian functional mapping g ,
- repair the damage done by using sparse approximations to GPs,
- support approximate Bayesian inference, and
- give better uncertainty estimates.



4 - Direct optimisation of EP energy

Approximate marginal likelihood given by EP:

$$\log p(\mathbf{Y}) \approx \mathcal{F}(\theta) = \phi(\theta) - \phi(\theta_{\text{prior}}) + \sum_{n=1}^N [\log \mathcal{Z}_n + \phi(\theta^n) - \phi(\theta)],$$

where $\log \mathcal{Z}_n = \log \int df q^n(f) p(y_n | f, \mathbf{x}_n)$; $\phi(\theta) = \int df p(f_{\neq \mathbf{u}} | \mathbf{u}) \exp[\Phi^T(\mathbf{u})\theta]$

7 - Comparison with Gaussian processes and Bayesian neural networks

Task: regression on 10 datasets from the UCI repository (8 have 20 train/test splits)
The largest dataset has **500k** datapoints and **90** dimensional inputs.

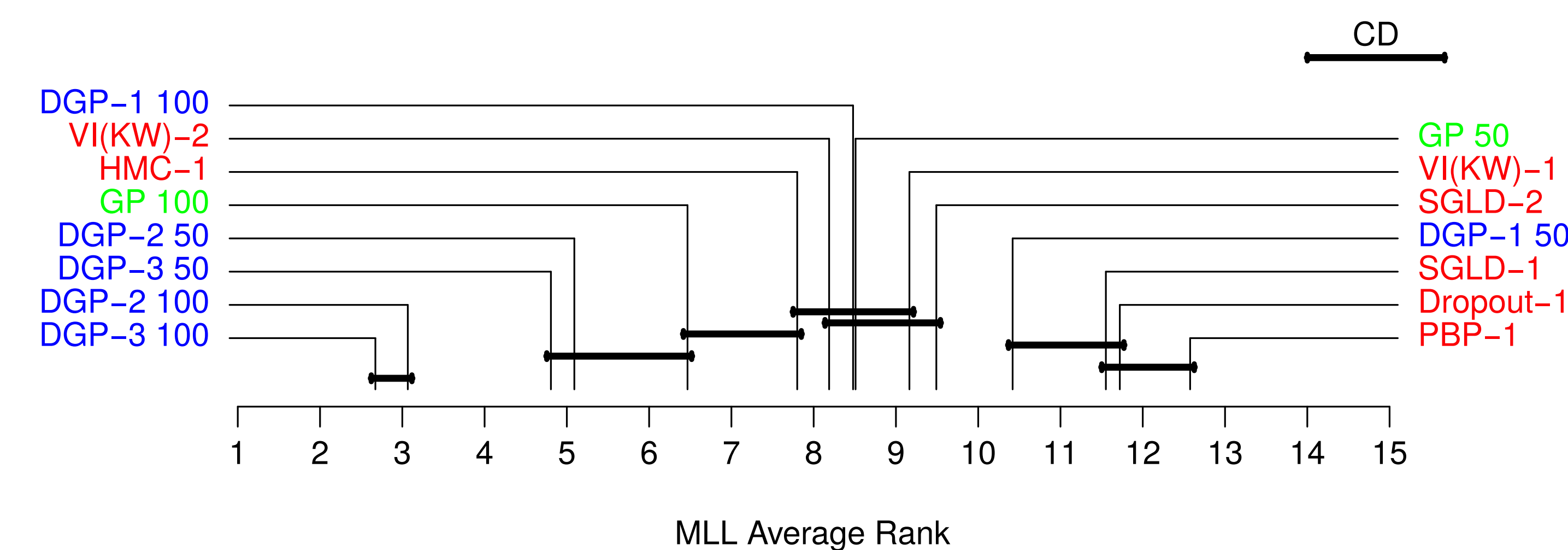


Figure : Rankings of all methods (lower is better).

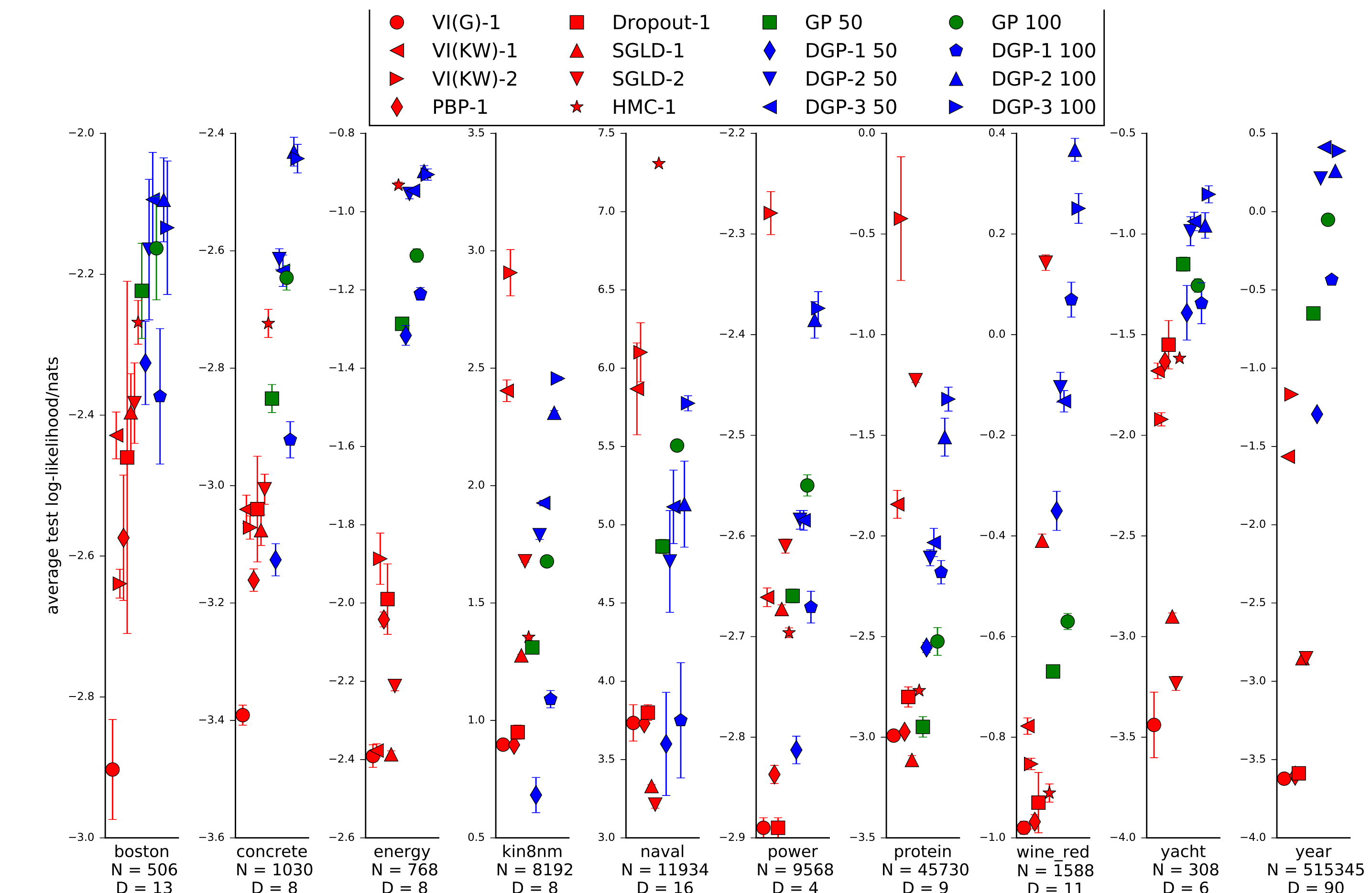


Figure : Average test log-likelihood (higher is better).

2 - A motivating example

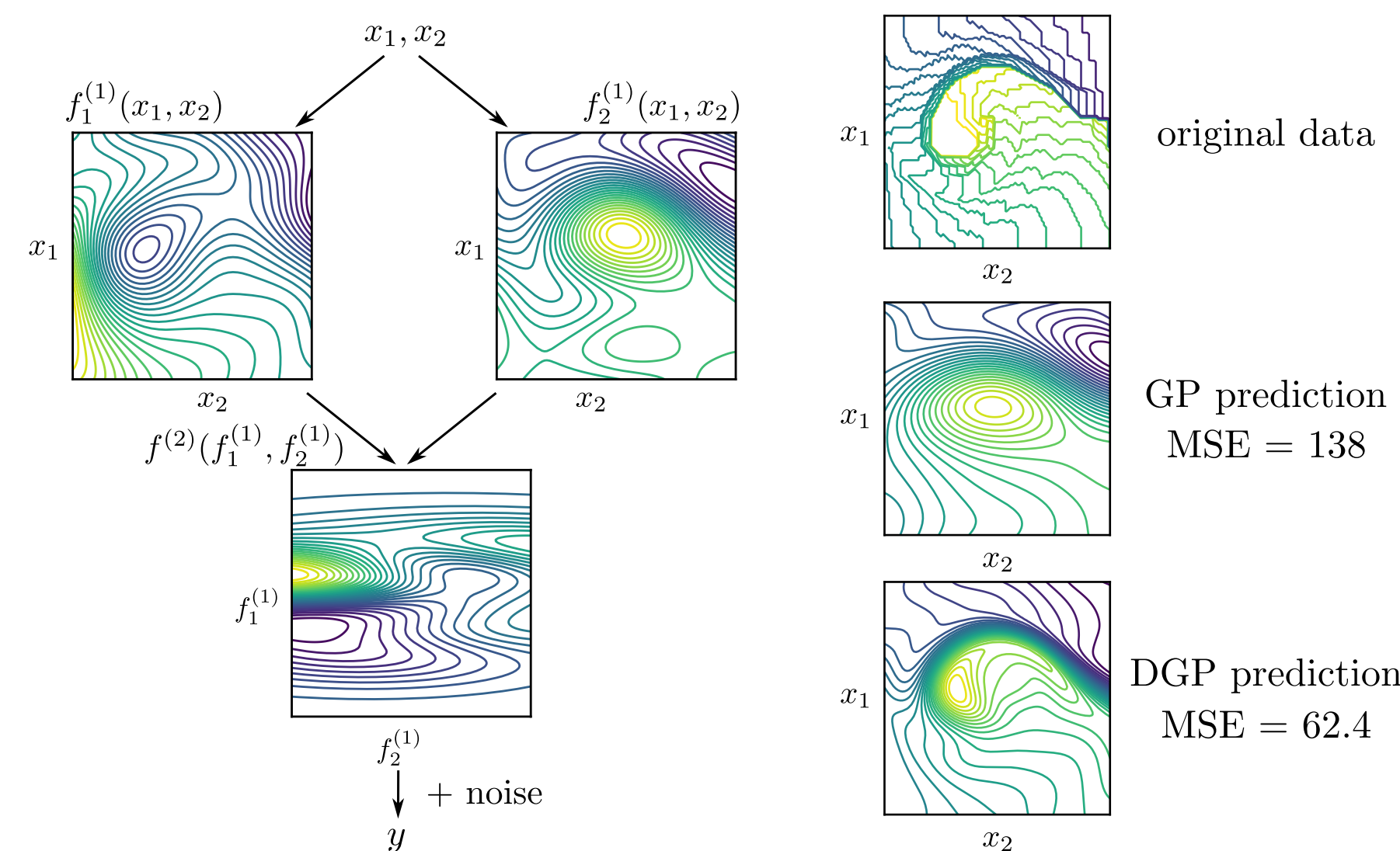


Figure : Fitting the value function of the 'mountain car' problem

5 - Approximate EP with tied approximate factors

The approximate posterior is parameterised by **pseudo-points**, \mathbf{u} and following [Li et al. 2015], we tie the factor approximations:

$$\begin{aligned} p(f|\mathbf{y}) &\propto p(f) \prod_{i=1}^3 t_i(f) \approx q(f) \propto p(f) \prod_{i=1}^3 \tilde{t}_i(\mathbf{u}) \\ p(f|\mathbf{y}) &\propto p(f) \prod_{i=1}^3 t_i(f) \approx q(f) \propto p(f) \prod_{i=1}^3 \tilde{t}^3(\mathbf{u}) \end{aligned}$$

We tie the factor approximations

6 - Nested Gaussian projections

Computing $\log \mathcal{Z}_n$ is challenging due to the uncertainty in the inputs of the GP predictive distribution. We approximate this by a recursive Gaussian projection and the projection for each layer is demonstrated in the figure below.

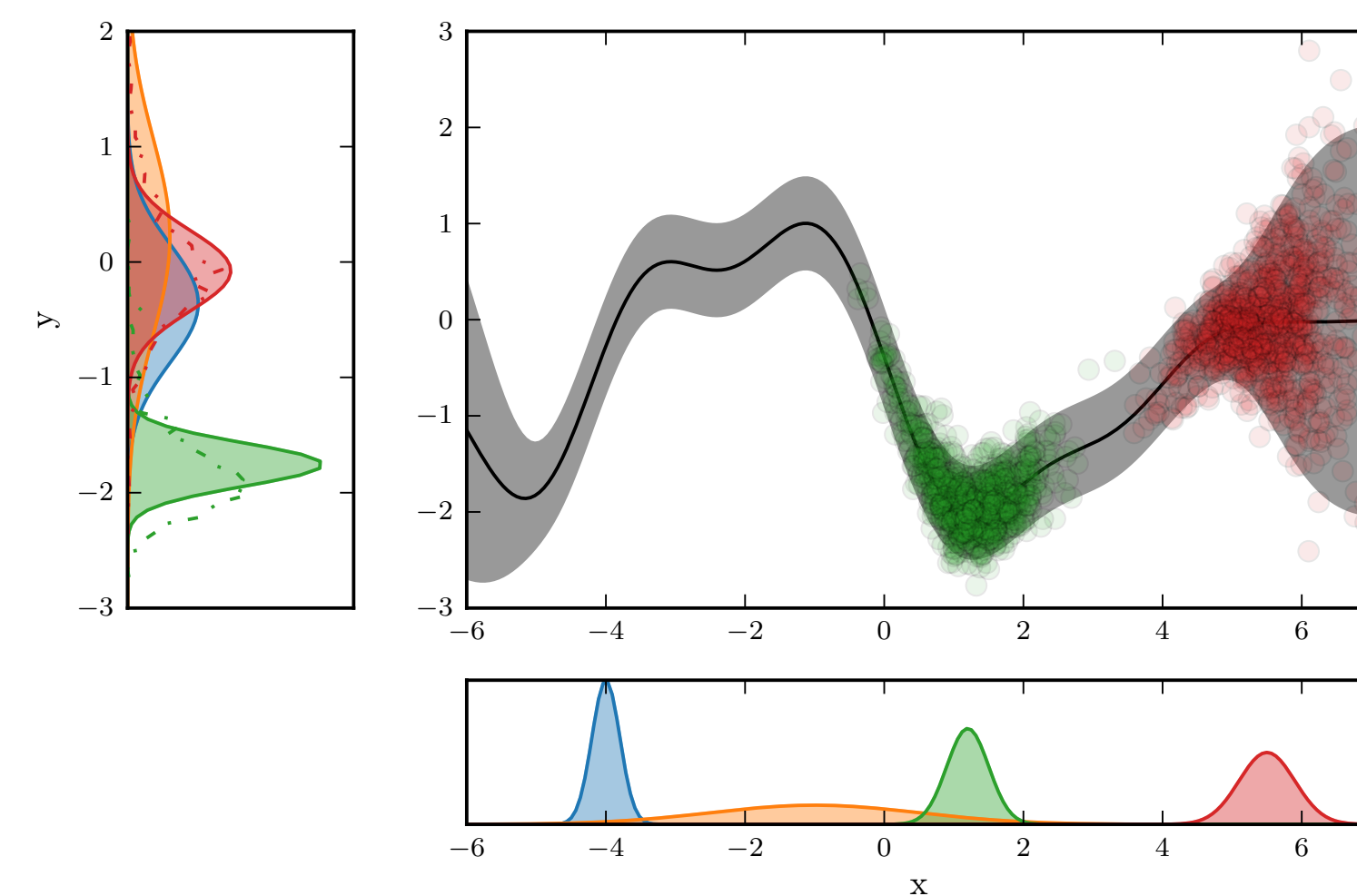
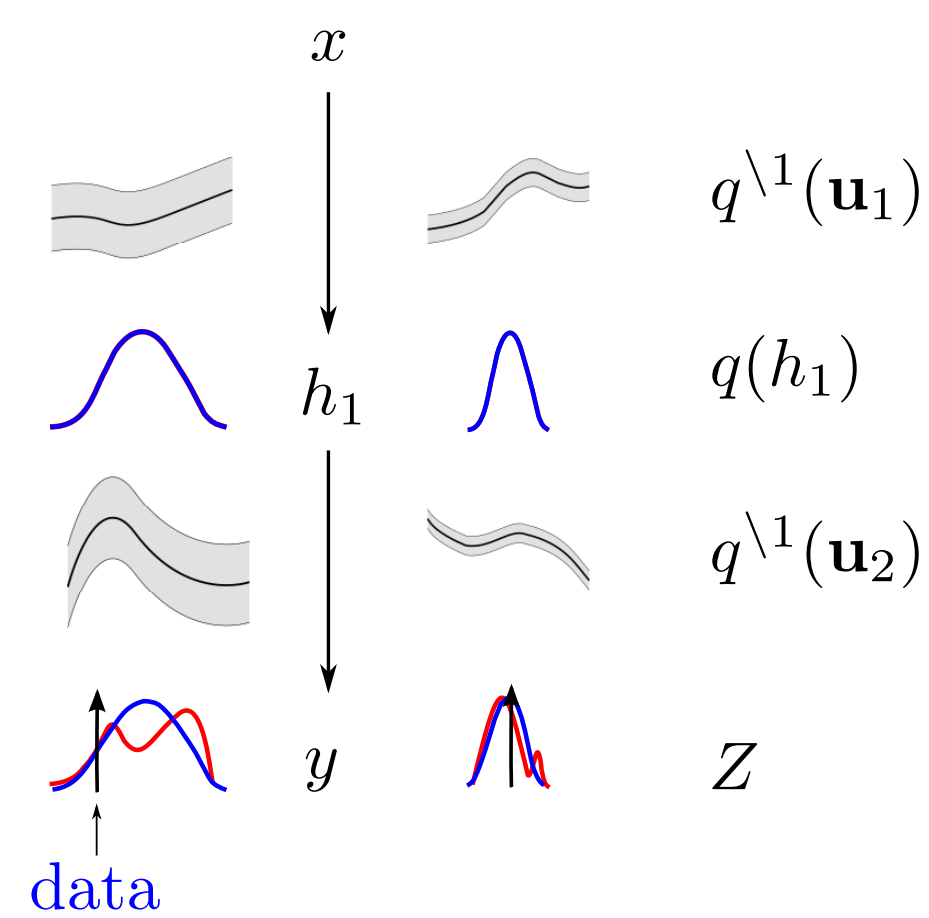


Figure : Gaussian projection: The bottom figure shows the distribution over the inputs of a layer, resulting in non-Gaussian output distributions (left, dash lines) which are then approximated by Gaussians.

The mean and variance of each approximate Gaussian can be computed analytically, for example, for an input distribution $q(h_1)$, the approximate output distribution is $q(h_2) = \mathcal{N}(h_2; m_2, v_2)$ where

$$\begin{aligned} m_2 &= E_{q(h_1)}[m_2|h_1] = E_{q(h_1)}[\mathbf{K}_{h_2, u_2} \mathbf{A}] \\ v_2 &= E_{q(h_1)}[v_2|h_1] + \text{var}_{q(h_1)}[m_2|h_1] \\ &= \sigma_2^2 + E_{q(h_1)}[\mathbf{K}_{h_2, h_2}] + \text{tr}(\mathbf{B} E_{q(h_1)}[\mathbf{K}_{u_2, h_2} \mathbf{K}_{h_2, u_2}]) - m_2^2 \end{aligned}$$

and $\mathbf{A} = \mathbf{K}_{u_2, u_2}^{-1} \mathbf{m}_2^{\setminus 1}$, $\mathbf{B} = \mathbf{K}_{u_2, u_2}^{-1} (\mathbf{V}_2^{\setminus 1} + \mathbf{m}_2^{\setminus 1} \mathbf{m}_2^{\setminus 1, T}) \mathbf{K}_{u_2, u_2}^{-1} - \mathbf{K}_{u_2, u_2}^{-1}$. We repeat the above step for each layer, compute $\log \mathcal{Z}$ at the last layer and find its gradient using **back-propagation**.



8 - Predicting the efficiency of organic photovoltaic molecules

Dataset: 50k/10k training/test points, 512-dim. binary input features. Need **calibrated error-bars** for active learning or Bayesian optimisation

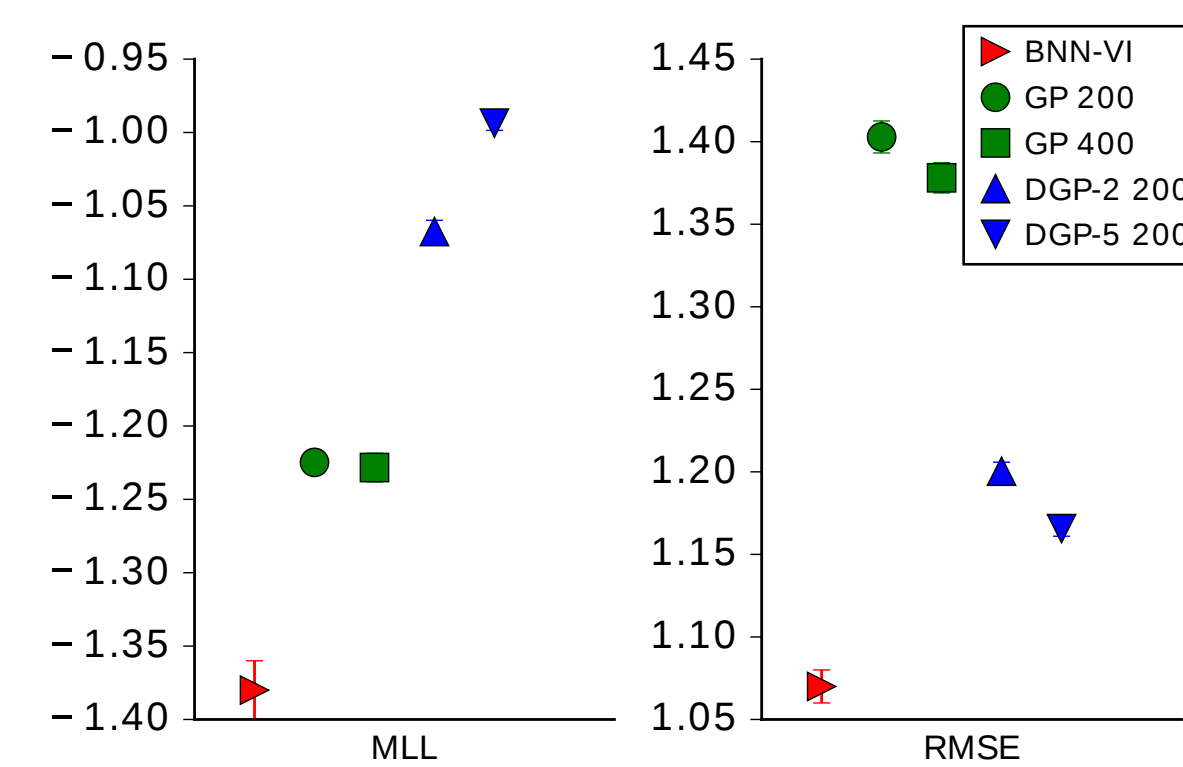


Figure : Average test log-likelihood [higher is better] and test error [lower is better].

Summary

- Deep GPs are state-of-the-art for regression
- Tying the factor approximations allows direct optimisation of the EP energy