



Mapping Gaussian Processes to Bayesian Neural Networks

Daniel Flam-Shepherd¹, James Requeima^{2,3}, David Duvenaud¹
1 University of Toronto, 2 University of Cambridge, 3 Invenia Labs



Priors in Function Space

- **Bayesian Neural Network Priors** are specified in parameter space. The implications of these priors in function space are hard to interpret.
- How do we incorporate prior knowledge about function properties in our prior?

Gaussian Processes

Gaussian Processes have a elegant mechanism for incorporating prior beliefs about the underlying function - the mean and covariance functions.

Kernel name:	Squared-exp (SE)	Periodic (Per)	Linear (Lin)
$k(x, x') =$	$\sigma_f^2 \exp\left(-\frac{(x-x')^2}{2\ell^2}\right)$	$\sigma_f^2 \exp\left(-\frac{2}{\ell^2} \sin^2\left(\pi \frac{x-x'}{p}\right)\right)$	$\sigma_f^2(x-c)(x'-c)$
Plot of $k(x, x')$:			
	$x - x'$	$x - x'$	x (with $x' = 1$)
Functions $f(x)$ sampled from GP prior:			
	x	x	x
Type of structure:	local variation	repeating structure	linear functions

Kernels can be combined using addition and multiplication to construct kernels with desired properties.

Lin \times Lin	SE \times Per	Lin \times SE	Lin \times Per
x (with $x' = 1$)	$x - x'$	x (with $x' = 1$)	x (with $x' = 1$)
quadratic functions	locally periodic	increasing variation	growing amplitude

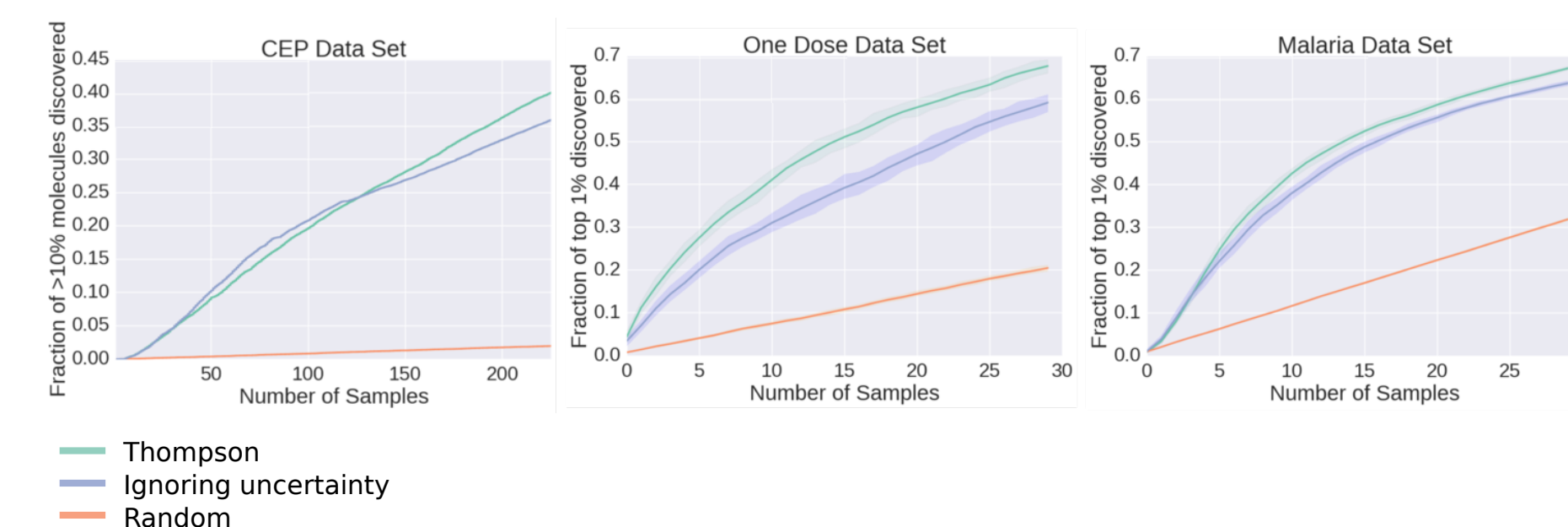
Mapping the Prior

We approximate the **KL divergence** between the Gaussian process $p_{GP}(\mathbf{f})$ and the BNN prior over functions $p_{BNN}(\mathbf{f})$.

$$\begin{aligned}\mathcal{L}_{p(\mathbf{X})}(\phi) &= \mathbb{E}_{p(\mathbf{X})}[\mathcal{KL}[p_{BNN}(\mathbf{f}(\mathbf{X})|\phi) | p_{GP}(\mathbf{f}(\mathbf{X}))]] \\ &= \mathbb{E}_{p(\mathbf{X})}[-\mathcal{H}[p_{BNN}(\mathbf{f}(\mathbf{X})|\phi)] \\ &\quad - \mathbb{E}_{p_{BNN}(\mathbf{f}|\phi)}[\log p_{GP}(\mathbf{f}(\mathbf{X}))]]\end{aligned}$$

The second term in this expectation can be approximated using Monte Carlo. The entropy term can be approximated 1.

Results: Bayesian Neural Networks



Data sets:

- CEP: Harvard Clean Energy Project data, 2.3M molecules.
- One-dose: percentage cell growth relative to control, 27,000 molecules.
- Malaria: drug concentration giving half max response, 19,000 molecules.

Batch sizes: 500 (CEP) and 200 (Malaria and One-dose).

Comparison with ϵ -greedy sampling

Average rank and standard errors:

Method	Rank
$\epsilon = 0.01$	3.42 \pm 0.28
$\epsilon = 0.025$	3.02 \pm 0.25
$\epsilon = 0.05$	2.86 \pm 0.23
$\epsilon = 0.075$	3.20 \pm 0.26
Thompson	2.51\pm0.20

Conclusions

We have proposed a batch BO method that

- runs in a parallel and distributed manner.
- can handle large batch sizes and large molecule libraries.
- is comparable to non-scalable approaches (parallel EI) in small problems with GPs.
- outperforms other alternative scalable approaches in large scale settings with Bayesian neural networks.