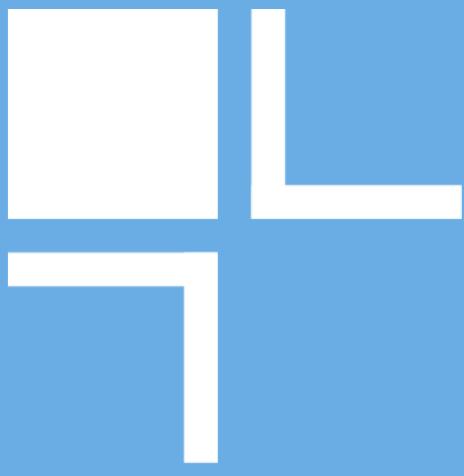


Mapping Gaussian Processes to Bayesian Neural Networks

Daniel Flam-Shepherd¹, James Requeima²³, David Duvenaud¹

¹ University of Toronto, ² University of Cambridge, ³ Invenia Labs



Priors in Function Space are Interpretable

- **Bayesian Neural Network Priors** are specified in parameter space. The implications of these priors in function space are hard to interpret.
- How to we incorporate prior knowledge about function properties in our prior?
- **Gaussian Processes** can elegantly incorporate prior beliefs about functions through 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 $k(x, x')$:			
Functions $f(x)$ sampled from GP prior:			
Type of structure:	local variation	repeating structure	linear functions

- Can we specify BNN priors using GP machinery?

Minimizing Divergence in Function Space

How can we find a prior on weights $p(w)$ that produces functions similar to a GP prior?

$$\phi^* = \underset{\phi}{\operatorname{argmin}} \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X})} \text{KL}[p_{\text{BNN}}(\mathbf{f}(\mathbf{X})|\phi) || p_{\text{GP}}(\mathbf{f}(\mathbf{X})]$$

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

We can estimate the likelihood with simple Monte Carlo, but how can we estimate the entropy of p_{BNN} ?

Estimating the Entropy of BNN Priors

1) Moment matching:

Approximate $p_{\text{BNN}}(\mathbf{f}(\mathbf{X})|\phi) \approx \mathcal{N}(\mathbf{f}|\mu_{\mathbf{f}(\mathbf{X})}, \Sigma_{\mathbf{f}(\mathbf{X})})$ using empirical mean and covariance. In this case,

$$\mathbb{H}[p_{\text{BNN}}(\mathbf{f}(\mathbf{X})|\phi)] \approx \frac{1}{2} \log |2\pi e \Sigma_{\mathbf{f}(\mathbf{X})}|$$

Reasonable since $p_{\text{BNN}}(\mathbf{f}|\phi) \rightarrow p_{\text{GP}}(\mathbf{f})$ over $\mathbf{X} \sim p(\mathbf{X})$.

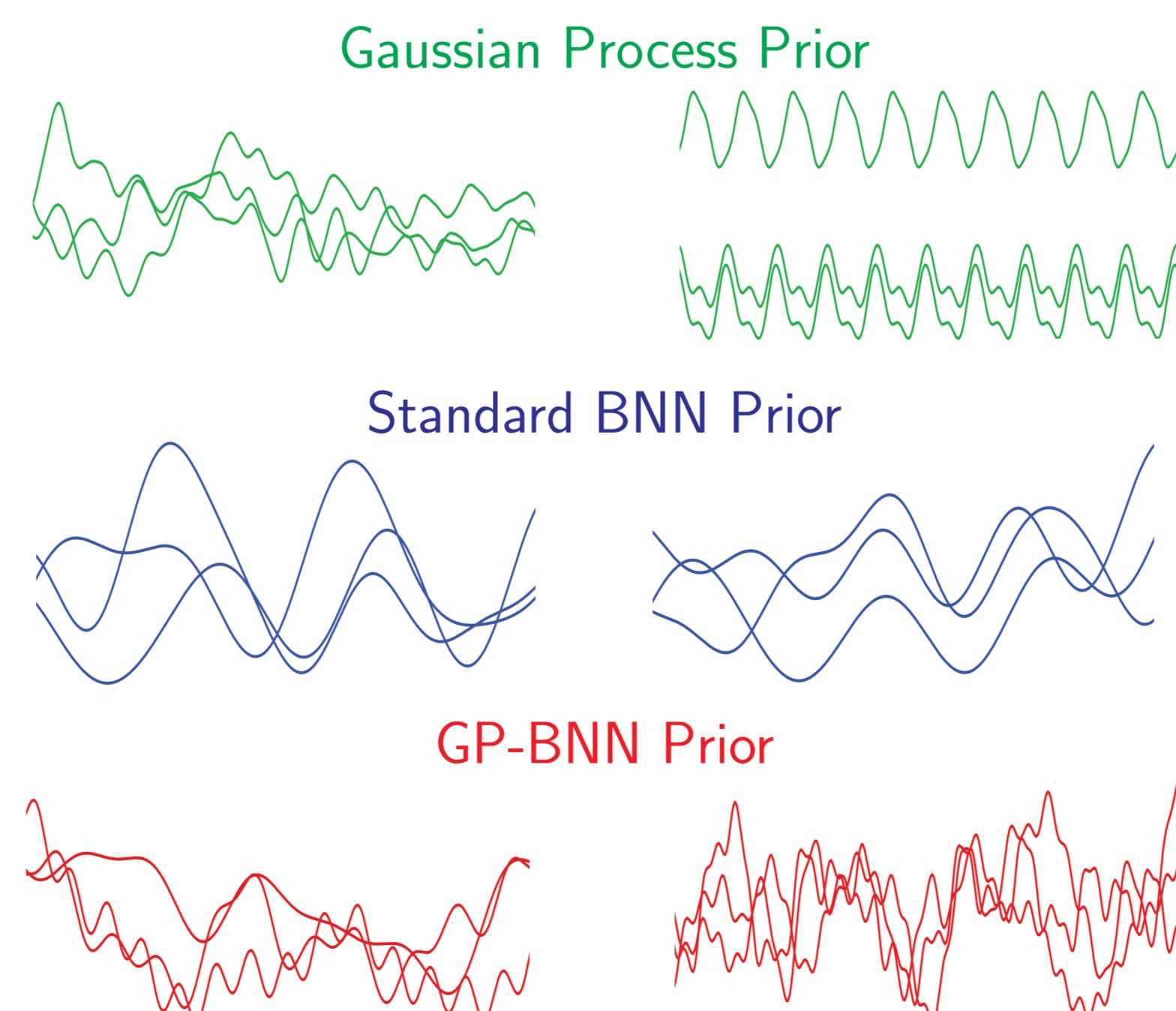
2) Nested Variational Bound:

Introduce noise $\mathbf{y} = \mathbf{f} + \epsilon$, and approximate $\log p(\mathbf{y}(\mathbf{X})|\phi)$ with a variational lower bound $\mathcal{L}(\boldsymbol{\lambda}, \mathbf{y}^{(s)}, \mathbf{X})$.

$$\begin{aligned} -\mathbb{H}[p_{\text{BNN}}(\mathbf{y}(\mathbf{X})|\phi)] &\approx \frac{1}{S} \sum_{s=1}^S \log p_{\text{BNN}}(\mathbf{y}^{(s)}(\mathbf{X})|\phi) \\ &\approx \frac{1}{S} \sum_{s=1}^S \max_{\boldsymbol{\lambda}} \mathcal{L}(\boldsymbol{\lambda}, \mathbf{y}^{(s)}, \mathbf{X}) \end{aligned}$$

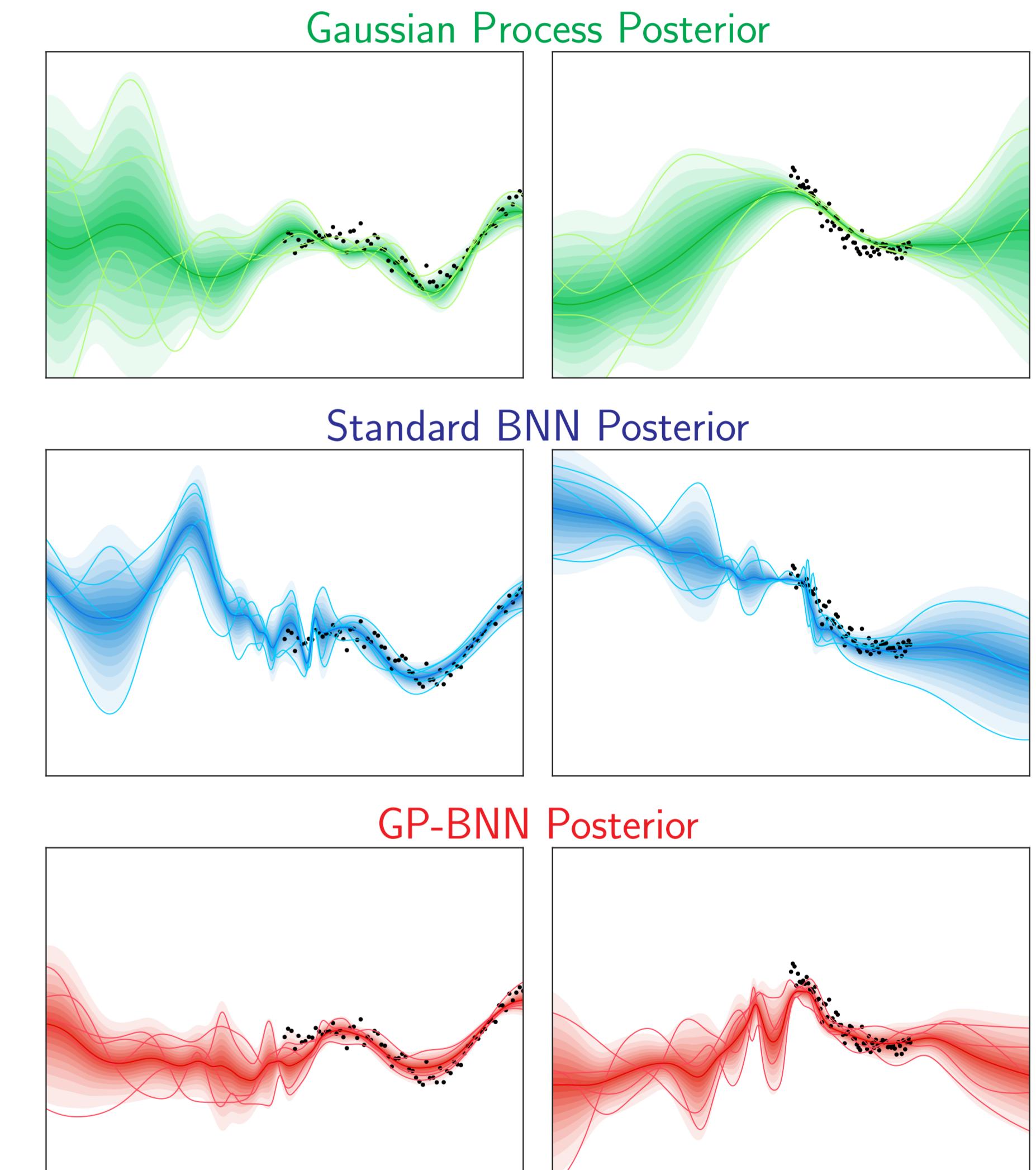
3) Early stopping to avoid mode collapse.

Samples from Approximate Priors



Posterior Results using GP Priors

We test our model (blue) on 2 different toy problems. All BNNs are 2 layer with rbf activation functions. The GP priors used RBF kernels.



Manifesto

- Specifying properties of functions should be the first decision, and computational architecture follows.
- Priors on functions are more interpretable than priors on parameters.