
Mapping Gaussian Process Priors to Bayesian Neural Networks

Daniel Flam-Shepherd
University of Toronto
daniel@utstat.utoronto.ca

James Requeima
University of Cambridge
jrr41@cam.ac.uk

David Duvenaud
University of Toronto
duvenaud@cs.toronto.edu

1 Introduction and Motivation

What defines a reasonable prior to use when forming and training Bayesian models? Recent work, such as [1] seek to improve priors over network parameters, by applying a horseshoe prior over preactivations of a Bayesian neural network (BNN) to effectively turn off weights that do not explain the data. Our goal with this work is to consider priors in function space as well.

It is difficult to incorporate meaningful prior information about functions to be modeled by BNNs since priors are generally specified over the network parameters. Often, normal distributions are placed over the weights for convenience and are interpreted as a bias toward less complex functions via smaller weights. Gaussian processes, on the other hand, have an elegant mechanism for incorporating prior beliefs about the underlying function - specifying the mean and covariance functions. However, Gaussian Processes have scalability limitations making Bayesian neural networks a more practical model in large data settings. In our work, we present an approach to specify a more principled prior for Bayesian Neural Networks that can leverage the well studied kernel design techniques from Gaussian process regression.

We consider matching the prior of a Bayesian neural network $p_{\text{BNN}}(\mathbf{f}|\phi)$ to the prior of a Gaussian process $p_{\text{GP}}(\mathbf{f})$ by minimizing their approximate KL divergence via some data distribution of interest $\mathbf{X} \sim p(\mathbf{X})$. We minimize this divergence with respect to the initial variational parameters ϕ of the proposal distribution $q(\mathbf{w}|\phi)$. These variational parameters $\phi^* = \{\mu_\phi^*, \log \sigma_\phi^*\}$ yield a prior on the BNN weights $p(\mathbf{w}|\phi^*) = \mathcal{N}(\mathbf{w}|\mu_\phi^*, \sigma_\phi^*)$. Then, variational inference allows us to perform approximate inference in our BNN using this more principled prior. We describe the implementation of both steps and demonstrate its success in the following sections.

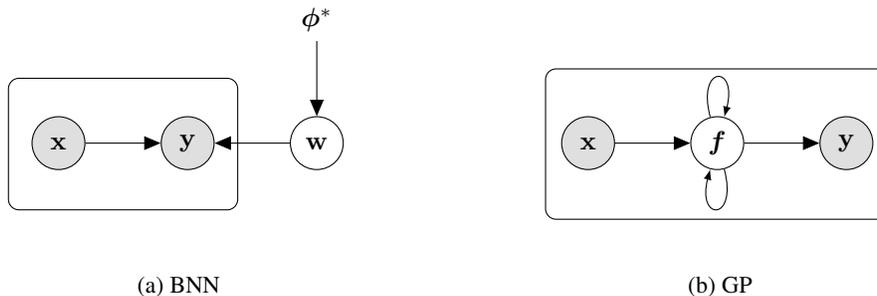


Figure 1: (a) and (b) display the graphical models of a BNN and GP.

2 Mapping a GP prior to BNN prior

In this section we describe the procedure used to minimize the KL divergence of the BNN prior distribution over functions $p_{\text{BNN}}(\mathbf{f}|\phi)$ and the GP prior distribution over functions $p_{\text{GP}}(\mathbf{f}) \equiv \mathcal{GP}(\mathbf{f}|\mathbf{0}, \mathbf{K})$ where \mathbf{K} is the Gram matrix of the kernel or covariance function chosen.

$$\mathcal{K}(\phi) = \mathbb{KL}[p_{\text{BNN}}(\mathbf{f}|\phi) | p_{\text{GP}}(\mathbf{f})] \quad (1)$$

$$= \int p_{\text{BNN}}(\mathbf{f}|\phi) \log \left[\frac{p_{\text{BNN}}(\mathbf{f}|\phi)}{p_{\text{GP}}(\mathbf{f})} \right] d\mathbf{f} \quad (2)$$

$$= -\mathbb{H}[p_{\text{BNN}}(\mathbf{f}|\phi)] - \mathbb{E}_{p_{\text{BNN}}(\mathbf{f}|\phi)}[\log p_{\text{GP}}(\mathbf{f})] \quad (3)$$

We define our first stochastic optimization objective $\mathcal{L}_{\mathbf{X}}(\phi)$ by approximating the KL divergence between these infinite dimensional distributions by taking expectations over some data distribution $p(\mathbf{X})$, explicitly : $\mathcal{L}_{\mathbf{X}}(\phi) \equiv \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X})}[\mathcal{K}(\phi)]$. This allows us to prioritize where in the input space we want $p_{\text{BNN}}(\mathbf{f}|\phi) \sim p_{\text{GP}}(\mathbf{f})$. We have :

$$\mathcal{L}_{\mathbf{X}}(\phi) = -\mathbb{E}_{\mathbf{X} \sim p(\mathbf{X})} \mathbb{H}[p_{\text{BNN}}(\mathbf{f}(\mathbf{X})|\phi)] - \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X})} \mathbb{E}_{p_{\text{BNN}}(\mathbf{f}(\mathbf{X})|\phi)}[\log p_{\text{GP}}(\mathbf{f}(\mathbf{X}))] \quad (4)$$

Both terms in equation (4) are intractable and force us to use some approximation. Since we are able to compute $p_{\text{GP}}(\mathbf{f}(\mathbf{X}))$ analytically, we can compute a Monte Carlo Estimate of the second term in (4):

$$\mathbb{E}_{\mathbf{X} \sim p(\mathbf{X})} \mathbb{E}_{p_{\text{BNN}}(\mathbf{f}|\phi)}[\log p_{\text{GP}}(\mathbf{f})] \approx \frac{1}{S} \sum_{s=1}^S \mathbb{E}[\log p_{\text{GP}}(\mathbf{f}^{(s)}(\mathbf{X}))]$$

where $\mathbf{f}^{(s)}(\mathbf{X}) \sim p_{\text{BNN}}(\mathbf{f}|\phi)$. For the entropy term in (4) we propose two estimation methods.

2.1 Estimate of entropy term using moment matching

The first estimate we make involves the approximation that the prior distribution over functions for the BNN is a multivariate Gaussian distribution that takes form $p_{\text{BNN}}(\mathbf{f}(\mathbf{X})|\phi) \sim \mathcal{N}(\mathbf{f}|\boldsymbol{\mu}_{\mathbf{f}(\mathbf{X})}, \boldsymbol{\Sigma}_{\mathbf{f}(\mathbf{X})})$. In this case, the entropy term can be evaluated analytically over $p(\mathbf{X})$ as $\mathbb{H}[p_{\text{BNN}}(\mathbf{f}(\mathbf{X})|\phi)] = \frac{1}{2} \log |2\pi e \boldsymbol{\Sigma}_{\mathbf{f}(\mathbf{X})}|$ where we use the sample mean and covariance for $\boldsymbol{\mu}_{\mathbf{f}(\mathbf{X})}$, $\boldsymbol{\Sigma}_{\mathbf{f}(\mathbf{X})}$ over $p(\mathbf{X})$. This is a reasonable approximation because the overall goal in this work is push $p_{\text{BNN}}(\mathbf{f}|\phi) \approx p_{\text{GP}}(\mathbf{f})$ over $p(\mathbf{X})$ and as $p_{\text{BNN}}(\mathbf{f}|\phi) \rightarrow p_{\text{GP}}(\mathbf{f})$ it will get more accurate. With this, the objective is

$$\mathcal{L}_{\mathbf{X}}(\phi) \approx -\frac{1}{2} \log |\boldsymbol{\Sigma}_{\mathbf{f}(\mathbf{X})}| - \frac{1}{S} \sum_{s=1}^S \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X})}[\log p_{\text{GP}}(\mathbf{f}^{(s)}(\mathbf{X}))] \quad (5)$$

We optimize this until convergence $\phi^* = \underset{\phi}{\operatorname{argmin}} \mathcal{L}_{\mathbf{X}}(\phi)$. The details are in Algorithm 1.

Algorithm 1 Optimization of the prior of the Bayesian neural net

- 1: **Initialize** $\phi = \{\boldsymbol{\mu}_{\phi}, \log \boldsymbol{\sigma}_{\phi}\}$
 - 2: **while** ϕ not converged **do**
 - 3: $\boldsymbol{\epsilon}^{(s)} \sim p(\boldsymbol{\epsilon}) = \mathcal{N}(\mathbf{0}, \mathbb{I})$ ▷ sample prior noise
 - 4: $\mathbf{w}^{(s)} \leftarrow t(\phi, \boldsymbol{\epsilon}) = \boldsymbol{\mu}_{\phi} + \boldsymbol{\sigma}_{\phi} \odot \boldsymbol{\epsilon}^{(s)}$ ▷ sample S weights $\mathbf{w} \sim q(\mathbf{w}|\phi)$
 - 5: $\mathbf{X} \leftarrow (\mathbf{x}_1, \dots, \mathbf{x}_n) \sim p(\mathbf{x}_1, \dots, \mathbf{x}_n)$ ▷ sample data
 - 6: $\mathbf{f}^{(s)}(\mathbf{X}) \leftarrow f(\mathbf{X}, \mathbf{w}^{(s)})$ ▷ sample S functions $\mathbf{f} \sim p_{\text{BNN}}(\mathbf{f}|\phi)$
 - 7: $\mathbf{g}_{\phi} \leftarrow \nabla_{\phi} \mathcal{L}_{\mathbf{X}}(\phi)$ ▷ compute gradients of the objective
 - 8: $\phi \leftarrow \text{adam}(\phi, \mathbf{g}_{\phi})$ ▷ update the parameters
 - 9: **Return** ϕ^*
-

The next section describes a formulation that allows to make an alternate approximation of the entropy term by considering the equivalent problem in observation space.

2.2 Estimate of entropy term by projecting into observation space

Consider the equivalent problem in observation space by adding noise $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbb{I})$ to our functions $\mathbf{y} = \mathbf{f} + \epsilon$ so that we are minimizing the KL divergence of $p(\mathbf{y}|\phi)$: the BNN prior distribution over noisy observations to the GP prior distribution over noisy observations $p_{\text{GP}}(\mathbf{y})$.

$$\mathcal{K}(\phi) = \mathbb{KL}[p(\mathbf{y}|\phi) | p_{\text{GP}}(\mathbf{y})] = -\mathbb{H}[p(\mathbf{y}|\phi)] - \mathbb{E}_{p(\mathbf{y}|\phi)}[\log p_{\text{GP}}(\mathbf{y})] \quad (6)$$

Now we use Monte Carlo estimates of both $\mathbb{H}[p(\mathbf{y}|\phi)]$ and $\mathbb{E}_{p(\mathbf{y}|\phi)}[\log p_{\text{GP}}(\mathbf{y}|\theta)]$:

$$\mathcal{K}(\phi) \approx \frac{1}{S} \sum_{s=1}^S \log p(\mathbf{y}^{(s)}|\phi) - \frac{1}{S} \sum_{s=1}^S \log p_{\text{GP}}(\mathbf{y}^{(s)}) \quad (7)$$

where we are sampling S times: $\mathbf{y}^{(s)} \sim p(\mathbf{y}|\phi)$. To estimate $\log p(\mathbf{y}^{(s)}|\phi)$ we do approximate inference on an intermediate proposal distribution $r(\mathbf{w}|\lambda)$, where $\lambda = \{\mu_\lambda, \log \sigma_\lambda\}$ are its variational parameters to be optimized. We then estimate $\log p(\mathbf{y}^{(s)}|\phi)$ using its ELBO $\mathcal{L}_{\mathcal{D}_s}(\lambda)$, optimized on data $\mathcal{D}_s = \{\mathbf{X}, \mathbf{y}^{(s)}\}$. This is similar to Hierarchical variational models by [2].

$$\log p(\mathbf{y}^{(s)}|\phi) \geq \mathcal{L}_{\mathcal{D}_s}(\lambda) \equiv \mathbb{E}_{r(\mathbf{w}|\lambda)} [\log p(\mathcal{D}_s|\mathbf{w}) + \log p(\mathbf{w}) - \log r(\mathbf{w}|\lambda)] \quad (8)$$

where $\log p(\mathcal{D}_s|\mathbf{w}) = \sum_i \log \mathcal{N}(\mathbf{y}_i^{(s)} | f(\mathbf{x}_i, \mathbf{w}), \sigma^2 \mathbb{I})$ defines a parametric model. Once again we take expectations $p(\mathbf{X})$ to define the objective $\mathcal{L}_{\mathbf{X}}(\phi, \lambda) \equiv \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X})} [\mathcal{K}(\phi)]$ by approximating the KL divergence between these infinite dimensional distributions. We obtain:

$$\mathcal{L}_{\mathbf{X}}(\phi, \lambda) = \min_{\phi} \frac{1}{S} \sum_{s=1}^S \max_{\lambda} \mathcal{L}_{\mathcal{D}_s}(\lambda) - \frac{1}{S} \sum_{s=1}^S \mathbb{E}_{\mathbf{X} \sim p(\mathbf{X})} [\log p_{\text{GP}}(\mathbf{y}^{(s)}(\mathbf{X}))] \quad (9)$$

We optimize (9) till convergence, this is described in the Algorithm 2.

Algorithm 2 Mapping a GP prior to a BNN prior in observation space

```

1: Initialize  $\phi = \{\mu_\phi, \log \sigma_\phi\}$ 
2: while  $\phi$  not converged do
3:    $\epsilon^{(s)} \sim p(\epsilon) = \mathcal{N}(\mathbf{0}, \mathbb{I})$  ▷ sample prior noise
4:    $\mathbf{w}^{(s)} \leftarrow \mathbf{t}(\phi, \epsilon^{(s)}) = \mu_\phi + \sigma_\phi \epsilon^{(s)}$  ▷ sample  $S$  weights  $\mathbf{w} \sim q(\mathbf{w}|\phi)$ 
5:    $\mathbf{X} \leftarrow (\mathbf{x}_1, \dots, \mathbf{x}_n) \sim p(\mathbf{x}_1, \dots, \mathbf{x}_n)$  ▷ sample from some data distribution
6:    $\mathbf{f}^{(s)} \leftarrow f(\mathbf{X}, \mathbf{w}^{(s)})$  ▷ sample  $S$  functions from the BNN
7:    $\mathbf{y}^{(s)} \leftarrow \mathbf{f}^{(s)}(\mathbf{X}) + \epsilon$  ▷ sample observations
8:   for each  $\mathcal{D}_s = \{\mathbf{X}, \mathbf{y}^{(s)}\}$  hold  $\phi$  constant and do
9:      $\lambda_s^* \leftarrow \underset{\lambda}{\operatorname{argmax}} \mathcal{L}_{\mathcal{D}_s}(\lambda)$  ▷ do approximate inference
10:     $\log p(\mathbf{y}^{(s)}) \leftarrow \mathcal{L}_{\mathcal{D}_s}(\lambda_s^*)$  ▷ estimate the marginal log likelihood
11:     $\mathbf{g}_\phi \leftarrow \nabla_{\phi} \mathcal{L}_{\mathbf{X}}(\phi, \lambda)$  ▷ compute gradients of the kl
12:     $\phi \leftarrow \operatorname{adam}(\phi, \mathbf{g}_\phi)$  ▷ update the parameters
13: Return  $\phi^*$ 

```

2.3 Early Stopping

Another rough approximation of (4) is to ignore the entropy term $\mathbb{H}[p_{\text{BNN}}(\mathbf{f}|\phi)]$ altogether, as it merely acts as a regularization term, so we avoid collapsing to a single high density value. To address this issue, we found that early stopping works well in practice and yielded our best results. The results presented in the following section use early stopping and no entropy term.

3 Experiments with different kernels and activation functions

We experiment with mapping GP priors (GPPs) with various covariance functions $k(x, x')$ to BNN priors that have various activation functions $a(x)$. We sample on order 10^3 data from a uniform data distribution $\mathbf{X} \sim p(\mathbf{X})$ about a symmetric interval $(-10, 10)$. We use automatic differentiation software : autograd [3, 4] to compute gradients of our objective, which is optimized using adam [5]. In all experiments we work directly in function space and found that reasonable results can be found even when ignoring the entropy term $\mathbb{H}[p(\mathbf{f})]$ altogether. Experiments are conducted on simple 1 or 2 layer neural networks. The prior proposal distribution on the BNN weights is a diagonal Gaussian in all cases. In the following subsections, plots have color scheme:

$$f(\mathbf{X}) \sim p_{\text{GP}}(f(\mathbf{X})), \quad \mathbf{f}(\mathbf{X}) \sim p_{\text{BNN}}(\mathbf{f}(\mathbf{X})|\phi^*), \quad \mathbf{f}(\mathbf{X}) \sim p_{\text{BNN}}(\mathbf{f}(\mathbf{X}))$$

3.1 Linear and Softplus Experiments

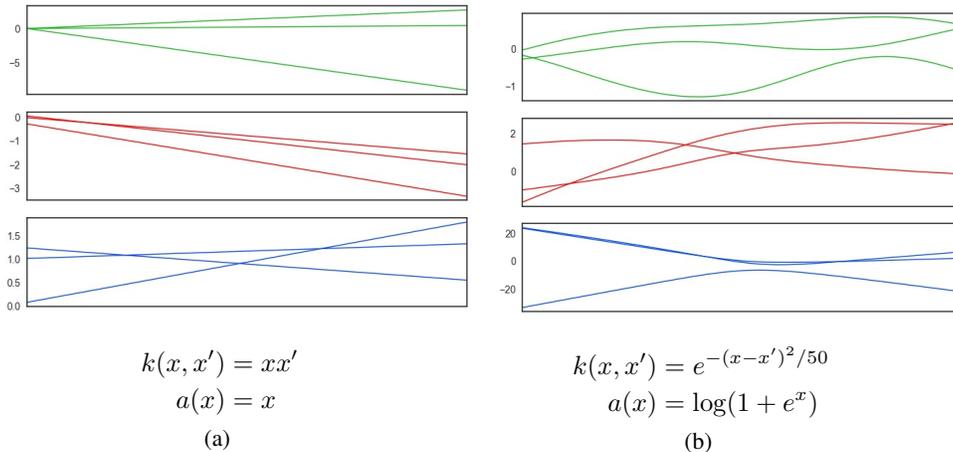


Figure 2: In (a) we map a GPP with a linear kernel to a linear model (Bayesian Linear regression), the optimized prior samples extend from the origin just as the GPP samples do. In (b) we map a GP prior with a RBF kernel to a BNN prior with a softplus activation function, the optimized samples have been mapped near the range of the GPP samples and almost as curvy as the GPP samples.

3.2 Experiments with rbf hyperbolic tan activations

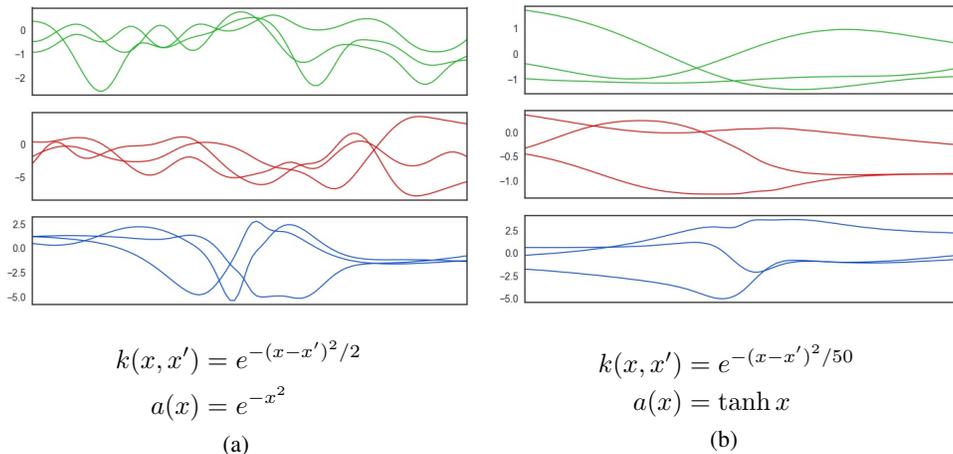


Figure 3: In (a) and (b) we map GP priors with a RBF kernel to BNNs with a RBF and tanh activation functions. It is clear that the optimized samples have taken on the curviness of the GPP samples.

3.3 Experiments with periodic kernels and activations

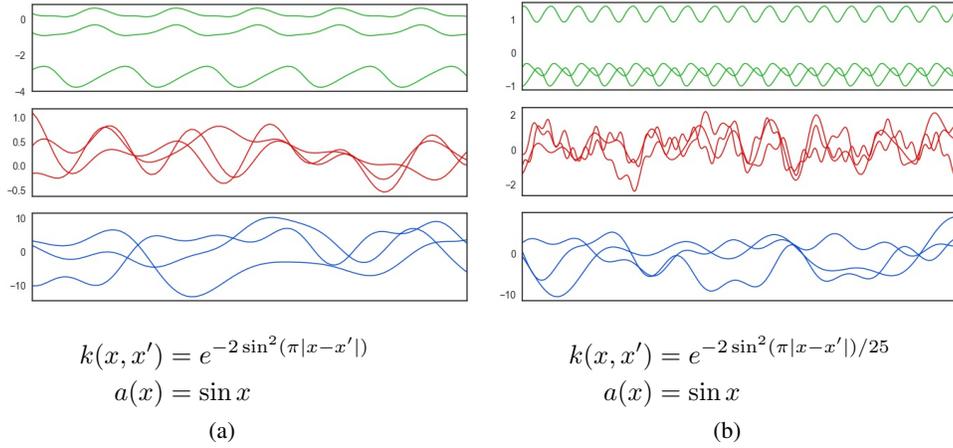


Figure 4: In (a) and (b) we map GPPs with periodic kernels to BNNs with sin activation functions. The periodic structure of the GPP samples has been mapped to the optimized BNN samples.

3.4 Experiments with other kernels

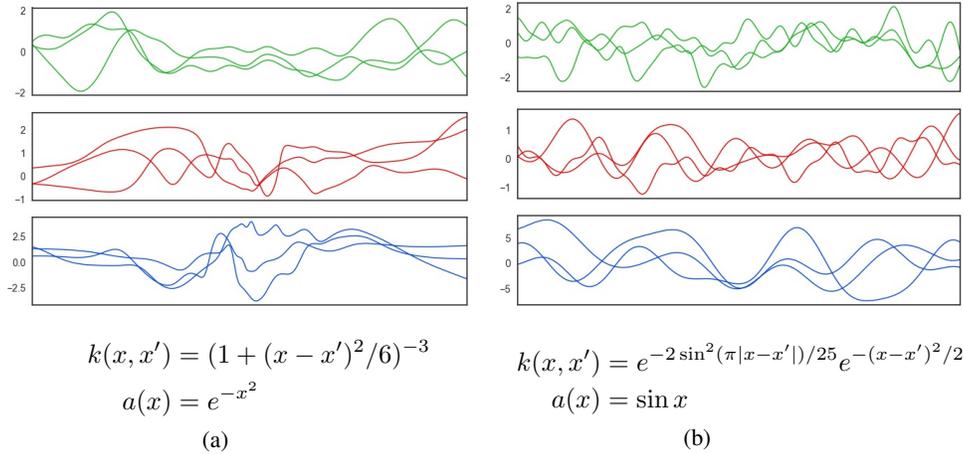


Figure 5: In (a) we map a GPP with a rational cubic kernel to a BNN prior, the optimized samples have taken on the "wigglyness" of the prior samples especially away from the origin (the middle). In (b) we map a GPP with a locally periodic kernel to a BNN prior, the optimized samples take on both the "local wigglyness" and the periodicity of the functions from the GPP

4 Inference using the optimized prior

Next we use the optimized parameters ϕ^* found by minimizing the kl of the $p_{\text{BNN}}(\mathbf{f}|\phi)$ to $p_{\text{GP}}(\mathbf{f})$ in a prior on the weights $p(\mathbf{w}|\phi^*) = \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}_{\phi^*}, \boldsymbol{\sigma}_{\phi^*})$. when maximizing the ELBO $\mathcal{L}_{\mathcal{D}}(\varphi) \leq \log p(\mathcal{D})$.

$$\mathcal{L}_{\mathcal{D}}(\varphi) = \mathbb{H}[q(\mathbf{w}|\varphi)] + \mathbb{E}_{q(\mathbf{w}|\varphi)} [\log p(\mathcal{D}|\mathbf{w}) + \log p(\mathbf{w}|\phi^*)] \quad (10)$$

Where $\mathcal{D} = \{\mathbf{y}, \mathbf{X}\}$ is the training data, $\varphi = \{\boldsymbol{\mu}_{\varphi}, \log \boldsymbol{\sigma}_{\varphi}\}$ are the variational parameters of the approximation $q(\mathbf{w}|\varphi) = \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}_{\varphi}, \boldsymbol{\sigma}_{\varphi})$ to the true posterior on the weights $p(\mathbf{w}|\mathcal{D})$. Which we can sample from using reparameterization trick [6, 7] $\mathbf{w}^{(\ell)} = \mathbf{t}(\varphi, \boldsymbol{\epsilon}) = \boldsymbol{\mu}_{\varphi} + \boldsymbol{\sigma}_{\varphi} \boldsymbol{\epsilon}^{(\ell)}$. Thus we can obtain unbiased stochastic gradients of the ELBO with respect to the variational parameters.

4.1 Experiments on toy data

We test our model (blue) on 2 different toy problems. We work directly in function space and use algorithm one, and we found that no entropy term estimate was necessary for reasonable results. We compare to a BNN (blue) using a standard normal prior distribution on the weights and to the posterior of a Gaussian process (green). The BNNs are trained using Bayes by Backprop [8]. For both Gaussian process priors, the kernel is the RBF covariance function. The observations y are sampled by passing some data sampled uniformly through a function $f(x)$ and adding some Gaussian noise $\epsilon \sim \mathcal{N}(0, \sigma^2)$. We use $\sigma = 1$. The networks are all 2 layer with rbf activations as well. In the following plots the previous color scheme also applies

$$f(\mathbf{X}) \sim p_{\text{GP}}(f|\mathcal{D}), \quad f(\mathbf{X}) \sim p_{\text{BNN}}(f(\mathbf{X})|\varphi^*, \phi^*), \quad f(\mathbf{X}) \sim p_{\text{BNN}}(f(\mathbf{X})|\varphi^*)$$

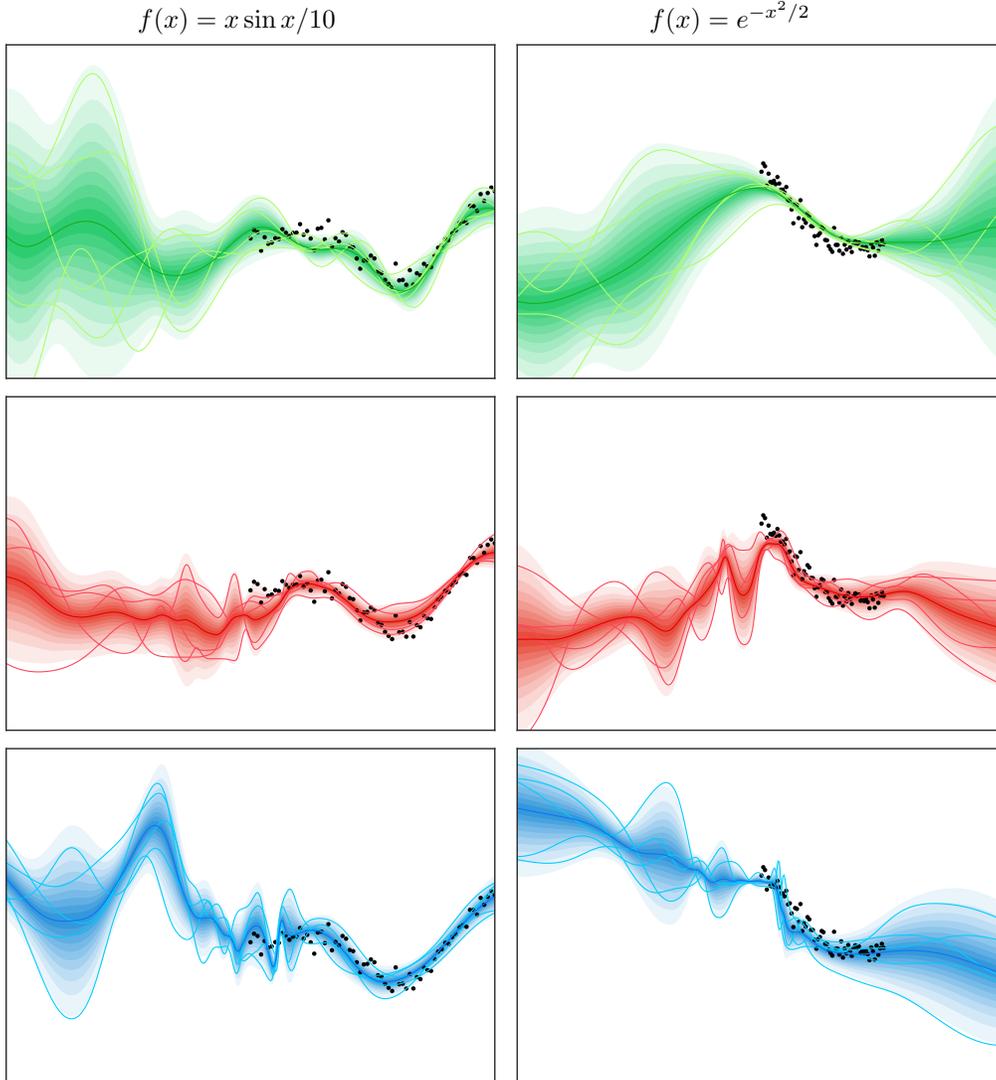


Figure 6: Plots of the 3 posteriors for 2 different functions. Different shades of color correspond to deciles of the predictive density $p(f(x_i))$. The samples from this density are a slightly different shade of green, red, and blue. The darkest is the mean. Notice that the mean and uncertainty bands of the posterior trained with the optimized prior closer resembles the uncertainty of the Gaussian process. This demonstrates that our method of approximating a GP prior via a BNN can successfully transfer characteristics of the GP posterior to the BNN posterior as well.

5 Limitations

In the previous sections we demonstrate the success of our proposed objective in completing the set out task of mapping the properties of a Gaussian process prior to a Bayesian neural network prior. However, there remain significant limitations of our method. We describe them below.

We found that we could not map GP priors to BNNs with certain activation functions, for example we could only map a GP prior with a linear kernel to a BNN with Relu activation functions at any capacity. In this case, the BNN could not, at any capacity, learn the smoothness or periodic properties of any more interesting GP prior.

As well, we could not map the properties of a GP prior with a periodic kernel to a BNN that did not have a periodic activation function. Our experiments focused on mapping properties of GP priors to BNNs with configurations that could reasonably learn those properties, for example periodic kernels to BNN priors with periodic activation functions.

We did not find an accurate estimator of the entropy term in our objective, as a result, training BNN priors to behave like GP priors was difficult and was negatively hampered by lack of a accurate regularization term. Furthermore, rough estimates of the entropy term did not stabilize or improve these difficulties. Also, working in observation space with a more accurate estimation of the entropy term introduced another set of difficulties brought on by the mini max objective. As a consequence, we found that ignoring the entropy term and performing early stopping worked best.

Often, without the entropy term, samples from the optimized BNN prior would take on the properties of the GP prior then gradually flatten and lose the desired properties. On the other hand, with the entropy term the samples would sometimes blow up. There was no stable regime where we could obtain consistent convergence. This would sometimes result in optimized samples which over emphasize the characteristics of the GP prior. An interesting example of this is displayed in Figure 7 where the red optimized samples take on an extremely local periodic behaviour later in training, while the samples of the GP prior being mapped do not have such a extreme periodic behavior.

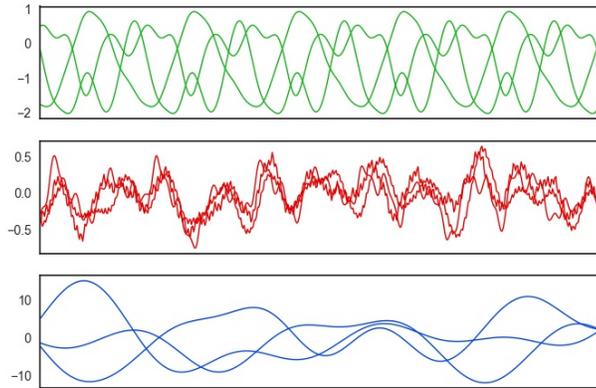


Figure 7: An example of extreme overtraining

6 Conclusions and Future Work

In this work, we formulate a method for mapping a Gaussian process prior to a Bayesian neural network prior. This is done in effort to transfer interesting properties of functions sampled from a GP prior to functions sampled from a Bayesian neural network prior. In this regard, we can implicitly think about priors in BNNs in function space rather than parameter space and overcome the limits that this entails. We demonstrate our method on a variety of kernels and activation functions, and gage its impact on posterior inference in simple toy examples. All together, while there are important described limitations we establish that our method works throughout these demonstrations and examples.

This work attempts to build a small bridge between the flexible models of Bayesian non parametrics and the more scalable models of Bayesian deep learning. In future we seek to continue to build upon this idea of specifying priors in function space. The goal is to building more expressive Bayesian models where the prior better captures necessary prior beliefs and is more useful for posterior inference.

Acknowledgements

The authors would like to thank Brian Ning and Guodong Zhang for helpful comments.

References

- [1] Soumya Ghosh and Finale Doshi-Velez. Model selection in bayesian neural networks via horseshoe priors. In *Neural Information Processing Systems*, 2016.
- [2] Rajesh Ranganath, Dustin Tran, and David M. Blei. Hierarchical variational models. *International Conference on Machine Learning*, 2016.
- [3] Dougal Maclaurin, David Duvenaud, Matthew Johnson, and Ryan P. Adams. Autograd: Reverse-mode differentiation of native python. 2015.
- [4] David Duvenaud and Ryan P. Adams. Black-box stochastic variational inference in five lines of python. *NIPS Workshop on Black-box Learning and Inference*, 2015.
- [5] D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. *International Conference on Learning Representations*, 2015.
- [6] D. P Kingma and M. Welling. Auto-Encoding Variational Bayes. *International Conference on Learning Representations*, 2014.
- [7] Danilo J Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *International Conference on Machine Learning*, 2014.
- [8] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. *International Conference on Machine Learning*, 2015.