

Walking on the Fiber: A Simple Geometric Approximation for Bayesian Neural Networks

Anonymous authors

Paper under double-blind review

Abstract

Bayesian Neural Networks provide a principled framework for uncertainty quantification by modeling the posterior distribution of network parameters. However, exact posterior inference is computationally intractable, and widely used approximations like the Laplace method struggle with scalability and posterior accuracy in modern deep networks. In this work, we revisit sampling techniques for posterior exploration, proposing a simple variation tailored to efficiently sample from the posterior in over-parameterized networks by leveraging the low-dimensional structure of loss minima. Building on this, we introduce a model that learns a deformation of the parameter space, enabling rapid posterior sampling without requiring iterative methods. Empirical results demonstrate that our approach achieves competitive posterior approximations with improved scalability compared to recent refinement techniques. These contributions provide a practical alternative for Bayesian inference in deep learning.

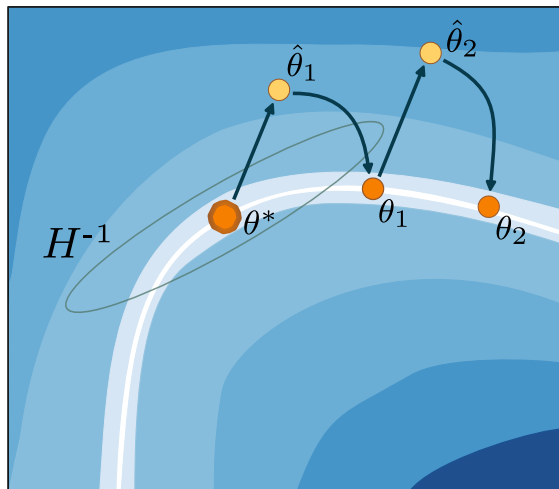


Figure 1: We propose a novel sampling scheme to approximate the posterior of a trained network. When the loss landscape of a trained network over its parameters (blue regions) is low on a very low-dimensional curve (white region), a simple Hessian (H^{-1}) fails in capturing its distribution. Our sampling scheme is composed of two iterative steps: randomly perturb a given solution (yellow points) then refine to minimize the loss function over the given dataset (orange points). This proposed scheme allows to estimate posteriors of arbitrary shapes and it is well suited when the space of solutions in a network are much lower dimensional than the parameter space.

1 Introduction

Bayesian Neural Networks (BNNs) have emerged as a principled framework for uncertainty quantification in deep learning by treating model parameters as random variables and inferring their posterior distribution. This Bayesian perspective is crucial for tasks requiring reliable decision-making under uncertainty, including medical diagnosis (Leibig et al., 2017), autonomous driving (Feng et al., 2018), and weather forecasting (Cofino et al., 2002). Despite their appeal, the high dimensionality and non-linearity of modern neural network parameter spaces render exact posterior inference intractable, requiring the development of efficient approximation techniques.

A widely used method to make a deterministic, already-trained neural network Bayesian is the Laplace approximation, which estimates the posterior distribution of the parameters as a Gaussian centered at the Maximum a Posteriori (MAP) estimate (MacKay, 1992; Daxberger et al., 2021). This approach is particularly convenient due to its simplicity and post-hoc applicability. However, the Laplace approximation faces significant challenges in modern deep networks. The computation and inversion of the Hessian scale poorly with network size (Martens, 2010), and its reliance on local curvature limits its ability to capture the complex, non-linear geometry of the posterior in over-parameterized models (Fort & Jastrzebski, 2019; Bergamin et al., 2024).

In contrast, sampling techniques, though less frequently used in recent years, can be surprisingly effective in exploring the posterior when the parameter space of loss minima resides on a much lower-dimensional manifold than the full parameter space (Garipov et al., 2018). By leveraging this property, sampling methods can efficiently explore the posterior with fewer samples, especially in over-parameterized settings. However, relying on sampling for posterior estimation can still be computationally expensive if every new inference requires a fresh sampling process. These limitations highlight the need for a method that combines the efficiency of sampling-based exploration with a more scalable and flexible posterior representation.

In this work, we propose a variation of a sampling method to explore the geometry of the loss landscape in the parameter space of neural networks. Our approach perturbs parameters around the MAP estimate using a set of drift directions and refines them with gradient updates, effectively maintaining computational efficiency as the network size increases. Using the data collected from our sampling process, we introduce a novel objective function to learn a structured latent representation of the posterior by deforming the original parameter space. This learned representation enables the direct generation of new parameter samples without relying on computationally expensive iterative sampling or restrictive variational inference approximations. We refer to these two components collectively as **MetricBNN**. Empirically, we demonstrate the efficacy of our method in estimating uncertainty from trained models in both regression and classification tasks. In particular, our approach produces better-calibrated uncertainty estimates than Gaussian approximations, especially on real-world datasets and high-dimensional tasks.

2 Related Work

Bayesian Neural Networks. BNNs provide a principled approach to uncertainty quantification by treating the network parameters θ as random variables and modeling their posterior distribution $p(\theta \mid \mathcal{D})$. However, exact inference is intractable due to the high dimensionality of parameter spaces in modern neural networks. Various methods have been proposed to approximate the posterior. Variational Inference (VI) approximates the posterior by optimizing a surrogate distribution, often in the form of a mean-field Gaussian (Graves, 2011; Blundell et al., 2015). Extensions to hierarchical and amortized variational methods have improved scalability but often struggle to capture complex posterior structures (Zhang et al., 2018). The Laplace Approximation (LA) defines the posterior locally around the MAP estimate using a Gaussian distribution (MacKay, 1998; Daxberger et al., 2021). Despite its simplicity and efficiency, it is limited by its reliance on local curvature and the Gaussian assumption. Sampling-based methods such as Stochastic Gradient Langevin Dynamics (SGLD) (Welling & Teh, 2011) and Hamiltonian Monte Carlo (HMC) (Neal, 2011) generate posterior samples by simulating stochastic dynamics. While these methods are more flexible, their computational cost often makes them impractical for large-scale neural networks.

Improving Posterior Samples. To address these limitations, recent works have sought to refine and extend the Laplace approximation by introducing more expressive approximate posteriors. Kristiadi et al. (2022) combines Laplace approximation with normalizing flows for a non-Gaussian posterior, refining the base Gaussian distribution. Similarly, Immer et al. (2021) refines Laplace approximation by leveraging Gaussian variational methods and Gaussian processes, improving linearized Laplace posterior accuracy. Miller et al. (2017) iteratively builds a mixture model to improve the posterior approximation by adding components to the variational distribution. Eschenhagen et al. (2021) combines multiple pre-trained models to form a weighted sum of Laplace approximations, improving posterior flexibility. Havasi et al. (2021) introduced auxiliary variables to locally refine mean-field variational approximations, achieving better fit in regions of interest. Bergamin et al. (2024) extends the Laplace approximation by leveraging Riemannian geometry to model the posterior distribution on a manifold, improving accuracy for non-linear loss landscapes.

In this paper, we reevaluate the efficiency of sampling methods for posterior estimation, particularly in the context of over-parameterized neural networks. Through empirical results, we demonstrate that our proposed simple sampling framework can outperform modern posterior refinement techniques in efficiency and accuracy. Furthermore, our novel use of an autoencoder moves beyond the Gaussian formulation of the posterior, enabling a flexible, easy-to-sample representation that significantly improves performance. This approach allows us to capture the non-linear geometry of the posterior while maintaining computational simplicity, addressing key limitations of both classical and modern refinement methods.

3 Preliminaries

BNNs provide a principled framework for quantifying uncertainty in neural network predictions by treating the model parameters as random variables. This section introduces the relevant background on BNNs, discusses the Laplace approximation for posterior estimation, and highlights its limitations.

Consider an independent and identically distributed (i.i.d.) dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ where $x_i \in \mathbb{R}^D$ and $y_i \in \mathbb{R}^C$. Let $f_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^C$ denote a parametric function (e.g., a neural network) with parameters $\theta \in \mathbb{R}^K$. The goal is to model the predictive distribution: $p(y' | x', \mathcal{D}) = \int p(y' | x', \theta) p(\theta | \mathcal{D}) d\theta$, where $p(\theta | \mathcal{D})$ is the posterior distribution of the parameters given the dataset. Bayesian inference provides a framework for estimating the posterior $p(\theta | \mathcal{D})$ using Bayes' theorem: $p(\theta | \mathcal{D}) = \frac{p(\mathcal{D} | \theta) p(\theta)}{p(\mathcal{D})}$, where $p(\mathcal{D} | \theta)$ is the likelihood of the data, $p(\theta)$ is the prior distribution over the parameters, and $p(\mathcal{D})$ is the marginal likelihood or evidence. In practice, computing the evidence $p(\mathcal{D})$ is often intractable, making direct posterior computation challenging.

3.1 Laplace Approximation

The Laplace approximation is a classic method for approximating the posterior distribution $p(\theta | \mathcal{D})$ using a Gaussian centered at the MAP estimate. Let $\mathcal{L}(\theta)$ denote the regularized negative log-likelihood:

$$\mathcal{L}(\theta) = -\log p(\mathcal{D} | \theta) - \alpha \log p(\theta), \quad (1)$$

where α is a regularization coefficient and depends on the choice of prior for the parameters. Assuming this to be the Normal distribution, the regularization can be rewritten as an Euclidean norm.

One of the main advantages of the Laplace approximation is that it allows to define a posterior distribution given an already fully trained network, *post-hoc*. When given a fully trained network we assume access to the set of parameters, θ^* , that minimize the loss function of Equation 1. This is the MAP solution. The exponential of this loss function is proportional to the posterior distribution as $p(\theta | \mathcal{D}) \propto p(\mathcal{D} | \theta) p(\theta)$. Laplace approximation methods propose to approximate the posterior distribution with a second-order Taylor expansion around the MAP of the loss function. This results in a Gaussian approximation of the parameters with the mean being the MAP solution and the covariance being the inverse of the Hessian computed in the MAP. The posterior can then be defined as:

$$q(\theta) = \mathcal{N}(\theta^*, H^{-1}), \quad (2)$$

where $H = \nabla^2 \mathcal{L}(\theta^*)$ is the Hessian of the loss function at θ^* .

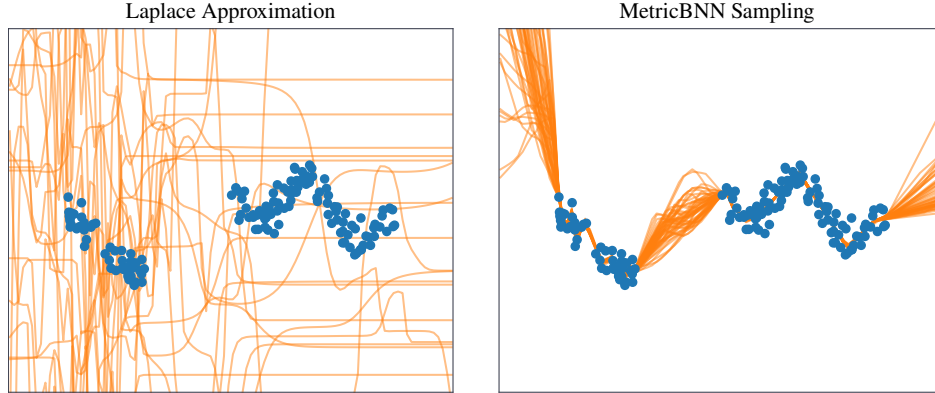


Figure 2: Posterior samples for Regression task. The blue points represent the dataset, the orange lines are samples from the estimated posteriors. The standard Laplace approximation fails in recovering the true parameter distribution. Our proposed sampling method correctly captures the uncertainty in the data gap.

3.2 Limitations of the Laplace Approximation

While widely used, the Laplace approximation has several limitations:

1. **Local Approximation:** It captures only the local curvature of the loss landscape around θ^* , ignoring the broader structure of the posterior, which can be highly non-Gaussian in high-dimensional spaces (Wilson & Izmailov, 2020).
2. **Scalability:** Computing and inverting the Hessian is computationally expensive for large-scale neural networks, limiting its applicability to small models (Martens, 2010).
3. **Positive-Definiteness:** In over-parameterized networks, the Hessian is often not positive definite, making it difficult to define a valid Gaussian approximation. Regularization or approximate methods are sometimes used to mitigate this issue, but these approaches can introduce biases (Daxberger et al., 2021).

The limitations of the Laplace approximation motivate the need for alternative methods that better capture the true posterior distribution. Specifically, the posterior for BNNs often lies on a complex, non-linear manifold in parameter space, which the Laplace approximation fails to represent. This motivates our approach, which combines efficient exploration of the parameter space with a flexible latent representation to better approximate the posterior.

4 Method

To address the limitations of the Laplace approximation and improve posterior estimation in BNNs, we propose *MetricBNN*, a two-step framework. The first step involves locally exploring the parameter space around the MAP estimate using a simple sampling method. The second step learns a latent representation to construct a flexible posterior distribution that captures the complex geometry of the parameter space.

4.1 Exploring Neighbor Solutions

Given a trained neural network with parameters θ^* that minimize the regularized loss function (Equation 1), the parameter space of deep networks is known to exhibit a high degree of redundancy due to overparameterization and reparametrization invariances Fort & Jastrzebski (2019); Garipov et al. (2018). These properties create a connected, lower-dimensional manifold of near-optimal solutions surrounding the MAP estimate θ^* . This phenomenon, often referred to as the *linear connectivity assumption*, suggests that different sets of parameters

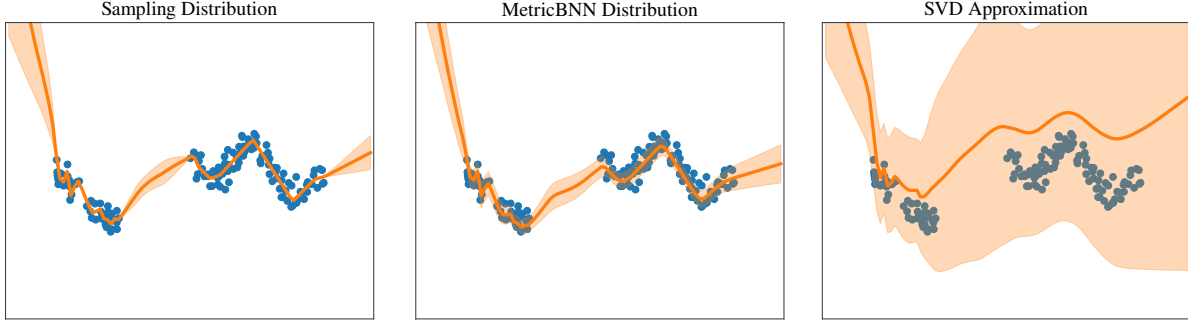


Figure 3: Estimated posterior for the regression task. The blue points represent the dataset and in orange the mean and standard deviation of posteriors sampled from the estimated distributions. A naive SVD approximation of the solutions found with the sampling scheme fails in representing correctly the true posterior. Our proposed MetricBNN posterior correctly approximates it.

can achieve similar performance, even when interpolating between them (Garipov et al., 2018; Fort & Jastrzebski, 2019; Brea et al., 2019).

In particular, empirical studies on the loss landscapes of neural networks have shown that minima are often connected by low-loss paths, forming smooth and structured regions in the parameter space (Garipov et al., 2018). This implies that the posterior distribution is not confined to a single mode but instead spans a broader, non-linear manifold. Capturing this geometric complexity is crucial for accurately modeling the posterior. By exploring neighboring solutions around θ^* , we aim to gather representative samples that reflect the underlying structure of this manifold, providing a richer and more accurate approximation of the posterior distribution.

Estimating the posterior distribution of the parameters amounts to identifying the distribution of these solutions. To achieve this, we propose collecting a set of such solutions using the following sampling technique (MetricBNN sampling):

1. **Initialization:** Start with N particles, each initialized at the MAP estimate θ^* , i.e., $\theta_{i,0}$.
2. **Random Drift Assignment:** Assign each particle a random drift vector d_i , sampled as a unit vector with uniform orientation in the parameter space.
3. **Iterative Exploration:** For T time steps, update each particle’s position as:
 - (a) **Drift:** Perturb the particle by adding the corresponding random drift.

$$\hat{\theta}_{i,t+1}^0 = \theta_{i,t} + \alpha \cdot d_i. \quad (3)$$

- (b) **Refinement:** After each drift update, refine the particle’s position using M steps of gradient descent to ensure alignment with the loss landscape.

$$\theta_{i,t+1}^{m+1} = \hat{\theta}_{i,t+1}^m - \eta \nabla_{\theta} \mathcal{L}(\hat{\theta}_{i,t+1}^m) \quad (4)$$

$$\theta_{i,t+1} = \hat{\theta}_{i,t+1}^M \quad (5)$$

This procedure generates a set of $N \times T$ viable solutions near θ^* . These samples represent a local exploration of the posterior distribution of interest, capturing the diversity of solutions around the MAP estimate. Figure 2(b) illustrates an example of these solutions for a two-dimensional regression problem.

While sampling methods have traditionally been considered computationally expensive for high-dimensional neural networks, leading to their relative neglect in favor of variational inference and other approximation techniques (MacKay, 1998; Blundell et al., 2015), we argue that these concerns warrant re-evaluation in light of the linear connectivity assumption. Specifically, the high redundancy and structured geometry of neural

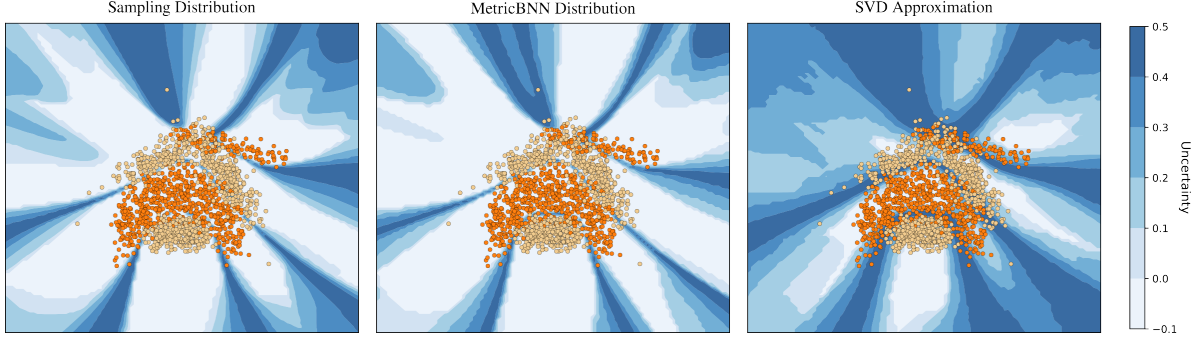


Figure 4: Estimated posterior for the Banana classification task. The orange and yellow points represent the data of the two classes in the classification dataset. In blue is the value of the uncertainty of the posteriors sampled from the estimated distributions.

network parameter spaces suggest that meaningful posterior exploration can be achieved through efficient sampling on a lower-dimensional, connected manifold of solutions (Garipov et al., 2018; Fort & Jastrzebski, 2019). This implies that the computational complexity of sampling methods may scale more favorably with the dimensionality of modern networks than previously assumed.

The empirical covariance matrix of the collected solutions provides an improved local Gaussian approximation of the posterior distribution. However, due to the extreme non-linearity of the solution manifold, a Gaussian approximation alone fails to accurately capture the true structure of the posterior distribution. Addressing this limitation requires moving beyond the Gaussian assumption.

4.2 Defining a Posterior

While the empirical covariance of the collected samples provides a basic Gaussian approximation, it fails to capture the non-linear geometry of the posterior, see Figure 3. Addressing this limitation requires moving beyond the Gaussian assumption. Given that the drift applied in the Drift step (3a) is small enough, we can assume that every element in between two sets of parameters is a viable solution. The curve defined by all the samples can, however, be arbitrarily non-linear. To overcome this, we propose to learn a representation that maps the parameters into a structured latent space where the collected trajectories of parameters are linearly distributed.

We define a latent representation through an autoencoder framework:

- $\varphi : \Theta \rightarrow Z$: An encoder that maps neural network parameters θ to a latent space $Z \subseteq \mathbb{R}^k$.
- $\varphi^{-1} : Z \rightarrow \Theta$: A decoder that reconstructs parameters θ from the latent space.

The training dataset for the autoencoder consists of the collected samples, i.e., D_θ . The goal is to learn a latent space where the posterior distribution is well-structured and linearly interpolable. We achieve this by optimizing the following loss function:

$$\mathcal{L} = \lambda_+ \mathcal{L}_+ + \lambda_- \mathcal{L}_- + \mathcal{L}_d, \quad (6)$$

with:

$$\mathcal{L}_+ = \mathbb{E}_{D_\theta} \left[\left(\|\varphi(\theta) - \varphi(\theta')\| - \frac{1}{T} \right)^2 \right], \quad (7)$$

$$\mathcal{L}_- = \mathbb{E}_{D_\theta} [-\log(\|\varphi(\theta) - \varphi(\theta'')\|)], \quad (8)$$

$$\mathcal{L}_d = \mathbb{E}_{D_\theta} [\|\varphi^{-1}(\varphi(\theta)) - \theta\|^2], \quad (9)$$

where θ and θ' are two sets of parameters of the same trajectory and successive time step while θ'' is another randomly sampled set of parameters. Both λ_+ and λ_- are scalar values. The role of these terms is as follows:

- \mathcal{L}_+ : encourages local distances between successive samples are preserved.
- \mathcal{L}_- : encourages these sets of parameters to stretch in the learned latent space by maximizing every pair-wise distance.
- \mathcal{L}_d : encourages reconstruction, mapping latent representations back to their original parameter space.

Using the trained autoencoder, we define the posterior distribution in the latent space Z , MetricBNN posterior. Each trajectory of samples is treated as independent, with a uniform probability assigned to each trajectory. Along each trajectory, we place a uniform probability between the MAP solution and the last sample of such trajectory, i.e., $\theta_{i,T}$. Sampling from the posterior involves: Sample a trajectory index i uniformly, Sample a scaling factor $\epsilon \sim \text{Uniform}(0, 1)$. Compute the parameter sample:

$$\theta = \varphi^{-1}(\varphi(\theta^*) + \epsilon \cdot (\varphi(\theta_{i,T}) - \varphi(\theta^*))). \quad (10)$$

This method captures the broader geometry of the posterior while maintaining computational efficiency. We provide a pseudo-code of the whole method in Appendix B

5 Experiments

In this section, we evaluate the ability of our proposed method to approximate the posterior efficiently. We first assess the quality of our sampling technique in capturing model uncertainty using simple regression and classification tasks. We then analyze the computational complexity of our approach as the network size increases. Finally, we present quantitative results on real-world datasets, evaluating the negative log-likelihood (NLL) and classification accuracy of our posterior distribution.

For all experiments, we compare our proposed method against the standard Laplace Approximation and the Riemannian Laplace Approximation proposed in Bergamin et al. (2024), both with and without linearization. Both approaches use the Gauss-Newton Hessian approximation and optimize the prior precision using marginal likelihood estimation. These methods, similarly to us, serve as post-hoc Bayesian inference techniques that estimate the parameter posterior without requiring modifications to the training process.

We conduct experiments across four different settings:

Simple Regression: We evaluate our method on the one-dimensional regression problem introduced in Snelson & Ghahramani (2005). The dataset consists of 200 points, with 50 held out to assess the model’s ability to estimate uncertainty. We use a fully connected neural network with three hidden layers of 32 units and ReLU activations.

Simple Classification: We test our method on the Banana dataset, a two-dimensional binary classification task with 5,300 points, of which 30% are reserved for testing. We use a fully connected network with two hidden layers of 6 units and Tanh activations.

Classification on UCI Datasets: To evaluate performance on real-world structured data, we experiment with six classification datasets from the UCI repository (Markelle et al., 2023). We use a fully connected network with two hidden layers of 32 units and ReLU activations.

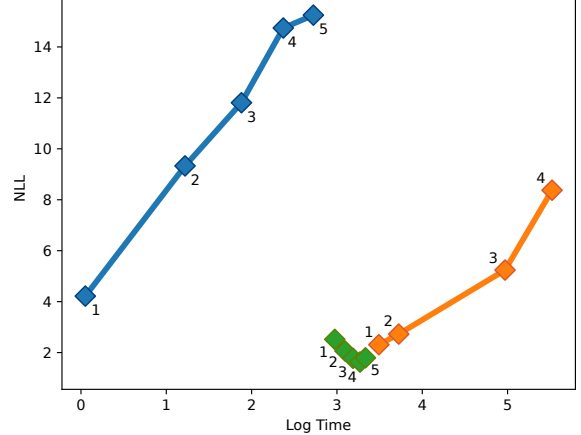


Figure 5: Trade-off between network size and computational complexity on the regression task. In blue is the Laplace Approximation, in orange is the Riemannian reformulation, in green is our proposed MetricBNN method. Each method is tested in its NLL performance (lower is better) on the test set with an increasing number of hidden layers in the network (numbers in the figure). Our proposed method does not considerably increase in computational complexity with an increase in the network size.

Table 1: Quantitative results on the UCI datasets, including Negative Log Likelihood (NLL), lower is better, and Computation Time (s).

	Model	Australian	Breast	Glass	Ionosphere	Vehicle	Waveform
NLL	LA	0.71 ± 0.06	0.73 ± 0.07	2.28 ± 0.28	0.72 ± 0.09	0.8 ± 0.16	0.88 ± 0.09
	Riem LA	0.54 ± 0.07	0.59 ± 0.1	1.42 ± 0.21	0.17 ± 0.02	0.65 ± 0.02	0.3 ± 0.01
	Lin LA	0.69 ± 0.04	0.61 ± 0.04	3.59 ± 1.17	0.42 ± 0.04	0.68 ± 0.01	0.4 ± 0.02
	Riem Lin LA	0.74 ± 0.05	2.32 ± 0.74	15.92 ± 0.05	13.5 ± 0.66	0.65 ± 0.01	0.38 ± 0.02
	MetricBNN	0.66 ± 0.05	0.57 ± 0.04	1.07 ± 0.07	0.17 ± 0.05	0.69 ± 0.02	0.3 ± 0.01
Log Time	LA	1.88 ± 0.98	1.54 ± 0.8	1.56 ± 0.67	1.88 ± 0.72	3.05 ± 1.95	2.68 ± 0.48
	Riem LA	4.36 ± 4.06	2.42 ± 1.61	3.69 ± 3.63	2.38 ± 1.61	3.33 ± 3.0	3.26 ± 2.45
	Lin LA	1.86 ± 1.0	1.52 ± 0.4	1.55 ± 0.48	1.89 ± 0.65	3.05 ± 1.95	2.68 ± 0.31
	Riem Lin LA	1.9 ± 1.04	1.63 ± 0.6	1.67 ± 0.54	1.93 ± 0.65	3.08 ± 1.42	2.69 ± 1.53
	MetricBNN	1.66 ± 0.34	1.61 ± 0.29	1.64 ± 0.32	1.76 ± 0.35	1.8 ± 0.36	1.82 ± 0.41

Image Classification: For high-dimensional problems, we experiment with the MNIST (LeCun, 1998) and FashionMNIST (Xiao et al., 2017) datasets. In line with standard practice, we use a shallow convolutional neural network with two convolutional layers followed by three fully connected layers with Tanh activations.

5.1 Scalability of the Proposed Sampling Method

Approximating the posterior distribution in neural networks is challenging due to the high-dimensional and non-linear nature of the parameter space. The standard Laplace approximation has been shown to struggle in these settings (Ritter et al., 2018; Lawrence, 2001), as it assumes a Gaussian posterior, which may not align well with the actual distribution. Our proposed sampling method provides a more flexible alternative, allowing for a tighter and better-calibrated posterior that is independent of the loss landscape’s curvature. Figure 2 compares our method with the Laplace approximation on the toy regression task, demonstrating a significantly improved uncertainty estimation.

Beyond accuracy, scalability is a key factor in Bayesian inference for deep networks. Many posterior approximation techniques rely on computing second-order derivatives, making them computationally expensive as the network size grows. To assess the computational trade-off, we analyze the inference time and NLL performance as the number of layers in the network increases. Figure 5 presents these results for our method, the standard Laplace approximation, and the Riemannian Laplace approximation (Bergamin et al., 2024). Our results show that while Hessian-based methods scale poorly, our sampling approach maintains a reasonable computational cost while preserving accuracy. Furthermore, our learned latent posterior model enables rapid parameter sampling, further reducing the overhead compared to iterative sampling.

Table 2: Quantitative results on image datasets, including Negative Log Likelihood (NLL), lower is better, and Computation Time (s).

	Model	MNIST	FMNIST
NLL	LA	2.42 ± 0.03	2.24 ± 0.08
	Riem LA	1.11 ± 0.11	1.2 ± 0.02
	Lin LA	0.96 ± 0.08	1.27 ± 0.13
	Riem Lin LA	0.62 ± 0.09	0.87 ± 0.07
	MetricBNN	0.13 ± 0.01	0.6 ± 0.1
Log Time	LA	3.69 ± 1.08	3.66 ± 2.93
	Riem LA	5.22 ± 3.16	5.01 ± 4.09
	Lin LA	3.64 ± 2.80	3.67 ± 2.96
	Riem Lin LA	3.90 ± 3.18	3.91 ± 3.29
	MetricBNN	3.61 ± 2.84	3.67 ± 2.86

5.2 Quality of the Proposed Geometric Posterior

A potential concern is whether the posterior obtained via our sampling method could be approximated using a naive covariance computation, similar to the Laplace approximation. However, as shown in Figure 3, this approach still does not yield accurate uncertainty estimates. By contrast, our latent posterior representation effectively captures the true posterior structure, demonstrating the benefits of learning a structured parameter space.

The same observations hold for classification tasks. Figure 4 illustrates the posterior obtained on the Banana dataset, where our method produces a more calibrated uncertainty estimate than Laplace-based approaches, particularly in high-uncertainty regions. This suggests that our method generalizes well beyond regression problems and provides a practical alternative for Bayesian classification.

5.3 Performance on UCI Regression and Image Classification

We further evaluate our approach to real-world structured data by testing it on standard UCI regression benchmarks. Table 1 reports NLL results, showing that our method remains competitive with Laplace-based refinements while being significantly more efficient. This demonstrates that our method offers a good balance between accuracy and computational cost.

To assess performance in high-dimensional settings, we train convolutional neural networks on MNIST and FashionMNIST. Table 2 presents classification accuracy and NLL results. As expected, performance differences become more pronounced as network complexity increases. While Laplace-based methods suffer from computational inefficiencies in these settings, our approach maintains strong performance while remaining scalable.

Our experiments highlight three key takeaways:

- Our sampling-based posterior exploration provides a more flexible and well-calibrated uncertainty estimate than the Laplace approximation, particularly as network size increases.
- The proposed latent posterior model enables efficient posterior sampling, avoiding the computational overhead of iterative methods.
- Our approach remains competitive in terms of NLL and classification accuracy on real-world datasets while being significantly more computationally efficient than Hessian-based refinements.

6 Conclusions

In this work, we proposed a simple variation of sampling-based techniques tailored to explore the posterior geometry of Bayesian Neural Networks efficiently, even in over-parameterized settings. By leveraging the low-dimensional structure of loss minima, our method achieves competitive posterior approximations while maintaining scalability as network size increases. Additionally, we introduced a model that learns a deformation of the parameter space based on the collected samples, enabling rapid posterior sampling without requiring iterative methods. Our empirical results demonstrate that this approach improves posterior accuracy and computational efficiency compared to recent refinement techniques. These contributions provide a practical and flexible framework for Bayesian inference in deep learning, offering new directions for scalable uncertainty quantification in complex models.

References

- Federico Bergamin, Pablo Moreno-Muñoz, Søren Hauberg, and Georgios Arvanitidis. Riemannian laplace approximations for bayesian neural networks. *Advances in Neural Information Processing Systems*, 36, 2024.
- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In *International conference on machine learning*, pp. 1613–1622. PMLR, 2015.
- Johanni Brea, Berfin Simsek, Bernd Illing, and Wulfram Gerstner. Weight-space symmetry in deep networks gives rise to permutation saddles, connected by equal-loss valleys across the loss landscape. *arXiv preprint arXiv:1907.02911*, 2019.
- Antonio S Cofino, Rafael Cano, Carmen Sordo, and Jose M Gutierrez. Bayesian networks for probabilistic weather prediction. In *15th European Conference on Artificial Intelligence (ECAI)*. Citeseer, 2002.

- Erik Daxberger, Agustinus Kristiadi, Alexander Immer, Runa Eschenhagen, Matthias Bauer, and Philipp Hennig. Laplace redux-effortless bayesian deep learning. *Advances in Neural Information Processing Systems*, 34:20089–20103, 2021.
- Runa Eschenhagen, Erik Daxberger, Philipp Hennig, and Agustinus Kristiadi. Mixtures of laplace approximations for improved post-hoc uncertainty in deep learning. *arXiv preprint arXiv:2111.03577*, 2021.
- Di Feng, Lars Rosenbaum, and Klaus Dietmayer. Towards safe autonomous driving: Capture uncertainty in the deep neural network for lidar 3d vehicle detection. In *2018 21st international conference on intelligent transportation systems (ITSC)*, pp. 3266–3273. IEEE, 2018.
- Stanislav Fort and Stanislaw Jastrzebski. Large scale structure of neural network loss landscapes. *Advances in Neural Information Processing Systems*, 32, 2019.
- Timur Garipov, Pavel Izmailov, Dmitrii Podoprikin, Dmitry P Vetrov, and Andrew G Wilson. Loss surfaces, mode connectivity, and fast ensembling of dnns. *Advances in neural information processing systems*, 31, 2018.
- Alex Graves. Practical variational inference for neural networks. *Advances in neural information processing systems*, 24, 2011.
- Marton Havasi, Jasper Snoek, Dustin Tran, Jonathan Gordon, and José Miguel Hernández-Lobato. Sampling the variational posterior with local refinement. *Entropy*, 23(11):1475, 2021.
- Alexander Immer, Maciej Korzepa, and Matthias Bauer. Improving predictions of bayesian neural nets via local linearization. In *International conference on artificial intelligence and statistics*, pp. 703–711. PMLR, 2021.
- Agustinus Kristiadi, Runa Eschenhagen, and Philipp Hennig. Posterior refinement improves sample efficiency in bayesian neural networks. *Advances in Neural Information Processing Systems*, 35:30333–30346, 2022.
- Neil David Lawrence. *Variational inference in probabilistic models*. PhD thesis, Citeseer, 2001.
- Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- Christian Leibig, Vaneeda Allken, Murat Seçkin Ayhan, Philipp Berens, and Siegfried Wahl. Leveraging uncertainty information from deep neural networks for disease detection. *Scientific reports*, 7(1):1–14, 2017.
- David JC MacKay. A practical bayesian framework for backpropagation networks. *Neural computation*, 4(3):448–472, 1992.
- David JC MacKay. Choice of basis for laplace approximation. *Machine learning*, 33:77–86, 1998.
- K Markelle, L Rachel, and N Kolby. Uci dataset. the uci machine learning repository, 2023.
- James Martens. Deep learning via hessian-free optimization. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, pp. 735–742, 2010.
- Andrew C Miller, Nicholas J Foti, and Ryan P Adams. Variational boosting: Iteratively refining posterior approximations. In *International Conference on Machine Learning*, pp. 2420–2429. PMLR, 2017.
- Radford M Neal. Mcmc using ensembles of states for problems with fast and slow variables such as gaussian process regression. *arXiv preprint arXiv:1101.0387*, 2011.
- Hippolyt Ritter, Aleksandar Botev, and David Barber. A scalable laplace approximation for neural networks. In *6th international conference on learning representations, ICLR 2018-conference track proceedings*, volume 6. International Conference on Representation Learning, 2018.
- Edward Snelson and Zoubin Ghahramani. Sparse gaussian processes using pseudo-inputs. *Advances in neural information processing systems*, 18, 2005.

- Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pp. 681–688. Citeseer, 2011.
- Andrew G Wilson and Pavel Izmailov. Bayesian deep learning and a probabilistic perspective of generalization. *Advances in neural information processing systems*, 33:4697–4708, 2020.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- Cheng Zhang, Judith Bütetage, Hedvig Kjellström, and Stephan Mandt. Advances in variational inference. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):2008–2026, 2018.

A Experimental Details

This section provides additional details regarding the experimental setup, including the models, hyperparameters, and computational settings.

A.1 Experimental Setup

We evaluate our method across four types of tasks: toy regression, toy classification, structured classification from the UCI repository, and high-dimensional image classification.

Toy Regression. We consider the *Snelson 1D regression dataset* (Snelson & Ghahramani, 2005), consisting of 200 points, with 50 held out for evaluating uncertainty estimation. We use a fully connected neural network with three hidden layers of 32 units and ReLU activations. The model is trained using the Adam optimizer with a learning rate of 1×10^{-3} , batch size of 200, and L2 regularization of 0.01 for 50,000 epochs.

For the sampling procedure of our approach, we used a drift-scaling factor of $\alpha = 0.001$ and performed $T = 1000$ sampling steps. At each step, we updated the parameters using $M = 10$ gradient refinement steps to improve posterior estimates. We initialized $N = 10$ particles, each following a perturbed trajectory in the parameter space, with a drift step scaled by an inner learning rate of $\eta = 0.001$. For the autoencoder-based posterior model, we projected the sampled parameter trajectories into a structured latent space of $k = 32$ dimensions. The autoencoder was trained using batch size 1024 for 1000 epochs, optimizing the contrastive loss with weight coefficients $\lambda_+ = 1.0$ and $\lambda_- = 1.0$. The architecture consists of an encoder with two fully connected hidden layers of 256 units and ReLU activations, and a decoder with a symmetric structure of two hidden layers with 256 units and ReLU activations. Training was performed using the Adam optimizer.

Toy Classification. For classification, we use the *Banana dataset*, a 2D binary classification problem with 5300 points, of which 30% are used for testing. The network consists of two hidden layers with 16 units and Tanh activations. The training procedure follows the regression setup, except with a smaller batch size (32) and 2500 training epochs. For our model, we use the same parameters as in the regression task except for $\alpha = 0.1$ and $T = 100$.

Structured Data (UCI Datasets). To evaluate our method in structured classification settings, we test on six datasets from the *UCI repository* (Markelle et al., 2023), using a fully connected architecture with two hidden layers of 32 units and ReLU activations. The training follows the same schedule as before, with a batch size of 32 and 1000 training epochs. Our method is estimated with the same parameters as in the classification experiment except for $T = 50$ steps, using $M = 10$ gradient refinements.

High-Dimensional Image Classification. For large-scale experiments, we train Bayesian neural networks on *MNIST* (LeCun, 1998) and *FashionMNIST* (Xiao et al., 2017). We use a shallow convolutional neural network consisting of two convolutional layers, followed by three fully connected layers, with Tanh activations. Training follows the same setup as UCI experiments, but with a reduced number of epochs (100). Our method follows the same parameters as in the UCI experiment.

A.2 Scalability and Computational Complexity

To conduct the experiment on the computational complexity depicted in Figure 5 we used a fully connected network with hidden layers from 1 to 5. Each layer with 15 units and Tanh as activation function. For each architecture, we measure the inference time and negative log-likelihood (NLL) performance, comparing our method with the standard Laplace approximation and Riemannian Laplace approximation (Bergamin et al., 2024) on the Regression experiment. All experiments were conducted on an NVIDIA RTX3080 GPU.

B Pseudocode

Algorithm 1 Proposed Posterior Approximation Method

Require: Trained neural network f_θ , dataset \mathcal{D} , drift scaling α , sampling steps T , gradient steps M , particles N , inner learning rate η , autoencoder latent dimension k , training epochs E

Ensure: Approximate posterior $q(\theta)$

```

1: Phase 1: Parameter Sampling
2: Initialize  $N$  particles at  $\theta^*$  (MAP estimate)
3: for each particle  $i = 1, \dots, N$  do
4:   Sample random drift direction  $d_i \leftarrow \frac{v}{\|v\|}$ ,  $v \sim \mathcal{N}(0, I)$ 
5:   for each step  $t = 1, \dots, T$  do
6:     Apply drift:  $\theta_{i,t} = \theta_{i,t-1} + \alpha d_i$ 
7:     for each gradient step  $m = 1, \dots, M$  do
8:       Refine using gradient descent:
9:        $\theta_{i,t} \leftarrow \theta_{i,t} - \eta \nabla_\theta \mathcal{L}(\theta_{i,t})$ 
10:    end for
11:  end for
12: end for
13: Store all sampled parameters  $\Theta = \{\theta_{i,t}\}$ 
14: Phase 2: Learning Latent Posterior Representation
15: Train an autoencoder  $\varphi : \Theta \rightarrow Z$  using:
16: for epoch  $e = 1, \dots, E$  do
17:   Sample minibatch of subsequent parameters  $\theta, \theta' \sim \Theta$ 
18:   Compute latent embeddings:  $z = \varphi(\theta)$  and  $z' = \varphi(\theta')$ 
19:   Compute negative parameters by reshuffling the batch:  $z''$ 
20:   Compute contrastive loss:
21:    $\mathcal{L}_+ = (\|z - z'\| - 1/T)^2$ 
22:    $\mathcal{L}_- = -\log(\|z - z''\|)$ 
23:    $\mathcal{L}_d = \|\varphi^{-1}(z) - \theta\|$ 
24:   Update autoencoder using  $\mathcal{L} = \lambda_+ \mathcal{L}_+ + \lambda_- \mathcal{L}_- + \lambda_d \mathcal{L}_d$ 
25: end for
26: Posterior Approximation and Sampling
27: Sample trajectory index  $j \sim \text{Categorical}(\{z_{\max,j}\}_{j=1}^N)$ 
28: Sample interpolation factor  $\epsilon \sim \mathcal{U}(0, 1)$ 
29: Compute latent posterior sample:  $z = z_{\text{MAP}} + \epsilon(z_{\max,j} - z_{\text{MAP}})$ 
30: Map to parameter space:  $\theta = \varphi^{-1}(z)$ 
31: return Posterior sample  $\theta$ 

```
