
On the Expressiveness of Approximate Inference in Bayesian Neural Networks

Andrew Y. K. Foong*
University of Cambridge
ykf21@cam.ac.uk

David R. Burt*
University of Cambridge
drb62@cam.ac.uk

Yingzhen Li
Microsoft Research
Yingzhen.Li@microsoft.com

Richard E. Turner
University of Cambridge
Microsoft Research
ret26@cam.ac.uk

Abstract

While Bayesian neural networks (BNNs) hold the promise of being flexible, well-calibrated statistical models, inference often requires approximations whose consequences are poorly understood. We study the quality of common variational methods in approximating the Bayesian predictive distribution. For single-hidden layer ReLU BNNs, we prove a fundamental limitation in *function-space* of two of the most commonly used distributions defined in *weight-space*: mean-field Gaussian and Monte Carlo dropout. We find there are simple cases where neither method can have substantially increased uncertainty in between well-separated regions of low uncertainty. We provide strong empirical evidence that exact inference does not have this pathology, hence it is due to the approximation and not the model. In contrast, for deep networks, we prove a universality result showing that there exist approximate posteriors in the above classes which provide flexible uncertainty estimates. However, we find empirically that pathologies of a similar form as in the single-hidden layer case can persist when performing variational inference in deeper networks. Our results motivate careful consideration of the implications of approximate inference methods in BNNs.

1 Introduction

Bayesian neural networks (BNNs) [30, 34] aim to combine the strong inductive biases and flexibility of neural networks (NNs) with the probabilistic framework for uncertainty quantification provided by Bayesian statistics. However, performing exact inference in BNNs is analytically intractable and requires approximations. A variety of scalable approximate inference techniques have been proposed, with mean-field variational inference (MFVI) [18, 6] and Monte Carlo dropout (MCDO) [14] among the most used methods. These methods have been successful in applications such as active learning and out-of-distribution detection [11, 37]. However, it is unclear to what extent the successes (and failures) of BNNs are attributable to the exact Bayesian predictive, rather than the peculiarities of the approximation method. From a Bayesian modelling perspective, it is therefore crucial to ask, *does the approximate predictive distribution retain the qualitative features of the exact predictive?*

Frequently, BNN approximations define a simple class of distributions over the model parameters, (an *approximating family*), and choose a member of this family as an approximation to the posterior.

*Equal contribution.

Both MFVI and MCDO follow this paradigm. For such a method to succeed, two criteria must be met:

Criterion 1 The approximating family must contain good approximations to the posterior.

Criterion 2 The method must then select a good approximate posterior within this family.

For nearly all tasks, the performance of a BNN only depends on the distribution over weights to the extent that it affects the distribution over predictions (i.e. in ‘function-space’). Hence for our purposes, a ‘good’ approximation is one that captures features of the exact posterior in function-space that are relevant to the task at hand. However, approximating families are often defined in weight-space for computational reasons. Evaluating **Criterion 1** therefore involves understanding how weight-space approximations translate to function-space, which is a non-trivial task for highly nonlinear models such as BNNs.

In this work we provide both theoretical and empirical analyses of the flexibility of the predictive mean and variance functions of approximate BNNs. Our major findings are:

1. For shallow BNNs, there exist simple situations where *no* mean-field Gaussian or MC dropout distribution can faithfully represent the exact posterior predictive uncertainty (**Criterion 1** is not satisfied). We prove in section 3 that in these instances the variance function of any fully-connected, single-hidden layer ReLU BNN using these families suffers a lack of ‘*in-between uncertainty*’: increased uncertainty in between well-separated regions of low uncertainty. This is especially problematic for lower-dimensional data where we may expect some datapoints to be in between others. Examples include spatio-temporal data, or Bayesian optimisation for hyperparameter search, where we frequently wish to make predictions in unobserved regions in between observed regions. We verify that the exact posterior predictive does not have this limitation; hence this pathology is attributable solely to the restrictiveness of the approximating family.

2. In section 4 we prove a universal approximation result showing that the mean and variance functions of deep approximate BNNs using mean-field Gaussian or MCDO distributions can uniformly approximate any continuous function and any continuous non-negative function respectively. However, it remains to be shown that appropriate predictive means and variances will be found when optimising the ELBO. To test this, we focus on the low-dimensional, small data regime where comparisons to references for the exact posterior such as the limiting GP [34, 25, 32] are easier to make. In section 4.2 we provide empirical evidence that in spite of its theoretical flexibility, VI in deep BNNs can still lead to distributions that suffer from similar pathologies to the shallow case, i.e. **Criterion 2** is not satisfied.

In section 5, we provide an active learning case study on a real-world dataset showing how in-between uncertainty can be a crucial feature of the posterior predictive. In this case, we provide evidence that although the inductive biases of the BNN model with exact inference can bring considerable benefits, these are lost when MFVI or MCDO are used. Code to reproduce our experiments can be found at <https://github.com/cambridge-mlg/expressiveness-approx-bnns>.

2 Background

Consider a regression dataset $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ with $\mathbf{x}_n \in \mathbb{R}^D$ and $y_n \in \mathbb{R}$. To define a BNN, we specify a prior distribution with density $p(\theta)$ over the NN parameters. Each parameter setting corresponds to a function $f_\theta : \mathbb{R}^D \rightarrow \mathbb{R}$. We specify a likelihood $p(\{y_n\}_{n=1}^N | \{\mathbf{x}_n\}_{n=1}^N, f_\theta)$ which describes the relationship between the observed data and the model parameters. The posterior distribution over parameters has density $p(\theta | \mathcal{D}) \propto p(\{y_n\}_{n=1}^N | \{\mathbf{x}_n\}_{n=1}^N, f_\theta)p(\theta)$. The posterior does not have a closed form and approximations must be made in order to make predictions.

2.1 Approximate Inference Methods

Many approximate inference algorithms define a parametric class of distributions, \mathcal{Q} , from which to select an approximation to the posterior. For BNNs, the distributions in \mathcal{Q} are defined over the model parameters θ . For example, \mathcal{Q} may be the set of all fully-factorised Gaussian distributions, in which case the variational parameters ϕ are a vector of means and variances. We denote this family as \mathcal{Q}_{FFG} . A density $q_\phi(\theta) \in \mathcal{Q}$ is then chosen to best approximate the exact posterior according to some criteria. Once q_ϕ is selected, predictions at a test point (\mathbf{x}_*, y_*) can be made by replacing the

expectation under the exact posterior by an expectation under the approximate posterior:

$$p(y_*|\mathbf{x}_*, \mathcal{D}) = \mathbb{E}_{p(\theta|\mathcal{D})} [p(y_*|\mathbf{x}_*, f_\theta)] \approx \mathbb{E}_{q_\phi(\theta)} [p(y_*|\mathbf{x}_*, f_\theta)] \approx \frac{1}{M} \sum_{m=1}^M p(y_*|\mathbf{x}_*, f_{\theta_m}), \quad (1)$$

where $\theta_m \sim q_\phi$ on the RHS of equation (1). Many approximate inference algorithms may share the same \mathcal{Q} , e.g. VI, the diagonal Laplace approximation [9], probabilistic backpropagation [16], stochastic expectation propagation [28], black-box alpha divergence minimisation [17], Rényi divergence VI [27], natural gradient VI [23] and functional variational BNNs [41] all frequently use \mathcal{Q}_{FFG} .

Mean-Field Variational Inference Variational inference [3, 22] is an approximate inference method that selects q_ϕ by minimising the KL divergence between the approximate and exact posterior [5]. This is equivalent to maximising an evidence lower bound (ELBO): $\mathcal{L}(\phi) = \sum_{n=1}^N \mathbb{E}_{q_\phi} [\log p(y_n|\mathbf{x}_n, f_\theta)] - \text{KL} [q_\phi(\theta)||p(\theta)]$. Most commonly in BNNs, q_ϕ is chosen from \mathcal{Q}_{FFG} . This is known as mean-field variational inference (MFVI).

Monte Carlo Dropout MCDO with ℓ_2 regularisation has been interpreted as VI [13]. Although the MCDO objective is not strictly an ELBO [20], we will sometimes refer to it as such. The variational family, $\mathcal{Q}_{\text{MCDO}}$, is the set of distributions determined by random variables of the form $\widehat{\mathbf{W}} := \mathbf{W} \text{diag}(\epsilon)$; where the weights \mathbf{W} are variational parameters and $\epsilon_i \stackrel{\text{i.i.d.}}{\sim} \text{Bern}(1-p)$, with p the dropout probability. Frequently, the first weight matrix \mathbf{W}_1 is deterministic (i.e. inputs are not dropped out) — we analyse this case in the main body and use $\mathcal{Q}_{\text{MCDO}}$ to refer to this family. There are fundamentally different considerations when \mathbf{W}_1 is also stochastic, addressed in appendix E.

2.2 BNN Priors and References for the Exact Posterior Predictive

In this paper, we examine how closely approximate BNN predictive distributions resemble exact inference. To make this comparison, a choice of BNN prior must be made. Common practice is to choose an independent $\mathcal{N}(0, 1)$ prior for all parameters, regardless of the size of the network. However, such priors are known to lead to extremely large variance in function space for wide or deep networks [34]. For example, choosing such a prior for a 4-hidden layer BNN with 50 neurons in each layer leads to a prior standard deviation of $\sim 10^3$ in function space at the origin. This is orders of magnitude too large for normalised data. It is conceivable that one may combine an unreasonable prior with poor approximate inference to obtain practically useful uncertainty estimates that bear little relation to the exact Bayesian predictive — we do not consider this case. Instead, we focus our study on the quality of approximate inference in models with moderate prior variances in function space.

There is a body of literature on BNN priors [34, 32, 39, 25] which shows how to select prior weight variances that lead to reasonable prior variances in function space, even as the width of the hidden layers tends to infinity. For a layer with N_{in} inputs, we choose independent $\mathcal{N}(0, \sigma_w^2/N_{\text{in}})$ priors for the weights, with σ_w^2 a constant. For regression with a Gaussian likelihood, as the width tends to infinity, both the prior and posterior of such a BNN converges to a Gaussian process (GP) [21, 34, 32]. It has been shown that even moderately wide BNNs closely resemble their corresponding infinite-width GP counterparts [32]. In this work, we use exact inference in the corresponding infinite-width limit GP and also ‘gold-standard’ Hamiltonian Monte Carlo (HMC) [35, 19] as references for the exact posterior.

3 Single-Hidden Layer Neural Networks

In this section, we prove that for single-hidden layer (IHL) ReLU BNNs, \mathcal{Q}_{FFG} and $\mathcal{Q}_{\text{MCDO}}$ are not expressive enough to satisfy **Criterion 1**. We identify limitations on the variance in function-space, $\mathbb{V}[f(\mathbf{x})]$, implied by these families. We show empirically that the exact posterior does not have these restrictions, implying that approximate inference does not qualitatively resemble the posterior.

Theorem 1 (Factorised Gaussian). *Consider any IHL fully-connected ReLU NN $f: \mathbb{R}^D \rightarrow \mathbb{R}$. Let x_d denote the d^{th} element of the input vector \mathbf{x} . Assume a fully factorised Gaussian distribution over the parameters. Consider any points $\mathbf{p}, \mathbf{q}, \mathbf{r} \in \mathbb{R}^D$ such that $\mathbf{r} \in \overline{\mathbf{pq}}$ and either: i. $\overline{\mathbf{pq}}$ contains $\mathbf{0}$ and \mathbf{r} is closer to $\mathbf{0}$ than both \mathbf{p} and \mathbf{q} , or ii. $\overline{\mathbf{pq}}$ is orthogonal to and intersects the plane $x_d = 0$, and \mathbf{r} is closer to the plane $x_d = 0$ than both \mathbf{p} and \mathbf{q} . Then $\mathbb{V}[f(\mathbf{r})] \leq \mathbb{V}[f(\mathbf{p})] + \mathbb{V}[f(\mathbf{q})]$.*

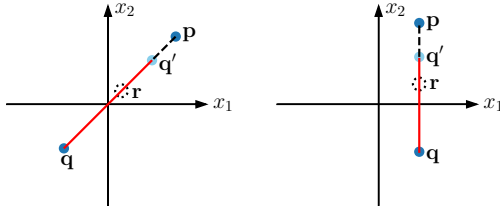


Figure 1: Illustration of the bounded regions in theorem 1, showing the input domain of a 1HL mean-field Gaussian BNN, for the case $\mathbf{x} \in \mathbb{R}^2$. Left (resp. Right): For any two points \mathbf{p} and \mathbf{q} such that the line joining them crosses the origin (resp. is orthogonal to and intersects a plane $x_d = 0$), the output variance at any point \mathbf{r} on the solid red portion of the line is upper bounded by $\mathbb{V}[f(\mathbf{p})] + \mathbb{V}[f(\mathbf{q})]$, illustrating condition (i) (resp. condition (ii)) of theorem 1. The bounded region extends from $\mathbf{q} = (q_1, q_2)$ to \mathbf{q}' , where $\mathbf{q}' = (-q_1, -q_2)$ (Left), or $\mathbf{q}' = (q_1, -q_2)$ (Right).

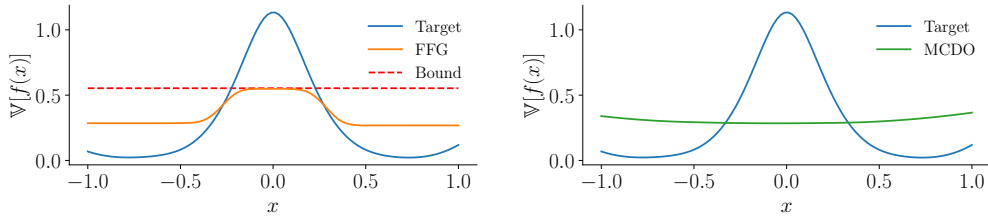


Figure 2: Results of *directly minimising the squared error in function space* between $\mathbb{V}[f(x)]$ (for a single-hidden layer NN) and a target variance function. Left: FFG distribution, Right: MCDO distribution. The bound for FFG distributions (red) applies on $[-1, 1]$ with $\mathbf{p} = -1, \mathbf{q} = 1$. The MCDO variance function is convex, and almost constant. The FFG and MCDO variance functions underestimate the target near the origin and overestimate it away from the origin.

Theorem 1 states that there are line segments (illustrated in figure 1) in input space such that the predictive variance on the line is bounded by the sum of the variance at the endpoints. Analogous but weaker bounds on higher dimensional sets in input space enclosed by these lines can be obtained as a corollary (see appendix B). Theorem 1 applies to any method using \mathcal{Q}_{FFG} , as listed in section 2.1. Although theorem 1 only bounds certain lines in input space, in appendix A we provide figures empirically showing that lines joining random points in input space suffer from similar behaviour. We provide similar results for MC dropout:

Theorem 2 (MC dropout). *Consider the same network architecture as in theorem 1. Assume an MC dropout distribution over the parameters, with inputs not dropped out. Then $\mathbb{V}[f(\mathbf{x})]$ is convex in \mathbf{x} .*

Remark 1. *In appendix E, we consider MC dropout with the inputs also dropped out. We prove that the variance at the origin is bounded by the maximum of the variance at any set of points containing the origin in their convex hull. This also applies to variational Gaussian dropout [24]. In the main body, we assume inputs are not dropped out.*

Theorem 2 implies the predictive variance on any line segment in input space is bounded by the maximum of the variance at its endpoints. Full proofs of theorems 1 and 2 are in appendix B. Theorems 1 and 2 show that there are simple cases where 1HL approximate BNNs using \mathcal{Q}_{FFG} and $\mathcal{Q}_{\text{MCDO}}$ cannot represent *in-between uncertainty*: i.e., increased uncertainty in between well separated regions of low uncertainty. As theorems 1 and 2 depend only on the approximating family, this cannot be fixed by improving the optimiser, regulariser or prior. Figure 2 shows a numerical verification of theorems 1 and 2. Since we are concerned with whether there are *any* distributions that show in-between uncertainty, we do not maximise the ELBO in this experiment (we consider ELBO maximisation in sections 3.2 and 4.2). Instead, we train 1HL networks of width 50 with \mathcal{Q}_{FFG} and $\mathcal{Q}_{\text{MCDO}}$ distributions to *directly* minimise the squared error between $\mathbb{V}[f(x)]$ and a pre-specified target variance function displaying in-between uncertainty. Full details are given in appendix E.3. Although theorems 1 and 2 apply only to 1HL BNNs, 1HL BNN regression tasks are a very common benchmark in the BNN literature [33, 43, 14, 16, 41], and have been used to assess different inference methods.

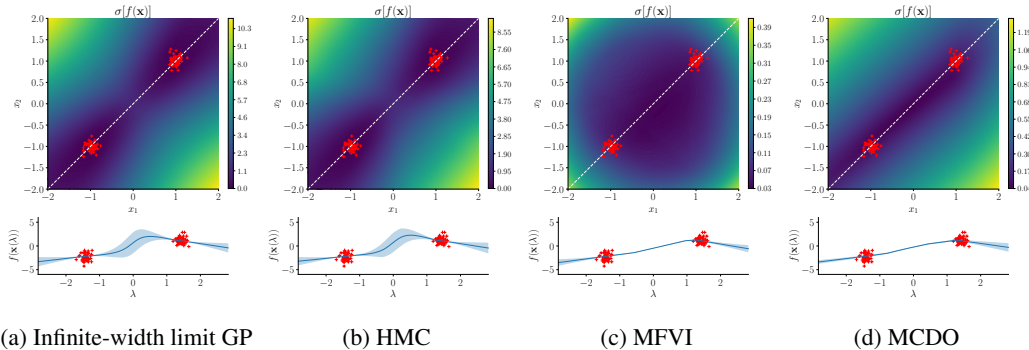


Figure 3: Regression on a 2D synthetic dataset (red crosses). The colour plots show the standard deviation of the output, $\sigma[f(\mathbf{x})]$, in 2D input space. The plots beneath show the mean with 2-standard deviation bars along the dashed white line (parameterised by λ). MFVI and MCDO are overconfident for $\lambda \in [-1, 1]$. Theorems 1 and 2 explain this: given the uncertainty is near zero at the data, there is *no* setting of the variational parameters that has variance much greater than zero in the line segment between them.

3.1 Intuition for Theorems 1 and 2

We now provide intuition for the proofs of theorems 1 and 2. Let θ_{in} be the parameters in the first layer. By the law of total variance, $\mathbb{V}[f(\mathbf{x})] = \mathbb{E}[\mathbb{V}[f(\mathbf{x})|\theta_{\text{in}}]] + \mathbb{V}[\mathbb{E}[f(\mathbf{x})|\theta_{\text{in}}]]$. For $\mathcal{Q}_{\text{MCDO}}$ the second term is 0 as θ_{in} is deterministic. Hence to prove theorem 2, it suffices to show the first term is convex. We have:

$$\mathbb{V}[f(\mathbf{x})|\theta_{\text{in}}] = \mathbb{V}\left[\sum_{i=1}^I w_i \psi(a_i(\mathbf{x}; \theta_{\text{in}})) + b \mid \theta_{\text{in}}\right] = \sum_{i=1}^I \mathbb{V}[w_i \psi(a_i(\mathbf{x}; \theta_{\text{in}}))^2] + \mathbb{V}[b], \quad (2)$$

where $\{w_i\}_{i=1}^I$ and b are the output weights and bias, $\psi(a) = \max(0, a)$, and $a_i(\mathbf{x}; \theta_{\text{in}})$ is the activation of the i^{th} neuron. Since $a_i(\mathbf{x}; \theta_{\text{in}})$ is affine in \mathbf{x} , $\psi(a_i(\mathbf{x}; \theta_{\text{in}}))^2$ is a ‘rectified quadratic’ in \mathbf{x} and therefore convex. This proves theorem 2. The same argument also applies to show that $\mathbb{V}[f(\mathbf{x})|\theta_{\text{in}}]$ is convex for \mathcal{Q}_{FFG} . To arrive at equation (2), we used that for \mathcal{Q}_{FFG} and $\mathcal{Q}_{\text{MCDO}}$, the output weights of each neuron are independent. Correlations between the weights could introduce negative covariance terms, leading to non-convex behaviour. Thus we see how *weight-space* factorisation assumptions can lead to *function-space* restrictions on the predictive uncertainty.

To complete the proof of theorem 1, we need to analyse $\mathbb{V}[\mathbb{E}[f(\mathbf{x})|\theta_{\text{in}}]]$. Because of the factorisation assumptions on the weights in the first layer, this term is a linear combination of the variances of each activation function. While these variances are not convex, in appendix B we show they satisfy restrictive conditions that imply bounds on arbitrary positive linear combinations of these functions.

3.2 Empirical Tests of Approximate Inference in Single-Hidden Layer BNNs

It is not immediately apparent that theorems 1 and 2 are problematic from the perspective of Bayesian inference. For example, even exact inference in a Bayesian linear regression model results in a convex predictive variance function. Here we provide strong evidence that, in contrast, the modelling assumptions of 1HL BNNs lead to *exact* posteriors that *do* show in-between uncertainty. Theorems 1 and 2 thus imply that it is *approximate* inference with \mathcal{Q}_{FFG} or $\mathcal{Q}_{\text{MCDO}}$ that fails to reflect this intuitively desirable property of the exact predictive, violating **Criterion 1**.

Figure 3 compares the predictive distributions obtained from MFVI and MCDO (here we optimise the ELBO for MFVI and the standard MCDO objective, in contrast with figure 2 — see appendix F for experimental details) with HMC and the limiting GP on a regression dataset consisting of two clusters of covariates. We use 1HL BNNs with 50 hidden units and ReLU activations. The HMC and limiting GP posteriors are almost indistinguishable, suggesting they both resemble the exact predictive. For these methods $\mathbb{V}[f(\mathbf{x})]$ is markedly larger near the origin than near the data. In contrast, MFVI and MCDO are as confident in between the data as they are near the data. This provides strong evidence that the lack of in-between uncertainty is not a feature of the BNN model or prior, but is caused by approximate inference.

4 Deeper Networks

Theorems 1 and 2 pose an important question: is the structural limitation observed in the 1HL case fundamental to \mathcal{Q}_{FFG} and $\mathcal{Q}_{\text{MCDO}}$ even in deeper networks, or can depth help these approximations satisfy **Criterion 1**? In theorem 3, we provide universality results for the mean and variance functions of approximate BNNs with at least two hidden layers using \mathcal{Q}_{FFG} and $\mathcal{Q}_{\text{MCDO}}$. As the predictive mean and variance often determine the performance of BNNs in regression applications, this provides theoretical evidence that approximate inference in *deep* BNNs satisfies **Criterion 1**.

Theorem 3 (Deeper networks). *Let g be any continuous function on a compact set $A \subset \mathbb{R}^D$, and h be any continuous, non-negative function on A . For any $\epsilon > 0$, for both \mathcal{Q}_{FFG} and $\mathcal{Q}_{\text{MCDO}}$ there exists a 2HL ReLU BNN such that $\sup_{\mathbf{x} \in A} |\mathbb{E}[f(\mathbf{x})] - g(\mathbf{x})| < \epsilon$ and $\sup_{\mathbf{x} \in A} |\mathbb{V}[f(\mathbf{x})] - h(\mathbf{x})| < \epsilon$.*

Remark 2. *If MC dropout is used with inputs also dropped out, the analogous statement to theorem 3 is false. In appendix E, we provide a counterexample that holds for arbitrarily deep networks and shows that if this is the case, $\mathbb{V}[f]$ cannot be made small at two points $\mathbf{x}_1, \mathbf{x}_2$ which have significantly different values of $\mathbb{E}[f(\mathbf{x}_1)]$ and $\mathbb{E}[f(\mathbf{x}_2)]$.*

Figure 4 shows the result of directly minimising the squared error between the network output mean and variance and a given target mean and variance function, using the same method and architecture as with the 1HL network in figure 2. In contrast to figure 2, the variances of both \mathcal{Q}_{FFG} and $\mathcal{Q}_{\text{MCDO}}$ are able to fit the target.

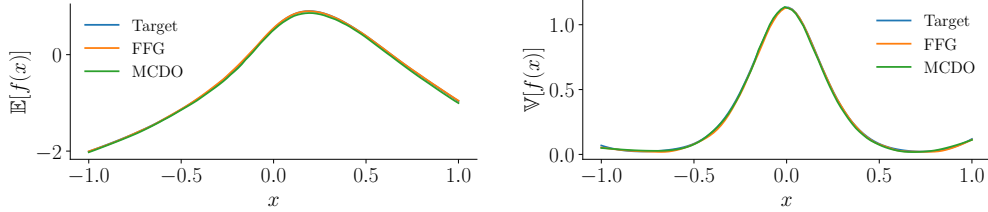


Figure 4: Results of minimising the squared error in function space between $\mathbb{E}[f(x)]$ and a target mean function (left), and between $\mathbb{V}[f(x)]$ and a target variance function (right), for a 2-hidden layer BNN with FFG and MCDO distributions.

While theorem 3 gives some cause for optimism for approximating family methods with deep BNNs, it only shows that the mean and variance of marginal distributions of the output are universal (it does not tell us about higher moments or covariances between outputs). Additionally, it does not say whether good distributions will actually be *found* by an optimiser when maximising the ELBO, i.e it does not address **Criterion 2**.

4.1 Proof Sketch of Theorem 3

To prove theorem 3 for \mathcal{Q}_{FFG} , we provide a construction that relies on the universal approximation theorem for deterministic NNs [26]. Consider a 2HL NN whose second hidden layer has two neurons, with activations a_1, a_2 . Let w_1, w_2 denote the weights connecting a_1, a_2 to the output, and b denote the output bias, such that the output $f(\mathbf{x}) = w_1\psi(a_1) + w_2\psi(a_2) + b$. In this construction, a_1 will be used to control the mean, and a_2 the variance, of the BNN output. By setting the variances in the first two linear layers to be sufficiently small, we can consider a_1 and a_2 to be essentially deterministic functions of \mathbf{x} . By the universal approximation theorem, a_1 and a_2 can approximate any continuous functions. Choose $a_1 \approx g(\mathbf{x}) - \min_{\mathbf{x}' \in A} g(\mathbf{x}')$ and $a_2 \approx \sqrt{h(\mathbf{x})}$. Choose $\mathbb{E}[b] = \min_{\mathbf{x}' \in A} g(\mathbf{x}')$, $\mathbb{V}[b] \approx 0$; $\mathbb{E}[w_1] = 1$, $\mathbb{V}[w_1] \approx 0$; and $\mathbb{E}[w_2] = 0$, $\mathbb{V}[w_2] = 1$. By linearity of expectation, the factorisation assumptions, and $a_1, a_2 \geq 0$:

$$\begin{aligned} \mathbb{E}[f(\mathbf{x})] &= \mathbb{E}[w_1\psi(a_1) + w_2\psi(a_2) + b] = \mathbb{E}[w_1] \mathbb{E}[\psi(a_1)] + \mathbb{E}[w_2] \mathbb{E}[\psi(a_2)] + \mathbb{E}[b] \\ &\approx g(\mathbf{x}) - \min_{\mathbf{x}' \in A} g(\mathbf{x}') + \min_{\mathbf{x}' \in A} g(\mathbf{x}') = g(\mathbf{x}), \end{aligned}$$

as desired. By the law of total variance, the variance of the network output is

$$\mathbb{V}[f(\mathbf{x})] = \mathbb{E}[\mathbb{V}[f(\mathbf{x})|a_1, a_2]] + \mathbb{V}[\mathbb{E}[f(\mathbf{x})|a_1, a_2]] \approx \mathbb{E}[\mathbb{V}[f(\mathbf{x})|a_1, a_2]] \approx \mathbb{E}[\psi(a_2)^2] + \mathbb{V}[b] \approx h(\mathbf{x}),$$

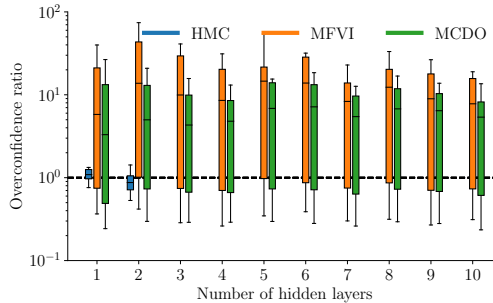


Figure 5: Box and whisker plots of the overconfidence ratios of HMC, MFVI and MCDO relative to exact inference in the corresponding infinite-width limit GP. The whiskers show the smallest and largest overconfidence ratios computed, and the the box extends from the lower to upper quartile values of the overconfidence ratios, with a line at the median. HMC is only run for 1 and 2 hidden layers due to difficulty ensuring convergence in larger models. We fix the BNN width to 50.

where we used that w_1, b are essentially deterministic and $\mathbb{V}[\mathbb{E}[f(\mathbf{x})|a_1, a_2]] \approx 0$ since a_1, a_2 are essentially deterministic. Also, we have that $\psi(a_2) \approx a_2$ since $a_2 \approx \sqrt{h(\mathbf{x})} \geq 0$. The approximations come from the standard universal function approximation theorem, and the variances of weights not being set exactly to 0 so that we remain in \mathcal{Q}_{CFG} . A rigorous proof, along with a proof for $\mathcal{Q}_{\text{MCDO}}$ with any dropout rate $p \in (0, 1)$, is given in appendix D. The proof for $\mathcal{Q}_{\text{MCDO}}$ uses a similar strategy, but is more involved as we cannot set individual weights to be essentially deterministic.

4.2 Empirical Tests of Approximate Inference in Deep BNNs

We now consider empirically whether the distributions found by optimising the ELBO with these families resemble the exact predictive distribution (**Criterion 2**). To do this, we define the ‘overconfidence ratio’ at an input \mathbf{x} as $\gamma(\mathbf{x}) = (\mathbb{V}_{\text{GP}}[f(\mathbf{x})]/\mathbb{V}_{q_\phi}[f(\mathbf{x})])^{1/2}$, where \mathbb{V}_{GP} is the predictive variance of exact inference in the infinite-width BNN. We then compute $\gamma(\mathbf{x})$ at 300 points $\{\mathbf{x}_n\}_{n=1}^{300}$ evenly spaced along the dashed white line joining the data clusters in figure 3, i.e., from $\mathbf{x} = (-1.2, -1.2)$ to $\mathbf{x} = (1.2, 1.2)$. We then create boxplots of the values $\{\gamma(\mathbf{x}_n)\}_{n=1}^{300}$ for varying BNN depths, shown in figure 5. Accurate inference should lead to similar uncertainty estimates to the limiting GP, i.e. the boxplot should be tightly centered around 1 (dashed line). For the 1HL and 2HL BNNs, the GP and HMC agree closely, suggesting both resemble the exact predictive. In contrast, MFVI and MCDO are often an order of magnitude overconfident ($\gamma(\mathbf{x}) > 1$) *between* the data clusters (upper tail of the boxplot) and somewhat underconfident ($\gamma(\mathbf{x}) < 1$) *at* the data clusters (lower tail of the boxplot). Increased depth does not alleviate this behaviour. See appendix F for experimental details and figures demonstrating this for different priors. In addition, in appendix A, we plot the uncertainty on line segments in between *random* clusters of data in a 5-dimensional input space, with similar results, showing that this phenomenon is not specific to the dataset from figure 3.

In light of theorem 3, it is perhaps surprising that VI fails to capture important properties of the predictive with deep networks. In appendix G we initialise the variational parameters such that the approximate predictive has mean and variance functions that closely match a reference predictive that exhibits in-between uncertainty. This is done by directly minimising the squared loss between the BNN mean and variance functions and the references. We find that proceeding to optimise the ELBO from this initialisation *still* leads to a lack of in-between uncertainty. This suggests that the objective function is at least partially at fault for the mismatch between the approximate and exact posteriors.

5 Case Study: Active Learning with BNNs

We now consider the impact of the pathologies described in sections 3 and 4 on active learning [40] on a real-world dataset, where the task is to use uncertainty information to intelligently select which points to label. Active learning with approximate BNNs has been considered in previous works, often showing improvements over random selection of datapoints [16, 15]. However, in cases when active

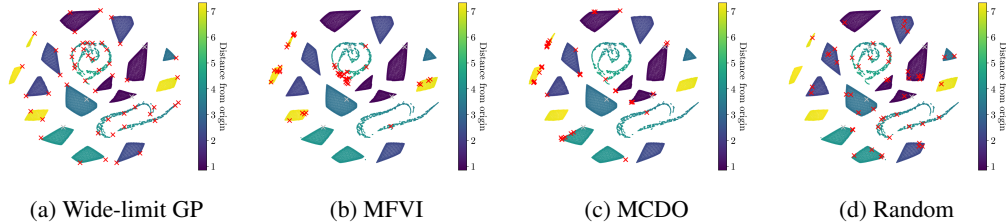


Figure 6: Points chosen during active learning in the 1HL case. Colours denote distance from the origin in 14-dimensional input space. Grey crosses (X) denote the five points randomly chosen as an initial training set. Red crosses (X) denote the 50 points selected by active learning. Both MFVI and MCDO entirely miss some clusters which are nearer the origin, and oversample certain clusters which are far from the origin, as might be expected of methods that struggle to represent in-between uncertainty. In contrast, the limiting GP samples the ‘corners’ of each cluster, without missing any entirely. Note that t-SNE does not preserve relative positions, so that clusters near the origin may appear on the ‘outside’ of the t-SNE plot.

learning fails, common metrics such as RMSE are insufficient to diagnose the causes. In particular, it is difficult to attribute the failure to the model or to poor approximate inference. In this section, we specifically analyse a dataset where we have observed active learning with approximate BNNs to fail — the Naval regression dataset [8], which is 14-dimensional and consists of 11,934 datapoints. We find via PCA that this dataset has most of its variance along a single direction. It hence may be especially problematic for methods that struggle with in-between uncertainty, as points are more likely to lie roughly in between others.

The main questions we address are: i) Is a lack of in-between uncertainty indicative of pathological behaviour in the 1HL case? In higher dimensional datasets such as Naval, it is not immediately apparent that theorems 1 and 2 are problematic, since the convex hull of the datapoints may have low volume in high dimensions. However, these theorems may be symptomatic of *related* problems with 1HL BNN uncertainty estimates. ii) Given theorem 3, will deeper approximate BNNs usefully reflect BNN modelling assumptions for active learning?

Experimental Set-up and Results We compare MFVI, MCDO and the limiting GP. We do not run HMC as this would take too long to wait for convergence at each iteration of active learning. We normalise the dataset to have zero mean and unit standard deviation in each dimension. 5 datapoints are chosen randomly as an initial active set, with the rest being the pool set. Models are trained on the active set, then the datapoint from the pool set that has the highest predictive variance is added to the active set, following Hernández-Lobato and Adams [16]. We train MFVI and MCDO for 20,000 iterations of ADAM at each step of active learning. This process is repeated 50 times. Table 1 shows the RMSE of each model on a held-out test set after this process, compared to a baseline where points are chosen randomly. Full details are in appendix H.1. Active learning significantly reduces RMSE for the GP compared to random selection, often by more than a factor of three. However it *increases* RMSE for 1HL MFVI and MCDO, and either increases it or does not significantly decrease it for deeper networks. The one exception is 3HL MCDO, where active performs about 10% better than random, which is still far less than the factor of three improvement suggested by exact inference in the infinite-width BNN.

Discussion In figure 6 we visualise the dataset using t-SNE [44]. The covariates of Naval are clustered, with points in the same cluster roughly the same distance from the origin. Since the

Table 1: Test RMSEs (± 1 standard error) after the 50th iteration of active learning, averaged over 20 random seeds. As the data is normalised, a method that predicts 0 will have an RMSE near 1.

	1 HL	2 HL	3 HL	4 HL
GP Active	0.04 \pm 0.00	0.04 \pm 0.00	0.04 \pm 0.00	0.05 \pm 0.00
GP Random	0.12 \pm 0.01	0.13 \pm 0.01	0.15 \pm 0.01	0.16 \pm 0.01
MFVI Active	0.94 \pm 0.11	0.46 \pm 0.04	0.35 \pm 0.03	0.31 \pm 0.02
MFVI Random	0.15 \pm 0.01	0.23 \pm 0.01	0.28 \pm 0.01	0.32 \pm 0.01
MCDO Active	0.69 \pm 0.04	0.36 \pm 0.02	0.38 \pm 0.02	0.45 \pm 0.02
MCDO Random	0.22 \pm 0.01	0.35 \pm 0.01	0.43 \pm 0.01	0.47 \pm 0.02

dataset is mean-centred, points closer to the origin are in a sense ‘in between’ others. We see that although the GP chooses points from every cluster during active learning, MFVI fails to select any points from many clusters — including all the clusters closest to the origin. It ignores points in the ‘inside’ and oversamples points on the ‘outside’, leading to a selection strategy worse than random. This behaviour is consistent with theorem 1. MCDO also fails to sample from many clusters; in appendix H.2 we show this is because it fails to reduce its uncertainty at clusters it has already heavily sampled. Interestingly, it sometimes chooses from clusters near the origin, even though its variance function is provably convex. This may be because the minimum of the variance function for MCDO is not centred at the origin, or because the variance has the shape of an elongated valley. In contrast, the GP seems to select the ‘corners’ of each cluster, which is intuitively efficient. The success of the infinite-width GP provides evidence that this BNN model combined with exact inference has desirable inductive biases for this task; it is *approximate* inference that has caused active learning to fail. In deeper networks, theorem 3 gives hope that the BNN predictive variance may be useful for active learning. While we find the problems are indeed less severe than in the IHL case, MFVI still oversamples points away from the origin compared to those near the origin (see appendix H.2).

6 Related Work

Concerns have been raised about the suitability of Q_{FFG} since early work on BNNs. MacKay [30, Figure 1] noted that a full-covariance Gaussian family was needed to obtain predictions with increased uncertainty away from data with the Laplace approximation, although no detailed explanation was provided. The desire to go beyond Q_{FFG} has motivated a great deal of research into more flexible approximating families [2, 29, 38]. However, to our knowledge, theorem 1 is the first theoretical result showing that Q_{FFG} can have a pathologically restrictive effect on BNN predictive uncertainties.

Recently, Farquhar et al. [10] argued that MFVI becomes a less restrictive approximation as depth increases in BNNs. However, they use different criteria to assess the quality of approximate inference. Farquhar et al. [10] use performance on classification benchmarks such as ImageNet and also the KL-divergence between certain Gaussian approximations to HMC samples in weight space to evaluate inference. In contrast, we investigate the resemblance between the function-space predictive distributions of the approximate and exact posteriors with the same prior, and focus on separating the effects of modelling and approximate inference. Additionally, we do not consider posterior tempering, and we use a different scaling for our BNN priors (see section 2.2). Recently, Wenzel et al. [46] also performed a study of the quality of approximate inference in deep BNNs. They focused on stochastic gradient Markov Chain Monte Carlo (SGMCMC) [45, 7, 47] in deep convolutional networks, and concluded that SGMCMC is accurate enough for inference, suggesting that the prior is at fault for poor performance. In contrast, we provide theoretical and empirical evidence for a specific pathology in VI in lower-dimensional regression problems, and demonstrate cases where BNN priors *do* encode useful inductive biases which are subsequently lost by approximate inference.

Osband et al. [36] note that the MCDO predictive distribution is invariant to duplicates of the data, and in the linear case predictive uncertainty does not decrease as dataset size increases (if the dropout rate and regulariser are fixed²). Theorem 2 shows that in the non-linear IHL case, predictive uncertainty in the MCDO posterior is fundamentally restricted even for datasets without repeated entries.

7 Conclusions

Principled approximate Bayesian inference involves defining a reasonable model, then finding an approximate posterior that retains the relevant properties of the exact posterior. We have shown that for BNNs, in-between uncertainty is a feature of the predictive that is often lost by variational inference. Although this is of greatest relevance for lower-dimensional regression tasks, the fact that MFVI and MCDO often fail these simple sanity checks indicates that these methods might generally have predictive distributions which are qualitatively different from the exact predictive. While BNNs have previously been shown to provide uncertainty estimates that are useful for a range of tasks, it remains an open question as to what extent this is attributable to a resemblance between the approximate and exact predictive posteriors.

²Note that for a fixed prior, the ‘KL condition’ [13, Section 3.2.3] requires the ℓ_2 regularisation constant to decrease with increasing dataset size.

Broader Impact

Bayesian approaches to deep learning problems are often proposed in situations where uncertainty estimation is critical. Often the justification given for this approach is the probabilistic framework of Bayesian inference. However, in cases where approximations are made, the quality of these approximations should also be taken into account. Our work illustrates that the uncertainty estimates given by approximate inference with commonly used algorithms may not qualitatively resemble the uncertainty estimates implied by Bayesian modelling assumptions. This may possibly have adverse consequences if Bayesian neural networks are used in safety-critical applications. Our work motivates a careful consideration of these situations.

Acknowledgments and Disclosure of Funding

We thank Wessel Bruinsma for the proof of lemma 5, and José Miguel Hernández-Lobato, Ross Clarke and Sebastian W. Ober for helpful discussions. AYKF gratefully acknowledges funding from the Trinity Hall Research Studentship and the George and Lilian Schiff Foundation. DRB gratefully acknowledges funding from the Herchel Smith Fellowship through Williams College, as well as the Qualcomm Innovation Fellowship. RET is supported by Google, Amazon, ARM, Improbable, EPSRC grants EP/M0269571 and EP/L000776/1, and the UKRI Centre for Doctoral Training in the Application of Artificial Intelligence to the study of Environmental Risks (AI4ER).

References

- [1] Milton Abramowitz and Irene A Stegun. *Handbook of mathematical functions: with formulas, graphs, and mathematical tables*, volume 55. Courier Corporation, 1965.
- [2] David Barber and Christopher M Bishop. Ensemble learning in Bayesian neural networks. *Neural networks and machine learning*, 168:215–237, 1998.
- [3] Matthew James Beal. *Variational algorithms for approximate Bayesian inference*. PhD thesis, University College London, 2003.
- [4] Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul Szerlip, Paul Horsfall, and Noah D. Goodman. Pyro: Deep universal probabilistic programming. *Journal of Machine Learning Research (JMLR)*, 2018.
- [5] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- [6] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. In *International Conference on Machine Learning (ICML)*, 2015.
- [7] Tianqi Chen, Emily Fox, and Carlos Guestrin. Stochastic gradient Hamiltonian Monte Carlo. In *International Conference on Machine Learning*, pages 1683–1691, 2014.
- [8] Andrea Coraddu, Luca Oneto, Alessandro Ghio, Stefano Savio, Davide Anguita, and Massimo Figari. Machine learning approaches for improving condition-based maintenance of naval propulsion plants. *Journal of Engineering for the Maritime Environment*, 2014.
- [9] John S Denker and Yann Lecun. Transforming neural-net output levels to probability distributions. In *Advances in Neural Information Processing Systems (NIPS)*, 1991.
- [10] Sebastian Farquhar, Lewis Smith, and Yarin Gal. Liberty or depth: Deep Bayesian neural nets do not need complex weight posterior approximations. *arXiv preprint arXiv:2002.03704*, 2020.
- [11] Angelos Filos, Sebastian Farquhar, Aidan N. Gomez, Tim G. J. Rudner, Zachary Kenton, Lewis Smith, Milad Alizadeh, Arnoud de Kroon, and Yarin Gal. Benchmarking Bayesian deep learning with diabetic retinopathy diagnosis. <https://github.com/OATML/bd1-benchmarks>, 2019.
- [12] Brendan J Frey and Geoffrey E Hinton. Variational learning in nonlinear Gaussian belief networks. *Neural Computation*, 11(1):193–213, 1999.

- [13] Yarin Gal. *Uncertainty in deep learning*. PhD thesis, University of Cambridge, 2016.
- [14] Yarin Gal and Zoubin Ghahramani. Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. In *International Conference on Machine Learning (ICML)*, 2016.
- [15] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep Bayesian active learning with image data. In *International Conference on Machine Learning (ICML)*, 2017.
- [16] José Miguel Hernández-Lobato and Ryan Adams. Probabilistic backpropagation for scalable learning of Bayesian neural networks. In *International Conference on Machine Learning (ICML)*, 2015.
- [17] José Miguel Hernández-Lobato, Yingzhen Li, Mark Rowland, Thang Bui, Daniel Hernández-Lobato, and Richard Turner. Black-box α -divergence minimization. In *International Conference on Machine Learning (ICML)*, 2016.
- [18] Geoffrey E Hinton and Drew Van Camp. Keeping the neural networks simple by minimizing the description length of the weights. In *Conference on Computational learning theory (COLT)*, 1993.
- [19] Matthew D Hoffman and Andrew Gelman. The No-U-Turn sampler: adaptively setting path lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research (JMLR)*, 15(1): 1593–1623, 2014.
- [20] Jiri Hron, Alex Matthews, and Zoubin Ghahramani. Variational Bayesian dropout: Pitfalls and fixes. In *International Conference on Machine Learning (ICML)*, 2018.
- [21] Jiri Hron, Yasaman Bahri, Roman Novak, Jeffrey Pennington, and Jascha Sohl-Dickstein. Exact posterior distributions of wide Bayesian neural networks. In *Uncertainty in deep learning Workshop, ICML.*, 2020.
- [22] Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An introduction to variational methods for graphical models. *Machine Learning*, 37(2):183–233, 1999.
- [23] Mohammad Emtiyaz Khan, Didrik Nielsen, Voot Tangkaratt, Wu Lin, Yarin Gal, and Akash Srivastava. Fast and scalable Bayesian deep learning by weight-perturbation in Adam. *International Conference on Machine Learning (ICML)*, 2018.
- [24] Durk P Kingma, Tim Salimans, and Max Welling. Variational dropout and the local reparameterization trick. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2575–2583, 2015.
- [25] Jaehoon Lee, Jascha Sohl-Dickstein, Jeffrey Pennington, Roman Novak, Sam Schoenholz, and Yasaman Bahri. Deep neural networks as Gaussian processes. In *International Conference on Learning Representations (ICLR)*, 2018.
- [26] Moshe Leshno, Vladimir Ya Lin, Allan Pinkus, and Shimon Schocken. Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. *Neural Networks*, 6(6):861–867, 1993.
- [27] Yingzhen Li and Richard E Turner. Rényi divergence variational inference. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1073–1081, 2016.
- [28] Yingzhen Li, José Miguel Hernández-Lobato, and Richard E Turner. Stochastic expectation propagation. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2323–2331, 2015.
- [29] Christos Louizos and Max Welling. Multiplicative normalizing flows for variational Bayesian neural networks. In *International Conference on Machine Learning (ICML)*, 2017.
- [30] David J. C. MacKay. A practical Bayesian framework for backpropagation networks. *Neural Computation*, 4(3):448–472, 1992.

- [31] Alexander G. de G. Matthews, Mark van der Wilk, Tom Nickson, Keisuke Fujii, Alexis Boukouvalas, Pablo León-Villagr a, Zoubin Ghahramani, and James Hensman. GPflow: A Gaussian process library using TensorFlow. *Journal of Machine Learning Research (JMLR)*, 18(40):1–6, 2017.
- [32] Alexander G. de G. Matthews, Jiri Hron, Mark Rowland, Richard E. Turner, and Zoubin Ghahramani. Gaussian process behaviour in wide deep neural networks. In *International Conference on Learning Representations (ICLR)*, 2018.
- [33] Jishnu Mukhoti, Pontus Stenertorp, and Yarin Gal. On the importance of strong baselines in Bayesian deep learning. *arXiv preprint arXiv:1811.09385*, 2018.
- [34] Radford M Neal. *Bayesian learning for neural networks*. PhD thesis, University of Toronto, 1995.
- [35] Radford M Neal et al. MCMC using Hamiltonian dynamics. *Handbook of Markov chain Monte Carlo*, 2(11):2, 2011.
- [36] Ian Osband, John Aslanides, and Albin Cassirer. Randomized prior functions for deep reinforcement learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 8617–8629, 2018.
- [37] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, D Sculley, Sebastian Nowozin, Joshua V Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty? Evaluating predictive uncertainty under dataset shift. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- [38] Hippolyt Ritter, Aleksandar Botev, and David Barber. A scalable Laplace approximation for neural networks. In *International Conference on Learning Representations (ICLR)*, 2018.
- [39] Samuel S Schoenholz, Justin Gilmer, Surya Ganguli, and Jascha Sohl-Dickstein. Deep information propagation. In *International Conference on Learning Representations (ICLR)*, 2017.
- [40] Burr Settles. Active learning literature survey. Technical report, University of Wisconsin-Madison Department of Computer Sciences, 2009.
- [41] Shengyang Sun, Guodong Zhang, Jiaxin Shi, and Roger Grosse. Functional variational Bayesian neural networks. In *International Conference on Learning Representations (ICLR)*, 2019.
- [42] Siddharth Swaroop, Cuong V Nguyen, Thang D Bui, and Richard E Turner. Improving and understanding variational continual learning. *arXiv preprint arXiv:1905.02099*, 2019.
- [43] Marcin B Tomczak, Siddharth Swaroop, and Richard E Turner. Neural network ensembles and variational inference revisited. In *1st Symposium on Advances in Approximate Bayesian Inference (AABI)*, 2018.
- [44] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research (JMLR)*, 9:2579–2605, 2008.
- [45] 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)*, pages 681–688, 2011.
- [46] Florian Wenzel, Kevin Roth, Bastiaan S Veeling, Jakub Świ atkowski, Linh Tran, Stephan Mandt, Jasper Snoek, Tim Salimans, Rodolphe Jenatton, and Sebastian Nowozin. How good is the Bayes posterior in deep neural networks really? In *International Conference on Machine Learning (ICML)*, 2020.
- [47] Ruqi Zhang, Chunyuan Li, Jianyi Zhang, Changyou Chen, and Andrew Gordon Wilson. Cyclical stochastic gradient MCMC for Bayesian deep learning. In *International Conference on Learning Representations (ICLR)*, 2020.

A In-between Uncertainty in Other Regions of Input Space

In this appendix, we show plots generated by placing two Gaussian clusters of data with centers randomly chosen on the sphere of radius \sqrt{D} , where $D = 5$ denotes input dimension. We generate synthetic data by sampling from a wide-limit Gaussian process. For each plot, we show the predictive mean and 2 standard deviations along the line segment in input space joining the centres of these two clusters. For all plots, we choose $\sigma_w = \sqrt{2}$, $\sigma_b = 1$, networks of width 50 and dropout probability of $p = 0.05$ for MCDO. We set the observation noise standard deviation to 0.01, which is the ground truth value used to generate the synthetic data. The initialisation of MFVI and MCDO is the same as discussed in appendix F.

In the 1HL case, Theorem 5 implies that MCDO’s predictive variance will be convex along any line, including the line plotted. In contrast, theorem 4 only applies to certain lines in input space, and does not bound the variance in these cases. However, we suspect that theorem 4 is indicative of a lack of in-between uncertainty in more general cases. Figure 7 shows that although that exact inference with the GP with the limiting BNN prior exhibits in-between uncertainty, this is lost by both MFVI and MCDO, even on general lines in input space. MFVI and MCDO are often *more* confident in between the data clusters than at the data clusters.

We next run the same experiment but with 3HL BNNs and their limiting GP. In this case theorem 6 implies that for sufficiently wide BNNs, there exist variational parameters that can approximate *any* predictive mean and standard deviation. However, in figure 8 we see that compared to exact inference in the limiting GP, MFVI and MCDO both underestimate in-between uncertainty — or sometimes show as large uncertainty *at* the data as *in between* the data.

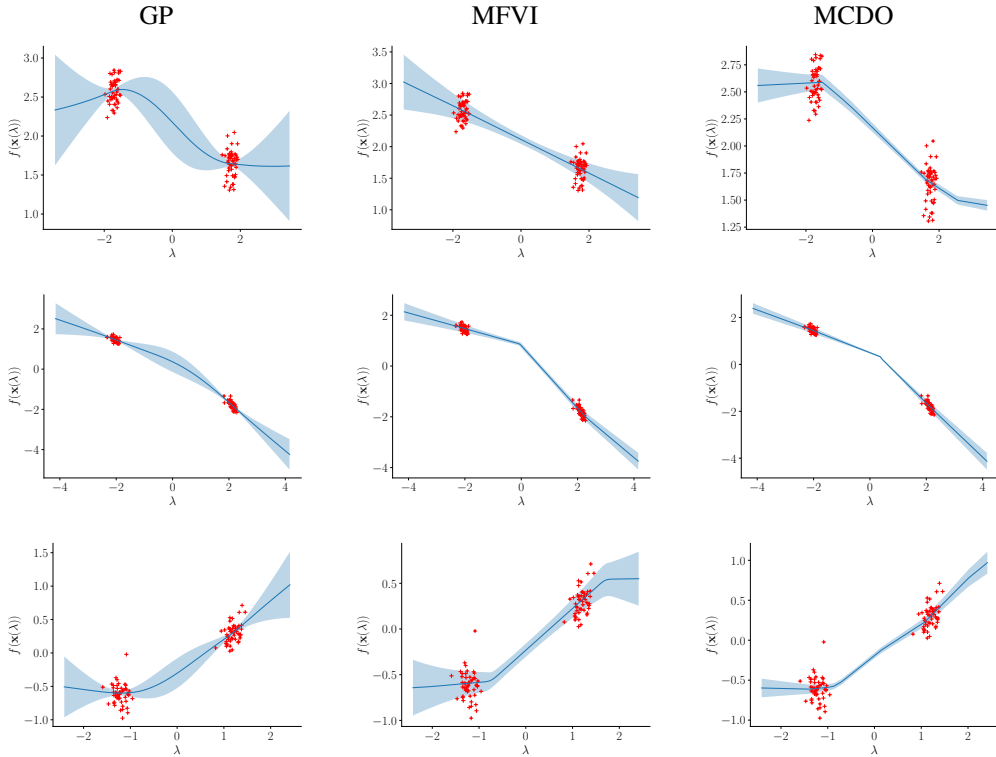


Figure 7: Mean and 2 standard deviation bars of the predictive distribution on lines joining random clusters of data, for 1HL BNNs. Each row represents the same random dataset. We also plot the projection of the 5-dimensional data onto this line segment as its λ -value, along with the output value of the data. Note that the data appears noisy, but this is due to the projection onto a lower-dimensional space.

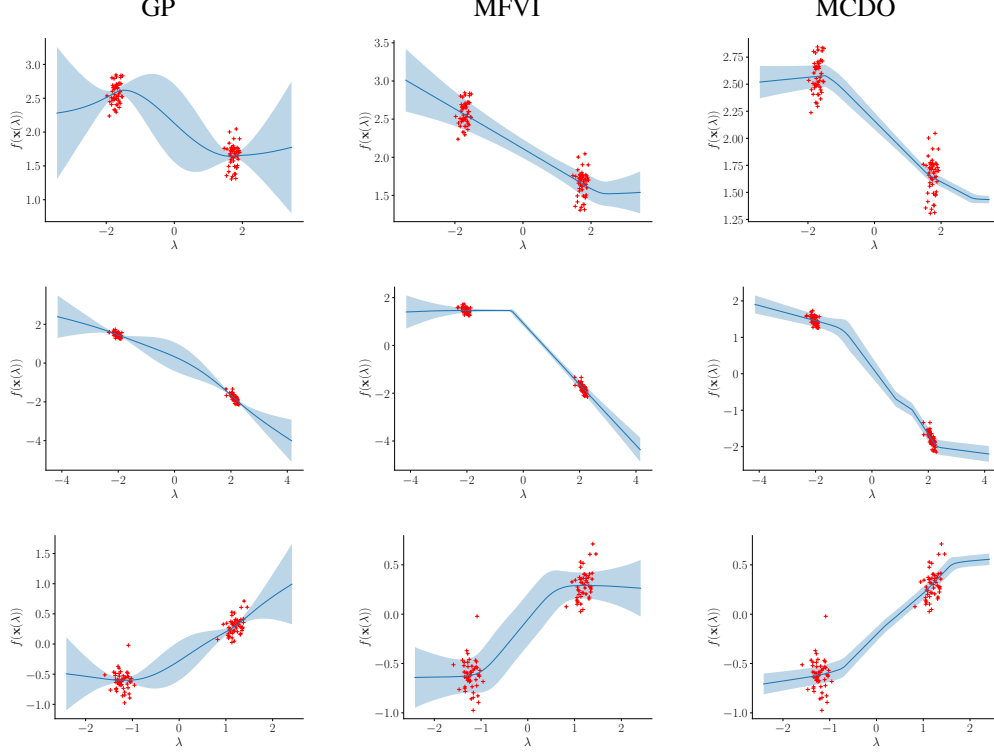


Figure 8: Same experimental set-up as in figure 7 for the 3HL case.

B General Statements and Proofs of Theorems 1 and 2

In section 3 we stated simplified versions of bounds concerning the variance of single-hidden layer networks with certain approximating families. The two main results we prove in this section are the following generalisations of theorems 1 and 2 respectively:

Theorem 4. Consider a single-hidden layer ReLU neural network mapping from $\mathbb{R}^D \rightarrow \mathbb{R}^K$ with $I \in \mathbb{N}$ hidden units. The corresponding mapping is given by $f^{(k)}(\mathbf{x}) = \sum_{i=1}^I w_{k,i} \psi\left(\sum_{d=1}^D u_{i,d} x_d + v_i\right) + b_k$ for $1 \leq k \leq K$, where $\psi(a) = \max(0, a)$. Suppose we have a distribution over network parameters with density of the form:

$$q(\mathbf{W}, \mathbf{b}, \mathbf{U}, \mathbf{v}) = \prod_{i=1}^I q_i(\mathbf{w}_i | \mathbf{U}, \mathbf{v}) q(\mathbf{b} | \mathbf{U}, \mathbf{v}) \prod_{i=1}^I \prod_{d=1}^D \mathcal{N}(u_{i,d}; \mu_{u_{i,d}}, \sigma_{u_{i,d}}^2) \prod_{i=1}^I \mathcal{N}(v_i; \mu_{v_i}, \sigma_{v_i}^2), \quad (3)$$

where $\mathbf{w}_i = \{w_{k,i}\}_{k=1}^K$ are the weights out of neuron i and $\mathbf{b} = \{b_k\}_{k=1}^K$ are the output biases, and $q_i(\mathbf{w}_i | \mathbf{U}, \mathbf{v})$ and $q(\mathbf{b} | \mathbf{U}, \mathbf{v})$ are arbitrary probability densities with finite first two moments. Consider a line in \mathbb{R}^D parameterised by $\mathbf{x}(\lambda)_d = \gamma_d \lambda + c_d$ for $\lambda \in \mathbb{R}$ such that $\gamma_d c_d = 0$ for $1 \leq d \leq D$. Then for any $\lambda_1 \leq 0 \leq \lambda_2$, and any λ_* such that $|\lambda_*| \leq \min(|\lambda_1|, |\lambda_2|)$,

$$\mathbb{V}[f^{(k)}(\mathbf{x}(\lambda_*))] \leq \mathbb{V}[f^{(k)}(\mathbf{x}(\lambda_1))] + \mathbb{V}[f^{(k)}(\mathbf{x}(\lambda_2))] \quad \text{for } 1 \leq k \leq K. \quad (4)$$

We now briefly show how the statement of theorem 4 in the main text can be deduced from this more general version. The fully factorised Gaussian family \mathcal{Q}_{FFG} is of the form in equation (3). It remains to show that both conditions *i.* and *ii.* imply that $\gamma_d c_d = 0$. Consider any line intersecting the origin (i.e. satisfying condition *ii.*). Such a line can be written in the form $\mathbf{x}(\lambda)_d = \gamma_d \lambda$ by choosing the origin to correspond to $\lambda = 0$. As $c_d = 0$ for all d , $\gamma_d c_d = 0$ for all d . In theorem 1 $\mathbf{p} = \mathbf{x}(\lambda_1)$ and $\mathbf{q} = \mathbf{x}(\lambda_2)$ are on opposite sides of the origin, hence the signs of λ_1 and λ_2 are opposite. Finally, the condition that $\mathbf{r} = \mathbf{x}(\lambda_*)$ is closer to the origin than both \mathbf{p} and \mathbf{q} is exactly that $|\lambda_*| \leq \min(|\lambda_1|, |\lambda_2|)$.

In order to verify condition ii), note that any line orthogonal to a hyperplane $x_{d'} = 0$ can be parameterised as $\mathbf{x}(\lambda)_d = \gamma_d \lambda + c_d$, where $\gamma_d = 0$ for $d \neq d'$ and $c_{d'} = 0$. Hence $\gamma_d c_d = 0$ for all d . The condition that the line segment $\overline{\mathbf{p}\mathbf{q}}$ intersects the plane, with $\mathbf{p} = \mathbf{x}(\lambda_1)$ and $\mathbf{q} = \mathbf{x}(\lambda_2)$ is exactly that the signs of λ_1 and λ_2 are opposite, and that $|\lambda_*| \leq \min(|\lambda_1|, |\lambda_2|)$.

As a corollary of theorem 4, we can obtain bounds on higher-dimensional objects than lines, such as on hypercubes. For instance, consider the case where $\mathbf{x} \in \mathbb{R}^2$. Let $\mathbf{p}, \mathbf{q}, \mathbf{r}, \mathbf{s}$ be the four corners of a rectangle centered the origin. For any point \mathbf{a} in the rectangle, we can upper bound $\mathbb{V}[f(\mathbf{a})]$ by the sum of the variances at the top and bottom edges of the rectangle. These in turn can be upper bounded by the variances at the corners of the rectangle. Hence we have that for any point \mathbf{a} in the rectangle, $\mathbb{V}[f(\mathbf{a})] \leq \mathbb{V}[f(\mathbf{p})] + \mathbb{V}[f(\mathbf{q})] + \mathbb{V}[f(\mathbf{r})] + \mathbb{V}[f(\mathbf{s})]$. Similarly the variance at any point in a hypercube centered at the origin can be bounded by the sum of the variances on its vertices, and we can obtain tighter bounds on diagonals and faces of the hypercube, by repeatedly applying theorem 4.

Theorem 5 (MC dropout). *Consider a single-hidden layer ReLU neural network mapping from $\mathbb{R}^D \rightarrow \mathbb{R}^K$ with $I \in \mathbb{N}$ hidden units. The corresponding mapping is given by $f^{(k)}(\mathbf{x}) = \sum_{i=1}^I w_{k,i} \psi \left(\sum_{d=1}^D u_{i,d} x_d + v_i \right) + b_k$ for $1 \leq k \leq K$, where $\psi(a) = \max(0, a)$. Assume \mathbf{U}, \mathbf{v} are set deterministically and*

$$q(\mathbf{W}, \mathbf{b}) = q(\mathbf{b}) \prod_{i=1}^I q_i(\mathbf{w}_i),$$

where $\mathbf{w}_i = \{w_{k,i}\}_{k=1}^K$ are the weights out of neuron i , $\mathbf{b} = \{b_k\}_{k=1}^K$ are the output biases and $q(\mathbf{b})$ and $q_i(\mathbf{w}_i)$ are arbitrary probability densities with finite first two moments. Then, $\mathbb{V}[f^{(k)}(\mathbf{x})]$ is convex in \mathbf{x} for $1 \leq k \leq K$.

Proof. The theorem follows immediately from lemma 1 since \mathbf{U} and \mathbf{v} are deterministic. \square

Remark 3. *Theorem 5 applies for any activation function ψ such that ψ^2 is convex. This is the only property of ψ used in lemma 1.*

B.1 Preliminary Lemmas

In order to prove theorems 1 and 2 we first collect a series of preliminary lemmas.

Lemma 1. *Assume a distribution for $\mathbf{W}, \mathbf{b}|\mathbf{U}, \mathbf{v}$ with density of the form*

$$q(\mathbf{W}, \mathbf{b}|\mathbf{U}, \mathbf{v}) = q(\mathbf{b}|\mathbf{U}, \mathbf{v}) \prod_i q_i(\mathbf{w}_i|\mathbf{U}, \mathbf{v}).$$

Then, $\mathbb{V}[f^{(k)}(\mathbf{x})|\mathbf{U}, \mathbf{v}]$ is a convex function of \mathbf{x} .

The proof of lemma 1 is in appendix C.1.

Lemma 2. *Consider the variance of a single neuron in the one dimensional case, with activation $a(x) \sim \mathcal{N}(\mu(x), \sigma^2(x))$, $\mu(x) = \mu_u x + \mu_v$ and $\sigma^2(x) = \sigma_u^2 x^2 + \sigma_v^2$. Let*

$$\mathcal{T}_1 = \{f \geq 0 : \forall 0 \leq b < a, f(a) \geq f(-a) \text{ and } f(b) \leq f(a)\}$$

and

$$\mathcal{T}_2 = \{f \geq 0 : \forall a < b \leq 0, f(a) \geq f(-a) \text{ and } f(b) \leq f(a)\}.$$

If $\mu_u \geq 0$, then $\mathbb{V}[\psi(a(x))] \in \mathcal{T}_1$. If $\mu_u \leq 0$, then $\mathbb{V}[\psi(a(x))] \in \mathcal{T}_2$.

The proof of lemma 2 is in appendix C.2.

Corollary 1 (Corollary of lemma 2). *Consider a line in \mathbb{R}^D parameterized by $[\mathbf{x}(\lambda)]_d = \gamma_d \lambda + c_d$ for $\lambda \in \mathbb{R}$ such that $\gamma_d c_d = 0$ for $1 \leq d \leq D$. Let $a(\mathbf{x}) := \sum_{d=1}^D u_d x_d + v$ with $\{u_d\}_{d=1}^D$ and v independent and Gaussian distributed. Then, $\mathbb{V}[\psi(a(\mathbf{x}(\lambda)))] \in \mathcal{T}_1 \cup \mathcal{T}_2$ (as a function of λ).*

Proof. The activation $a(\mathbf{x}(\lambda))$ is a linear combination of Gaussian random variables, and is therefore Gaussian distributed. Moreover the mean is linear in λ . The variance of $a(\mathbf{x}(\lambda))$ is given by:

$$\begin{aligned}\mathbb{V}[a(\mathbf{x}(\lambda))] &= \sum_{d=1}^D \mathbb{V}[u_d](\gamma_d \lambda + c_d)^2 + \mathbb{V}[v] \\ &= \sum_{d=1}^D \sigma_{u_d}^2 (\gamma_d \lambda + c_d)^2 + \sigma_v^2 \\ &= \lambda^2 \left(\sum_{d=1}^D \sigma_{u_d}^2 \gamma_d^2 \right) + 2\lambda \left(\sum_{d=1}^D \sigma_{u_d}^2 \gamma_d c_d \right) + \left(\sum_{d=1}^D \sigma_{u_d}^2 c_d^2 + \sigma_v^2 \right) \\ &= \lambda^2 \left(\sum_{d=1}^D \sigma_{u_d}^2 \gamma_d^2 \right) + \left(\sum_{d=1}^D \sigma_{u_d}^2 c_d^2 + \sigma_v^2 \right).\end{aligned}$$

Defining $\sigma_u^2 = \sum_{d=1}^D \sigma_{u_d}^2 \gamma_d^2$ and $\sigma_v^2 = \sum_{d=1}^D \sigma_{u_d}^2 c_d^2 + \sigma_v^2$, the corollary follows from lemma 2. \square

Lemma 3. *Let \mathcal{C} be the set of convex functions from $\mathbb{R} \rightarrow [0, \infty)$. Fix any $a < 0 < b$ and c such that $|c| \leq \min(|a|, |b|)$. Then any function f that can be written as a linear combination of functions in $\mathcal{T}_1 \cup \mathcal{T}_2 \cup \mathcal{C}$ with non-negative weights satisfies, $f(c) \leq f(a) + f(b)$.*

The proof of lemma 3 can be found in appendix C.3.

B.2 Proof of Theorem 1

Having collected the necessary preliminary lemmas we now prove theorem 1.

Proof of Theorem 1. By the law of total variance,

$$\mathbb{V}[f^{(k)}(\mathbf{x})] = \mathbb{E}[\mathbb{V}[f^{(k)}(\mathbf{x})|\mathbf{U}, \mathbf{v}]] + \mathbb{V}[\mathbb{E}[f^{(k)}(\mathbf{x})|\mathbf{U}, \mathbf{v}]].$$

Using lemma 1, $\mathbb{V}[f^{(k)}(\mathbf{x})|\mathbf{U}, \mathbf{v}]$ is convex as a function of \mathbf{x} . As the expectation of a convex function is convex, the first term is a convex function of \mathbf{x} . For the second term we have

$$\mathbb{E}[f^{(k)}(\mathbf{x})|\mathbf{U}, \mathbf{v}] = \mathbb{E} \left[\sum_{i=1}^I w_{k,i} \psi(a_i) + b_k \middle| \mathbf{U}, \mathbf{v} \right] = \sum_{i=1}^I \mu_{w_{k,i}} \psi(a_i) + \mu_{b_k},$$

where $\mu_{w_{k,i}} := \mathbb{E}[w_{k,i}]$, $\mu_{b_k} := \mathbb{E}[b_k]$. In the second line we used linearity of expectation and that conditioned on (\mathbf{U}, \mathbf{v}) , the a_i are deterministic. Next,

$$\mathbb{V}[\mathbb{E}[f^{(k)}(\mathbf{x})|\mathbf{U}, \mathbf{v}]] = \mathbb{V} \left[\sum_{i=1}^I \mu_{w_{k,i}} \psi(a_i) + \mu_{b_k} \right] = \sum_{i=1}^I \mu_{w_{k,i}}^2 \mathbb{V}[\psi(a_i)], \quad (5)$$

since the a_i are independent of each other.

Consider a line in \mathbb{R}^D parameterised by $[\mathbf{x}(\lambda)]_d = \gamma_d \lambda + c_d$ for $\lambda \in \mathbb{R}$ such that $\gamma_d c_d = 0$ for $1 \leq d \leq D$.

By corollary 1, $\mathbb{V}[\psi(a_i(\mathbf{x}(\lambda)))] \in \mathcal{T}_1 \cup \mathcal{T}_2$ (as a function of λ). Since $\mathbb{V}[f^{(k)}(\mathbf{x})|\mathbf{U}, \mathbf{v}]$ is convex as a function of \mathbf{x} , it is also convex as a function of λ . We have written $\mathbb{V}[f^{(k)}(\mathbf{x}(\lambda))]$ in the form assumed in lemma 3, completing the proof. \square

C Proof of Lemmas

In this section we prove the preliminary lemmas stated in appendix B.1.

C.1 Proof of Lemma 1

Proof. We assume a distribution for the network weights such that:

$$q(\mathbf{W}, \mathbf{b}|\mathbf{U}, \mathbf{v}) = q(\mathbf{b}|\mathbf{U}, \mathbf{v}) \prod_{i=1}^I q_i(\mathbf{w}_i|\mathbf{U}, \mathbf{v}).$$

By this factorisation assumption, the outgoing weights from each neuron are conditionally independent. This means the conditional variance of the output under this distribution can be written

$$\mathbb{V}[f^{(k)}(\mathbf{x})|\mathbf{U}, \mathbf{v}] = \sum_i \mathbb{V}[w_{k,i}|\mathbf{U}, \mathbf{v}]\psi(a_i)^2 + \mathbb{V}[b_k|\mathbf{U}, \mathbf{v}]. \quad (6)$$

with $a_i := a_i(\mathbf{x}) = \sum_{d=1}^D u_{i,d}x_d + v_i$.

Since $\mathbb{V}[f^{(k)}(\mathbf{x})|\mathbf{U}, \mathbf{v}]$ is a linear combination of the $\psi(a_i)^2$ with non-negative weights (plus a constant), to prove convexity it suffices to show that each $\psi(a_i)^2$ is convex as a function of \mathbf{x} . $\psi(a_i)^2$ is convex as a function of a_i , since it is 0 for $a_i \leq 0$ and a_i^2 for $a_i > 0$. To show that it is convex as a function of \mathbf{x} , we write

$$\begin{aligned} \psi(a_i(t\mathbf{x}_1 + (1-t)\mathbf{x}_2))^2 &= \psi\left(\sum_d u_{i,d}(t[\mathbf{x}_1]_d + (1-t)[\mathbf{x}_2]_d) + v_i\right)^2 \\ &= \psi\left(t\left(\sum_d u_{i,d}[\mathbf{x}_1]_d + v_i\right) + (1-t)\left(\sum_d u_{i,d}[\mathbf{x}_2]_d + v_i\right)\right)^2 \\ &\leq t\psi\left(\sum_d u_{i,d}[\mathbf{x}_1]_d + v_i\right)^2 + (1-t)\psi\left(\sum_d u_{i,d}[\mathbf{x}_2]_d + v_i\right)^2 \\ &= t\psi(a_i(\mathbf{x}_1))^2 + (1-t)\psi(a_i(\mathbf{x}_2))^2. \end{aligned}$$

The inequality uses convexity of $\psi(a)$ as a function of a . □

C.2 Proof of Lemma 2

Throughout, we assume σ_u, σ_v and μ_v are fixed and suppress dependence on these parameters. Let $v_{\mu_u}(x) := \mathbb{V}[\psi(a(x))]$ where the variance is taken with respect to a distribution with parameter μ_u . Then, $v_{\mu_u}(x) = v_{-\mu_u}(-x)$ since $\mu(x)$ and $\sigma^2(x)$ are unchanged by the transformation $\mu_u, x \rightarrow -\mu_u, -x$.

Suppose $v_{\mu_u} \in \mathcal{T}_1$ for $\mu_u > 0$, then for $x \leq 0$,

$$v_{-\mu_u}(x) = v_{\mu_u}(-x) \geq v_{\mu_u}(x) = v_{-\mu_u}(-x),$$

and for $x < y \leq 0$,

$$v_{-\mu_u}(y) = v_{\mu_u}(-y) \leq v_{\mu_u}(-x) = v_{-\mu_u}(x).$$

In words, if $v_{\mu_u} \in \mathcal{T}_1$ then $v_{-\mu_u} \in \mathcal{T}_2$. It therefore suffices to consider the case when $\mu_u \geq 0$.

We first show that if $x \geq 0$, $v_{\mu_u}(x) \geq v_{\mu_u}(-x)$. Henceforth, we assume $\mu_u \geq 0$ is fixed and suppress it notationally. From Frey and Hinton [12],

$$v(x) = \sigma(x)^2 \alpha(r(x)), \quad (7)$$

Here $r(x) = \mu(x)/\sigma(x)$. We define $h(r) = N(r) + r\Phi(r)$, where N is the standard Gaussian pdf, Φ is the standard Gaussian cdf. We define $\alpha(r) = \Phi(r) + rh(r) - h(r)^2$.

As $\sigma(x)^2 = \sigma(-x)^2$, it suffices to show $\alpha(r(x)) \geq \alpha(r(-x))$ for $x > 0$. To show this, we first show that $r(x) \geq r(-x)$ for $x > 0$, then show that $\alpha(r)$ is monotonically increasing.

$$r(x) = \mu(x)/\sigma(x) = \mu(-x)/\sigma(-x) + 2\mu_u x/\sigma(-x) \geq \mu(-x)/\sigma(-x) = r(-x).$$

The inequality uses that both μ_u and x are non-negative. It remains to show that $\alpha(r)$ is monotonically increasing. A straightforward calculation shows that,

$$\alpha'(r) = 2h(r)(1 - \Phi(r)).$$

As $1 - \Phi(r) > 0$, we must show $h(r) \geq 0$. We have $\lim_{r \rightarrow -\infty} h(r) = 0$ and $h'(r) = \Phi(r) > 0$, implying $h(r) > 0$. We conclude $\alpha'(r) > 0$ for all r , showing that $v_{\mu_u}(x) \geq v_{\mu_u}(-x)$ for $x \geq 0$.

To complete the proof, we must show that $v(x)$ is monotonically increasing for $x \geq 0$. As $\sigma(x)^2$ is increasing as a function of x and $\alpha(r)$ is increasing as a function of r , $v(x)$ is increasing as a function of x whenever $r(x)$ is increasing as a function of x . As $r'(x) = \frac{\sigma_v^2 \mu_u - \sigma_u^2 \mu_v x}{\sigma(x)^3}$, this completes the proof if $\sigma_v^2 \mu_u - \sigma_u^2 \mu_v x \geq 0$. In particular, we need only consider cases when $\mu_v > 0$. In this case, we write,

$$v(x) = \mu(x)^2 \beta(r(x)) \quad (8)$$

where $\beta(r) = \alpha(r)/r^2$. Also in this region, we have the inequality,

$$r'(x)\sigma(x) = \frac{\sigma_v^2 \mu_u - \sigma_u^2 \mu_v x}{\sigma_u^2 x^2 + \sigma_v^2} \leq \frac{\sigma_v^2 \mu_u}{\sigma_u^2 x^2 + \sigma_v^2} \leq \frac{\sigma_v^2 \mu_u}{\sigma_v^2} = \mu_u,$$

which leads to $r'(x) \leq \mu_u/\sigma(x)$.

Differentiating equation (8),

$$\begin{aligned} v'(x) &= 2\mu_u \mu(x) \beta(r(x)) + \mu(x)^2 \left(\frac{\sigma_v^2 \mu_u - \sigma_u^2 \mu_v x}{\sigma(x)^3} \right) \beta'(r(x)) \\ &\geq 2\mu_u \mu(x) \left(\beta(r(x)) + \frac{1}{2} r(x) \beta'(r(x)) \right). \end{aligned}$$

The inequality uses that $r(x) > 0$, so that by lemma 4, $\beta'(r(x)) < 0$. It suffices to show that $\beta(r) + \frac{1}{2} r \beta'(r) > 0$ for $r > 0$.

$$\beta(r) + \frac{1}{2} r \beta'(r) = \beta(r) + \frac{1}{2} r \frac{d}{dr} \left(\frac{\alpha(r)}{r^2} \right) = \frac{\alpha(r)}{r^2} + \frac{1}{2} r \frac{\alpha'(r)r^2 - 2r\alpha(r)}{r^4} = \frac{\alpha'(r)}{2r} \geq 0.$$

We conclude that $v'(x) \geq 0$ for $x \geq 0$, implying that $v(x)$ is monotonically increasing in this region. This completes the proof that $v_{\mu_u}(x) \in \mathcal{T}_1$ for $\mu_u > 0$.

Lemma 4. For β defined as in the proof of lemma 2 and for $r > 0$, $\beta'(r) < 0$

Proof. For $r \neq 0$, $\beta'(r) = (-2\Phi(r) + 2N(r)^2 + 2N(r)\Phi(r))/r^3$. As $r > 0$,

$$\beta'(r) \leq 0 \Leftrightarrow I(r) := -\Phi(r) + N(r)^2 + N(r)r\Phi(r) \leq 0.$$

Rearranging [1, 7.1.13] yields:

$$1 - \frac{2}{r + \sqrt{r^2 + 8/\pi}} N(r) \leq \Phi(r) < 1 - \frac{2}{r + \sqrt{r^2 + 4}} N(r). \quad (9)$$

for $r \geq 0$.

$$\begin{aligned} I(r) &= -\Phi(r) + N(r)^2 + rN(r)\Phi(r) \\ &\leq -\Phi(r) + N(r)^2 + rN(r) \left(1 - \frac{2}{r + \sqrt{r^2 + 4}} N(r) \right) \\ &\leq -1 + \frac{2}{r + \sqrt{r^2 + 8/\pi}} N(r) + N(r)^2 + rN(r) \left(1 - \frac{2}{r + \sqrt{r^2 + 4}} N(r) \right) \\ &= -1 + \frac{2}{r + \sqrt{r^2 + 8/\pi}} N(r) + rN(r) + N(r)^2 \left(1 - \frac{2r}{r + \sqrt{r^2 + 4}} \right) \end{aligned} \quad (10)$$

We now make use of numerous crude bounds which hold for $r > 0$:

1. $N(r) \leq 1/\sqrt{2\pi}$,
2. $\frac{2}{r + \sqrt{r^2 + 8/\pi}} \leq \sqrt{\pi/2}$,
3. $rN(r) \leq 1/\sqrt{2\pi e}$

$$4. \frac{2r}{r+\sqrt{r^2+4}} \geq 0.$$

Plugging these into equation (10),

$$I(r) \leq -1 + \frac{\sqrt{\pi/2}}{\sqrt{2\pi}} + \frac{1}{\sqrt{2\pi e}} + \frac{1}{2\pi} = -\frac{1}{2} + \frac{1}{\sqrt{2\pi e}} + \frac{1}{2\pi} \approx -0.098 < 0. \quad \square$$

C.3 Proof of Lemma 3

Proof. Recall that

$$\mathcal{T}_1 = \{f \geq 0 : \forall 0 \leq b < a, f(a) \geq f(-a) \text{ and } f(b) \leq f(a)\}$$

and

$$\mathcal{T}_2 = \{f \geq 0 : \forall a < b \leq 0, f(a) \geq f(-a) \text{ and } f(b) \leq f(a)\}.$$

First, note that $\mathcal{T}_1, \mathcal{T}_2$ and the set of non-negative convex functions, \mathcal{C} are all closed under addition and positive scalar multiplication. We can therefore write f as a sum of three functions, $f(x) = t_1(x) + t_2(x) + s(x)$ with $t_1 \in \mathcal{T}_1, t_2 \in \mathcal{T}_2$ and $s \in \mathcal{C}$. We prove the case when $a \leq c \leq 0 \leq -c \leq b$. The case $a \leq -c \leq 0 \leq c \leq b$ follows a symmetric argument.

$$\begin{aligned} f(c) &= t_1(c) + t_2(c) + s(c) \quad (\text{def.}) \\ &\leq t_1(c) + t_2(a) + s(c) \quad (\text{second condition for } \mathcal{T}_2) \\ &\leq t_1(-c) + t_2(a) + s(c) \quad (\text{first condition for } \mathcal{T}_1) \\ &\leq t_1(b) + t_2(a) + s(c) \quad (\text{second condition for } \mathcal{T}_1) \\ &\leq t_1(b) + t_2(a) + \max(s(a), s(b)) \quad (s \text{ convex}) \\ &\leq t_1(b) + t_2(a) + s(a) + s(b) \\ &\leq t_1(a) + t_1(b) + t_2(a) + t_2(b) + s(a) + s(b) \quad (\text{non-negativity}) \\ &= f(a) + f(b). \end{aligned} \quad \square$$

D Proof of Theorem 3

We now restate and prove Theorem 3 from the main body:

Theorem 6. *Let $A \subset \mathbb{R}^D$ be compact, and let $C(A)$ be the space of continuous functions on A to \mathbb{R} . Similarly, let $C^+(A)$ be the space of continuous functions on A to $\mathbb{R}_{\geq 0}$. Then for any $g \in C(A)$ and $h \in C^+(A)$, and any $\epsilon > 0$, for both the mean-field Gaussian and MC dropout families, there exists a 2-hidden layer ReLU NN such that*

$$\sup_{\mathbf{x} \in A} |\mathbb{E}[f(\mathbf{x})] - g(\mathbf{x})| < \epsilon \quad \text{and} \quad \sup_{\mathbf{x} \in A} |\mathbb{V}[f(\mathbf{x})] - h(\mathbf{x})| < \epsilon,$$

where $f(\mathbf{x})$ is the (stochastic) output of the network.

Our proof will make use of the standard universal approximation theorem for deterministic NNs as given in Leshno et al. [26]:

Theorem 7 (Universal approximation for deterministic NNs). *Let $\psi(a) = \max(0, a)$. Then for every $g \in C(\mathbb{R}^D)$ and every compact set $A \subset \mathbb{R}^D$, for any $\epsilon > 0$ there exists a function $f \in S$ such that $\|g - f\|_\infty < \epsilon$. Here*

$$S = \left\{ \sum_{i=1}^I w_i \psi \left(\sum_{d=1}^D u_{i,d} x_d + v_i \right) : I \in \mathbb{N}, w_i, u_{i,d}, v_i \in \mathbb{R} \right\}.$$

We first prove a useful lemma.

Lemma 5. *Let $\psi(a) = \max(0, a)$. Let a be a random variable with finite first two moments. Then $\mathbb{V}[\psi(a)] \leq \mathbb{V}[a]$.*

Proof. For all $x, y \in \mathbb{R}$, we have $|x - y|^2 \geq |\psi(x) - \psi(y)|^2$. Consider two i.i.d. copies of any random variable with finite first two moments, denoted a_1 and a_2 . Then

$$\begin{aligned} \mathbb{V}[a_1] &= \mathbb{E}[a_1^2] - \mathbb{E}[a_1]^2 = \frac{1}{2} \mathbb{E}[a_1^2 + a_2^2 - 2a_1 a_2] = \frac{1}{2} \mathbb{E}[|a_1 - a_2|^2] \geq \frac{1}{2} \mathbb{E}[|\psi(a_1) - \psi(a_2)|^2] \\ &= \mathbb{V}[\psi(a_1)]. \end{aligned} \quad \square$$

D.1 Proof of Theorem 3 for \mathcal{Q}_{FFG}

We prove theorem 6 for the fully-factorised Gaussian approximating family. We begin by proving results about 1HL networks within this family. The overall goal of these results is lemma 8, which informally says that for any set of mean parameters for the weights, we can find a setting of the standard deviations of the weights, such that the mean output of the network is close to the output of the deterministic network, with weights equal to the mean parameters. Our proof of this proceeds in 3 parts: First, in lemma 9, we show that by making the standard deviation parameters sufficiently small, we can ensure that the variance of the output of the network is uniformly small on some compact set A . Next, in lemma 7, we show that again by choosing the standard deviation sufficiently small, we can show that most of the sample functions of the 1HL network are close to the function that would be obtained by using the mean parameters. Finally, in the proof of lemma 8, we use Chebyshev's inequality and the triangle inequality to conclude that the mean of the network must also be close to the function defined by the mean parameters.

These networks will be used to construct the desired 2HL network.

Notation Consider a 1HL ReLU NN with input $\mathbf{x} \in \mathbb{R}^D$ and output $\mathbf{f} \in \mathbb{R}^K$. Let the network have I hidden units and be parameterised by input weights $U \in \mathbb{R}^{I \times D}$, input biases $v \in \mathbb{R}^I$, output weights $W \in \mathbb{R}^{K \times I}$ and output biases $b \in \mathbb{R}^K$. Let $\theta = (U, v, W, b)$. Denote the k^{th} output of the network by $f_{\theta}^{(k)}(\mathbf{x})$. Consider a factorised Gaussian distribution over the parameters θ in the network. Let the means of the Gaussians be denoted $\boldsymbol{\mu} = (\mu_U, \mu_v, \mu_W, \mu_b)$, where e.g. μ_U is a matrix whose elements are the means of U . Each mean is always taken to be $\in \mathbb{R}$. Let the standard deviations be denoted $\boldsymbol{\sigma} = (\sigma_U, \sigma_v, \sigma_W, \sigma_b)$. Each standard deviation is always taken to be $\in \mathbb{R}_{>0}$.

The following lemma states that we can make the output of a 1HL BNN have low variance by setting the standard deviation of the weights to be small.

Lemma 6. *Let $A \subset \mathbb{R}^D$ be a compact set and $f_{\theta}^{(k)}(\mathbf{x})$ be the k^{th} output of a 1HL ReLU NN with a mean-field Gaussian distribution mapping from $A \rightarrow \mathbb{R}$. Fix any $\boldsymbol{\mu}$ and any $\epsilon > 0$. Let all the standard deviations in $\boldsymbol{\sigma}$ be equal to a shared constant $\sigma > 0$. Then there exists $\sigma' > 0$ such that for all $\sigma < \sigma'$ and for all $\mathbf{x} \in A$, $\mathbb{V}[\psi(f_{\theta}^{(k)}(\mathbf{x}))] < \epsilon$ for all $1 \leq k \leq K$.*

Proof. Define $a_i = \sum_{d=1}^D u_{i,d}x_d + v_i$, so that $f_{\theta}^{(k)}(\mathbf{x}) = \sum_{i=1}^I w_{k,i}\psi(a_i) + b_k$. Then

$$\begin{aligned} \mathbb{V}[f_{\theta}^{(k)}(\mathbf{x})] &= \mathbb{V}\left[\sum_{i=1}^I w_{k,i}\psi(a_i)\right] + \sigma^2 \\ &= \sum_{i=1}^I \sum_{j=1}^I \text{Cov}(w_{k,i}\psi(a_i), w_{k,j}\psi(a_j)) + \sigma^2 \\ &\leq \sum_{i=1}^I \sum_{j=1}^I |\text{Cov}(w_{k,i}\psi(a_i), w_{k,j}\psi(a_j))| + \sigma^2 \\ &\leq \sum_{i=1}^I \sum_{j=1}^I \sqrt{\mathbb{V}[w_{k,i}\psi(a_i)]\mathbb{V}[w_{k,j}\psi(a_j)]} + \sigma^2, \end{aligned}$$

where the final line follows from the Cauchy–Schwarz inequality. We now analyse each of the constituent terms. Since $w_{k,i}$ and $\psi(a_i)$ are independent,

$$\mathbb{V}[w_{k,i}\psi(a_i)] = \mu_{w_{k,i}}^2 \mathbb{V}[\psi(a_i)] + \mathbb{E}[\psi(a_i)]^2 \sigma^2 + \sigma^2 \mathbb{V}[\psi(a_i)].$$

As A is compact, it is bounded, so there exists an M such that $|x_d| \leq M$ for all $1 \leq d \leq D$. Using lemma 5, and the mean-field assumptions,

$$\mathbb{V}[\psi(a_i)] \leq \mathbb{V}[a_i] = \sigma^2 \left(\sum_{d=1}^D x_d^2 + 1 \right) \leq \sigma^2 (DM^2 + 1).$$

Since a_i is a linear combination of Gaussian random variables, we have that $a_i \sim \mathcal{N}(\mu_{a_i}, \sigma_{a_i}^2)$, where $\mu_{a_i} = \sum_{d=1}^D \mu_{u_{i,d}} x_d + \mu_{v_i}$ and $\sigma_{a_i}^2 = \sigma^2 \left(\sum_{d=1}^D x_d^2 + 1 \right)$. Therefore, we have that [12]

$$\begin{aligned} \mathbb{E}[\psi(a_i)]^2 &= \left(\mu_{a_i} \Phi\left(\frac{\mu_{a_i}}{\sigma_{a_i}}\right) + \sigma_{a_i} N\left(\frac{\mu_{a_i}}{\sigma_{a_i}}\right) \right)^2 \leq \left(|\mu_{a_i}| \Phi\left(\frac{\mu_{a_i}}{\sigma_{a_i}}\right) + \sigma_{a_i} N\left(\frac{\mu_{a_i}}{\sigma_{a_i}}\right) \right)^2 \\ &\leq \left(|\mu_{a_i}| + \frac{\sigma_{a_i}}{\sqrt{2\pi}} \right)^2. \end{aligned}$$

We can then upper bound $\mathbb{V}[w_{k,i}\psi(a_i)]$ as follows:

$$\begin{aligned} \mathbb{V}[w_{k,i}\psi(a_i)] &\leq \mu_{w_{k,i}}^2 \sigma^2 (DM^2 + 1) + \left(|\mu_{a_i}| + \frac{\sigma_{a_i}}{\sqrt{2\pi}} \right)^2 \sigma^2 + \sigma^4 (DM^2 + 1) \\ &\leq \mu_{w_{k,i}}^2 \sigma^2 (DM^2 + 1) + \left(M \sum_{d=1}^D |\mu_{u_{i,d}}| + |\mu_{v_i}| + \frac{\sqrt{\sigma^2 (M^2 D + 1)}}{\sqrt{2\pi}} \right)^2 \sigma^2 + \sigma^4 (DM^2 + 1) \\ &:= v_{k,i}(\sigma). \end{aligned}$$

The second inequality follows since A is compact and we have $|\mu_{a_i}| \leq M \sum_{d=1}^D |\mu_{u_{i,d}}| + |\mu_{v_i}|$. Note that the upper bound $v_{k,i}(\sigma)$ is continuous and monotonically increasing in σ , and $v_{k,i}(0) = 0$. We can then upper bound the variance of the output:

$$\mathbb{V}[f_{\theta}^{(k)}(\mathbf{x})] \leq \sum_{i=1}^I \sum_{j=1}^I \sqrt{v_{k,i}(\sigma) v_{k,j}(\sigma)} + \sigma^2.$$

We then choose σ' such that for all $1 \leq k \leq K$ and for all $1 \leq i \leq I$, $v_{k,i}(\sigma') < \frac{\epsilon}{2I^2}$, and such that $\sigma'^2 < \frac{\epsilon}{2}$. Then

$$\mathbb{V}[f_{\theta}^{(k)}(\mathbf{x})] \leq I^2 \frac{\epsilon}{2I^2} + \sigma'^2 < \epsilon$$

for $1 \leq k \leq K$. Finally, applying lemma 5, we have $\mathbb{V}[\psi(f_{\theta}^{(k)}(\mathbf{x}))] < \epsilon$ for $1 \leq k \leq K$. \square

The following lemma states that by setting the standard deviation of the weights to be sufficiently small, we can with high probability make the sampled BNN output close to the BNN output evaluated at the mean parameters.

Lemma 7. *Let $A \subset \mathbb{R}^D$ be any compact set. Fix any $\boldsymbol{\mu}$ and any $\epsilon, \delta > 0$. Let all the standard deviations in $\boldsymbol{\sigma}$ be equal to a shared constant $\sigma > 0$. Then there exists $\sigma' > 0$ such that for all $\sigma < \sigma'$, and for any $\mathbf{x} \in A$,*

$$\Pr\left(|\psi(f_{\boldsymbol{\mu}}^{(k)}(\mathbf{x})) - \psi(f_{\theta}^{(k)}(\mathbf{x}))| > \epsilon\right) < \delta$$

for all $1 \leq k \leq K$.

Proof. Let $\theta \in \mathbb{R}^P$. We first note that $\psi(f_{\theta}^{(k)}(\mathbf{x}))$ is continuous as a function from $A \times \mathbb{R}^P \rightarrow \mathbb{R}$, under the metric topology induced by the Euclidean metric on $A \times \mathbb{R}^P$. Next, define a ball in parameter space

$$B_{\gamma} = \{\theta : \|\theta - \boldsymbol{\mu}\|_2 < \gamma\}.$$

Consider the closed ball of unit radius around $\boldsymbol{\mu}$, \bar{B}_1 . Note that \bar{B}_1 is compact, and therefore $A \times \bar{B}_1$ is compact as a product of compact spaces.

Since a continuous map from a compact metric space to another metric space is uniformly continuous, given $\epsilon > 0$, there exists a $0 < \tau < 1$ such that for all pairs $(\mathbf{x}_1, \theta_1), (\mathbf{x}_2, \theta_2) \in A \times \bar{B}_1$ such that $d((\mathbf{x}_1, \theta_1), (\mathbf{x}_2, \theta_2)) < \tau$, $|\psi(f_{\theta_1}^{(k)}(\mathbf{x}_1)) - \psi(f_{\theta_2}^{(k)}(\mathbf{x}_2))| < \epsilon$. Here $d(\cdot, \cdot)$ is the Euclidean metric on $A \times \mathbb{R}^P$. Since this is true for all $1 \leq k \leq K$, we can find a $0 < \tau < 1$ such that $|\psi(f_{\theta_1}^{(k)}(\mathbf{x}_1)) - \psi(f_{\theta_2}^{(k)}(\mathbf{x}_2))| < \epsilon$ holds for all k simultaneously, by taking the minimum of the τ over k .

Now choose $\sigma' > 0$ such that for all $\sigma < \sigma'$, $\Pr(\theta \in B_\tau) > 1 - \delta$. This event implies $d((\mathbf{x}, \theta), (\mathbf{x}, \boldsymbol{\mu})) = \|\theta - \boldsymbol{\mu}\|_2 < \tau$. Furthermore, $\theta \in \bar{B}_1$, since $\tau < 1$. Hence $|\psi(f_{\boldsymbol{\mu}}^{(k)}(\mathbf{x})) - \psi(f_{\theta}^{(k)}(\mathbf{x}))| < \epsilon$ holds for all $1 \leq k \leq K$. \square

The following lemma shows that for 1HL networks, we can make $\mathbb{E}[\psi(f_{\theta}^{(k)})]$ (the mean BNN output) close to $\psi(f_{\boldsymbol{\mu}}^{(k)})$ (the BNN output evaluated at the mean parameter settings) by choosing the standard deviation of the weights to be sufficiently small.

Lemma 8. *Let $A \subset \mathbb{R}^D$ be any compact set. Then, for any $\epsilon > 0$ and any $\boldsymbol{\mu}$, there exists a $\sigma_1 > 0$ such that for any shared standard deviation $\sigma < \sigma_1$,*

$$\left\| \mathbb{E}[\psi(f_{\theta}^{(k)})] - \psi(f_{\boldsymbol{\mu}}^{(k)}) \right\|_{\infty} < \epsilon$$

for all $1 \leq k \leq K$.

Proof. For all $\mathbf{x} \in A$ and any θ^* , by the triangle inequality

$$\left| \mathbb{E}[\psi(f_{\theta}^{(k)}(\mathbf{x}))] - \psi(f_{\boldsymbol{\mu}}^{(k)}(\mathbf{x})) \right| \leq \left| \mathbb{E}[\psi(f_{\theta}^{(k)}(\mathbf{x}))] - \psi(f_{\theta^*}^{(k)}(\mathbf{x})) \right| + \left| \psi(f_{\boldsymbol{\mu}}^{(k)}(\mathbf{x})) - \psi(f_{\theta^*}^{(k)}(\mathbf{x})) \right|.$$

Applying lemma 7 with $\epsilon' = \epsilon/2$ and $\delta = 1/4$, we can find a σ' such that for all $\sigma < \sigma'$, $\left| \psi(f_{\boldsymbol{\mu}}^{(k)}(\mathbf{x})) - \psi(f_{\theta^*}^{(k)}(\mathbf{x})) \right| \leq \epsilon/2$ with probability at least $3/4$. By lemma 6, we can find a σ'' such that for all $\sigma < \sigma''$, $\mathbb{V}[\psi(f_{\theta}^{(k)}(\mathbf{x}))] < \frac{\epsilon^2}{16K}$. Choose $0 < \sigma < \min(\sigma', \sigma'')$. We can apply Chebyshev's inequality to each random variable $\psi(f_{\theta}^{(k)}(\mathbf{x}))$,

$$\Pr \left[\left| \psi(f_{\theta}^{(k)}(\mathbf{x})) - \mathbb{E}[\psi(f_{\theta}^{(k)}(\mathbf{x}))] \right| > \epsilon/2 \right] < \frac{1}{4K}.$$

Applying the union bound, the probability that $|\psi(f_{\theta}^{(k)}(\mathbf{x})) - \mathbb{E}[\psi(f_{\theta}^{(k)}(\mathbf{x}))]| \leq \epsilon/2$ for all k simultaneously is at least $3/4$. Therefore, for any \mathbf{x} we can find a θ^* such that $|\psi(f_{\theta^*}^{(k)}(\mathbf{x})) - \mathbb{E}[\psi(f_{\theta}^{(k)}(\mathbf{x}))]| \leq \epsilon/2$ and $|\psi(f_{\boldsymbol{\mu}}^{(k)}(\mathbf{x})) - \psi(f_{\theta^*}^{(k)}(\mathbf{x}))| \leq \epsilon/2$ simultaneously because both events occur with probability at least $1/2$ and therefore have a non-empty intersection. Therefore for all \mathbf{x} and all k

$$\left| \mathbb{E}[\psi(f_{\theta}^{(k)}(\mathbf{x}))] - \psi(f_{\boldsymbol{\mu}}^{(k)}(\mathbf{x})) \right| \leq \epsilon. \quad \square$$

We can now complete the proof of theorem 3 for \mathcal{Q}_{FFG} .

Proof of theorem 6. Consider the case of a 2-hidden layer ReLU Bayesian neural network with 2 units in the second hidden layer. Denote the inputs to these units as $f_{\theta}^{(1)}(\mathbf{x})$ and $f_{\theta}^{(2)}(\mathbf{x})$ respectively, where θ are the parameters in the bottom two weight matrices and biases of the network. The output of the network can then be written as,

$$f(\mathbf{x}) = s_1 \psi(f_{\theta}^{(1)}(\mathbf{x})) + s_2 \psi(f_{\theta}^{(2)}(\mathbf{x})) + t \quad (11)$$

where the s_i are the weights in the final layer and t is the bias. Taking expectations on both sides,

$$\mathbb{E}[f(\mathbf{x})] = \mathbb{E}[s_1 \psi(f_{\theta}^{(1)}(\mathbf{x}))] + \mathbb{E}[s_2 \psi(f_{\theta}^{(2)}(\mathbf{x}))] + \mathbb{E}[t]$$

Choose $\mu_{s_1} = 1, \mu_{s_2} = 0$, and note that s_1 is independent of θ by the mean field assumption. Then

$$\mathbb{E}[f(\mathbf{x})] = \mathbb{E}[\psi(f_{\theta}^{(1)}(\mathbf{x}))] + \mathbb{E}[t]. \quad (12)$$

Define $\mu_t = -\min_{\mathbf{x}' \in A} g(\mathbf{x}')$ (as A is compact and g is continuous, this minimum is well-defined). Define $\tilde{g}(\mathbf{x}) \geq 0$ to be $g(\mathbf{x}) - \min_{\mathbf{x}' \in A} g(\mathbf{x}')$. By the universal approximation theorem (theorem 7)

we can find a setting of the mean parameters, $\boldsymbol{\mu}$ in the first two layers (i.e. excluding the parameters of the distributions on s_1, s_2 and t) such that

$$\|f_{\boldsymbol{\mu}}^{(1)} - \tilde{g}\|_{\infty} < \epsilon/2 \quad \text{and} \quad \|f_{\boldsymbol{\mu}}^{(2)} - \sqrt{h}\|_{\infty} < \epsilon/2.$$

This can be done by splitting the neurons in the first hidden layer into two sets, where the first and second set are responsible for $f^{(1)}, f^{(2)}$ respectively, and the weights from each set to the output of the other set are zero. Since $\tilde{g}(\mathbf{x}) > 0$, applying the ReLU can only make $f^{(1)}$ closer to \tilde{g} . Hence $\|\psi(f_{\boldsymbol{\mu}}^{(1)}) - \tilde{g}\|_{\infty} < \epsilon/2$.

By lemma 8, we can find a $\sigma_1 > 0$ for this $\boldsymbol{\mu}$ such that when the standard deviations in the first two layers are set to any shared constant $\sigma < \sigma_1$,

$$\left\| \mathbb{E} \left[\psi(f_{\theta}^{(1)}) \right] - \psi(f_{\boldsymbol{\mu}}^{(1)}) \right\|_{\infty} < \epsilon/2.$$

By the triangle inequality, $\left\| \mathbb{E} \left[\psi(f_{\theta}^{(1)}) \right] - \tilde{g} \right\|_{\infty} < \epsilon$. Combining with equation (12), it follows that the expectation can approximate any continuous function g .

We now consider the variance of equation (11).

$$\begin{aligned} \mathbb{V}[f(\mathbf{x})] &= \mathbb{V}[s_1\psi(f_{\theta}^{(1)}(\mathbf{x})) + s_2\psi(f_{\theta}^{(2)}(\mathbf{x}))] + \mathbb{V}[t] \\ &= \mathbb{V}[s_1\psi(f_{\theta}^{(1)}(\mathbf{x}))] + \mathbb{V}[s_2\psi(f_{\theta}^{(2)}(\mathbf{x}))] + 2\text{Cov}(s_1\psi(f_{\theta}^{(1)}(\mathbf{x})), s_2\psi(f_{\theta}^{(2)}(\mathbf{x}))) + \sigma_t^2. \end{aligned}$$

Choose $\sigma_t^2 = \epsilon$. We now consider $\mathbb{V}[s_1\psi(f_{\theta}^{(1)}(\mathbf{x}))]$. As s_1 is independent of θ ,

$$\mathbb{V}[s_1\psi(f_{\theta}^{(1)}(\mathbf{x}))] = \mu_{s_1}^2 \mathbb{V}[\psi(f_{\theta}^{(1)}(\mathbf{x}))] + \sigma_{s_1}^2 \mathbb{E} \left[\psi(f_{\theta}^{(1)}(\mathbf{x})) \right]^2 + \mathbb{V}[\psi(f_{\theta}^{(1)}(\mathbf{x}))] \sigma_{s_1}^2.$$

Recall $\mu_{s_1} = 1$ and choose $\sigma_{s_1}^2 = \min \left(1, \epsilon / \left(\max_{x \in A} \mathbb{E} \left[\psi(f_{\theta}^{(1)}(\mathbf{x})) \right]^2 \right) \right)$, then

$$\mathbb{V}[s_1\psi(f_{\theta}^{(1)}(\mathbf{x}))] \leq 2\mathbb{V}[\psi(f_{\theta}^{(1)}(\mathbf{x}))] + \epsilon.$$

By lemma 6, we can find a σ_2 such that for any $\sigma < \sigma_2$, $\mathbb{V}[\psi(f_{\theta}^{(1)}(\mathbf{x}))] \leq \epsilon$. For any such σ , $\mathbb{V}[s_1\psi(f_{\theta}^{(1)}(\mathbf{x}))] \leq 3\epsilon$.

We now choose $\sigma_{s_2}^2 = 1$ and consider

$$\begin{aligned} \mathbb{V}[s_2\psi(f_{\theta}^{(2)}(\mathbf{x}))] &= \mu_{s_2}^2 \mathbb{V}[\psi(f_{\theta}^{(2)}(\mathbf{x}))] + \sigma_{s_2}^2 \mathbb{E} \left[\psi(f_{\theta}^{(2)}(\mathbf{x})) \right]^2 + \sigma_{s_2}^2 \mathbb{V}[\psi(f_{\theta}^{(2)}(\mathbf{x}))] \\ &= \mathbb{E} \left[\psi(f_{\theta}^{(2)}(\mathbf{x})) \right]^2 + \mathbb{V}[\psi(f_{\theta}^{(2)}(\mathbf{x}))]. \end{aligned}$$

By lemma 6, we can find a σ_3 such that for any $\sigma < \sigma_3$, $\mathbb{V}[\psi(f_{\theta}^{(2)}(\mathbf{x}))] < \epsilon$.

By the universal function approximator theorem (theorem 7) we can find a setting of the mean parameters, $\boldsymbol{\mu}$ in the first two layers such that $\|f_{\boldsymbol{\mu}}^{(2)} - \sqrt{h}\|_{\infty} < \epsilon/2$. Since $\sqrt{h(\mathbf{x})} > 0$, the ReLU can only make $f^{(2)}$ closer to \sqrt{h} , $\|\psi(f_{\boldsymbol{\mu}}^{(2)}) - \sqrt{h}\|_{\infty} < \epsilon/2$.

By lemma 8, we can find a setting of σ for this $\boldsymbol{\mu}$ such that

$$\left\| \mathbb{E} \left[\psi(f_{\theta}^{(2)}) \right] - \psi(f_{\boldsymbol{\mu}}^{(2)}) \right\|_{\infty} < \epsilon/2.$$

By the triangle inequality,

$$\left\| \mathbb{E} \left[\psi(f_{\theta}^{(2)}) \right] - \sqrt{h} \right\|_{\infty} < \epsilon.$$

This implies,

$$\begin{aligned} \left\| \mathbb{E} \left[\psi(f_{\theta}^{(2)}) \right]^2 - h \right\|_{\infty} &= \left\| \left(\mathbb{E} \left[\psi(f_{\theta}^{(2)}) \right] - \sqrt{h} \right) \left(\mathbb{E} \left[\psi(f_{\theta}^{(2)}) \right] + \sqrt{h} \right) \right\|_{\infty} \\ &\leq \epsilon \left\| \mathbb{E} \left[\psi(f_{\theta}^{(2)}) \right] + \sqrt{h} \right\|_{\infty} \\ &\leq \epsilon(2\|\sqrt{h}\|_{\infty} + \epsilon) \end{aligned}$$

We therefore have,

$$\begin{aligned}\|\mathbb{V}[f] - h\|_\infty &\leq E(\epsilon) + 2\text{Cov}(s_1\psi(f_\theta^{(1)}(\mathbf{x})), s_2\psi(f_\theta^{(2)}(\mathbf{x}))) \\ &\leq E(\epsilon) + 2\sqrt{\mathbb{V}[s_1\psi(f_\theta^{(1)}(\mathbf{x}))]\mathbb{V}[s_2\psi(f_\theta^{(2)}(\mathbf{x}))]} \\ &\leq E(\epsilon) + C\sqrt{\epsilon}\end{aligned}$$

where the first inequality is Cauchy-Schwarz, and $E(\epsilon)$ is a function that tends to zero with ϵ and C is a constant. The theorem follows by choosing $\sigma < \min\{\sigma_1, \sigma_2, \sigma_3\}$. \square

The construction in our proof used a 2HL BNN with only two neurons in the second hidden layer. The construction still works for wider hidden layers, by setting the unused neurons to have zero mean and sufficiently small variance.

An analogous statement to theorem 3 for networks with more than two hidden layers can be proved inductively: applying theorem 3 for 2HL BNNs we can choose the variance to be uniformly small, thus satisfying the condition stated in lemma 6. The proof of lemma 7 applies equally for the output of 2HL BNNs. The rest of the proof then follows as stated.

D.2 Proof of theorem 6 for MCDO

In order to prove the universality result for deep dropout, we first prove two lemmas about 1HL dropout networks. The following lemma states that the mean of a 1HL dropout network is a universal function approximator, while its variance can simultaneously be made arbitrarily small.

Lemma 9. *Consider any $\epsilon > 0$ and any continuous function, m mapping from a compact subset, A of \mathbb{R}^D to \mathbb{R} . Then there exists a (random) ReLU neural network of the form*

$$f(x) = \sum_{i=1}^I w_i \gamma_i \psi \left(\sum_{d=1}^D u_{i,d} x_d + v_i \right) + b$$

with $\gamma_i \stackrel{\text{i.i.d.}}{\sim} \text{Bern}(1-p)$ such that $\|\mathbb{E}[f] - m\|_\infty < \epsilon$ and $\|\mathbb{V}[f]\|_\infty \leq \epsilon$.

Proof. By the universal approximation theorem Leshno et al. [26], there exists a J and 1HL network of the form,

$$g(x) = \sum_{j=1}^J \tilde{w}_j \psi \left(\sum_{d=1}^D \tilde{u}_{j,d} x_d + v_j \right) + b,$$

such that $\|g - m\|_\infty \leq \epsilon$. Define the dropout network,

$$f^{(1)}(x) = \sum_{j=1}^J \frac{\tilde{w}_j}{1-p} \psi \left(\sum_{d=1}^D \tilde{u}_{j,d} x_d + v_j \right) + b.$$

Then $\mathbb{E}[f^{(1)}] = g$, so that $\|\mathbb{E}[f^{(1)}] - m\|_\infty \leq \epsilon$. Let $S = \|\mathbb{V}[f^{(1)}]\|_\infty < \infty$.

Define $f = \frac{1}{L} \sum_{\ell=1}^L f^{(1,\ell)}$ where each $f^{(1,\ell)}$ is an independent realisation of $f^{(1)}$. Then $\mathbb{E}[f] = g$ and $\mathbb{V}[f] = \frac{\mathbb{V}[f^{(1)}]}{\sqrt{L}} \leq \frac{S}{\sqrt{L}}$. f can be realised by a dropout network by combining L copies of $f^{(1)}$ together with identical weights within each copy and 0 weights connecting the various copies. Choosing $L = (S/\epsilon)^2$ completes the proof. \square

The following lemma states that the mean of the MCDO network can approximate any continuous positive function, after application of the ReLU non-linearity.

Lemma 10. *Given a positive mean function m with $0 < \delta \leq \|m\|_\infty \leq \Delta$ and a stochastic process f such that $\|\mathbb{E}[f] - m\|_\infty \leq \epsilon \leq \delta$ and $\|\mathbb{V}[f]\|_\infty \leq \epsilon$,*

$$\|\mathbb{E}[\psi(f)] - m\|_\infty \leq \epsilon + \frac{\sqrt{\epsilon^2 + \epsilon(\Delta + \epsilon)^2}}{\delta - \epsilon} = \mathcal{O}(\Delta\sqrt{\epsilon}/(\delta - \epsilon))$$

and $\|\mathbb{V}[\psi(f)]\|_\infty \leq \epsilon$. In the big-O notation, we assume Δ is bounded below by a constant and ϵ, δ are bounded above by a constant.

Proof. The bound $\|\mathbb{V}[\psi(f)]\|_\infty \leq \epsilon$ follows from lemma 5. We consider the expectation of $\psi(f(\mathbf{x}))$ for some arbitrary fixed \mathbf{x} ,

$$\begin{aligned} |\mathbb{E}[\psi(f(\mathbf{x}))] - m(\mathbf{x})| &= |\mathbb{E}[f(\mathbf{x})] - m(\mathbf{x}) - \mathbb{E}[\min(0, f(\mathbf{x}))]| \\ &\leq |\mathbb{E}[f(\mathbf{x})] - m(\mathbf{x})| + |\mathbb{E}[\min(0, f(\mathbf{x}))]| \\ &\leq \epsilon + |\mathbb{E}[\min(0, f(\mathbf{x}))]|. \end{aligned}$$

We therefore bound $|\mathbb{E}[\min(0, f(\mathbf{x}))]|$.

$$|\mathbb{E}[\min(0, f(\mathbf{x}))]| = |\mathbb{E}[f(\mathbf{x})\mathbf{1}\{\mathbf{x} : f(\mathbf{x}) < 0\}]| \leq \sqrt{\mathbb{E}[f(\mathbf{x})^2] \Pr(f(\mathbf{x}) < 0)}.$$

The inequality uses Cauchy-Schwarz, that the square of an indicator function is itself and reinterprets the expectation of an indicator function as a probability. We bound the two terms on the RHS separately.

$$\mathbb{E}[f(\mathbf{x})^2] = \mathbb{V}[f(\mathbf{x})] + \mathbb{E}[f(\mathbf{x})]^2 \leq \epsilon + \mathbb{E}[f(\mathbf{x})]^2 \leq \epsilon + (m(\mathbf{x}) + \epsilon)^2 \leq \epsilon + (\Delta + \epsilon)^2$$

We use Chebyshev's inequality to bound the probability $f(\mathbf{x}) < 0$,

$$\begin{aligned} \Pr(f(\mathbf{x}) < 0) &\leq \Pr(|f(\mathbf{x}) - \mathbb{E}[f(\mathbf{x})]| > m(\mathbf{x}) - \epsilon) \\ &\leq \frac{\mathbb{V}[f(\mathbf{x})]}{(m(\mathbf{x}) - \epsilon)^2} \\ &\leq \frac{\epsilon}{(m(\mathbf{x}) - \epsilon)^2} \\ &\leq \frac{\epsilon}{(\delta - \epsilon)^2}. \quad \square \end{aligned}$$

Having collected the necessary lemmas, we provide a construction that proves theorem 6.

Proof of theorem 6. Consider a 2HL dropout NN. Let the pre-activations in the first hidden layer be collectively denoted \mathbf{a}_1 , and the random dropout masks by ϵ_1 . Let the second hidden layer have $I + 2$ hidden units. Let \odot denote the elementwise product of two vectors of the same length. Define the pre-activations of two of the second hidden layer units by $a_v = \mathbf{w}_v^\top(\epsilon_1 \odot \psi(\mathbf{a}_1))$, i.e. both these hidden units have identical weight vectors \mathbf{w}_v and dropout masks, and are hence the same random variable. Similarly, let the remaining I second hidden layer pre-activations be defined by $a_m = \mathbf{w}_m^\top(\epsilon_1 \odot \psi(\mathbf{a}_1))$, again all being the same random variable. Furthermore, let $(\mathbf{w}_v)_i = 0$ whenever $(\mathbf{w}_m)_i \neq 0$ and vice versa, so that the first hidden layer neurons that influence a_v and those that influence a_m form disjoint sets. Then the output of the 2HL network is:

$$f = \epsilon_a w_{2,a} \psi(a_v) + \epsilon_b w_{2,b} \psi(a_v) + \sum_{i=1}^I \epsilon_i w_{2,i} \psi(a_m) + b_2,$$

where $\epsilon_a, \epsilon_b, \{\epsilon_i\}_{i=1}^I$ are the final layer dropout masks and $\{w_{2,i}\}_{i=1}^I, b_2$ are the final layer weights and bias.

We now make the choices $w_{2,a} = 1, w_{2,b} = -1, w_{2,i} = \alpha$, where $\alpha I = 1/(1-p)$. Then $\mathbb{E}[f] = \mathbb{E}[\psi(a_m)] + b_2$. Let $b_2 = \min_{\mathbf{x} \in A} g - \delta$, where $\delta > 0$ and the min exists due to compactness of A . Define $g' = g - b_2$. Since a_m is just the output of a single-hidden layer dropout network, for any $\gamma' > 0$ we can use lemma 9 to choose $\|\mathbb{E}[a_m] - g'\|_\infty < \gamma'$ and $\|\mathbb{V}[a_m]\|_\infty < \gamma'$. Since g' is bounded below by δ and bounded above by some $\Delta \in \mathbb{R}$ (by continuity of g and compactness of A), we can then apply lemma 10 to obtain $\|\mathbb{E}[a_m] - g'\|_\infty = \mathcal{O}(\Delta\sqrt{\epsilon'}/(\delta - \epsilon'))$ and $\|\mathbb{V}[\psi(a_m)]\|_\infty < \gamma'$. We can use this to bound the error in the mean of the 2HL network output:

$$\|\mathbb{E}[f] - g\|_\infty = \|\mathbb{E}[\psi(a_m)] + b_2 - g\|_\infty = \|\mathbb{E}[\psi(a_m)] - g'\|_\infty = \mathcal{O}(\Delta\sqrt{\gamma'}/(\delta - \gamma')).$$

We can choose γ' to depend on δ, Δ such that $\|\mathbb{E}[f] - g\|_\infty < \gamma$, proving the first part of the theorem. Next, calculating the variance,

$$\mathbb{V}[f] = \mathbb{V}\left[(\epsilon_a - \epsilon_b)\psi(a_v) + \alpha\psi(a_m) \sum_{i=1}^I \epsilon_i\right] = \mathbb{V}[(\epsilon_a - \epsilon_b)\psi(a_v)] + \alpha^2 \mathbb{V}\left[\psi(a_m) \sum_{i=1}^I \epsilon_i\right]. \quad (13)$$

Next we show that by taking I sufficiently large, we can make the second term arbitrarily small. We have,

$$\begin{aligned}\mathbb{V}\left[\psi(a_m)\sum_{i=1}^I\epsilon_i\right] &= \mathbb{V}[\psi(a_m)]\mathbb{V}\left[\sum_{i=1}^I\epsilon_i\right] + \mathbb{V}[\psi(a_m)]\mathbb{E}\left[\sum_{i=1}^I\epsilon_i\right]^2 + \mathbb{V}\left[\sum_{i=1}^I\epsilon_i\right]\mathbb{E}[\psi(a_m)]^2 \\ &= \mathbb{V}[\psi(a_m)]Ip(1-p) + \mathbb{V}[\psi(a_m)]I^2(1-p)^2 + Ip(1-p)\mathbb{E}[\psi(a_m)]^2 \\ &\leq \gamma'Ip(1-p) + \gamma'I^2(1-p)^2 + Ip(1-p)\mathbb{E}[\psi(a_m)]^2\end{aligned}$$

The first two of these three terms can be made arbitrarily small by choosing γ' sufficiently small. The third term, upon multiplying by α^2 , becomes

$$\alpha^2Ip(1-p)\mathbb{E}[\psi(a_m)]^2 = \frac{p}{I(1-p)}\mathbb{E}[\psi(a_m)]^2,$$

which can also be made arbitrarily small by choosing $I \in \mathbb{N}$ sufficiently large. We now show that the first term in equation (13) can well approximate our target variance function h .

$$\mathbb{V}[(\epsilon_a - \epsilon_b)\psi(a_v)] = \mathbb{V}[\epsilon_a - \epsilon_b]\mathbb{V}[\psi(a_v)] + \mathbb{V}[\epsilon_a - \epsilon_b]\mathbb{E}[\psi(a_v)]^2 + \mathbb{V}[\psi(a_v)]\mathbb{E}[\epsilon_a - \epsilon_b]^2 \quad (14)$$

$$= 2p(1-p)\mathbb{V}[\psi(a_v)] + 2p(1-p)\mathbb{E}[\psi(a_v)]^2 \quad (15)$$

Define

$$h' = \sqrt{\frac{h}{2p(1-p)}} + \delta',$$

for some $\delta' > 0$. Again applying lemma 9 (which we can do independently of the choice of a_m since neurons influencing a_v and a_m are disjoint), for any $\gamma'' > 0$ we can choose $\|\mathbb{E}[a_v] - h'\|_\infty < \gamma''$ and $\|\mathbb{V}[a_v]\|_\infty < \gamma''$. The first term in equation (15) can be made arbitrarily small by choosing γ'' small enough. We can again apply lemma 10 so that $\|\mathbb{E}[\psi(a_v)] - h'\|_\infty = \mathcal{O}(\Delta'\sqrt{\gamma''}/(\delta' - \gamma''))$. We then bound the difference between the second term in equation (15) and our target variance function:

$$\left\|2p(1-p)\mathbb{E}[\psi(a_v)]^2 - h\right\|_\infty \leq \left\|\sqrt{2p(1-p)}\mathbb{E}[\psi(a_v)] + \sqrt{h}\right\|_\infty \left\|\sqrt{2p(1-p)}\mathbb{E}[\psi(a_v)] - \sqrt{h}\right\|_\infty \quad (16)$$

$$\leq \left(\left\|2\sqrt{h}\right\|_\infty + \left\|\sqrt{2p(1-p)}\mathbb{E}[\psi(a_v)] - \sqrt{h}\right\|_\infty\right) \left\|\sqrt{2p(1-p)}\mathbb{E}[\psi(a_v)] - \sqrt{h}\right\|_\infty \quad (17)$$

where equation (16) follows from sub-multiplicativity of the infinity norm. Expanding the second term in equation (17),

$$\begin{aligned}\left\|\sqrt{2p(1-p)}\mathbb{E}[\psi(a_v)] - \sqrt{h}\right\|_\infty &= \sqrt{2p(1-p)}\|\mathbb{E}[\psi(a_v)] - h' + \delta'\|_\infty \\ &= \mathcal{O}(\delta' + \Delta'\sqrt{\gamma''}/(\delta' - \gamma''))\end{aligned}$$

By first choosing δ' sufficiently small, and then choosing γ'' depending on δ' , we can make this error term arbitrarily small. Since all the other contributions to $\mathbb{V}[f]$ were made arbitrarily small, this allows us to set $\|\mathbb{V}[f] - h\| < \gamma$, for any $\gamma > 0$, completing the proof. \square

In order to provide an analogous construction for MCDO BNNs with more than 2 hidden layers, we note that the above proof only requires a BNN output with a universal mean function and an arbitrarily small variance function in lemma 9. Instead of a 1HL network, we can apply theorem 3 to construct a 2 or more hidden layer network to provide these mean and variance functions. The rest of the proof then follows as in the 2HL case.

E Dropout With Inputs Dropped Out

The behaviour of MC dropout with inputs dropped out is somewhat different, both theoretically and empirically, from the case when inputs are not dropped out as discussed in the main body.

E.1 Single-Hidden Layer Networks

In the single-hidden layer case, the variance is no longer convex as a function of \mathbf{x} . On the other hand, this approximating family still struggles to represent in-between uncertainty:

Theorem 8 (MC dropout, dropped out inputs). *Consider a single-hidden layer ReLU neural network mapping from $\mathbb{R}^D \rightarrow \mathbb{R}^K$ with $I \in \mathbb{N}$ hidden units. The corresponding mapping is given by $f^{(k)}(\mathbf{x}) = \sum_{i=1}^I w_{k,i} \psi\left(\sum_{d=1}^D u_{i,d} x_d + v_i\right) + b_k$ for $1 \leq k \leq K$, where $\psi(a) = \max(0, a)$. Assume \mathbf{v} is set deterministically and*

$$q(\mathbf{W}, \mathbf{b}, \mathbf{U}) = q(\mathbf{U})q(\mathbf{b}|\mathbf{U}) \prod_i q_i(\mathbf{w}_i|\mathbf{U}),$$

where $\mathbf{w}_i = \{w_{k,i}\}_{k=1}^K$ are the weights out of neuron i , $\mathbf{b} = \{b_k\}_{k=1}^K$ are the output biases and $q(\mathbf{U})$, $q(\mathbf{b}|\mathbf{U})$ and $q_i(\mathbf{w}_i|\mathbf{U})$ are arbitrary probability densities with finite first two moments. Then, for any finite set of points $\mathcal{S} \subset \mathbb{R}^D$ such that $\mathbf{0}$ is in the convex hull of \mathcal{S} ,

$$\mathbb{V}[f^{(k)}(\mathbf{0})] \leq \max_{\mathbf{s} \in \mathcal{S}} \left\{ \mathbb{V}[f^{(k)}(\mathbf{s})] \right\} \quad \text{for } 1 \leq k \leq K. \quad (18)$$

In order to prove theorem 8 we use the following simple lemma,

Lemma 11. *Let $f : \mathbb{R}^D \rightarrow \mathbb{R}$ be a convex function and consider a finite set of points $\mathcal{S} \subset \mathbb{R}^D$. Then for any point \mathbf{r} in the convex hull of \mathcal{S} , $f(\mathbf{r}) \leq \max_{\mathbf{s} \in \mathcal{S}} \{f(\mathbf{s})\}$.*

Proof of lemma 11. Let $\{\mathbf{s}_n\}_{n=1}^N = \mathcal{S}_N \subset \mathbb{R}^D$. We proceed by induction. The lemma is true for $N = 2$ by the definition of convexity. Assume it is true for N . Let $\text{Conv}(\mathcal{S}_{N+1})$ denote the convex hull of \mathcal{S}_{N+1} . Consider a point $\mathbf{r}_{N+1} \in \text{Conv}(\mathcal{S}_{N+1})$. Then

$$f(\mathbf{r}_{N+1}) = f\left(\sum_{n=1}^{N+1} \alpha_n \mathbf{s}_n\right) \quad (19)$$

with $\sum_{n=1}^{N+1} \alpha_n = 1$ and $\alpha_n \geq 0$ for $1 \leq n \leq N+1$. We can write

$$f(\mathbf{r}_{N+1}) = f\left(\left(\sum_{n=1}^N \alpha_n\right) \mathbf{t}_N + \alpha_{N+1} \mathbf{s}_{N+1}\right) \leq \max\{f(\mathbf{t}_N), f(\mathbf{s}_{N+1})\} \quad (20)$$

where $\mathbf{t}_N := \sum_{n=1}^N \alpha_n \mathbf{s}_n / \sum_{n=1}^N \alpha_n$, and we have used the convexity of f . By the induction assumption, $f(\mathbf{t}_N) \leq \max_{\mathbf{s} \in \mathcal{S}_N} \{f(\mathbf{s})\}$, since $\mathbf{t}_N \in \text{Conv}(\mathcal{S}_N)$. Combining this with equation (20) completes the proof. \square

Proof of theorem 8. By the law of total variance,

$$\mathbb{V}[f^{(k)}(\mathbf{x})] = \mathbb{E}[\mathbb{V}[f^{(k)}(\mathbf{x})|\mathbf{U}]] + \mathbb{V}[\mathbb{E}[f^{(k)}(\mathbf{x})|\mathbf{U}]].$$

Using lemma 1, $\mathbb{V}[f^{(k)}(\mathbf{x})|\mathbf{U}]$ is convex as a function of \mathbf{x} . As the expectation of a convex function is convex, the first term is a convex function of \mathbf{x} . This implies

$$\mathbb{E}[\mathbb{V}[f^{(k)}(\mathbf{0})|\mathbf{U}]] \leq \max_{\mathbf{s} \in \mathcal{S}} \left\{ \mathbb{E}[\mathbb{V}[f^{(k)}(\mathbf{s})|\mathbf{U}]] \right\},$$

by lemma 11. $\mathbb{V}[\mathbb{E}[f^{(k)}(\mathbf{x})|\mathbf{U}]]$ is non-negative everywhere. As the output of the first layer is independent of the matrix \mathbf{U} at $\mathbf{x} = \mathbf{0}$, $\mathbb{E}[f^{(k)}(\mathbf{0})|\mathbf{U}]$ is deterministic. So $\mathbb{V}[\mathbb{E}[f^{(k)}(\mathbf{0})|\mathbf{U}]] = 0$, completing the proof. \square

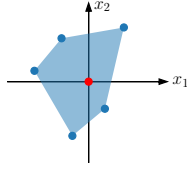


Figure 9: Schematic illustration of the bound in theorem 2, showing the input domain of a single-hidden layer MC dropout BNN, for the case $\mathbf{x} \in \mathbb{R}^2$. The convex hull (in blue) of the blue points contains the origin. Hence the variance at the origin cannot exceed the variance at any of the blue points.

E.2 Deep Networks

In the case when the network has several hidden layers, dropout with inputs dropped defines a posterior with somewhat strange properties, as observed in Gal [13, Section 4.2.1]. In particular, in D dimensions, a typical sample function from the approximate posterior will be constant as a function of roughly pD of the input dimensions. However, which dimensions it is constant along depends on the particular sample. This behaviour is unlikely to be shared by the exact posterior. We are able to exploit this type of behaviour to show that if inputs are dropped out, there are simple combinations of mean and variance functions that cannot be simultaneously approximated by the corresponding approximating family.

Proposition 1. *Consider f the (stochastic) output of an MC dropout network of arbitrary depth with inputs dropped out. For any $x, x' \in \mathbb{R}$ such that $\mathbb{V}[f(x)], \mathbb{V}[f(x')] < \epsilon^2$, $|\mathbb{E}[f(x)] - \mathbb{E}[f(x')]| \leq 2\epsilon\sqrt{2/p}$.*

Proof. With probability p , the input is dropped out, so $\Pr(f(x) = f(x')) \geq p$. We apply Chebyshev's inequality giving the bounds,

$$\Pr(|f(x) - \mathbb{E}[f(x)]| \leq r\epsilon) \geq 1 - 1/r^2 \quad \text{and} \quad \Pr(|f(x') - \mathbb{E}[f(x')]| \leq r\epsilon) \geq 1 - 1/r^2.$$

for any $r > 0$. Choose $r = \sqrt{2/p} + \delta$ for any $\delta > 0$, then there exists a realisation of the dropout network such that $|f(x) - \mathbb{E}[f(x)]| \leq r\epsilon$, $|f(x') - \mathbb{E}[f(x')]| \leq r\epsilon$ and $f(x) = f(x')$ simultaneously. Consequently,

$$\begin{aligned} |\mathbb{E}[f(x)] - \mathbb{E}[f(x')]| &= |\mathbb{E}[f(x)] - f(x) + f(x) - \mathbb{E}[f(x')]| \\ &= |\mathbb{E}[f(x)] - f(x) + f(x') - \mathbb{E}[f(x')]| \\ &\leq |\mathbb{E}[f(x)] - f(x)| + |f(x') - \mathbb{E}[f(x')]| \\ &\leq 2r\epsilon = 2\epsilon\sqrt{2/p} + 2\epsilon\delta. \end{aligned}$$

Taking the limit as $\delta \rightarrow 0$ completes the proof. \square

In other words we can bound the difference in the mean output at two points in terms of the uncertainty at those points and the dropout probability.

In $D > 1$ dimensions, we can get similarly tight bounds on lines parallel to a coordinate axis: for \mathbf{x}, \mathbf{x}' on such a line $\Pr(f(\mathbf{x}) = f(\mathbf{x}')) \geq p$ still holds. If the dimension on which \mathbf{x} and \mathbf{x}' differ is dropped out $f(\mathbf{x}) = f(\mathbf{x}')$.

Alternatively in D dimensions for arbitrary $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^D$, $\Pr(f(\mathbf{x}) = f(\mathbf{x}')) \geq p^D$. This comes from noting that with probability p^D the output of the network is a constant function. However, we note this bound becomes exponentially weak as the input dimension increases.

E.3 Details of Experiments Minimising Squared Loss

We generated a dataset that consisted of two separated clusters in one dimension. We then fit a Gaussian process to the dataset and computed the predictive mean and variance on a one-dimensional

grid of $N = 40$ points, call these point X . Let $\mu(X) \in \mathbb{R}^N$ denote the mean of the GP regression at these points $\sigma^2(X) \in \mathbb{R}^N$ denote its variance. We define a loss function as

$$\mathcal{L}(\phi) = \|\mathbb{E}_{q_\phi}[f(X)] - \mu(X)\|_2^2 + \|\mathbb{V}_{q_\phi}[f(X)] - \sigma^2(X)\|_2^2.$$

The expectation and variance are Monte Carlo estimated using 128 samples. We use full batch optimisation with ADAM with learning rate 1×10^{-3} for 50,000 iterations. A dropout rate of 0.05 is used for MCDO. Weights and biases are initialized at the prior for MFVI.

F Details and Additional Figures for Section 4.2

In this appendix, we provide details of the protocol used to generate figure 5.

F.1 Experimental Details

Data: We consider the dataset from figure 3 with $\mathbf{x} \in \overrightarrow{\mathbf{pq}}$, where $\mathbf{p} = (-1.2, -1.2)$ and $\mathbf{q} = (1.2, 1.2)$ i.e. the line segment between and including the two data clusters. We evaluate the overconfidence ratio on a discretisation of $\overrightarrow{\mathbf{pq}}$.

Choosing the Prior: For each depth a fully-connected ReLU network with 50 hidden units per layer is used. The prior mean for all parameters is chosen to be 0. The prior standard deviation for the bias parameters is chosen as $\sigma_b = 1$ for all experiments. In figure 5, the prior weight standard deviation is selected so that the prior standard deviation in function space at the region containing data is approximately constant. In particular, let σ_w/\sqrt{H} be the prior standard deviation of each weight, where H is the number of inputs to the weight matrix. We choose $\sigma_w = \{4, 3, 2.25, 2, 2, 1.9, 1.75, 1.75, 1.7, 1.65\}$ for depths 1-10 respectively, which ensures the prior standard deviations (of both the GP and the BNN) in function space at the points $(1, 1)$ and $(-1, -1)$ (the centres of the data clusters) are between 10 and 15. Choosing a fixed σ_w such as $\sigma_w = 4$ for all depths would have caused the prior standard deviation in function space to grow unreasonably large with increasing depth; see Schoenholz et al. [39]. All models used a fixed Gaussian likelihood with standard deviation 0.1.

Fitting the GP: The Gaussian process was implemented using GPFlow [31] with the infinite-width ReLU BNN kernel implemented following [25]. All hyperparameters were fixed and exact inference was performed using Cholesky decomposition.

Fitting MFVI: We initialize the standard deviations of weights to be small and train for many epochs, following Tomczak et al. [43], Swaroop et al. [42] who found this led to good predictive performance. The weight means in each weight matrix were initialised by sampling from $\mathcal{N}(0, 1/\sqrt{2n_{\text{out}}})$, where n_{out} is the number of outputs of the weight matrix. The weight standard deviations were all initialised to a very small value of 1×10^{-5} , (we tried a larger initialization with weight standard deviations initialized to $1 \times 10^{-2.5}$ and found no significant difference). Bias means were initialised to zero, with the variances initialised to the same small value as the weight variance. 100,000 iterations of full batch training on the dataset were performed using ADAM with a learning rate of 1×10^{-3} . The ELBO was estimated using 32 Monte Carlo samples during training. The local reparameterisation trick was used [24]. The predictive distribution at test time was estimated using 500 samples from the approximate posterior.

Fitting MCDO: The weights and biases were initialised using the default `torch.linear` initialization. The dropout rate was fixed at $p = 0.05$. The ℓ^2 regularisation parameter was set following Gal [13, Section 3.2.3] for the given prior, in such a way that the ‘KL condition’ is met, in the interpretation of dropout training as variational inference. 100,000 iterations of full batch training on the dataset were performed using ADAM with a learning rate of 1×10^{-3} . The dropout objective was estimated using 32 Monte Carlo samples during training. The predictive distribution at test time was estimated using 500 samples from the approximate posterior.

Fitting HMC: For HMC on the 1HL BNN, 250,000 samples of HMC were taken using the NUTS implementation in Pyro [19, 4] after 10,000 warmup steps. For the 2HL case, 1,000,000 samples of HMC were taken after 20,000 warmup steps. We set the maximum tree depth in NUTS to 5, and adapt the step size and mass matrix during warmup.

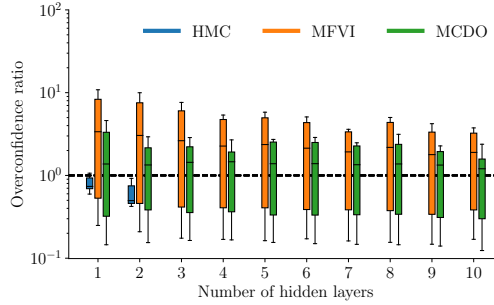


Figure 10: Boxplots of the overconfidence ratios of HMC, MFVI and MCDO relative to exact inference in an infinite width BNN (GP) with $\sigma_w = \sqrt{2}$.

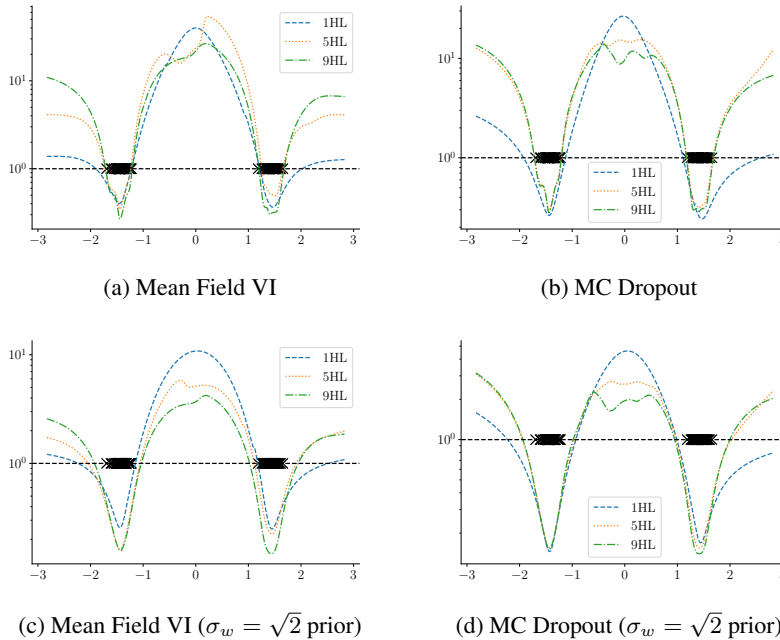


Figure 11: Plots of the overconfidence ratio against λ (where λ is defined as in figure 3) for several depths of neural networks with $\sigma_w = 4, 2, 1.7$ for 1, 5 and 9 hidden layers respectively (top), $\sigma_w = \sqrt{2}$ for all depths (bottom). Projections of the datapoints onto the diagonal slice between the clusters are shown as black crosses (X). We see that both MCDO and MFVI are overconfident (> 1) in between data, and underconfident (< 1) at the locations where we have observed data, relative to the GP reference.

F.2 Additional Figures

In order to assess the robustness of our findings to different choices of prior, we also consider the same experiment with $\sigma_w = \sqrt{2}$ for all depths. We choose this prior as it leads to similar variances in function space as depth increases [39]. We note that the variance of this prior is significantly smaller than the variance of the prior in the previous setting. The corresponding box plot is shown in figure 10. With this prior both methods tend to be less over-confident between data clusters, but more underconfident at the data, especially as depth increases (see figure 11).

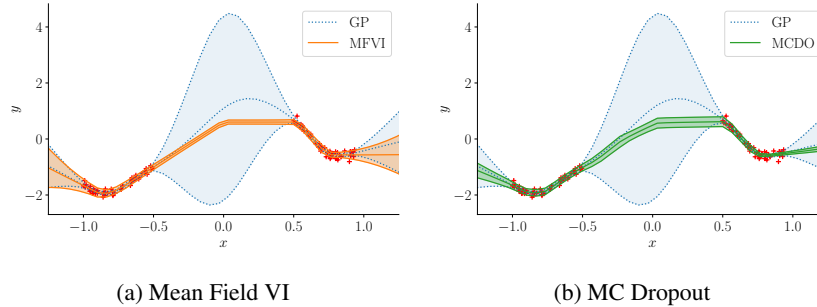


Figure 12: Mean and error bars (± 2 standard deviations) for the GP and the BNN with each inference scheme, trained on the data shown by the red crosses. The inference algorithms were initialised by first minimising the squared error to the reference GP mean and variance, and then running the respective inference algorithm.

G Initialisation of VI

In order to assess whether the variational objective (ELBO) or optimisation is primarily responsible for the lack of in-between uncertainty when performing MFVI and MCDO, we considered the effect of initialisation on the quality of the posterior obtained after variational inference. In order to find setting of the weights so that the posterior distribution in function space was close to the exact posterior in function space, we initialised the weights of the network by training the network using mean squared loss between the mean and variance functions of a reference posterior and the approximate posterior (as in figure 4). The reference posterior was obtained by fitting the limiting GP on the dataset (shown in crosses). We used these weights as an initialisation for variational inference. The noise variance was fixed to the true noise variance that generated the data. The data itself was sampled from the limiting GP prior, so that the model should be able to fit the data well.

Two-hidden layer MFVI and MCDO networks were used, with 50 hidden units in both layers. The solution found by minimising squared loss for 50,000 iterations between the mean and variance functions and a reference posterior may lead to distributions over weights such that the KL from these distributions to the prior is high. This can lead to very high values of the variational objective function. To alleviate this behaviour, we gradually interpolate between the squared-error loss and the variational objective, by taking convex combinations of the losses. Call the function space squared loss L_1 and the standard variational objective L_2 . Then after the first 50,000 iterations of using L_1 , we train for 10,000 iterations using $.9L_1 + .1L_2$, 10,000 iterations using $.8L_1 + .2L_2$ and so on until we are only training using L_2 . We then train for 100,000 iterations using just L_2 , to ensure the variational objective has converged. Figure 12 shows that the obtained posterior still lacks in-between uncertainty, providing evidence that this may be due to the objective function rather than overfitting.

H Details and Additional Plots for Active Learning

H.1 Experimental Setup

We use the same initialisation as in appendix F. As the dataset has low noise, we use a homoskedastic Gaussian noise model with a fixed standard deviation of 0.01 for all models. We used the ADAM optimiser with learning rate 1×10^{-3} for 20,000 epochs to optimise both MFVI and MCDO. We perform full batch training. All BNNs are retrained from scratch after the acquisition of each point from the pool set. We used 32 Monte Carlo samples from q_ϕ to estimate the objective function for both MFVI and MCDO. All networks had 50 neurons in each hidden layer. The prior for all BNNs and the GP was chosen to have $\sigma_w = \sqrt{2}$, $\sigma_b = 1$ (see appendix F for definitions). $\sigma_w = \sqrt{2}$ was chosen so that the prior in function space has a stable variance as depth increases [39]. The dropout probability was set at $p = 0.05$ for all MCDO networks. The dropout ℓ_2 regularisation was chosen to match the ‘KL condition’ as stated in Gal [13, Section 3.2.3]. The results are averaged over 20 random initialisations and selections of the 5 initial points in the active set. For MFVI and MCDO, the predictive distribution at test time and the predictive variances used for active learning were

estimated using 500 samples from the approximate posterior. The parameter initialisations are the same as those in appendix F.

H.2 Additional Figures

Figure 13 shows the points chosen by deep BNNs. Again the GP chooses points from every cluster, and seems to focus on the ‘corners’ of each cluster. MFVI samples from more clusters than the 1HL case, but still comparatively oversamples clusters further from the origin, and undersamples those near the origin. MCDO has a more spread out choice of points than the 1HL case, but still fails to obtain significantly better RMSE than random.

Figures 14 and 15 show the predictive uncertainty of 1HL models at the beginning and end of active learning. All models significantly reduce their uncertainty around clusters that have been heavily sampled, except for MCDO. This causes MCDO to repeatedly sample near locations that have already been labelled, in contrast to the GP. Note also that MFVI is most confident at clusters near the origin that have never been sampled, and least confident at clusters far from the origin that have already been heavily sampled.

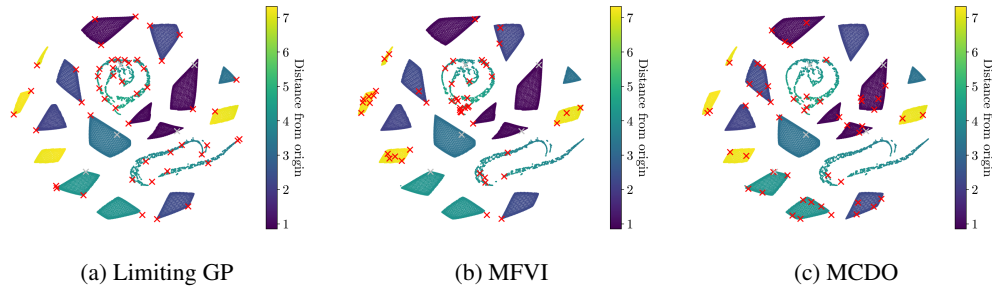


Figure 13: Points chosen during active learning in the 3HL case. Colours denote distance from the origin in 14-dimensional input space. Grey crosses (×) denote the five points randomly chosen as an initial training set. Red crosses (×) denote the 50 points selected by active learning. Again, the GP samples the corners of each cluster, and MFVI oversamples clusters far from the origin.

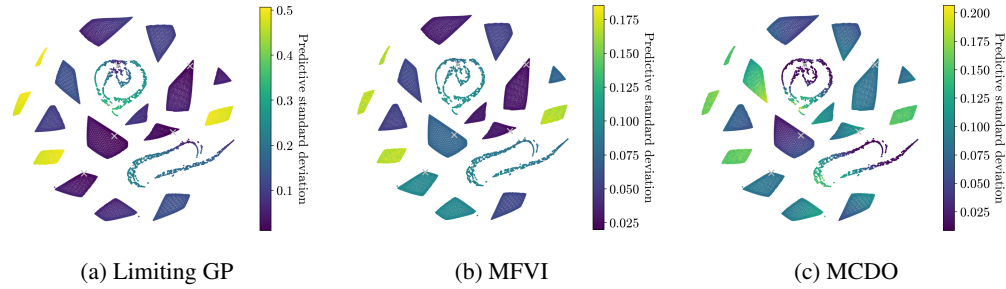


Figure 14: Colours denote *predictive uncertainties* in the 1HL case, at the beginning of active learning. As the noise standard deviation was fixed to 0.01 for all models, changes in the predictive standard deviation reflect model uncertainty. Grey crosses (×) denote the five points randomly chosen as an initial training set.

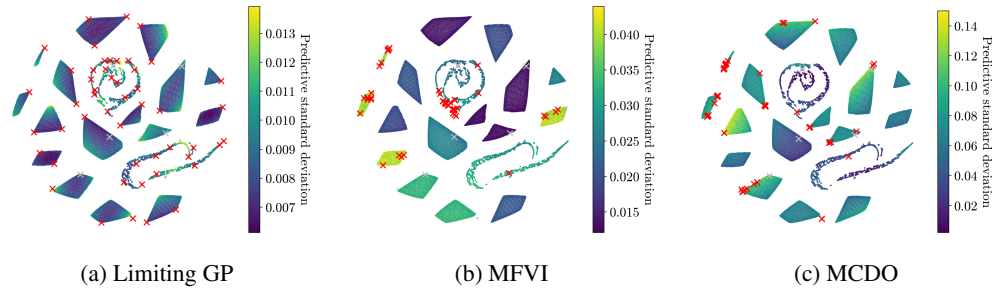


Figure 15: Colours denote *predictive uncertainties* in the 1HL case, after 50 iterations of active learning. As the noise standard deviation was fixed to 0.01 for all models, changes in the predictive standard deviation reflect model uncertainty. Grey crosses (×) denote the five points randomly chosen as an initial training set. Red crosses (×) denote the 50 points selected by active learning. Note how, compared to figure 14, the GP has reduced its uncertainty near points it has observed, and is most uncertain on corners opposite those points. In contrast, for both MFVI and MCDO, the network is still uncertain around regions it has already collected points from, leading it to oversample those clusters and undersample others.