

---

# Uncertainty in Neural Networks: Approximately Bayesian Ensembling

---

Tim Pearce<sup>1</sup>, Felix Leibfried<sup>2</sup>, Alexandra Brintrup<sup>1</sup>, Mohamed Zaki<sup>1</sup>, Andy Neely<sup>1</sup>  
<sup>1</sup>University of Cambridge, <sup>2</sup>PROWLER.io

## Abstract

Understanding the uncertainty of a neural network’s (NN) predictions is essential for many purposes. The Bayesian framework provides a principled approach to this, however applying it to NNs is challenging due to large numbers of parameters and data. Ensembling NNs provides an easily implementable, scalable method for uncertainty quantification, however, it has been criticised for not being Bayesian. This work proposes one modification to the usual process that we argue does result in approximate Bayesian inference; regularising parameters about values drawn from a distribution which can be set equal to the prior. A theoretical analysis of the procedure in a simplified setting suggests the recovered posterior is centred correctly but tends to have underestimated marginal variance, and overestimated correlation. However, two conditions can lead to exact recovery. We argue that these conditions are partially present in NNs. Empirical evaluations demonstrate it has an advantage over standard ensembling, and is competitive with variational methods.

Interactive demo: [teapearce.github.io](https://teapearce.github.io).

to guard against adversarial examples (Smith and Gal, 2018; Sünderhauf et al., 2018)

A principled approach to modelling uncertainty is provided by the Bayesian framework. Bayesian Neural Networks (BNNs) model the parameters of a NN as probability distributions computed via Bayes rule (MacKay, 1992). Whilst appealing, the large number of parameters and data points used with modern NNs renders many Bayesian inference techniques that work well in small-scale settings infeasible, e.g. MCMC methods.

Ensembling provides an alternative way to estimate uncertainty: it aggregates the estimates of multiple individual NNs, trained from different initialisations and sometimes on noisy versions of the training data. The variance of the ensemble’s predictions may be interpreted as its uncertainty. The intuition is simple: predictions converge to similar results where data has been observed, and will be diverse elsewhere. The chief attraction is that the method scales well to large parameter and data settings, with each individual NN implemented in precisely the usual way.

Whilst ensembling has proven empirically successful (Tibshirani, 1996; Lakshminarayanan et al., 2017; Osband et al., 2016), the absence of connection to Bayesian methodology has drawn critics and inhibited uptake, e.g. Gal (2016) [p. 27].

## 1 Introduction

Neural networks (NNs) are the current dominant force within machine learning, however, quantifying the uncertainty of their predictions is a challenge. This is important for many real-world applications (Bishop, 1994) as well as in auxiliary ways; to drive exploration in reinforcement learning (RL), for active learning, and

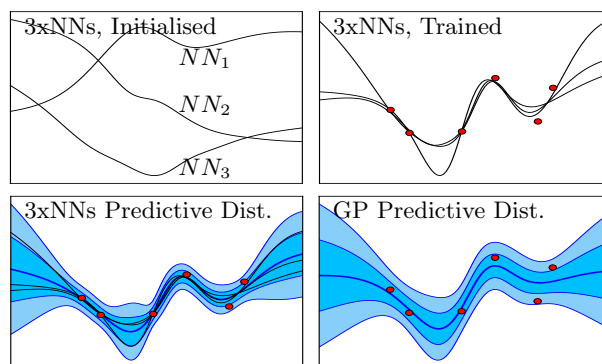


Figure 1: An ensemble of NNs, starting from different initialisations and trained with the proposed modification, produce a predictive distribution approximating that of a GP. This improves with number of NNs.

### 1.1 Contribution

This paper proposes, analyses and tests one modification to the usual NN ensembling process, with the purpose of examining how closely the resulting procedure aligns with Bayesian inference. The modification regularises parameters about values drawn from an anchor distribution, which can be set to be equal to the prior distribution. We name this procedure *anchored ensembling* - see figure 1 for an illustration. This falls into a family of little known Bayesian inference methods, *randomised MAP sampling* (RMS) (section 2.2).

Our first contributions do not specifically consider NNs; we derive an abstracted version of RMS in parameter space rather than output space. (This abstraction later allows us to propose RMS for classification tasks for the first time.) Under the assumption that the joint parameter likelihood and prior obey a multivariate normal distribution, we show that it is always possible to design an RMS procedure to recover the true posterior.

This design requires knowing the parameter likelihood covariance a priori, which is infeasible except in the simplest models. We propose a workaround that results in an approximation of the posterior. In general this approximation has correct mean but underestimated variance and overestimated correlation. However, two conditions lead to an exact recovery: 1) perfectly correlated parameters, 2) parameters whose marginal likelihood variance is infinite ('extrapolation parameters').

We proceed by considering the applicability of RMS to NNs. We discuss the appropriateness of assumptions used in the theoretical analysis, and argue that the two conditions leading to exact recovery of the posterior are partially present in NNs. We postulate this as the reason that predictive posteriors produced by anchored ensembling appear very similar to those by exact Bayesian methods in figures 4, 6, 7 & 8.

The performance of anchored ensembling is assessed experimentally on regression, image classification, sentiment analysis and RL tasks. It provides an advantage over standard ensembling procedures, and is competitive with variational methods.

## 2 Background

### 2.1 Bayesian Neural Networks

A variety of methods have been developed to perform Bayesian inference in NNs. Variational inference (VI) has received much attention both explicitly modelling parameters with distributions (Graves, 2011; Hernández-Lobato and Adams, 2015) and also implicitly through noisy optimisation procedures - MC Dropout (Gal and Ghahramani, 2015), Vadam (Khan

et al., 2018). Correlations between parameters are often ignored - mean-field VI (MFVI).

Other inference methods include: Hamiltonian Monte Carlo (HMC), a MCMC variant which provides 'gold standard' inference but at limited scalability (Neal, 1997); The Laplace method fits a multivariate normal distribution to the posterior (Ritter et al., 2018). Whilst ensembling is generally seen as a non-Bayesian alternative, Duvenaud et al. (2016) interpreted it, with early stopping, as approximate inference. Aside from *doing* Bayesian inference, recent works have begun exploring prior design in BNNs, e.g. Pearce et al. (2019).

BNNs of infinite width converge to GPs (Neal, 1997). Analytical kernels exist for NNs with certain activation functions, including sigmoidal (Error Function, ERF) (Williams, 1996), Rectified Linear Unit (ReLU) (Cho and Saul, 2009), and leaky ReLU (Tsuchida et al., 2018). Whilst GPs scale superlinearly with data (though see (Wang et al., 2019)), they provide a convenient method for doing exact inference on small problems. In this paper we use these GPs as 'ground truth' predictive distributions to compare to wide NNs. In section 5, we benchmark the ReLU GP on UCI datasets.

### 2.2 Randomised MAP Sampling

Recent work in the Bayesian community, and independently in the RL community, has begun to explore a novel approach to Bayesian inference. Roughly speaking, it exploits the fact that adding a regularisation term to a loss function returns a maximum a posteriori (MAP) parameter estimate - a point estimate of the Bayesian posterior. Injecting noise into this loss, either to targets or regularisation term, and sampling repeatedly (i.e. ensembling), produces a *distribution* of MAP solutions mimicking that of the true posterior. This can be an efficient method to sample from high-dimensional posteriors (Gu et al., 2007; Chen and Oliver, 2012; Bardsley et al., 2014).

Whilst it is possible to specify a noise injection that produces exact inference in linear regression, there is difficulty in transferring this idea to more complex settings, such as NNs or classification. Directly applying the noise distribution from linear regression to NNs has had some empirical success, despite not reproducing the true posterior (Lu and Van Roy, 2017; Osband et al., 2018) (section 3.2). A more accurate, though more computationally demanding solution, is to wrap the optimisation step in an MCMC procedure (Bardsley, 2012; Bardsley et al., 2014).

These works have been proposed under several names including randomise-then-optimize, randomised prior functions, and ensemble sampling. We refer to this family of procedures *randomised MAP sampling* (RMS).

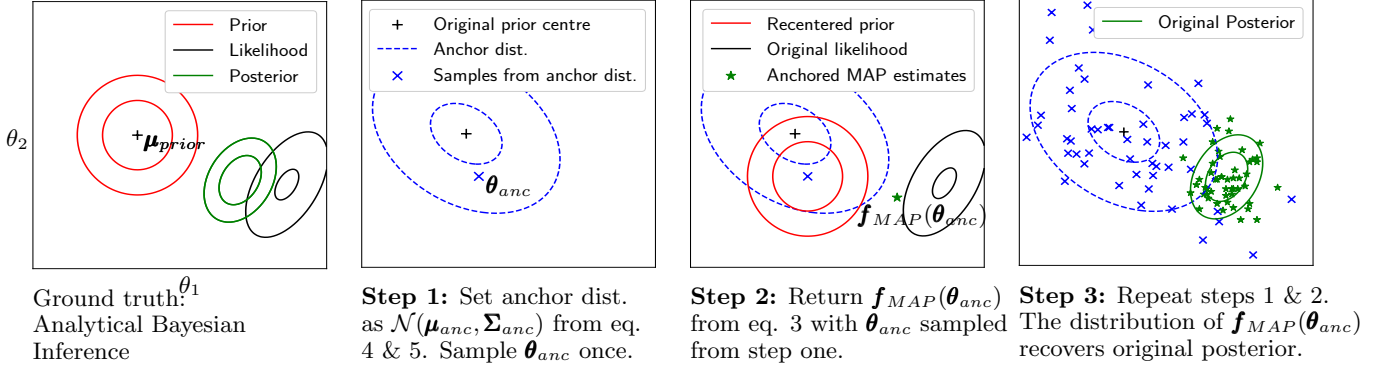


Figure 2: Demonstration of (exact) RMS in a 2D parameter space.

### 3 RMS Theoretical Results

This section presents several novel results. We first derive a general form of RMS by analysing the procedure in parameter space, using the simplifying assumption that both prior and parameter likelihood are multivariate normal distributions. This is an abstraction compared to previous works. Appendix A contains definitions and proofs in full.

If the parameter likelihood covariance is known a priori, we show how RMS can be designed to recover the true posterior. In general, this will not be known, and we propose a practical workaround requiring knowledge only of the prior distribution.

This workaround no longer guarantees exact recovery of the posterior. We derive results specifying in what ways the estimated RMS posterior is in general biased, including underestimated marginal variance, and overestimated correlation coefficient. We discover two special conditions that lead to an exact recovery.

The appropriateness of the normal assumption in non-linear models for general data likelihoods will be discussed in section 4, when we consider applying this RMS scheme with workaround to NNs.

#### 3.1 Parameter-Space Derivation

Consider multivariate normal prior and parameter likelihood distributions,  $P(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_{\text{prior}}, \boldsymbol{\Sigma}_{\text{prior}})$ ,  $P_{\boldsymbol{\theta}}(\mathcal{D}|\boldsymbol{\theta}) \propto \mathcal{N}(\boldsymbol{\mu}_{\text{like}}, \boldsymbol{\Sigma}_{\text{like}})$ . We make a distinction between two forms of likelihood: *data likelihood*, which is defined on the domain of the target variable, and *parameter likelihood*, which is specified in parameter space. (See definition 1.)

From Bayes rule the posterior, also normal, is,

$$\mathcal{N}(\boldsymbol{\mu}_{\text{post}}, \boldsymbol{\Sigma}_{\text{post}}) \propto \mathcal{N}(\boldsymbol{\mu}_{\text{prior}}, \boldsymbol{\Sigma}_{\text{prior}}) \mathcal{N}(\boldsymbol{\mu}_{\text{like}}, \boldsymbol{\Sigma}_{\text{like}}) \quad (1)$$

The MAP solution is simply  $\boldsymbol{\theta}_{\text{MAP}} = \boldsymbol{\mu}_{\text{post}}$ ,

$$\boldsymbol{\theta}_{\text{MAP}} = \boldsymbol{\Sigma}_{\text{post}} \boldsymbol{\Sigma}_{\text{like}}^{-1} \boldsymbol{\mu}_{\text{like}} + \boldsymbol{\Sigma}_{\text{post}} \boldsymbol{\Sigma}_{\text{prior}}^{-1} \boldsymbol{\mu}_{\text{prior}}, \quad (2)$$

where  $\boldsymbol{\Sigma}_{\text{post}} = (\boldsymbol{\Sigma}_{\text{like}}^{-1} + \boldsymbol{\Sigma}_{\text{prior}}^{-1})^{-1}$ . In RMS we assume availability of a mechanism for returning  $\boldsymbol{\theta}_{\text{MAP}}$ , and are interested in injecting noise into eq. 2 so that a *distribution* of  $\boldsymbol{\theta}_{\text{MAP}}$  solutions are produced, matching the true posterior distribution.

A practical choice of noise source is the mean of the prior,  $\boldsymbol{\mu}_{\text{prior}}$ , since a modeller has full control over this value. Let us replace  $\boldsymbol{\mu}_{\text{prior}}$  with some noisy random variable,  $\boldsymbol{\theta}_{\text{anc}}$ . This is the same place as a hyperprior over  $\boldsymbol{\mu}_{\text{prior}}$ , though with a subtly different role. Denote  $\mathbf{f}_{\text{MAP}}(\boldsymbol{\theta}_{\text{anc}})$  a function that takes as input  $\boldsymbol{\theta}_{\text{anc}}$  and returns the resulting MAP estimate,

$$\mathbf{f}_{\text{MAP}}(\boldsymbol{\theta}_{\text{anc}}) = \boldsymbol{\Sigma}_{\text{post}} \boldsymbol{\Sigma}_{\text{like}}^{-1} \boldsymbol{\mu}_{\text{like}} + \boldsymbol{\Sigma}_{\text{post}} \boldsymbol{\Sigma}_{\text{prior}}^{-1} \boldsymbol{\theta}_{\text{anc}}. \quad (3)$$

Accuracy of this procedure hinges on selection of an appropriate distribution for  $\boldsymbol{\theta}_{\text{anc}}$ , which we term the *anchor distribution*. The distribution that will produce the true posterior can be found by setting  $\mathbb{E}[\mathbf{f}_{\text{MAP}}(\boldsymbol{\theta}_{\text{anc}})] = \boldsymbol{\mu}_{\text{post}}$  and  $\text{Var}[\mathbf{f}_{\text{MAP}}(\boldsymbol{\theta}_{\text{anc}})] = \boldsymbol{\Sigma}_{\text{post}}$ .

**Theorem 1.** *In order that,  $P(\mathbf{f}_{\text{MAP}}(\boldsymbol{\theta}_{\text{anc}})) = P(\boldsymbol{\theta}|\mathcal{D})$ , the required distribution of  $\boldsymbol{\theta}_{\text{anc}}$  is also multivariate normal,  $P(\boldsymbol{\theta}_{\text{anc}}) = \mathcal{N}(\boldsymbol{\mu}_{\text{anc}}, \boldsymbol{\Sigma}_{\text{anc}})$ , where,*

$$\boldsymbol{\mu}_{\text{anc}} = \boldsymbol{\mu}_{\text{prior}} \quad (4)$$

$$\boldsymbol{\Sigma}_{\text{anc}} = \boldsymbol{\Sigma}_{\text{prior}} + \boldsymbol{\Sigma}_{\text{prior}} \boldsymbol{\Sigma}_{\text{like}}^{-1} \boldsymbol{\Sigma}_{\text{prior}}. \quad (5)$$

Figure 2 provides a demonstration of the RMS algorithm in 2D parameter space.

#### 3.2 Comparison to Prior Work

Previous work on RMS (Lu and Van Roy, 2017; Osband et al., 2019) was motivated via linear regression. Noting that the MAP solution is given by,

$$\boldsymbol{\theta}_{\text{MAP}} = \left( \frac{1}{\sigma_{\epsilon}^2} \mathbf{X}^T \mathbf{X} + \boldsymbol{\Sigma}_{\text{prior}}^{-1} \right)^{-1} \left( \frac{1}{\sigma_{\epsilon}^2} \mathbf{X}^T \mathbf{y} + \boldsymbol{\Sigma}_{\text{prior}}^{-1} \boldsymbol{\mu}_{\text{prior}} \right), \quad (6)$$

these works added Gaussian noise to  $\boldsymbol{\mu}_{\text{prior}}$ , *in addition* to adding noise to  $\mathbf{y}$ , either by additive Gaussian noise or bootstrapping. Eq. 6 is a special case of our own derivation, substituting  $\boldsymbol{\Sigma}_{\text{like}}^{-1} = 1/\sigma_{\epsilon}^2 \mathbf{X}^T \mathbf{X}$  into eq. 2.

### 3.3 Practical Workaround: General Case

The previous section showed how to design an RMS procedure that will precisely recover the true Bayesian posterior. Unfortunately, in eq. 5 one must specify the parameter likelihood covariance in order to set the anchor distribution. For most models, this is infeasible.

A practical workaround is to simply ignore the second term in eq. 5 and set  $\Sigma_{anc} := \Sigma_{prior}$ . Using RMS with this anchor distribution will not generally lead to an exact recovery of the true posterior, however the resulting distribution can be considered an approximation of it. Corollary 1.1 derives this RMS approximate posterior in terms of the true posterior.

**Corollary 1.1.** *Set  $\mu_{anc} := \mu_{prior}$  and  $\Sigma_{anc} := \Sigma_{prior}$ . The RMS approximate posterior is  $P(\mathbf{f}_{MAP}(\theta_{anc})) = \mathcal{N}(\mu_{post}, \Sigma_{post} \Sigma_{prior}^{-1} \Sigma_{post})$ .*

*Proof sketch.* This follows similar working to theorem 1, but instead of enforcing  $\mathbb{E}[\mathbf{f}_{MAP}(\theta_{anc})] = \mu_{post}$ ,  $\text{Var}[\mathbf{f}_{MAP}(\theta_{anc})] = \Sigma_{post}$  and solving for  $\mu_{anc}, \Sigma_{anc}$ , we now enforce  $\mu_{anc} := \mu_{prior}$ ,  $\Sigma_{anc} := \Sigma_{prior}$  and solve for  $\mathbb{E}[\mathbf{f}_{MAP}(\theta_{anc})]$ ,  $\text{Var}[\mathbf{f}_{MAP}(\theta_{anc})]$ .

This corollary shows that the means of the two distributions are aligned, although the covariances are not. Next we state several results quantifying how the RMS approximate posterior covariance differs compared to the true posterior covariance. All results assume multivariate normal prior and parameter likelihood. They can be observed in figure 3 (A).

**Lemma 1.1.** *When  $\mu_{anc} := \mu_{prior}$ ,  $\Sigma_{anc} := \Sigma_{prior}$  the RMS approximate posterior will in general underestimate the marginal variance compared to the true posterior,  $\text{Var}[\mathbf{f}_{MAP}(\theta_{anc})] < \text{Var}[\theta|\mathcal{D}]$ .*

*Proof sketch.*  $\Sigma_{post} \Sigma_{prior}^{-1} \Sigma_{post}$  can be rearranged as  $\Sigma_{post} - \Sigma_{post} \Sigma_{like}^{-1} \Sigma_{post}$ . The second term will be pos-

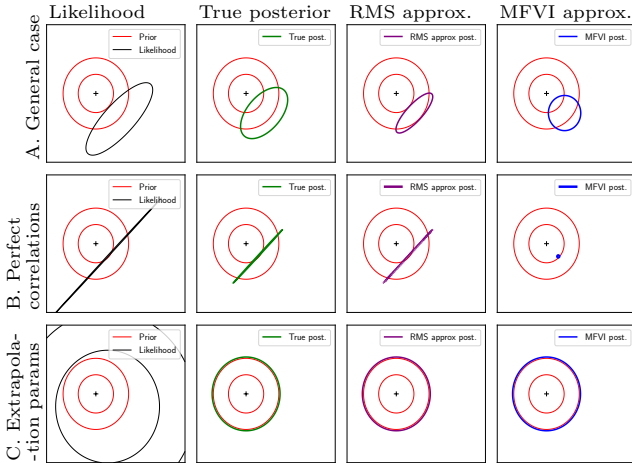


Figure 3: Examples of the RMS approximate posterior when  $\mu_{anc} := \mu_{prior}$ ,  $\Sigma_{anc} := \Sigma_{prior}$ . MFVI also shown.

itive definite, so the diagonal entry is positive, and hence,  $\text{diag}(\Sigma_{post} - \Sigma_{post} \Sigma_{like}^{-1} \Sigma_{post})_i < \text{diag}(\Sigma_{post})_i$ .

**Lemma 1.2.** *Additionally assume the prior is isotropic. When  $\mu_{anc} := \mu_{prior}$ ,  $\Sigma_{anc} := \Sigma_{prior}$  the eigenvectors (or ‘orientation’) of the RMS approximate posterior equal those of the true posterior.*

*Proof sketch.*  $\Sigma_{post} \Sigma_{prior}^{-1} \Sigma_{post} = 1/\sigma_{prior}^2 \Sigma_{post}^2$ . Squaring a matrix only modifies eigenvalues not eigenvectors. As does multiplying by a constant.

**Theorem 2.** *Additionally assume the prior is isotropic. For a two parameter model, when  $\mu_{anc} := \mu_{prior}$ ,  $\Sigma_{anc} := \Sigma_{prior}$ , the RMS approximate posterior will in general overestimate the magnitude of the true posterior parameter correlation coefficient,  $|\rho|$ . If  $|\rho| = 1$ , then it will recover it precisely.*

*Proof sketch.* We compute the individual entries resulting from the required  $2 \times 2$  matrix multiplications.

We were unable to generalise theorem 2 beyond a two parameter model, but numerical examples (appendix B.1) suggest that it holds for higher dimensionality.

### 3.4 Practical Workaround: Special Cases

Having described the covariance bias that in general will be present in the RMS approximate posterior, we now give two special conditions under which there is no bias, and the true posterior is exactly recovered. Illustrations of these cases are shown in figure 3 (B, C).

**Theorem 3.** *For extrapolation parameters (def. 2 - parameters which do not affect data likelihood but may affect new predictions) of a model, setting  $\mu_{anc} := \mu_{prior}$ ,  $\Sigma_{anc} := \Sigma_{prior}$ , means the marginal RMS approximate posterior equals that of the marginal true posterior.*

*Proof sketch.* We show that the required matrix multiplications,  $\Sigma_{post} \Sigma_{prior}^{-1} \Sigma_{post}$ , do not affect rows corresponding to extrapolation parameters.

**Theorem 4.** *Set  $\mu_{anc} := \mu_{prior}$ ,  $\Sigma_{anc} := \Sigma_{prior}$ . The RMS approximate posterior will exactly equal the true posterior,  $\Sigma_{post}$ , when all eigenvalues of a scaled version of  $\Sigma_{post}$  (scaled such that the prior equals the identity matrix) are equal to either 0 or 1. This corresponds to posteriors that are a mixture of perfectly correlated and perfectly uncorrelated parameters.*

*Proof sketch.* We are searching for solutions to  $\Sigma_{post} = \Sigma_{post} \Sigma_{prior}^{-1} \Sigma_{post}$ . Applying a scaling,  $\Sigma'_{post} = \Sigma_{prior}^{-1/2} \Sigma_{post} \Sigma_{prior}^{-1/2}$ , results in a slightly simpler equation to find a solution to,  $\Sigma'_{post} = \Sigma_{post}^2$ . Results for idempotent matrices tell us that if  $\Sigma'_{post}$  is singular with all eigenvalues equal to 0 or 1, this will be a solution.

To provide intuition behind theorem 4 consider a two

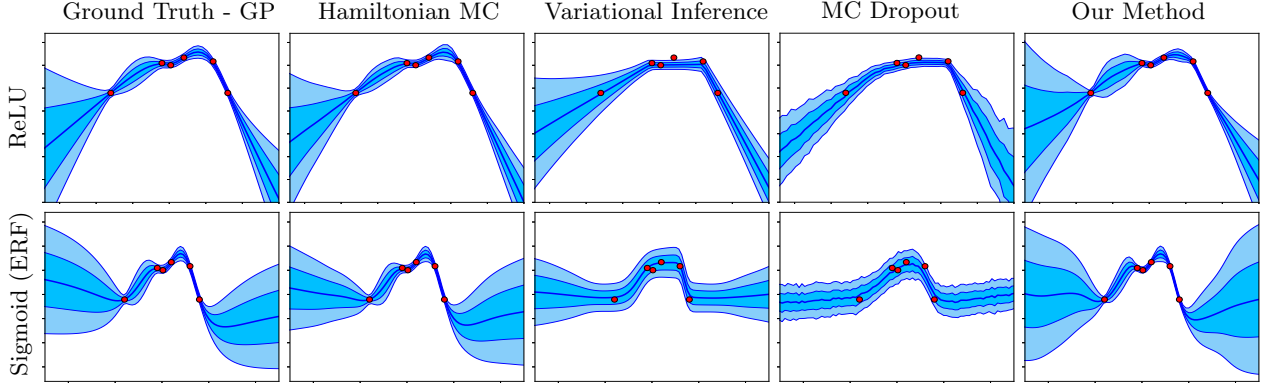


Figure 4: Predictive distributions produced by various inference methods (columns) with varying activation functions (rows) in single-layer NNs on a toy regression task.

parameter model. If parameters are perfectly correlated, the effect on the data likelihood of an increase in the first can be exactly compensated for by a change in the second. If the region over which this applies is large relative to the prior, the likelihood is a line of negligible width. This leads to a posterior of negligible width spanning the prior. Examples in appendix B.1 show what combinations of parameters this holds for.

This section’s proofs show that if these two conditions exist, RMS makes a precise recovery. In practise, one would expect to see an increasingly accurate RMS approximation as these conditions are approached.

## 4 RMS for Neural Networks

We now apply RMS with practical workaround to NNs. We will refer to this as ‘anchored ensembling’.

First, we define the NN loss function to be optimised that corresponds to RMS. We then discuss the validity of the RMS procedure in the context of NNs, given the assumptions made. Finally we consider some matters arising in implementation of the scheme. Appendix, algorithm 1 details the full procedure.

### 4.1 Loss Function

Consider a NN containing parameters,  $\theta$ , making predictions,  $\hat{y}$ , with  $H$  hidden nodes and  $N$  data points. If the prior is given by  $P(\theta) = \mathcal{N}(\mu_{prior}, \Sigma_{prior})$ , maximising the following returns MAP parameter estimates. (See appendix A.1 for the standard derivation.)

$$\theta_{MAP} = \operatorname{argmax}_{\theta} \log(P_{\mathcal{D}}(\mathcal{D}|\theta)) - \frac{1}{2} \|\Sigma_{prior}^{-1/2} \cdot (\theta - \mu_{prior})\|_2^2 \quad (7)$$

When  $\mu_{prior} = \mathbf{0}$ , this is standard L2 regularisation. In order to apply RMS we instead replace  $\mu_{prior}$  with some random variable  $\theta_{anc}$ . To use the practical form of RMS, we will draw  $\theta_{anc} \sim \mathcal{N}(\mu_{prior}, \Sigma_{prior})$ .

Conveniently, no parametric form of data likelihood has yet been specified. For a regression task assuming homoskedastic Gaussian noise of variance  $\sigma_{\epsilon}^2$ , MAP estimates are found by minimising,

$$Loss_j = \frac{1}{N} \|\mathbf{y} - \hat{\mathbf{y}}_j\|_2^2 + \frac{1}{N} \|\mathbf{\Gamma}^{1/2} \cdot (\theta_j - \theta_{anc,j})\|_2^2. \quad (8)$$

We have defined a diagonal regularisation matrix,  $\mathbf{\Gamma}$ , where the  $i^{th}$  diagonal element is the ratio of data noise of the target variable to prior variance for parameter  $\theta_i$ ,  $\operatorname{diag}(\mathbf{\Gamma})_i = \sigma_{\epsilon}^2 / \sigma_{prior_i}^2$ . Note a subscript has been introduced,  $j \in \{1 \dots M\}$ , with the view of an ensemble of  $M$  NNs, each with a distinct draw of  $\theta_{anc}$ .

For classification tasks, cross entropy is normally maximised, which assumes a multinomial data likelihood,

$$= -\frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C y_{n,c} \log \hat{y}_{n,c,j} + \frac{1}{N} \|\mathbf{\Gamma}^{1/2} \cdot (\theta_j - \theta_{anc,j})\|_2^2, \quad (9)$$

where  $y_c$  is the label for class  $c \in \{1 \dots C\}$ . Here,  $\operatorname{diag}(\mathbf{\Gamma})_i = 1/2\sigma_{prior_i}^2$ .

### 4.2 Validity of RMS in NNs

Theory derived to motivate and analyse RMS assumed a simplified setting of multivariate normal parameter likelihoods. This section discusses this assumption, then considers the prevalence of special conditions (section 3.4) that would lead to a close approximation of the true posterior.

#### 4.2.1 Normal Distribution

Earlier proofs assumed parameter likelihoods follow a multivariate normal distribution. We provide two justifications for using this assumption in NNs.

1) Other approximate Bayesian methods incorporate similar assumptions into their methodologies. MFVI commonly fits a factorised normal distribution to the

posterior. The Laplace approximation fits a multivariate normal distribution to the mode of a MAP solution.

2) In figure 5 we visualise conditional parameter likelihoods for actual NNs trained on regression and classification tasks. After training, a parameter is randomly selected, and all others are frozen. The chosen parameter is varied over a small range and the data likelihood calculated at each point. Hence conditional distributions may be plotted. The plots suggest that thinking of local modes as approximately normally distributed is not unreasonable for the purpose of analysis.

This justifies modelling a single mode of the parameter likelihood as multivariate normal. However, the parameter space of a NN is likely to contain many such modes, with each member of an anchored ensemble ending up at a different one. We believe that many of these modes would be exchangeable, for example arising from parameter symmetries. In this case we believe that MAP solutions would also be exchangeable.

Empirically we did not observe this multimodality being problematic - plots such as figure 8 show predictive posteriors with low bias compared to the true posterior.

### 4.2.2 Presence of Special Cases

Setting the anchor distribution equal to the prior leads to an RMS approximate posterior that, in general, has underestimated variance and overestimated correlation.

Figures 4, 6 & 8 show predictive distributions for anchored ensembles that very closely approximate the true Bayesian posterior, with little sign of bias. This demands an answer to why, rather than if, anchored ensembling performs such accurate inference in these examples. We believe the reason is the presence of the two special conditions that can lead to exact recovery.

It should be straightforward to see that extrapolation parameters (definition 2) exist in the figures. Many hidden nodes will be dead across the range which contains data. Their corresponding final layer weight then has no effect on the data likelihood, but they do affect predictions outside of the training data.

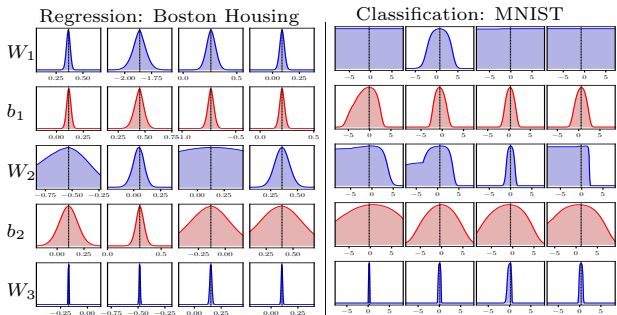


Figure 5: Empirical plots of conditional likelihoods for 4 randomly sampled parameters in two-layer NNs.

It is more difficult to see that perfect correlations also exist, and we provide a numerical example illustrating this in appendix B.3. Essentially it relies on two hidden nodes becoming live in between the same two data points. The associated final layer weights are then perfectly correlated. Whether these special conditions exist beyond fully-connected NNs is something tested indirectly in later experiments with CNNs.

One obvious way to further encourage these conditions is to increase the width of the NN, creating more parameters and an increasing probability of strong correlations. See also a study of multicollinearity in NNs (Cheng et al., 2018) [7.1].

### 4.3 Implementation Practicalities

*How many NNs to use in an RMS ensemble?* A large number of samples (and therefore NNs) would be required to fully capture the posterior parameter distributions. By contrast, if one thinks of each NN as an iid sample from a posterior *predictive* distribution, a much smaller number are required, given output dimensionality is typically small. Note this is unaffected by input dimension. Our experiments in section 5 used 5-10 NNs per ensemble, delivering good performance on tasks ranging from 1-10 outputs. See also figure 8. This results in anchored ensembles scaling by  $O(MN)$ .

*Should the NNs be initialised at anchor points?* It is convenient to draw parameter initialisations from the anchor distribution, and regularise directly around these initialised values, however, we found decoupling initialisations from anchor points benefited experiments.

## 5 Experiments

This section shares high-level findings from experiments. Further details and hyperparameter settings are given in appendix E. Appendix C additionally includes two RL experiments; one testing uncertainty-aware agents for model-free RL, and one applying anchored ensembles to noisy environments for model-based RL. Code is available online ([github/TeaPearce](https://github.com/TeaPearce)). Also see our [interactive demo](#).

### 5.1 Qualitative Tests

We first examine anchored ensembles on toy problems to gain intuition about its behaviour compared to popular approximate inference and ensembling methods.

Figure 4 compares popular Bayesian inference methods in single-layer NNs for ReLU and sigmoidal nonlinearities. GP and HMC produce ‘gold standard’ Bayesian inference, and we judge the remaining methods, which are scalable approximations, to them. Both

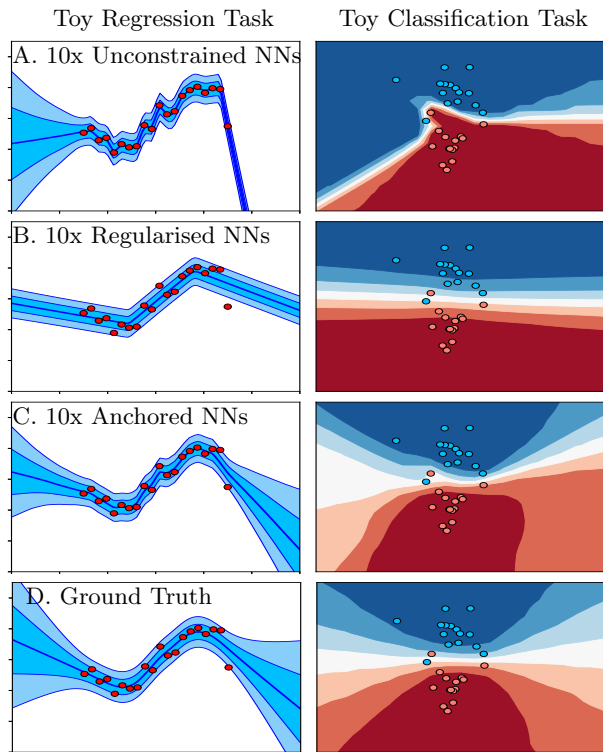


Figure 6: Comparison of NN ensemble loss choices.

MFVI (with a factorised normal distribution) and MC dropout do a poor job of capturing interpolated uncertainty. This is a symptom of the posterior approximation ignoring parameter correlations - see also figure 3 which shows MFVI failing to capture correlations in the posterior. This was explored in Foong et al. (2019).

Figure 6 contrasts anchored ensembles trained on eq. 8 & 9, with NN ensembles using standard loss functions, either with no regularisation term (‘unconstrained’,  $\Gamma = \mathbf{0}$ ), or regularised around zero (‘regularised’,  $\theta_{anc,j} = \mathbf{0}$ ). Regularised produces poor results since it encourages all NNs to the same single solution and diversity is reduced. Unconstrained is also inappropriate - although it produces diversity, no notion of prior is maintained

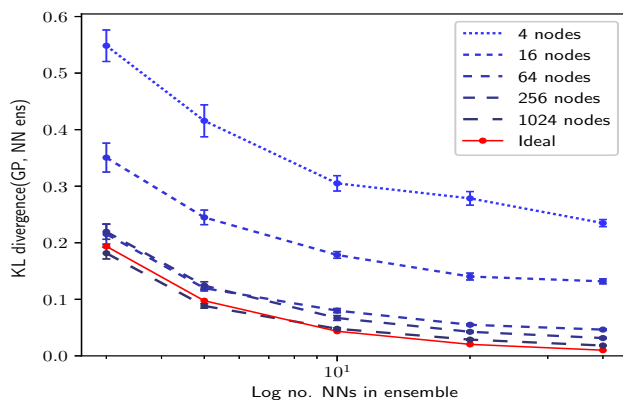


Figure 7: Difference in predictive distributions of an anchored ensemble and a ReLU GP as a function of width and number of NNs. Mean  $\pm 1$  standard error.

and it overfits the data.

Figure 8 shows the predictive distribution improving with number of NNs compared to a ReLU GP, however it appears a small residual difference remains.

## 5.2 Convergence Behaviour

To assess how precisely anchored ensembling performs Bayesian inference on a real dataset, we compared its predictive distribution with that of an exact method (ReLU GP) on the Boston housing dataset. Figure 7 quantifies the difference when varying both the width of the NN, and number of NNs in the ensemble. KL divergence between the two predictive distributions was measured and found to decrease as both NN width and number of NNs was increased. As in figure 8 a small amount of residual difference remains even for 40xNNs of 1,024 nodes.

## 5.3 UCI Regression Benchmarks

In order to compare anchored ensembles against popular approximate inference methods, we used a standard BNN benchmark. This assesses uncertainty quality for UCI regression tasks on data drawn from the same distribution as the training data (Hernández-Lobato and Adams, 2015). We also implemented the ReLU GP to assess the performance limit on these datasets.

Table 1 lists our results. We include results reported for Deep Ensembles (Lakshminarayanan et al., 2017), which is considered the state-of-the-art ensemble method. Appendix C.3 provides a full comparison with other approximate Bayesian methods including Probabilistic Backpropagation, MC Dropout, and Stochastic Gradient HMC.

Ordering results according to the level of estimated data noise,  $\hat{\sigma}_\epsilon^2$ , shows a clear pattern - anchored ensembles perform best in datasets with low data noise, surpassing both Deep Ensembles and all approximate inference methods listed in appendix C.3. This may be due to an increased importance of interpolation uncertainty when data noise is low, which anchored ensembles models well. On other datasets, the method is also competitive (the Deep Ensemble implementation used additional complexity to capture heteroskedastic uncertainty and has an advantage on higher data noise datasets).

## 5.4 Out-of-Distribution Classification

We now test on classification tasks, for out-of-distribution (OOD) data, with complex NN architectures, and compare against other ensemble methods.

An uncertainty-aware NN should make predictions of decreasing confidence as it is asked to predict on data

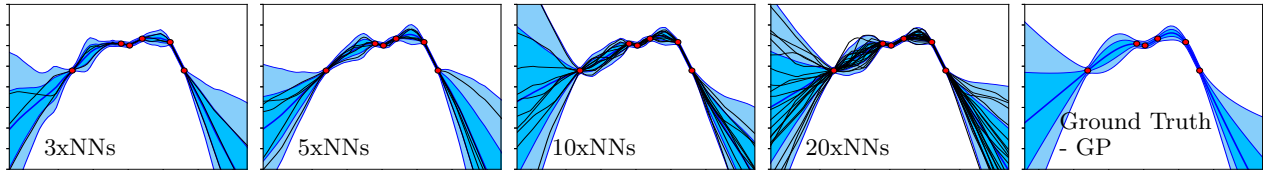


Figure 8: The predictive distribution of an anchored ensemble approaches that of a ReLU GP.

 Table 1: NLL regression benchmark results. See appendix C for RMSE and variants of our method. Mean  $\pm 1$  standard error.

	$\hat{\sigma}_\epsilon^2$	Deep Ens. <i>State-Of-Art</i>	Anch. Ens. <i>Our Method</i>	ReLU GP <sup>1</sup> <i>Gold Standard</i>
High Epistemic Uncertainty				
Energy	1e-7	1.38 $\pm$ 0.22	<b>0.96 <math>\pm</math> 0.13</b>	0.86 $\pm$ 0.02
Naval	1e-7	-5.63 $\pm$ 0.05	<b>-7.17 <math>\pm</math> 0.03</b>	-10.05 $\pm$ 0.02
Yacht	1e-7	1.18 $\pm$ 0.21	<b>0.37 <math>\pm</math> 0.08</b>	0.49 $\pm$ 0.07
Equal Epistemic & Aleatoric Uncertainty				
Kin8nm	0.02	-1.20 $\pm$ 0.02	<b>-1.09 <math>\pm</math> 0.01</b>	-1.22 $\pm$ 0.01
Power	0.05	<b>2.79 <math>\pm</math> 0.04</b>	2.83 $\pm$ 0.01	2.80 $\pm$ 0.01
Concrete	0.05	3.06 $\pm$ 0.18	<b>2.97 <math>\pm</math> 0.02</b>	2.96 $\pm$ 0.02
Boston	0.08	<b>2.41 <math>\pm</math> 0.25</b>	2.52 $\pm$ 0.05	2.45 $\pm$ 0.05
High Aleatoric Uncertainty				
Protein	0.5	<b>2.83 <math>\pm</math> 0.02</b>	2.89 $\pm$ 0.01	*2.88 $\pm$ 0.00
Wine	0.5	<b>0.94 <math>\pm</math> 0.12</b>	<b>0.95 <math>\pm</math> 0.01</b>	0.92 $\pm$ 0.01
Song	0.7	<b>3.35 <math>\pm</math> NA</b>	3.60 $\pm$ NA	**3.62 $\pm$ NA

<sup>1</sup> For comparison only (not a scalable method). \* Trained on 10,000 rows of data. \*\* Trained on 20,000 rows of data, tested on 5,000 data points.

further from the distribution seen during training. To test this, we report the proportion of high confidence predictions (defined as a softmax output class being  $\geq 90\%$ ) made by various ensemble systems - unconstrained, regularised, and anchored (as in section 5.1).

We trained on three different datasets, using a NN architecture appropriate to each: 1) Fashion MNIST image classification; 3 fully-connected layers of 100 hidden nodes. 2) IMDb movie review sentiment classification; embedding + 1D convolution + fully-connected layer. 3) CIFAR-10 image classification; convolutional NN (CNN) similar to VGG-13 (9 million parameters).

The confidence of predictions on novel data categories not seen during training was assessed. Table 2 shows example OOD images shown to the NNs trained on CIFAR-10. Edge refers to two CIFAR classes held out during training (ships, dogs). Appendix E provides OOD examples for other datasets.

The three tables show similar patterns. Whilst all methods predict with similar confidence on the training data, confidence differs greatly for other data categories, with anchored ensembles generally producing the most conservative predictions. This gap increases for data drawn further from the training distribution. Encouragingly, we observe similar (though less extreme) behaviour to that in the toy examples of figure 6.

Table 2: Proportion of predictions that were high confidence on out-of-distribution data, e.g. a single regularised NN trained on CIFAR-10 made high confidence predictions 54% of the time when asked to predict on MNIST. Mean over five runs (three for CIFAR).

CIFAR-10 Image Classification, VGG-13 CNN								
Train	— Edge —	Fashion	MNIST	Scramble	Invert	Noise		
Accuracy	Train	Edge	Fashion	MNIST	Scramble	Invert	Noise	
1xNNs Reg.	81.6%	0.671	0.466	0.440	0.540	0.459	0.324	0.948
5xNNs Uncons.	85.0%	0.607	0.330	0.208	0.275	0.175	0.209	0.380
5xNNs Reg.	86.1%	0.594	0.296	0.219	0.188	<b>0.106</b>	0.153	0.598
5xNNs Anch.	85.6%	0.567	<b>0.258</b>	<b>0.184</b>	<b>0.149</b>	0.134	<b>0.136</b>	<b>0.118</b>
10xNNs Anch.	86.0%	0.549	0.256	0.119	0.145	0.122	0.124	0.161

IMDb Text Sentiment Classification, Embedding+CNN

	Accuracy	Train	Reuters	Rand. 1	Rand. 2	Rand. 3
1xNNs Reg.	85.3%	0.637	0.119	0.153	0.211	0.326
5xNNs Uncons.	89.1%	0.670	0.102	0.141	0.100	0.075
5xNNs Reg.	87.1%	0.612	0.051	0.091	0.076	0.055
5xNNs Anch.	87.7%	0.603	<b>0.049</b>	<b>0.075</b>	<b>0.061</b>	<b>0.009</b>

Fashion MNIST Image Classification, Fully-Connected NN

	Accuracy	Train	Edge	CIFAR	MNIST	Distort	Noise
1xNN Reg.	86.8%	0.660	0.584	0.143	0.160	0.429	0.364
5xNNs Uncons.	89.0%	0.733	0.581	0.301	0.104	0.364	0.045
5xNNs Reg.	87.8%	0.634	<b>0.429</b>	0.115	0.072	0.342	0.143
5xNNs Anch.	88.0%	0.631	0.452	<b>0.065</b>	<b>0.041</b>	<b>0.246</b>	<b>0.006</b>

## 6 Conclusion

This paper proposed, analysed, and tested a modification to the usual NN ensembling process that results in approximate Bayesian inference - regularising parameters around values drawn from a prior distribution.

Under simplifying assumptions, we derived an abstracted form of RMS motivating this. We analysed a practical RMS variant to understand the bias of its approximate posterior. Two special conditions were shown to lead to recovery of the true posterior: perfectly correlated parameters and extrapolation parameters. We discussed the validity of applying RMS to NNs, arguing that these two special conditions are partially present in NNs.

On regression benchmarking experiments, state-of-the-art performance was achieved on 3/10 datasets - outperforming popular approximate inference methods. On image and text classification tasks, anchored ensembles were shown to be more robust than alternative ensemble methods.



## Acknowledgements

Thanks to all anonymous reviewers for their helpful comments and suggestions. The lead author was funded through EPSRC (EP/N509620/1) and partially accommodated by the Alan Turing Institute. The model-based RL experiments were run during an internship at PROWLER.io. Thanks to Nicolas Anastassacos for collaborating on an early version of the paper, and Ayman Boustati and Ahmed Al-Ali for helpful discussions.

## References

- Bardsley, J. M. (2012). MCMC-based image reconstruction with uncertainty quantification. *SIAM Journal on Scientific Computing*, 34(3):1316–1332.
- Bardsley, J. M., Solonen, A., Haario, H., and Laine, M. (2014). Randomize-Then-Optimize: A Method for Sampling from Posterior Distributions in Nonlinear Inverse Problems. *SIAM Journal on Scientific Computing*, 36(4).
- Bishop, C. (1994). Novelty detection and neural network validation. *IEEE Proceedings - Vision, Image, and Signal Processing*, 141(4):217.
- Chen, Y. and Oliver, D. S. (2012). Ensemble Randomized Maximum Likelihood Method as an Iterative Ensemble Smoother. *International Association for Mathematical Geosciences*, (D):1–26.
- Cheng, X., Khomtchouk, B., Matloff, N., and Mohanty, P. (2018). Polynomial Regression As an Alternative to Neural Nets. In *arXiv:1806.06850*.
- Cho, Y. and Saul, L. K. (2009). Kernel Methods for Deep Learning. In *Advances in Neural Information Processing Systems 22*.
- Chua, K., Calandra, R., McAllister, R., and Levine, S. (2018). Deep Reinforcement Learning in a Handful of Trials using Probabilistic Dynamics Models. In *NeurIPS*.
- Dearden, R., Friedman, N., and Russell, S. (1998). Bayesian Q-learning. In *American Association for Artificial Intelligence (AAAI)*.
- Duvenaud, D., Maclaurin, D., and Adams, R. P. (2016). Early Stopping as Nonparametric Variational Inference. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 51, pages 1070–1077.
- Foong, A. Y. K., Li, Y., Hernández-Lobato, J. M., and Turner, R. E. (2019). ‘In-Between’ Uncertainty in Bayesian Neural Networks. In *Workshop on Uncertainty and Robustness in Deep Learning, ICML*.
- Gal, Y. (2016). *Uncertainty in Deep Learning*. PhD thesis.
- Gal, Y. and Ghahramani, Z. (2015). Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. In *Proceedings of the 33rd International Conference on Machine Learning*.
- Graves, A. (2011). Practical Variational Inference for Neural Networks. *Advances in Neural Information Processing Systems*, pages 1–9.
- Gu, Y., Oliver, D. S., and Oklahoma, U. (2007). An Iterative Ensemble Kalman Filter for Multiphase Fluid Flow Data Assimilation. *SPE Journal* 12(4).
- Hernández-Lobato, J. M. and Adams, R. P. (2015). Probabilistic Backpropagation for Scalable Learning of Bayesian Neural Networks. In *Proceedings of the 32nd International Conference on Machine Learning*.
- Khan, M. E., Nielsen, D., Tangkaratt, V., Lin, W., Gal, Y., and Srivastava, A. (2018). Fast and Scalable Bayesian Deep Learning by Weight-Perturbation in Adam. In *ICML*.
- Lakshminarayanan, B., Pritzel, A., and Blundell, C. (2017). Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles. In *31st Conference on Neural Information Processing Systems*.
- Lu, X. and Van Roy, B. (2017). Ensemble Sampling. In *31st Conference on Neural Information Processing Systems*.
- MacKay, D. J. C. (1992). A Practical Bayesian Framework for Backpropagation Networks. *Neural Computation*, 4(3):448–472.
- Mukhoti, J., Stenatorp, P., and Gal, Y. (2018). On the Importance of Strong Baselines in Bayesian Deep Learning. In *Bayesian Deep Learning Workshop, Neural Information Processing Systems (NeurIPS)*, pages 1–4.
- Neal, R. M. (1997). *Bayesian Learning for Neural Networks*. PhD thesis.
- Osband, I., Aslanides, J., and Cassirer, A. (2018). Randomized Prior Functions for Deep Reinforcement Learning. In *32nd Conference on Neural Information Processing Systems (NIPS 2018)*.
- Osband, I., Blundell, C., Pritzel, A., and Van Roy, B. (2016). Deep Exploration via Bootstrapped DQN. In *Advances in neural information processing systems*, pages 1–18.
- Osband, I., Russo, D., Wen, Z., and Van Roy, B. (2019). Deep Exploration via Randomized Value Functions. *JMLR*.

Pearce, T., Tsuchida, R., Zaki, M., Brintrup, A., and Neely, A. (2019). Expressive Priors in Bayesian Neural Networks: Kernel Combinations and Periodic Functions. In *Uncertainty in Artificial Intelligence (UAI)*.

Pedersen, M. S., Baxter, B., Templeton, B., Rishøj, C., Theobald, D. L., Hoegh-rasmussen, E., Casteel, G., Gao, J. B., Dedecius, K., Strim, K., Christiansen, L., Hansen, L. K., Wilkinson, L., He, L., Bar, M., Winther, O., Sakov, P., Hattinger, S., Petersen, K. B., and Rishøj, C. (2008). The Matrix Cookbook. *Matrix*, M:1–71.

Ritter, H., Botev, A., and Barber, D. (2018). A Scalable Laplace Approximation for Neural Networks. In *ICLR*, pages 1–15.

Smith, L. and Gal, Y. (2018). Understanding Measures of Uncertainty for Adversarial Example Detection. In *Uncertainty in Artificial Intelligence (UAI)*.

Sünderhauf, N., Brock, O., Scheirer, W., Hadsell, R., Fox, D., Leitner, J., Upcroft, B., Abbeel, P., Burgard, W., Milford, M., and Corke, P. (2018). The limits and potentials of deep learning for robotics. *The International Journal of Robotics Research*, 37:405–420.

Tibshirani, R. (1996). A Comparison of Some Error Estimates for Neural Network Models. *Neural Computation*, 8:152–163.

Tsuchida, R., Roosta-Khorasani, F., and Gallagher, M. (2018). Invariance of Weight Distributions in Rectified MLPs. In *Proceedings of the 35th International Conference on Machine Learning*.

Wang, K. A., Pleiss, G., Gardner, J. R., Tyree, S., Weinberger, K. Q., and Wilson, A. G. (2019). Exact Gaussian Processes on a Million Data Points. In *Neural Information Processing Systems*.

Williams, C. K. I. (1996). Computing with infinite networks. In *Advances in Neural Information Processing Systems 9*.

# Appendix to Uncertainty in Neural Networks: Approximately Bayesian Ensembling

## A Proofs

**Definition 1.** Data likelihood and parameter likelihood

We take care to define two versions of the likelihood, one in output space,  $P_{\mathcal{D}}(\mathcal{D}|\boldsymbol{\theta})$  (data likelihood), and one in parameter space,  $P_{\boldsymbol{\theta}}(\mathcal{D}|\boldsymbol{\theta})$  (parameter likelihood). Both return the same values given some data set  $\mathcal{D}$  and parameter values  $\boldsymbol{\theta}$ , and hence are exchangeable, but their forms are subtly different.

The data likelihood,  $P_{\mathcal{D}}(\mathcal{D}|\boldsymbol{\theta})$ , is defined on the output domain. Typically the log of this,  $\log(P_{\mathcal{D}}(\mathcal{D}|\boldsymbol{\theta}))$ , might be optimised as the cross entropy loss or (negative) mean squared error.

In contrast,  $P_{\boldsymbol{\theta}}(\mathcal{D}|\boldsymbol{\theta})$  defines a likelihood function in the parameter domain.

### Illustrative Example

Consider a linear regression model with dataset,  $\mathcal{D}$ , consisting of tuples,  $\{\mathbf{x}, y\}$ ; a vector of predictor variables  $\mathbf{x} \in \mathbb{R}^p$ , predicting a single scalar  $y \in \mathbb{R}$ . If the model is of the form,  $\boldsymbol{\theta}^T \mathbf{x}$ , a Gaussian data likelihood with variance  $\sigma_{\epsilon}^2$  on the output might be assumed.

This leads to a data likelihood for the target,  $y$ ,

$$P_{\mathcal{D}}(\mathcal{D}|\boldsymbol{\theta}) = \mathcal{N}(y|\boldsymbol{\theta}^T \mathbf{x}, \sigma_{\epsilon}^2). \quad (10)$$

For this linear model, the corresponding parameter likelihood is a multivariate normal distribution,

$$P_{\boldsymbol{\theta}}(\mathcal{D}|\boldsymbol{\theta}) \propto \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\mu}_{like}, \boldsymbol{\Sigma}_{like}). \quad (11)$$

where,  $\boldsymbol{\mu}_{like} \in \mathbb{R}^p$  &  $\boldsymbol{\Sigma}_{like} \in \mathbb{R}^{p \times p}$ , can be found analytically. They are implicitly functions of the dataset,  $\mathcal{D}$ , although to lighten notation we do not write this. Subsequently we also drop the explicit referral to  $\boldsymbol{\theta}$ .

Note that whilst both the data and parameter likelihood follow a normal distribution, they are defined in different domains.

The correspondence between a Gaussian data likelihood and multivariate normal parameter likelihood is only exact for a linear regression model. For non-linear models with Gaussian data likelihoods, and other data likelihoods, the parameter likelihood is not in general multivariate normal. Nevertheless it can be convenient to model it as such.

**Standard Result 1.** Product of two multivariate Gaussians (§8.1.8, The Matrix Cookbook, 2008)

$$\mathcal{N}(\boldsymbol{\mu}_{like}, \boldsymbol{\Sigma}_{like})\mathcal{N}(\boldsymbol{\mu}_{prior}, \boldsymbol{\Sigma}_{prior}) \propto \mathcal{N}(\boldsymbol{\mu}_{post}, \boldsymbol{\Sigma}_{post}) \quad (12)$$

$$\boldsymbol{\Sigma}_{post} = (\boldsymbol{\Sigma}_{prior}^{-1} + \boldsymbol{\Sigma}_{like}^{-1})^{-1}, \quad (13)$$

$$\boldsymbol{\mu}_{post} = \boldsymbol{\Sigma}_{post}\boldsymbol{\Sigma}_{prior}^{-1}\boldsymbol{\mu}_{prior} + \boldsymbol{\Sigma}_{post}\boldsymbol{\Sigma}_{like}^{-1}\boldsymbol{\mu}_{like}. \quad (14)$$

**Standard Result 2.** Affine transform of a normal random variable (§8.1.4, The Matrix Cookbook, 2008)

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad (15)$$

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}, \quad (16)$$

$$\mathbf{y} \sim \mathcal{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T). \quad (17)$$

**Theorem 1.** Assume that a model’s parameter likelihood follows a multivariate normal distribution,  $P_{\theta}(\mathcal{D}|\theta) \propto \mathcal{N}(\boldsymbol{\mu}_{like}, \boldsymbol{\Sigma}_{like})$ , and the prior also,  $P(\theta) = \mathcal{N}(\boldsymbol{\mu}_{prior}, \boldsymbol{\Sigma}_{prior})$ . The posterior is then also multivariate normal,  $P(\theta|\mathcal{D}) = \mathcal{N}(\boldsymbol{\mu}_{post}, \boldsymbol{\Sigma}_{post})$ .

Further assume availability of some function which returns MAP parameter estimates taking as input the location of the prior centre,  $\mathbf{f}_{MAP}(\theta_{anc})$ . In order that,  $P(\mathbf{f}_{MAP}(\theta_{anc})) = P(\theta|\mathcal{D})$ , then the required distribution of  $\theta_{anc}$  is also multivariate normal,  $P(\theta_{anc}) = \mathcal{N}(\boldsymbol{\mu}_{anc}, \boldsymbol{\Sigma}_{anc})$ , where,  $\boldsymbol{\mu}_{anc} = \boldsymbol{\mu}_{prior}$ , and,  $\boldsymbol{\Sigma}_{anc} = \boldsymbol{\Sigma}_{prior} + \boldsymbol{\Sigma}_{prior}\boldsymbol{\Sigma}_{like}^{-1}\boldsymbol{\Sigma}_{prior}$ .

*Proof.* Consider a model’s parameters having a multivariate normal prior,

$$P(\theta) = \mathcal{N}(\boldsymbol{\mu}_{prior}, \boldsymbol{\Sigma}_{prior}), \quad (18)$$

where,  $\theta \in \mathbb{R}^p$ ,  $\boldsymbol{\mu}_{prior} \in \mathbb{R}^p$ ,  $\boldsymbol{\Sigma}_{prior} \in \mathbb{R}^{p \times p}$ .

This theorem makes the assumption that the form of the parameter likelihood (def. 1) is multivariate normal,

$$P_{\theta}(\mathcal{D}|\theta) \propto \mathcal{N}(\boldsymbol{\mu}_{like}, \boldsymbol{\Sigma}_{like}) \quad (19)$$

where,  $\boldsymbol{\mu}_{like} \in \mathbb{R}^p$ ,  $\boldsymbol{\Sigma}_{like} \in \mathbb{R}^{p \times p}$ . Here  $\propto$  is used since it is not a true probability distribution in  $\theta$  so need not sum to 1.

The posterior is calculated by Bayes rule. Recalling that data likelihood and parameter likelihood are exchangeable (def. 1), and using Standard Result 1,

$$P(\theta|\mathcal{D}) = \frac{P_{\mathcal{D}}(\mathcal{D}|\theta)P(\theta)}{P(\mathcal{D})} = \frac{P_{\theta}(\mathcal{D}|\theta)P(\theta)}{P(\mathcal{D})} \propto \mathcal{N}(\boldsymbol{\mu}_{like}, \boldsymbol{\Sigma}_{like})\mathcal{N}(\boldsymbol{\mu}_{prior}, \boldsymbol{\Sigma}_{prior}) \propto \mathcal{N}(\boldsymbol{\mu}_{post}, \boldsymbol{\Sigma}_{post}), \quad (20)$$

where,  $\boldsymbol{\mu}_{post}$  &  $\boldsymbol{\Sigma}_{post}$  are given by eq. 14 & 13.

We introduce a further distribution, termed ‘anchor distribution’, which we enforce as multivariate normal,

$$P(\theta_{anc}) = \mathcal{N}(\boldsymbol{\mu}_{anc}, \boldsymbol{\Sigma}_{anc}). \quad (21)$$

It will be used as described in the main text (see figure 2 and algorithm 1) so that samples are drawn from the anchor distribution,  $\theta_{anc} \sim P(\theta_{anc})$ , with a prior then recentred at each sample, denoted  $P_{anc}(\theta)$ ,

$$P_{anc}(\theta) = \mathcal{N}(\theta_{anc}, \boldsymbol{\Sigma}_{prior}). \quad (22)$$

Note that this anchor distribution is in the same position as a hyperprior on  $\boldsymbol{\mu}_{prior}$ , but will have a subtly different role.  $\boldsymbol{\Sigma}_{prior}$  is unchanged from eq. 18,

Denote  $\mathbf{f}_{MAP}(\theta_{anc})$  as the MAP estimates given this recentred prior and the original likelihood from eq. 19.

$$\mathbf{f}_{MAP}(\theta_{anc}) := \operatorname{argmax}_{\theta} P_{anc}(\theta)P_{\theta}(\mathcal{D}|\theta) \quad (23)$$

In order to prove the theorem, three things regarding  $\mathbf{f}_{MAP}(\theta_{anc})$  must be shown:

1. Its distribution is multivariate normal - denote mean and covariance  $\boldsymbol{\mu}_{post}^{\text{RMS}}, \boldsymbol{\Sigma}_{post}^{\text{RMS}}$ ,

$$P(\mathbf{f}_{MAP}(\theta_{anc})) = \mathcal{N}(\boldsymbol{\mu}_{post}^{\text{RMS}}, \boldsymbol{\Sigma}_{post}^{\text{RMS}}), \quad (24)$$

2. That  $\boldsymbol{\mu}_{anc}$  &  $\boldsymbol{\Sigma}_{anc}$  can be selected in such a way that the mean of the distribution is equal to that of the original posterior

$$\boldsymbol{\mu}_{post}^{\text{RMS}} = \boldsymbol{\mu}_{post}, \quad (25)$$

3. And also so that the covariance of the distribution is equal to that of the original posterior

$$\Sigma_{post}^{RMS} = \Sigma_{post}. \quad (26)$$

For a multivariate normal distribution, the MAP solution is simply equal to the mean of the posterior,  $\boldsymbol{\mu}_{post}$ . For the typical case this is given by eq. 14. In our procedure, the location of the prior mean has been replaced by  $\boldsymbol{\theta}_{anc}$ , so the MAP solution is given by,

$$\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc}) = \Sigma_{post}\Sigma_{prior}^{-1}\boldsymbol{\theta}_{anc} + \Sigma_{post}\Sigma_{like}^{-1}\boldsymbol{\mu}_{like} \quad (27)$$

$$= \mathbf{A}_1\boldsymbol{\theta}_{anc} + \mathbf{b}_1 \quad (28)$$

where two constants have been defined for convenience,

$$\mathbf{A}_1 = \Sigma_{post}\Sigma_{prior}^{-1} \quad (29)$$

$$\mathbf{b}_1 = \Sigma_{post}\Sigma_{like}^{-1}\boldsymbol{\mu}_{like}, \quad (30)$$

which is the same form as eq. 16. Hence, from Standard Result 2, if  $\boldsymbol{\theta}_{anc}$  is normally distributed,  $\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})$  **will also be normally distributed**.

Regarding the mean of  $\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})$ , we have,

$$\mathbb{E}[\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})] = \mathbb{E}[\mathbf{A}_1\boldsymbol{\theta}_{anc} + \mathbf{b}_1] \quad (31)$$

$$= \mathbf{A}_1\mathbb{E}[\boldsymbol{\theta}_{anc}] + \mathbf{b}_1. \quad (32)$$

By choosing the anchor distribution to be centred about the original prior,  $\mathbb{E}[\boldsymbol{\theta}_{anc}] = \boldsymbol{\mu}_{prior}$ , we have,

$$= \mathbf{A}_1\boldsymbol{\mu}_{prior} + \mathbf{b}_1 \quad (33)$$

$$= \Sigma_{post}\Sigma_{prior}^{-1}\boldsymbol{\mu}_{prior} + \Sigma_{post}\Sigma_{like}^{-1}\boldsymbol{\mu}_{like}, \quad (34)$$

This is consistent with eq. 14 and proves that **the means of the distributions are aligned when  $\boldsymbol{\mu}_{anc} = \boldsymbol{\mu}_{prior}$** .

Finally we consider the variance of  $\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})$ , which we wish to equal  $\Sigma_{post}$  by choosing  $\Sigma_{anc}$ . Using the form from eq. 28 and applying Standard Result 2,

$$\text{Var}[\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})] = \text{Var}[\mathbf{A}_1\boldsymbol{\theta}_{anc} + \mathbf{b}_1] \quad (35)$$

$$= \mathbf{A}_1\text{Var}[\boldsymbol{\theta}_{anc}]\mathbf{A}_1^T \quad (36)$$

$$= \mathbf{A}_1\Sigma_{anc}\mathbf{A}_1^T \quad (37)$$

We require  $\text{Var}[\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})] = \Sigma_{post}$ .

$$\Sigma_{post} = \mathbf{A}_1\Sigma_{anc}\mathbf{A}_1^T. \quad (38)$$

Note that transposes of covariance matrices may be ignored since they are symmetric.

$$\Sigma_{anc} = \mathbf{A}_1^{-1}\Sigma_{post}\mathbf{A}_1^{-1T} \quad (39)$$

$$= (\Sigma_{post}\Sigma_{prior}^{-1})^{-1}\Sigma_{post}(\Sigma_{prior}^{-1}\Sigma_{post})^{-1} \quad (40)$$

$$= \Sigma_{prior}\Sigma_{post}^{-1}\Sigma_{post}\Sigma_{post}^{-1}\Sigma_{prior} \quad (41)$$

$$= \Sigma_{prior}\Sigma_{post}^{-1}\Sigma_{prior} \quad (42)$$

$$= \Sigma_{prior}(\Sigma_{prior}^{-1} + \Sigma_{like}^{-1})\Sigma_{prior} \quad (43)$$

$$= \Sigma_{prior} + \Sigma_{prior}\Sigma_{like}^{-1}\Sigma_{prior}. \quad (44)$$

This proves that **the covariances of the two distributions are aligned when  $\Sigma_{anc} = \Sigma_{prior} + \Sigma_{prior}\Sigma_{like}^{-1}\Sigma_{prior}$** .

□

**Corollary 1.1.** *Following from theorem 1 (and under the same assumptions), set  $\boldsymbol{\mu}_{anc} := \boldsymbol{\mu}_{prior}$  and  $\boldsymbol{\Sigma}_{anc} := \boldsymbol{\Sigma}_{prior}$ . The RMS approximate posterior is  $P(\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})) = \mathcal{N}(\boldsymbol{\mu}_{post}, \boldsymbol{\Sigma}_{post}\boldsymbol{\Sigma}_{prior}^{-1}\boldsymbol{\Sigma}_{post})$ .*

*Proof.* Independent of the choice of anchor distribution covariance  $\boldsymbol{\Sigma}_{anc}$ , theorem 1 demonstrated that the resulting posterior,  $P(\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc}))$  is normally distributed, with mean equal to that of the true posterior  $\boldsymbol{\mu}_{post}$ .

To discover the covariance of the resulting distribution,  $\text{Var}[\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})]$ , we take eq. 37 and simply set,  $\boldsymbol{\Sigma}_{anc} := \boldsymbol{\Sigma}_{prior}$ .

$$\text{Var}[\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})] = \mathbf{A}_1\boldsymbol{\Sigma}_{anc}\mathbf{A}_1^T \quad (45)$$

$$= \mathbf{A}_1\boldsymbol{\Sigma}_{prior}\mathbf{A}_1^T \quad (46)$$

$$= \boldsymbol{\Sigma}_{post}\boldsymbol{\Sigma}_{prior}^{-1}\boldsymbol{\Sigma}_{prior}\boldsymbol{\Sigma}_{prior}^{-1}\boldsymbol{\Sigma}_{post} \quad (47)$$

$$= \boldsymbol{\Sigma}_{post}\boldsymbol{\Sigma}_{prior}^{-1}\boldsymbol{\Sigma}_{post} \quad (48)$$

□

**Lemma 1.1.** *Following from corollary 1.1 (and under the same assumptions), when  $\boldsymbol{\mu}_{anc} := \boldsymbol{\mu}_{prior}$ ,  $\boldsymbol{\Sigma}_{anc} := \boldsymbol{\Sigma}_{prior}$ , the RMS approximate posterior will in general underestimate the marginal variance compared to the true posterior,  $\text{Var}[\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})] < \text{Var}[\boldsymbol{\theta}|\mathcal{D}]$ .*

*Proof.* We consider the marginal posterior of a single parameter,  $\theta := \boldsymbol{\theta}_i$ , again assuming multivariate normal prior and parameter likelihood. First consider the following rearrangement of eq. 48, beginning by noting,  $\boldsymbol{\Sigma}_{prior}^{-1} = \boldsymbol{\Sigma}_{post}^{-1} - \boldsymbol{\Sigma}_{like}^{-1}$ .

$$\boldsymbol{\Sigma}_{post}\boldsymbol{\Sigma}_{prior}^{-1}\boldsymbol{\Sigma}_{post} = \boldsymbol{\Sigma}_{post}(\boldsymbol{\Sigma}_{post}^{-1} - \boldsymbol{\Sigma}_{like}^{-1})\boldsymbol{\Sigma}_{post} \quad (49)$$

$$= (\mathbb{I} - \boldsymbol{\Sigma}_{post}\boldsymbol{\Sigma}_{like}^{-1})\boldsymbol{\Sigma}_{post} \quad (50)$$

$$= \boldsymbol{\Sigma}_{post} - \boldsymbol{\Sigma}_{post}\boldsymbol{\Sigma}_{like}^{-1}\boldsymbol{\Sigma}_{post} \quad (51)$$

To show that RMS generally underestimates the marginal variance, it must hold that diagonal elements of the true posterior covariance matrix are greater than or equal to the same diagonal element of the RMS posterior.

$$\text{Var}[\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})] < \text{Var}[\boldsymbol{\theta}|\mathcal{D}] \quad (52)$$

$$\text{diag}(\boldsymbol{\Sigma}_{post}\boldsymbol{\Sigma}_{prior}^{-1}\boldsymbol{\Sigma}_{post})_i < \text{diag}(\boldsymbol{\Sigma}_{post})_i \quad (53)$$

substituting in the diagonal of the rearrangement in eq. 51,

$$\text{diag}(\boldsymbol{\Sigma}_{post})_i - \text{diag}(\boldsymbol{\Sigma}_{post}\boldsymbol{\Sigma}_{like}^{-1}\boldsymbol{\Sigma}_{post})_i < \text{diag}(\boldsymbol{\Sigma}_{post})_i \quad (54)$$

We know that  $ABA^T$  is positive definite if  $A, B$  are positive definite, and also that the inverse of a positive definite matrix is positive definite (§9.6.4, §9.6.10, The Matrix Cookbook, 2008). The diagonal of a positive definite matrix is positive. Hence,  $\text{diag}(\boldsymbol{\Sigma}_{post}\boldsymbol{\Sigma}_{like}^{-1}\boldsymbol{\Sigma}_{post})_i > 0$ , and we have shown that  $\text{Var}[\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})] < \text{Var}[\boldsymbol{\theta}|\mathcal{D}]$ .

□

**Lemma 1.2.** *This lemma follows from corollary 1.1. Again parameter likelihood and prior are assumed normally distributed. The prior is additionally assumed isotropic. When  $\boldsymbol{\mu}_{anc} := \boldsymbol{\mu}_{prior}$ ,  $\boldsymbol{\Sigma}_{anc} := \boldsymbol{\Sigma}_{prior}$  the eigenvectors (or ‘orientation’) of the RMS approximate posterior equal those of the true posterior.*

*Proof.* From eq. 48,  $\text{Var}[\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})] = \boldsymbol{\Sigma}_{post}\boldsymbol{\Sigma}_{prior}^{-1}\boldsymbol{\Sigma}_{post}$ . If the prior is isotropic,  $\boldsymbol{\Sigma}_{prior} = \sigma_{prior}^2\mathbb{I}$ , then,  $\text{Var}[\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})] = 1/\sigma_{prior}^2\boldsymbol{\Sigma}_{post}$ . Hence the prior only scales the eigenvalues, and doesn’t affect the eigenvectors. (Note that this won’t be the case for non-isotropic  $\boldsymbol{\Sigma}_{prior}$ .)

Consider some matrix  $A$  and a specific eigenvalue  $\lambda_i$  and eigenvector  $\mathbf{v}_i$  so that,  $A\mathbf{v}_i = \lambda_i\mathbf{v}_i$ . It then follows that if  $A$  is squared,  $A^2\mathbf{v}_i = A(A\mathbf{v}_i) = \lambda_i A\mathbf{v}_i = \lambda_i^2\mathbf{v}_i$ . Hence eigenvalues are squared but eigenvectors are unaffected. This applies to the transformation  $\boldsymbol{\Sigma}_{post}^2$ .

Hence both the square and the multiplication of prior covariance,  $1/\sigma_{prior}^2\boldsymbol{\Sigma}_{post}^2$ , do not modify the original eigenvectors of  $\boldsymbol{\Sigma}_{post}$  and its orientation is unaffected.

□

**Theorem 2.** *For a two parameter model with normally distributed parameter likelihood and isotropic prior, the RMS approximate posterior will in general overestimate the magnitude of the true posterior parameter correlation coefficient,  $|\rho|$ . However, if  $|\rho| = 1$ , then it will recover it precisely. We set  $\boldsymbol{\mu}_{anc} := \boldsymbol{\mu}_{prior}$ ,  $\boldsymbol{\Sigma}_{anc} := \boldsymbol{\Sigma}_{prior}$ .*

*Proof.* From corollary 1.1, we have that the RMS approximate posterior is given by  $\boldsymbol{\Sigma}_{post} \boldsymbol{\Sigma}_{prior}^{-1} \boldsymbol{\Sigma}_{post}$ . Let  $\boldsymbol{\Sigma}_{prior} := \sigma_{prior}^2 \mathbb{I}$ , the RMS approximate posterior is then given by  $1/\sigma_{prior}^2 \boldsymbol{\Sigma}_{post}^2$ . Denote the true posterior covariance as the following  $2 \times 2$  matrix.

$$\boldsymbol{\Sigma}_{post} = \begin{bmatrix} a & b \\ b & c \end{bmatrix} \quad (55)$$

For general covariance matrices, the correlation coefficient,  $\rho$ , can be found by solving,  $b = \rho\sqrt{ac}$ .

Our RMS approximate posterior is given as follows.

$$1/\sigma_{prior}^2 \boldsymbol{\Sigma}_{post}^2 = 1/\sigma_{prior}^2 \begin{bmatrix} a & b \\ b & c \end{bmatrix} \begin{bmatrix} a & b \\ b & c \end{bmatrix} \quad (56)$$

$$= 1/\sigma_{prior}^2 \begin{bmatrix} a^2 + b^2 & ab + bc \\ ab + bc & b^2 + c^2 \end{bmatrix} \quad (57)$$

The correlation coefficient here, denoted  $\rho_{RMS}$ , is found by solving,  $ab + bc = \rho_{RMS} \sqrt{(a^2 + b^2)(b^2 + c^2)}$ .

To prove the correlation coefficient is generally overestimated, we must show that  $\rho_{RMS}^2 > \rho^2$  when  $\rho^2 < 1$ .

$$\rho_{RMS}^2 > \rho^2 \quad (58)$$

$$\frac{(ab + bc)^2}{(a^2 + b^2)(b^2 + c^2)} > \frac{b^2}{ac} \quad (59)$$

$$(ab + bc)^2 ac > b^2(a^2 + b^2)(b^2 + c^2) \quad (60)$$

$$(a + c)^2 ac > (a^2 + b^2)(b^2 + c^2) \quad (61)$$

$$a^3c + ac^3 + 2a^2c^2 > a^2b^2 + a^2c^2 + b^4 + b^2c^2 \quad (62)$$

$$a^3c + ac^3 + a^2c^2 > a^2b^2 + b^4 + b^2c^2 \quad (63)$$

We now note that from the set up of the proof,  $b^2 = \rho^2 ac$ . Since  $\rho^2 < 1$ , we have,  $b^2/ac < 1 \implies b^2 < ac$ . We can use this to provide an upper bound on the right hand side of eq. 63.

$$a^2b^2 + b^4 + b^2c^2 < a^2(ac) + (ac)^2 + (ac)c^2 \quad (64)$$

$$< a^3c + a^2c^2 + ac^3 \quad (65)$$

Coincidentally, this is precisely the inequality in eq. 63, that we were proving.

Alternatively, if  $|\rho| = 1 \implies b^2 = ac$ , and  $\rho_{RMS}^2 = \rho^2$ . □

**Definition 2.** Extrapolation Parameters

We define extrapolation parameters as model parameters which have no effect on the data likelihood of a training dataset, but which nevertheless could influence model predictions made on a new data point.

**Illustrative Example**

Consider a fully-connected NN trained on the MNIST digit dataset. Further consider a preprocessing such that the pixel values by default are set to 0, and where they contain part of the digit take values (0, 1].

Certain pixels may be zero across the entire training dataset, such as those in the corners of the image. First layer weights connected to such pixels will never receive input across the whole training dataset. Hence, the values of these parameter weights have no effect on the data likelihood. However, the weights would still influence predictions for some test image containing values for these pixels. Hence, these are named extrapolation parameters, since they influence extrapolation properties of the model.

The top row of figure 5 empirically shows examples of flat likelihoods for precisely these types of weights on MNIST.

**Theorem 3.** For extrapolation parameters (definition 2) of a model, setting  $\boldsymbol{\mu}_{anc} := \boldsymbol{\mu}_{prior}$ ,  $\boldsymbol{\Sigma}_{anc} := \boldsymbol{\Sigma}_{prior}$ , means the marginal RMS approximate posterior equals that of the marginal true posterior. This holds for any distributional form of parameter likelihood.

*Proof.* Extrapolation parameters (definition 2) do not have any effect on the data likelihood, therefore their parameter likelihoods are flat. This means that their marginal posterior equals their marginal prior.

This results in a posterior covariance matrix structure as follows, where parameter  $i$  is an extrapolation parameter (here shown in the first row for convenience), so has marginal variance equal to the prior variance and is uncorrelated with all other parameters.

$$\boldsymbol{\Sigma}_{post} = \begin{bmatrix} \sigma_{prior,i}^2 & 0 & \dots & 0 \\ 0 & a_{22} & \dots & a_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{D2} & \dots & a_{DD} \end{bmatrix}$$

From corollary 1.1,  $\text{Var}[\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})] = \boldsymbol{\Sigma}_{post} \boldsymbol{\Sigma}_{prior}^{-1} \boldsymbol{\Sigma}_{post}$ .

$$\begin{aligned} \text{Var}[\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})] &= \begin{bmatrix} \sigma_{prior,i}^2 & 0 & \dots & 0 \\ 0 & a_{22} & \dots & a_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{D2} & \dots & a_{DD} \end{bmatrix} \begin{bmatrix} 1/\sigma_{prior,i}^2 & 0 & \dots & 0 \\ 0 & b_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & b_{DD} \end{bmatrix} \begin{bmatrix} \sigma_{prior,i}^2 & 0 & \dots & 0 \\ 0 & a_{22} & \dots & a_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{D2} & \dots & a_{DD} \end{bmatrix} \\ &= \begin{bmatrix} \sigma_{prior,i}^2 & 0 & \dots & 0 \\ 0 & c_{22} & \dots & c_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & c_{D2} & \dots & c_{DD} \end{bmatrix} \end{aligned}$$

This shows that the marginal variance of the RMS approximate posterior equals that of the true posterior,  $\text{Var}[\mathbf{f}_{MAP}(\boldsymbol{\theta}_{anc})]_{i,i} = [\boldsymbol{\Sigma}_{post}]_{i,i}$ , for extrapolation parameters. Note that the values of the rest of the covariance matrices ( $a$ 's,  $b$ 's,  $c$ 's) are irrelevant since these have no effect on the marginals of interest.

This proof did not assume any specific distributional form of parameter likelihood, only that it is flat for these extrapolation parameters. □



**Theorem 4.** Set  $\boldsymbol{\mu}_{anc} := \boldsymbol{\mu}_{prior}$ ,  $\boldsymbol{\Sigma}_{anc} := \boldsymbol{\Sigma}_{prior}$ . The RMS approximate posterior will exactly equal the true posterior,  $\boldsymbol{\Sigma}_{post}$ , when all eigenvalues of a scaled version of  $\boldsymbol{\Sigma}_{post}$  (scaled such that the prior equals the identity matrix) are equal to either 0 or 1. This corresponds to posteriors that are a mixture of perfectly correlated and perfectly uncorrelated parameters.

*Proof.* We will initially consider a scaled version of the parameter space. This conveniently allows standard results for idempotent matrices to apply to the posterior covariance. A reverse scaling is subsequently applied to show that results hold for the original unscaled version. Finally, we articulate arguments allowing relaxation of the distributional assumptions.

Corollary 1.1 showed that the RMS approximate posterior is normally distributed and centered at the true posterior mean, but with modified variance,  $\mathcal{N}(\boldsymbol{\mu}_{post}, \boldsymbol{\Sigma}_{post} \boldsymbol{\Sigma}_{prior}^{-1} \boldsymbol{\Sigma}_{post})$ . This proof requires specifying conditions that allow  $\boldsymbol{\Sigma}_{post} \boldsymbol{\Sigma}_{prior}^{-1} \boldsymbol{\Sigma}_{post} = \boldsymbol{\Sigma}_{post}$  to hold.

One solution is given by  $\boldsymbol{\Sigma}_{post} = \boldsymbol{\Sigma}_{prior}$ , and is a trivial extension of theorem 3. Here we consider alternative solutions.

The two inputs into the inference process are the likelihood and prior covariances. Consider a scaling  $\boldsymbol{\Sigma}'_{like} := \boldsymbol{\Sigma}_{prior}^{-1/2} \boldsymbol{\Sigma}_{like} \boldsymbol{\Sigma}_{prior}^{-1/2}$  and  $\boldsymbol{\Sigma}'_{prior} := \boldsymbol{\Sigma}_{prior}^{-1/2} \boldsymbol{\Sigma}_{prior} \boldsymbol{\Sigma}_{prior}^{-1/2} = \mathbb{I}$ . The posterior for this scaled version will be denoted by  $\boldsymbol{\Sigma}'_{post}$ , and is given as follows.

$$\boldsymbol{\Sigma}'_{post} = (\boldsymbol{\Sigma}'_{like} + \boldsymbol{\Sigma}'_{prior})^{-1} \tag{66}$$

$$= (\boldsymbol{\Sigma}_{prior}^{1/2} \boldsymbol{\Sigma}_{like} \boldsymbol{\Sigma}_{prior}^{1/2} + \boldsymbol{\Sigma}_{prior}^{1/2} \boldsymbol{\Sigma}_{prior}^{-1} \boldsymbol{\Sigma}_{prior}^{1/2})^{-1} \tag{67}$$

$$= \boldsymbol{\Sigma}_{prior}^{-1/2} (\boldsymbol{\Sigma}_{like} + \boldsymbol{\Sigma}_{prior})^{-1} \boldsymbol{\Sigma}_{prior}^{-1/2} \tag{68}$$

$$= \boldsymbol{\Sigma}_{prior}^{-1/2} \boldsymbol{\Sigma}_{post} \boldsymbol{\Sigma}_{prior}^{-1/2} \tag{69}$$

Hence, unsurprisingly the same scaling applies to the posterior covariance,  $\boldsymbol{\Sigma}'_{post} = \boldsymbol{\Sigma}_{prior}^{-1/2} \boldsymbol{\Sigma}_{post} \boldsymbol{\Sigma}_{prior}^{-1/2}$ .

We now consider conditions under which  $\boldsymbol{\Sigma}'_{post} \boldsymbol{\Sigma}'_{prior} \boldsymbol{\Sigma}'_{post} = \boldsymbol{\Sigma}'_{post}$  holds. From our choice of rescaling, we have that  $\boldsymbol{\Sigma}'_{prior} = \mathbb{I}$ . So we require that  $\boldsymbol{\Sigma}'_{post} = \boldsymbol{\Sigma}'_{post}$ .

This conveniently allows use of results for idempotent matrices - defined as a square matrix,  $A$ , for which  $A^2 = A$ . Aside from the case when  $A = \mathbb{I}$  (which corresponds to  $\boldsymbol{\Sigma}_{post} = \boldsymbol{\Sigma}_{prior}$ ), a matrix is idempotent if and only if it is singular and all eigenvalues are 0 or 1.

In order that,  $\boldsymbol{\Sigma}'_{post} \boldsymbol{\Sigma}'_{prior} \boldsymbol{\Sigma}'_{post} = \boldsymbol{\Sigma}'_{post}$ , it is therefore sufficient that our scaled posterior,  $\boldsymbol{\Sigma}'_{post}$ , is singular with all eigenvalues 0 or 1. Any possible permutation is allowed.

Naturally, applying a reverse scaling recovers the original parameter space,  $\boldsymbol{\Sigma}_{post} = \boldsymbol{\Sigma}_{prior}^{1/2} \boldsymbol{\Sigma}'_{post} \boldsymbol{\Sigma}_{prior}^{1/2}$ .

**Remark.** To summarise, we have shown that provided the RMS approximate posterior equals the true posterior in the scaled space, it will also be equal in the original unscaled space. In order for this equality to hold, eigenvalues must be 0 or 1 in the scaled space.

See section B.1.2 for numerical examples in a three parameter model when this condition holds.

□

### A.1 MAP solution and regularisation interpretation

For completeness, we write out the MAP solution for the case of normally distributed prior, and data likelihoods often used in regression and classification. From this derives our interpretation of the regularisation matrix,  $\mathbf{\Gamma}$ .

$$\begin{aligned}\boldsymbol{\theta}_{MAP} &= \operatorname{argmax}_{\boldsymbol{\theta}} P(\boldsymbol{\theta}|\mathcal{D}) \\ &= \operatorname{argmax}_{\boldsymbol{\theta}} P_{\mathcal{D}}(\mathcal{D}|\boldsymbol{\theta})P(\boldsymbol{\theta}) \\ &= \operatorname{argmax}_{\boldsymbol{\theta}} \log(P_{\mathcal{D}}(\mathcal{D}|\boldsymbol{\theta})) + \log(P(\boldsymbol{\theta}))\end{aligned}$$

If prior is normally distributed,  $P(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,

$$\begin{aligned}&= \operatorname{argmax}_{\boldsymbol{\theta}} \log(P_{\mathcal{D}}(\mathcal{D}|\boldsymbol{\theta})) - \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\theta} - \boldsymbol{\mu}) + \text{const.} \\ &= \operatorname{argmax}_{\boldsymbol{\theta}} \log(P_{\mathcal{D}}(\mathcal{D}|\boldsymbol{\theta})) - \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\theta} - \boldsymbol{\mu}).\end{aligned}$$

Typically in BNNs the prior covariance is chosen as diagonal. This is sometimes set as isotropic,  $\boldsymbol{\Sigma} = \lambda \mathbb{I}$ , but here we will keep it in matrix form (but assuming it is diagonal) so that different prior variances can be assigned to different layer weights.

$$= \operatorname{argmax}_{\boldsymbol{\theta}} \log(P_{\mathcal{D}}(\mathcal{D}|\boldsymbol{\theta})) - \frac{1}{2} \|\boldsymbol{\Sigma}^{-1/2} \cdot (\boldsymbol{\theta} - \boldsymbol{\mu})\|_2^2$$

In the case that the prior mean is zero  $\boldsymbol{\mu} = \mathbf{0}$ ,

$$= \operatorname{argmax}_{\boldsymbol{\theta}} \log(P_{\mathcal{D}}(\mathcal{D}|\boldsymbol{\theta})) - \frac{1}{2} \|\boldsymbol{\Sigma}^{-1/2} \cdot \boldsymbol{\theta}\|_2^2.$$

One is free to choose any suitable expression for  $\log(P_{\mathcal{D}}(\mathcal{D}|\boldsymbol{\theta}))$ . Next we describe the resulting forms for common choices of log likelihood in regression and classification tasks.

#### Regression

For regression, a common choice is that the NN predicts the mean of the function,  $\hat{\mathbf{y}}$ , and there is additive noise on the true targets  $\mathbf{y}$ ,  $P_{\mathcal{D}}(\mathcal{D}|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{y}|\hat{\mathbf{y}}, \sigma_{\epsilon}^2)$

$$\begin{aligned}\boldsymbol{\theta}_{MAP} &= \operatorname{argmax}_{\boldsymbol{\theta}} - \frac{1}{2\sigma_{\epsilon}^2} \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2 + \text{const.} - \frac{1}{2} \|\boldsymbol{\Sigma}^{-1/2} \cdot \boldsymbol{\theta}\|_2^2 \\ &= \operatorname{argmax}_{\boldsymbol{\theta}} - \frac{1}{2\sigma_{\epsilon}^2} \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2 - \frac{1}{2} \|\boldsymbol{\Sigma}^{-1/2} \cdot \boldsymbol{\theta}\|_2^2\end{aligned}$$

Generally the mean squared error is minimised,

$$= \operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{N} \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2 + \frac{1}{N} \|\sigma_{\epsilon} \boldsymbol{\Sigma}^{-1/2} \cdot \boldsymbol{\theta}\|_2^2$$

More compactly, we can define  $\mathbf{\Gamma} := \sigma_{\epsilon}^2 \boldsymbol{\Sigma}^{-1}$ , as a diagonal matrix with,  $\operatorname{diag}(\mathbf{\Gamma})_i = \sigma_{\epsilon}^2 / \sigma_{\text{prior},i}^2$ ,

$$= \operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{N} \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2 + \frac{1}{N} \|\mathbf{\Gamma}^{1/2} \cdot \boldsymbol{\theta}\|_2^2$$

#### Classification

The data likelihood is commonly chosen as a multinomial distribution,  $P_{\mathcal{D}}(\mathcal{D}|\boldsymbol{\theta}) \propto \prod_{n=1}^N \prod_{c=1}^C \hat{y}_{n,c}^{y_{n,c}}$ , for  $C$  classes, and  $N$  data points, where  $\hat{y} \in [0, 1]$  denotes predicted probability, and  $y_{n,c} \in \{0, 1\}$  the true targets.

$$\boldsymbol{\theta}_{MAP} = \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{n=1}^N \sum_{c=1}^C y_{n,c} \log(\hat{y}_{n,c}) - \frac{1}{2} \|\boldsymbol{\Sigma}^{-1/2} \cdot \boldsymbol{\theta}\|_2^2$$

Cross entropy is typically minimised,

$$= \operatorname{argmin}_{\boldsymbol{\theta}} - \sum_{n=1}^N \sum_{c=1}^C y_{n,c} \log(\hat{y}_{n,c}) + \frac{1}{2} \|\boldsymbol{\Sigma}^{-1/2} \cdot \boldsymbol{\theta}\|_2^2$$

Here we can simply define  $\mathbf{\Gamma} := \frac{1}{2} \boldsymbol{\Sigma}^{-1}$ , with  $\operatorname{diag}(\mathbf{\Gamma})_i = 1/2\sigma_{\text{prior},i}^2$ .

## B Numerical Examples

### B.1 Numerical Examples of Proofs

In this section we print covariance matrices that illustrate theoretical results numerically.

#### B.1.1 General case: Example of lemma 1.1, 1.2, theorem 2

For general  $\Sigma_{post}$ , RMS will return a posterior with underestimated marginal variances and overestimated correlations. Since the prior is isotropic, the orientation (eigenvectors) of the RMS approximate posterior will be unchanged. Here a three parameter is shown. Note  $\Sigma_{post}\Sigma_{prior}^{-1}\Sigma_{post}$  represents the RMS approximate posterior.

$$\Sigma_{prior} = \begin{bmatrix} 2.0 & 0.0 & 0.0 \\ 0.0 & 2.0 & 0.0 \\ 0.0 & 0.0 & 2.0 \end{bmatrix} \quad \Sigma_{like} = \begin{bmatrix} 2.0 & 0.707 & 0.283 \\ 0.707 & 1.0 & 0.4 \\ 0.283 & 0.4 & 1.0 \end{bmatrix}$$

$$\Sigma_{post} = \begin{bmatrix} 0.953 & 0.238 & 0.067 \\ 0.238 & 0.589 & 0.166 \\ 0.067 & 0.166 & 0.638 \end{bmatrix} \quad \Sigma_{post}\Sigma_{prior}^{-1}\Sigma_{post} = \begin{bmatrix} 0.485 & 0.189 & 0.073 \\ 0.189 & 0.215 & 0.11 \\ 0.073 & 0.11 & 0.22 \end{bmatrix}$$

Note,  $\Sigma_{post_{i,i}} > \Sigma_{post}\Sigma_{prior}^{-1}\Sigma_{post}_{i,i}, \forall i$ .

$$\text{Correlation}(\Sigma_{post}) = \begin{bmatrix} 1.0 & 0.317 & 0.086 \\ 0.317 & 1.0 & 0.27 \\ 0.086 & 0.27 & 1.0 \end{bmatrix} \quad \text{Correlation}(\Sigma_{post}\Sigma_{prior}^{-1}\Sigma_{post}) = \begin{bmatrix} 1.0 & 0.585 & 0.224 \\ 0.585 & 1.0 & 0.504 \\ 0.224 & 0.504 & 1.0 \end{bmatrix}$$

Note,  $\text{Correlation}(\Sigma_{post})_{i,j} < \text{Correlation}(\Sigma_{post}\Sigma_{prior}^{-1}\Sigma_{post})_{i,j}, \forall i \neq j$ .

Eigenvalues and eigenvectors of  $\Sigma_{post}$ ,

$$\lambda = 1.1101, \mathbf{v} = [-0.8352 \quad -0.471 \quad -0.284]$$

$$\lambda = 0.6667, \mathbf{v} = [-0.465 \quad 0.3288 \quad 0.822]$$

$$\lambda = 0.4032, \mathbf{v} = [0.2938 \quad -0.8186 \quad 0.4936]$$

Eigenvalues and eigenvectors of  $\Sigma_{post}\Sigma_{prior}^{-1}\Sigma_{post}$ ,

$$\lambda = 0.6162, \mathbf{v} = [-0.8352 \quad -0.471 \quad -0.284]$$

$$\lambda = 0.2222, \mathbf{v} = [-0.465 \quad 0.3288 \quad 0.822]$$

$$\lambda = 0.0813, \mathbf{v} = [0.2938 \quad -0.8186 \quad 0.4936]$$

#### B.1.2 Special case: Examples of theorem 3, 4

We again print out covariance matrices for a three parameter model. Firstly we provide an example where two parameters are perfectly correlated, and one has no effect on the likelihood. We print unscaled and scaled versions. Note that all eigenvalues of the scaled posterior,  $\Sigma'_{post}$ , are either 0 or 1.

$$\Sigma_{post} = \begin{bmatrix} 1.0 & 1.0 & 0.0 \\ 1.0 & 1.0 & 0.0 \\ 0.0 & 0.0 & 2.0 \end{bmatrix} \quad \Sigma_{prior} = \begin{bmatrix} 2.0 & 0.0 & 0.0 \\ 0.0 & 2.0 & 0.0 \\ 0.0 & 0.0 & 2.0 \end{bmatrix} \quad \Sigma_{post}\Sigma_{prior}^{-1}\Sigma_{post} = \begin{bmatrix} 1.0 & 1.0 & 0.0 \\ 1.0 & 1.0 & 0.0 \\ 0.0 & 0.0 & 2.0 \end{bmatrix}$$

$$\Sigma'_{post} = \begin{bmatrix} 0.5 & 0.5 & 0.0 \\ 0.5 & 0.5 & 0.0 \\ 0.0 & 0.0 & 1.0 \end{bmatrix} \quad \Sigma'_{prior} = \begin{bmatrix} 1.0 & 0.0 & 0.0 \\ 0.0 & 1.0 & 0.0 \\ 0.0 & 0.0 & 1.0 \end{bmatrix} \quad \Sigma'_{post}\Sigma'^{-1}_{prior}\Sigma'_{post} = \begin{bmatrix} 0.5 & 0.5 & 0.0 \\ 0.5 & 0.5 & 0.0 \\ 0.0 & 0.0 & 1.0 \end{bmatrix}$$

Eigenvalues and eigenvectors of  $\Sigma_{post} = \Sigma_{post}\Sigma_{prior}^{-1}\Sigma_{post}$ ,

$$\begin{aligned}\lambda &= 2.0, \mathbf{v} = [0.7071 \quad 0.7071 \quad 0.] \\ \lambda &= 0.0, \mathbf{v} = [-0.7071 \quad 0.7071 \quad 0.] \\ \lambda &= 2.0, \mathbf{v} = [0. \quad 0. \quad 1.] \end{aligned}$$

Eigenvalues and eigenvectors of  $\Sigma'_{post} = \Sigma'_{post}\Sigma'_{prior}^{-1}\Sigma'_{post}$ ,

$$\begin{aligned}\lambda &= 1.0, \mathbf{v} = [0.707 \quad 0.707 \quad 0.] \\ \lambda &= 0.0, \mathbf{v} = [-0.707 \quad 0.707 \quad 0.] \\ \lambda &= 1.0, \mathbf{v} = [0. \quad 0. \quad 1.] \end{aligned}$$

Following is the same set up as the previous case, but now all parameters are perfectly correlated. Eigenvalues of the scaled posterior,  $\Sigma'_{post}$ , are again either 0 or 1.

$$\begin{aligned}\Sigma_{post} &= \begin{bmatrix} 0.667 & 0.667 & 0.667 \\ 0.667 & 0.667 & 0.667 \\ 0.667 & 0.667 & 0.667 \end{bmatrix} & \Sigma_{prior} &= \begin{bmatrix} 2.0 & 0.0 & 0.0 \\ 0.0 & 2.0 & 0.0 \\ 0.0 & 0.0 & 2.0 \end{bmatrix} & \Sigma_{post}\Sigma_{prior}^{-1}\Sigma_{post} &= \begin{bmatrix} 0.667 & 0.667 & 0.667 \\ 0.667 & 0.667 & 0.667 \\ 0.667 & 0.667 & 0.667 \end{bmatrix} \\ \Sigma'_{post} &= \begin{bmatrix} 0.333 & 0.333 & 0.333 \\ 0.333 & 0.333 & 0.333 \\ 0.333 & 0.333 & 0.333 \end{bmatrix} & \Sigma'_{prior} &= \begin{bmatrix} 1.0 & 0.0 & 0.0 \\ 0.0 & 1.0 & 0.0 \\ 0.0 & 0.0 & 1.0 \end{bmatrix} & \Sigma'_{post}\Sigma'_{prior}^{-1}\Sigma'_{post} &= \begin{bmatrix} 0.333 & 0.333 & 0.333 \\ 0.333 & 0.333 & 0.333 \\ 0.333 & 0.333 & 0.333 \end{bmatrix} \end{aligned}$$

Eigenvalues and eigenvectors of  $\Sigma_{post} = \Sigma_{post}\Sigma_{prior}^{-1}\Sigma_{post}$ ,

$$\begin{aligned}\lambda &= 2.0, \mathbf{v} = [-0.5774 \quad -0.5774 \quad -0.5774] \\ \lambda &= 0.0, \mathbf{v} = [-0. \quad -0.7071 \quad 0.7071] \\ \lambda &= 0.0, \mathbf{v} = [-0.6667 \quad -0.0749 \quad 0.7416] \end{aligned}$$

Eigenvalues and eigenvectors of  $\Sigma'_{post} = \Sigma'_{post}\Sigma'_{prior}^{-1}\Sigma'_{post}$ ,

$$\begin{aligned}\lambda &= 1.0, \mathbf{v} = [0.577 \quad 0.577 \quad 0.577] \\ \lambda &= 0.0, \mathbf{v} = [0. \quad -0.707 \quad 0.707] \\ \lambda &= 0.0, \mathbf{v} = [-0.521 \quad -0.284 \quad 0.805] \end{aligned}$$

Now we detail an example of a non-isometric prior.

$$\begin{aligned}\Sigma_{post} &= \begin{bmatrix} 1.818 & 0.0 & 1.818 \\ 0.0 & 2.0 & 0.0 \\ 1.818 & 0.0 & 1.818 \end{bmatrix} & \Sigma_{prior} &= \begin{bmatrix} 20.0 & 0.0 & 0.0 \\ 0.0 & 2.0 & 0.0 \\ 0.0 & 0.0 & 2.0 \end{bmatrix} & \Sigma_{post}\Sigma_{prior}^{-1}\Sigma_{post} &= \begin{bmatrix} 1.818 & 0.0 & 1.818 \\ 0.0 & 2.0 & 0.0 \\ 1.818 & 0.0 & 1.818 \end{bmatrix} \\ \Sigma'_{post} &= \begin{bmatrix} 0.091 & 0.0 & 0.287 \\ 0.0 & 1.0 & 0.0 \\ 0.287 & 0.0 & 0.909 \end{bmatrix} & \Sigma'_{prior} &= \begin{bmatrix} 1.0 & 0.0 & 0.0 \\ 0.0 & 1.0 & 0.0 \\ 0.0 & 0.0 & 1.0 \end{bmatrix} & \Sigma'_{post}\Sigma'_{prior}^{-1}\Sigma'_{post} &= \begin{bmatrix} 0.091 & 0.0 & 0.287 \\ 0.0 & 1.0 & 0.0 \\ 0.287 & 0.0 & 0.909 \end{bmatrix} \end{aligned}$$

Eigenvalues and eigenvectors of  $\Sigma_{post} = \Sigma_{post} \Sigma_{prior}^{-1} \Sigma_{post}$ ,

$$\begin{aligned} \lambda &= 3.636, \mathbf{v} = [0.707 \quad 0. \quad 0.707] \\ \lambda &= 0.0, \mathbf{v} = [-0.707 \quad 0. \quad 0.707] \\ \lambda &= 2.0, \mathbf{v} = [0. \quad 1. \quad 0.] \end{aligned}$$

Eigenvalues and eigenvectors of  $\Sigma'_{post} = \Sigma'_{post} \Sigma'_{prior}^{-1} \Sigma'_{post}$ ,

$$\begin{aligned} \lambda &= 0.0, \mathbf{v} = [-0.953 \quad 0. \quad 0.302] \\ \lambda &= 1.0, \mathbf{v} = [-0.302 \quad 0. \quad -0.953] \\ \lambda &= 1.0, \mathbf{v} = [0. \quad 1. \quad 0.] \end{aligned}$$

## B.2 Mixtures of Parameter Types

Here, we provide examples of a five parameter model containing a mixture of perfectly correlated, partially correlated, and extrapolation parameters.

First we consider distinct blocks of perfectly and partially correlated parameters, as well as one extrapolation parameter. In this situation both the perfectly correlated block and the extrapolation parameter posterior are recovered exactly. The RMS approximate posterior of the partially correlated block is biased as per the general case.

$$\Sigma_{post} = \begin{bmatrix} 1.0 & 1.0 & 0.0 & 0.0 & 0.0 \\ 1.0 & 1.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.5 & 0.2 & 0.0 \\ 0.0 & 0.0 & 0.2 & 0.8 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 2.0 \end{bmatrix} \quad \Sigma_{prior} = \begin{bmatrix} 2 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 & 2 \end{bmatrix} \quad \Sigma_{post} \Sigma_{prior}^{-1} \Sigma_{post} = \begin{bmatrix} 1.0 & 1.0 & 0.0 & 0.0 & 0.0 \\ 1.0 & 1.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.145 & 0.13 & 0.0 \\ 0.0 & 0.0 & 0.13 & 0.34 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 2.0 \end{bmatrix}$$

Secondly the perfectly correlated block overlaps with the partially correlated block. In this scenario, a small amount of bias is introduced on the perfectly correlated block, but not in terms of the correlation. The extrapolation parameter is unaffected.

$$\Sigma_{post} = \begin{bmatrix} 1.0 & 1.0 & 0.1 & 0.2 & 0.0 \\ 1.0 & 1.0 & 0.1 & 0.2 & 0.0 \\ 0.1 & 0.1 & 0.5 & 0.2 & 0.0 \\ 0.2 & 0.2 & 0.2 & 0.8 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 2.0 \end{bmatrix} \quad \Sigma_{prior} = \begin{bmatrix} 2 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 & 2 \end{bmatrix} \quad \Sigma_{post} \Sigma_{prior}^{-1} \Sigma_{post} = \begin{bmatrix} 1.03 & 1.03 & 0.15 & 0.29 & 0.0 \\ 1.03 & 1.03 & 0.15 & 0.29 & 0.0 \\ 0.15 & 0.15 & 0.16 & 0.15 & 0.0 \\ 0.29 & 0.29 & 0.15 & 0.38 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 2.0 \end{bmatrix}$$

### B.3 Example of Perfectly Correlated Parameters in a Neural Network

Section 4.2.2 stated that perfectly correlated parameters can exist in a NN. Here, we provide a concrete example of such a case, and show that anchored ensembling *does* recover the true posterior for these parameters when  $\Sigma_{anc} = \Sigma_{prior}$ .

We consider finding the posterior of final layer weights in a small single layer ReLU NN of two hidden nodes for a regression problem with two data points. (Choosing the final layer weights for our analysis allows analytical equations associated with linear regression to be used, simplifying our analysis, though these results would also apply to the first layer weights and biases.)

We design the problem such that the point where both hidden nodes becomes greater than zero (the elbow points of ReLU units) falls in between the two data points, and the active half of the output is also shared, so that the final layer weights are perfectly correlated.

Figure 9 illustrates our set up. Data points and NN parameters are as follows,

$$\mathcal{D} = \{x_1 = -5, y_1 = 0; x_2 = 5, y_2 = 0\}, \sigma_\epsilon^2 = 0.01,$$

$$\mathbf{W}_1 = [-0.8, -0.4], \mathbf{b}_1 = [-1, 0.1], \hat{y} = \mathbf{W}_2^T \max(x\mathbf{W}_1 + \mathbf{b}_1, 0)$$

We set prior means to zero, with isotropic covariance according to  $1/H$ ,

$$\Sigma_{prior} = \begin{bmatrix} 0.5 & 0.0 \\ 0.0 & 0.5 \end{bmatrix}$$

We print out matrices of interest,

$$\Sigma_{like} = \begin{bmatrix} -6.89e + 13 & 9.85e + 13 \\ 9.85e + 13 & -1.40e + 14 \end{bmatrix}$$

$$\Sigma_{like}^{-1} = \begin{bmatrix} 90.0 & 63.0 \\ 63.0 & 44.1 \end{bmatrix}$$

$$\Sigma_{post} = \begin{bmatrix} 0.169 & -0.231 \\ -0.231 & 0.338 \end{bmatrix}$$

$$\Sigma_{post}\Sigma_{prior}^{-1}\Sigma_{post} = \begin{bmatrix} 0.166 & -0.235 \\ -0.235 & 0.338 \end{bmatrix}$$

The anchored ensembling posterior,  $\Sigma_{post}\Sigma_{prior}^{-1}\Sigma_{post}$ , provides a close approximation of the true posterior covariance (and would be exact discounting numerical rounding issues). Note the similarity of the posterior in figure 9 (middle) with the perfect correlations example shown in figure 3 (B).

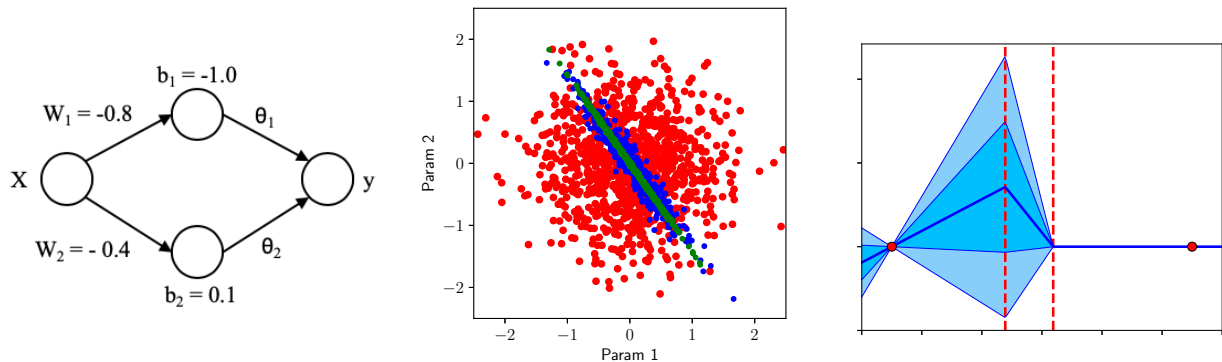


Figure 9: Left: Single layer NN of two hidden nodes. Middle: Draws in parameter space for prior (red), analytical posterior (blue) and anchored posterior (green). Right: Posterior predictive distribution - dashed red lines are elbows of ReLU units.

## C Further Results

### C.1 Uncertainty-Aware Model-Free Reinforcement Learning

An anchored ensemble of 5xNNs, each with two hidden layers, was trained to complete a discretised version of FetchPush - an agent controls a robotic arm, with rewards received when a randomly placed cube is pushed to a goal.

We used Bayesian Q-learning (Dearden et al., 1998), similar to regular Q-learning, but with Q-values modelled as *distributions* rather than point estimates - the wider the distribution, the less certain the agent. This is beneficial both to drive the exploration/exploitation process via Thompson sampling, and for identifying OOD examples.

Figure 10 shows the agent’s awareness of its uncertainty. After training for 40,000 episodes, its confidence over actions was plotted for three scenarios: A) Cube and goal are in positions often encountered during training, the agent has learnt that it must move the arm left - the narrow distributions with significantly different means reflect its confidence in this. B) The goal has already been achieved - narrow overlapping distributions with higher means. C) A peculiar goal position that has never been encountered - the broad distributions over all actions reflect its high uncertainty.

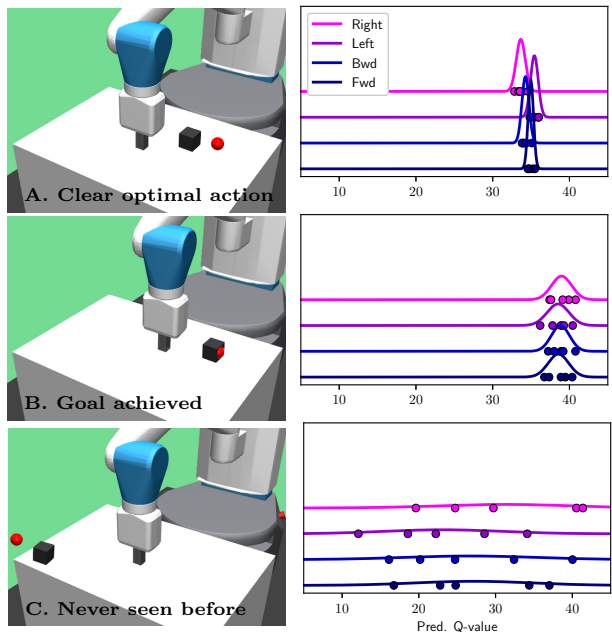


Figure 10: Anchored ensembling creates uncertainty-aware agents.

### C.2 Model-Based Reinforcement: Learning in Noisy Environments

We tested the benefit of using an anchored ensemble in noisy RL environments. We modified the classic cartpole swingup environment such that different levels of stochastic noise could be added to the future state; consider a state action pair given by,  $s, a$ , and a noisy state,  $\bar{s}$  such that  $P(\bar{s}_{t+1}|a_t, s_t) = \mathcal{N}(s_{t+1}, \sigma_\epsilon^2)$ .

The value of  $\sigma_\epsilon^2$  was given three settings: low, medium and high,  $\sigma_\epsilon^2 \in \{0.001, 0.002, 0.005\}$ . The task was learnt using a model-based RL approach similar to the heteroskedastic ensembles used by Chua et al. (2018). Fully-connected three-layer NNs learnt to predict the dynamics of the environment given some state and action. Planning was performed using the cross-entropy method, rolling out for a horizon of 25 steps, with 10 particles.

Figure 11 (A, B, C) shows learning curves for the three noise levels. All ensemble techniques perform similarly in the low noise setting. As noise is increased the overall performance of all methods drops. Anchored ensembles are most resilient, followed by the unconstrained ensemble. In panel D, we compared against an unconstrained ensemble employing early stopping - this corresponds to the proposal that applying early stopping to an ensemble can produce approximate inference (Duvenaud et al., 2016). Whilst careful tuning did offer some improvement over the default setting, a performance gap to anchored ensembling remained.

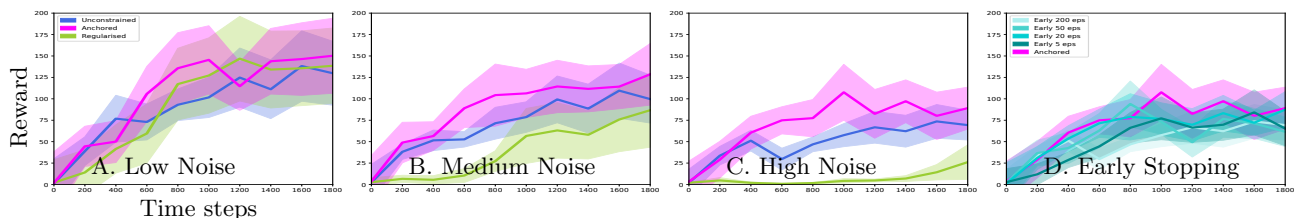


Figure 11: Anchored ensembling creates robust MBRL agents in noisy environments. Mean and standard error over five runs.

### C.3 Regression Benchmarking

Tables 3 & 4 show all experiments run on the regression benchmarking datasets. The below discussion focuses on NLL results in table 4.

ERF GP refers to the equivalent GP for an infinite width, single-layer BNN with ERF activations. It was tuned and implemented as for the ReLU GP. We were interested to discover how different activation functions would affect uncertainty estimates. In general the ReLU GP performed better than the ERF GP, with some exceptions, such as for Wine. The target variable for Wine is ordinal, containing five factors, it is therefore understandable that the ReLU GP, which extrapolates linearly, is at a slight disadvantage.

10x 50 NNs refers to an anchored ensemble of ten NNs with 50 hidden nodes. We find that these results fall in between the 5x 50 NNs and the ReLU GP. This agrees with the convergence analysis done in section 5.2.

We also implemented an anchored ensemble of five two-layer NNs, 5x 50-50 NNs. Even with minimal hyperparameter tuning (section E) we found an extra layer gave a performance boost over the 5x 50 NNs. We expect with more careful tuning this margin would increase.

Single 50 NN refers to a single regularised NN, of one hidden layer with 50 hidden nodes, for which we used a constant value of predictive variance. Although this performs poorly in several cases, e.g. Boston and Yacht, the results are surprisingly close to those achieved by both our method and Deep Ensembles, even surpassing them on the Energy dataset. A method outputting constant predictive variance should not perform well in experiments designed to test uncertainty quantification, and this raises questions over the validity of the benchmarks.

Table 5 compares anchored ensembles against results reported for other methods.

Table 3: Variants of our method on benchmark regression datasets, RMSE.

	$N$	$D$	RMSE					
			ReLU GP	ERF GP	5x 50 NNs	10x 50 NNs	5x 50-50 NNs	Single 50 NN
Boston	506	13	2.86 ± 0.16	2.94 ± 0.18	3.09 ± 0.17	3.09 ± 0.17	3.00 ± 0.18	3.40 ± 0.20
Concrete	1,030	8	4.88 ± 0.13	5.21 ± 0.12	4.87 ± 0.11	4.73 ± 0.11	4.75 ± 0.12	5.17 ± 0.13
Energy	768	8	0.60 ± 0.02	0.78 ± 0.03	0.35 ± 0.01	0.34 ± 0.01	0.40 ± 0.01	0.40 ± 0.01
Kin8nm	8,192	8	0.07 ± 0.00	0.08 ± 0.00	0.07 ± 0.00	0.07 ± 0.00	0.06 ± 0.00	0.07 ± 0.00
Naval	11,934	16	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
Power	9,568	4	3.97 ± 0.04	3.94 ± 0.04	4.07 ± 0.04	4.07 ± 0.04	4.03 ± 0.04	4.23 ± 0.04
Protein	45,730	9	4.34 ± 0.02	4.23 ± 0.02	4.36 ± 0.02	4.34 ± 0.02	4.23 ± 0.02	4.56 ± 0.02
Wine	1,599	11	0.61 ± 0.01	0.60 ± 0.01	0.63 ± 0.01	0.62 ± 0.01	0.62 ± 0.01	0.64 ± 0.01
Yacht	308	6	0.60 ± 0.08	1.48 ± 0.15	0.57 ± 0.05	0.54 ± 0.05	0.85 ± 0.08	0.81 ± 0.07
Song Year	515,345	90	9.01 ± NA	8.90 ± NA	8.82 ± NA	8.82 ± NA	8.66 ± NA	8.77 ± NA

Table 4: Variants of our method on benchmark regression datasets, NLL.

	$\hat{\sigma}_\epsilon^2$	NLL					
		ReLU GP	ERF GP	5x 50 NNs	10x 50 NNs	5x 50-50 NNs	Single 50 NN
Boston	0.08	2.45 ± 0.05	2.46 ± 0.05	2.52 ± 0.05	2.50 ± 0.05	2.50 ± 0.07	2.70 ± 0.05
Concrete	0.05	2.96 ± 0.02	3.06 ± 0.02	2.97 ± 0.02	2.94 ± 0.02	2.94 ± 0.02	3.08 ± 0.03
Energy	1e-7	0.86 ± 0.02	1.06 ± 0.03	0.96 ± 0.13	0.52 ± 0.06	0.61 ± 0.07	0.57 ± 0.03
Kin8nm	0.02	-1.22 ± 0.01	-1.17 ± 0.00	-1.09 ± 0.01	-1.16 ± 0.01	-1.25 ± 0.01	-1.17 ± 0.01
Naval	1e-7	-10.05 ± 0.02	-9.66 ± 0.04	-7.17 ± 0.03	-7.29 ± 0.02	-7.08 ± 0.13	-6.58 ± 0.04
Power	0.05	2.80 ± 0.01	2.79 ± 0.01	2.83 ± 0.01	2.83 ± 0.01	2.82 ± 0.01	2.86 ± 0.01
Protein	0.5	2.88 ± 0.00	2.86 ± 0.00	2.89 ± 0.01	2.88 ± 0.01	2.86 ± 0.01	2.94 ± 0.00
Wine	0.5	0.92 ± 0.01	0.91 ± 0.01	0.95 ± 0.01	0.94 ± 0.01	0.94 ± 0.01	0.97 ± 0.01
Yacht	1e-7	0.49 ± 0.07	1.50 ± 0.13	0.37 ± 0.08	0.18 ± 0.03	0.04 ± 0.08	1.50 ± 0.02
Song Year	0.7	3.62 ± NA	3.61 ± NA	3.60 ± NA	3.60 ± NA	3.57 ± NA	3.59 ± NA



Table 5: Comparison against inference methods on UCI benchmark regression datasets, log likelihood. Adapted from Mukhoti et al. (2018).

	Log Likelihood ( <i>not negative</i> )							
	Anch. Ens.	Drop conv.	Drop tune	VMG	HS-BNN	PBP-MV	SGHMC tune	SGHMC adap.
Boston	-2.52 ± 0.05	<b>-2.40 ± 0.04</b>	<b>-2.40 ± 0.04</b>	-2.46 ± 0.09	-2.54 ± 0.15	-2.54 ± 0.08	-2.49 ± 0.15	-2.54 ± 0.04
Concrete	-2.97 ± 0.02	-2.97 ± 0.02	<b>-2.93 ± 0.02</b>	-3.01 ± 0.03	-3.09 ± 0.06	-3.04 ± 0.03	-4.17 ± 0.72	-3.38 ± 0.24
Energy	<b>-0.96 ± 0.13</b>	-1.72 ± 0.01	-1.21 ± 0.01	-1.06 ± 0.03	-2.66 ± 0.13	-1.01 ± 0.01	---	---
Kim8mm	1.09 ± 0.01	0.97 ± 0.00	1.14 ± 0.01	1.10 ± 0.01	1.12 ± 0.03	<b>1.28 ± 0.01</b>	---	---
Naval	<b>7.17 ± 0.03</b>	3.91 ± 0.01	4.45 ± 0.00	2.46 ± 0.00	5.52 ± 0.10	4.85 ± 0.06	---	---
Power	-2.83 ± 0.01	-2.79 ± 0.01	-2.80 ± 0.01	-2.82 ± 0.01	-2.81 ± 0.03	<b>-2.78 ± 0.01</b>	---	---
Protein	-2.89 ± 0.01	-2.87 ± 0.00	-2.87 ± 0.00	-2.84 ± 0.00	-2.89 ± 0.00	<b>-2.77 ± 0.01</b>	---	---
Wine	-0.95 ± 0.01	<b>-0.92 ± 0.01</b>	-0.93 ± 0.01	-0.95 ± 0.01	-0.95 ± 0.05	-0.97 ± 0.01	-1.29 ± 0.28	-1.04 ± 0.17
Yacht	<b>-0.37 ± 0.08</b>	-1.38 ± 0.01	-1.25 ± 0.01	-1.30 ± 0.02	-2.33 ± 0.01	-1.64 ± 0.02	-1.75 ± 0.19	<b>-1.10 ± 0.08</b>

### C.4 Fashion MNIST

Table 6 provides a breakdown of results from the OOD classification test in section 5.4 for fashion MNIST. Also included are results for entropy, where high entropy represents high uncertainty. These correlated strongly with the proportion metrics, which was true across all three OOD experiments.

Table 6: Fashion MNIST results: proportion of predictions made with  $\geq 90\%$  probability, and entropy of predicted categorical distribution. Also shown is relative advantage (percentage change) for each method compared to anchored ensembles. Averaged over five runs/random seeds, mean  $\pm$  1 standard error. Best result in blue.

	—Edge Cases—			—Out-of-distribution—		—Natural Adversarial—			—Pure Adversarial—	
	Train	Sneaker	Trouser	CIFAR	MNIST	Rotate	Flip	Invert	Noise	Sparse
reg 1xNN	0.660 ± 0.006	0.739 ± 0.056	0.429 ± 0.047	0.143 ± 0.008	0.160 ± 0.007	0.609 ± 0.007	0.330 ± 0.009	0.349 ± 0.015	0.271 ± 0.007	0.456 ± 0.006
free 5xNN	0.733 ± 0.001	0.781 ± 0.015	0.380 ± 0.030	0.301 ± 0.013	0.104 ± 0.010	0.571 ± 0.011	0.300 ± 0.011	0.222 ± 0.052	0.042 ± 0.005	0.048 ± 0.003
reg 5xNN	0.634 ± 0.002	0.589 ± 0.054	<b>0.269 ± 0.020</b>	0.115 ± 0.004	0.072 ± 0.007	0.556 ± 0.007	0.256 ± 0.012	0.213 ± 0.002	0.112 ± 0.005	0.174 ± 0.005
anc 5xNN	0.631 ± 0.002	<b>0.578 ± 0.049</b>	0.325 ± 0.037	<b>0.065 ± 0.002</b>	<b>0.041 ± 0.002</b>	<b>0.497 ± 0.003</b>	<b>0.215 ± 0.005</b>	<b>0.025 ± 0.010</b>	<b>0.006 ± 0.001</b>	<b>0.006 ± 0.001</b>
Proportion Relative Advantage										
1xNN Reg. to 5xNN Anch.	-4.4%	-21.8%	-24.2%	-54.5%	-74.4%	-18.4%	-34.8%	-92.8%	-97.8%	-98.7%
5xNN Uncons. to 5xNN Anch.	-13.9%	-26.0%	-14.5%	-78.4%	-60.6%	-13.0%	-28.3%	-88.7%	-85.7%	-87.5%
5xNN Reg. to 5xNN Anch.	-0.5%	-1.9%	20.8%	-43.5%	-43.1%	-10.6%	-16.0%	-88.3%	-94.6%	-96.6%
Entropy (larger better)										
1xNN Reg.	0.328 ± 0.005	0.253 ± 0.043	0.575 ± 0.050	1.176 ± 0.010	0.984 ± 0.015	0.484 ± 0.008	0.713 ± 0.009	0.836 ± 0.035	0.808 ± 0.010	0.580 ± 0.008
5xNN Uncons.	0.230 ± 0.001	0.161 ± 0.010	0.535 ± 0.021	0.688 ± 0.009	1.016 ± 0.021	0.453 ± 0.011	0.685 ± 0.011	0.573 ± 0.037	1.036 ± 0.014	0.992 ± 0.012
5xNN Reg.	0.352 ± 0.001	<b>0.365 ± 0.039</b>	<b>0.707 ± 0.019</b>	1.239 ± 0.009	1.161 ± 0.012	0.564 ± 0.008	0.807 ± 0.017	1.014 ± 0.014	1.048 ± 0.008	0.919 ± 0.009
5xNN Anch.	0.349 ± 0.001	0.327 ± 0.034	0.623 ± 0.042	<b>1.251 ± 0.011</b>	<b>1.295 ± 0.013</b>	<b>0.624 ± 0.006</b>	<b>0.868 ± 0.002</b>	<b>1.098 ± 0.035</b>	<b>1.238 ± 0.013</b>	<b>1.191 ± 0.014</b>
Entropy Relative Advantage										
1xNN Reg. to 5xNN Anch.	6.4%	29.2%	8.3%	6.4%	31.6%	28.9%	21.7%	31.3%	53.2%	105.3%
5xNN Uncons. to 5xNN Anch.	51.7%	103.1%	16.4%	81.8%	27.5%	37.7%	26.7%	91.6%	19.5%	20.1%
5xNN Reg. to 5xNN Anch.	-0.9%	-10.4%	-11.9%	1.0%	11.5%	10.6%	7.6%	8.3%	18.1%	29.6%

## D Additional Material

### D.1 Algorithms

---

**Algorithm 1** Implementing anchored ensembles of NNs

---

**Input:** Training data,  $\mathbf{X}$  &  $\mathbf{Y}$ , test data point,  $\mathbf{x}^*$ , prior mean and covariance,  $\boldsymbol{\mu}_{prior}$ ,  $\boldsymbol{\Sigma}_{prior}$ , ensemble size,  $M$ , data noise variance estimate,  $\hat{\sigma}_\epsilon^2$  (regression only).

**Output:** Estimate of mean and variance,  $\hat{\mathbf{y}}$ ,  $\hat{\sigma}_y^2$  for regression, or class probabilities,  $\hat{\mathbf{y}}$  for classification.

*# Set regularisation matrix*

$\boldsymbol{\Gamma} \leftarrow \hat{\sigma}_\epsilon^2 \boldsymbol{\Sigma}_{prior}^{-1}$  (regression)   OR    $\boldsymbol{\Gamma} \leftarrow \frac{1}{2} \boldsymbol{\Sigma}_{prior}^{-1}$  (classification)

*# Create ensemble*

$\boldsymbol{\mu}_{anc} \leftarrow \boldsymbol{\mu}_{prior}$ ,  $\boldsymbol{\Sigma}_{anc} \leftarrow \boldsymbol{\Sigma}_{prior}$

**for**  $j = 1$  **to**  $M$

$\boldsymbol{\theta}_{anc,j} \sim \mathcal{N}(\boldsymbol{\mu}_{anc}, \boldsymbol{\Sigma}_{anc})$  *# Sample anchor points*

$NN_j.create(\boldsymbol{\Gamma}, \boldsymbol{\theta}_{anc,j})$  *# Create custom regulariser*

$NN_j.initialise()$  *# Initialisations independent of  $\boldsymbol{\theta}_{anc,j}$*

*# Train ensemble*

**for**  $j = 1$  **to**  $M$

$NN_j.train(\mathbf{X}, \mathbf{Y})$ , loss in eq. 8 (regression) or eq. 9 (classification) or eq. 7 (custom)

*# Predict with ensemble*

**for**  $j = 1$  **to**  $M$

$\hat{\mathbf{y}}_j \leftarrow NN_j.predict(\mathbf{x}^*)$

*# Regression - combine ensemble estimates*

$\hat{\mathbf{y}} = \frac{1}{M} \sum_{j=1}^M \hat{\mathbf{y}}_j$ , *# Mean prediction*

$\hat{\sigma}_{model}^2 = \frac{1}{M-1} \sum_{j=1}^M (\hat{\mathbf{y}}_j - \hat{\mathbf{y}})^2$  *# Epistemic var.*

$\hat{\sigma}_y^2 = \hat{\sigma}_{model}^2 + \hat{\sigma}_\epsilon^2$  *# Total var. = epistemic + data noise*

*# Classification - combine ensemble estimates*

$\hat{\mathbf{y}} = \frac{1}{M} \sum_{j=1}^M \hat{\mathbf{y}}_j$ , *# Average softmax output*

$\hat{\sigma}_y^2 = \text{None}$  *# N/A for classification*

**return**  $\hat{\mathbf{y}}$ ,  $\hat{\sigma}_y^2$

---

## E Experimental Details

### E.1 Introduction to Anchored Ensembles

Experimental details for figure 1 are as follows.

Six randomly generated data points were used.

Hyperparameters: activation = ERF,  $\sigma_\epsilon^2 = 0.003$ ,  $b_1$  variance = 1,  $W_1$  variance = 1,  $H = 100$ ,  $M = 3$  (number of ensembles), optimiser = adam, epochs = 400, learning rate = 0.005.

### E.2 Panel of Inference Methods

Experimental details for figure 4 are as follows.

Same six data points were used for all methods and activation functions, generated by  $y = x \sin(5x)$ , evaluated at, [-0.8, -0.1, 0.02, 0.2, 0.6, 0.8].

Hyperparameters:  $b_1$  variance = 10,  $W_1$  variance = 10,  $H = 100$ ,  $M = 10$ , epochs= 4,000,  $\sigma_\epsilon^2 = 0.001$ , leaky ReLU  $\alpha = 0.2$ , optimiser = adam, MC Dropout probability = 0.4, MC Dropout samples = 200, HMC step size = 0.001, HMC no. steps = 150, HMC burn in = 500, HMC total samples = 1000, HMC predict samples = 50, VI predict samples = 50, VI iterations = 2000, VI gradient samples = 200.

### E.3 Ensembling Loss Functions

Experimental details for figure 6 are as follows.

#### E.3.1 Regression

Generated  $\mathbf{X}$  by sampling 20 points linearly spaced from the interval [-1.5, 1.5],  $y = \sin(2x) + \epsilon$  with  $\epsilon \sim \mathcal{N}(0, 0.2^2)$ . The  $y$  value corresponding to the largest  $x$  value was shifted -0.4 to produce a slight outlier.

Sub-plot A was trained via mean square error, B was regularised, C was anchored. D shows a ReLU GP.

Hyperparameters: activation = ReLU,  $\sigma_\epsilon^2 = 0.08$ ,  $b_1$  variance = 10,  $W_1$  variance = 10,  $H = 1000$ , optimiser = adam, epochs = 2,000, learning rate = 0.003,  $M = 10$ , hidden layers = 1.

#### E.3.2 Classification

Generated  $\mathbf{X}$  using sklearn’s ‘make blobs’ function, n samples = 30.

Sub-plot A was trained via cross entropy, B was regularised, C was anchored. D shows inference with HMC.

Hyperparameters: activation = ReLU,  $b_1$  variance = 15/2,  $W_1$  variance = 15/2,  $b_2$  variance = 1/50,  $W_2$  variance = 1/50,  $W_3$  variance = 10/50,  $H = 50$ , optimiser = adam, epochs = 100, learning rate = 0.001,  $M = 10$ , hidden layers = 2.

### E.4 1-D Convergence Plots

Experimental details for figure 8 are as follows.

Data as in section E.2 was used, with  $M = [3,5,10,20]$ .

Hyperparameters: activation = ReLU,  $\sigma_\epsilon^2 = 0.001$ ,  $b_1$  variance = 20,  $W_1$  variance = 20,  $H = 100$ , optimiser = adam, epochs = 4,000, learning rate = 0.005.

### E.5 KL Convergence Results

Experimental details for figure 7 are as follows.

Training was done on 50% of the data, with KL computed over the other 50%. Results were averaged over ten runs. The ‘ideal’ line shows the metric when posterior samples from the GP itself, rather than anchored NNs,

were used.

The Boston Housing dataset was used, with 50% of data used for training, and testing on the other 50%.

Hyperparameters: activation = ReLU,  $\sigma_\epsilon^2 = 0.1$ ,  $b_1$  variance = 2,  $W_1$  variance = 2,  $H = [4, 16, 64, 256, 1024]$ ,  $M = [3, 5, 10, 20, 40]$ , optimiser = adam, no. runs = 10, epochs = 1,000, learning rate = 0.001 when  $H < 20$  else learning rate = 0.0002.

## E.6 Regression Benchmarking Experiments

We complied with the established protocol (Hernández-Lobato and Adams, 2015). Single-layer NNs of 50 nodes were used, experiments repeated 20 times with random train/test splits of 90%/10%. The larger Protein and Song datasets allow 100 node NNs, and were repeated five and one time respectively.

The hyperparameter tuning process and final settings for experiments in table 1, 3 & 4 are as follows.

### E.6.1 Hyperparameter Tuning

Hyperparameter tuning was done on a single train/validation split of 80%/20%. We found it convenient to begin by tuning data noise variance and prior variances. We restricted the prior variance search space by enforcing,  $\sigma_{W_1}^2 = \sigma_{b_1}^2 / D$ , and  $\sigma_{W_2}^2 = 1/H$ . We therefore had only two hyperparameters to optimise initially:  $\sigma_{b_1}^2$  and  $\sigma_\epsilon^2$ . We did this with the GP model, using grid search, maximising marginal log likelihood over the training portion, and minimising NLL of the validation portion. For the larger datasets, when inference over the 80% training portion was too slow, we reduced the training split to 2,000 data points.

Hyperparameters for priors and data noise estimates were shared between the GP and anchored ensembles. Hyperparameters requiring tuning specifically for anchored ensembles were batch size, learning rate, number of epochs and decay rate. This was done on the same 80%/20% split used to select data noise and prior variance. We used random search, directed by our knowledge of the optimisation process (e.g. a lower learning rate requires more epochs to converge), minimising NLL on the validation portion.

We did not retune hyperparameters from scratch for the double layer NN (5x 50-50 NNs). We used settings as for the single-layer NNs (5x 50 NNs), but divided learning rate by 4, and multiplied epochs by 1.5.

For the single regularised NN with constant noise, we again used hyperparameters as for the single-layer ensemble (5x 50 NNs), tuning only the constant amount of variance to be added on the same 80%/20% split.

### E.6.2 Hyperparameter Settings

Table 7 provides the key hyperparameters used. The adam optimiser was used for all experiments. ReLU activations were used for all except the ERF GP (prior variance was separately tuned for this, values aren't given in the table).

Table 7: Hyperparameters used for regression benchmark results.

	$N$	Batch Size	Learn Rate	$\hat{\sigma}_\epsilon^2$	$b_1$ variance	$W_1$ variance	No. Epochs	Decay Rate	Single NN var.
Boston	506	64	0.05	0.06	10	0.77	3000	0.995	0.45
Concrete	1,030	64	0.05	0.05	40	5.00	2000	0.997	0.28
Energy	768	64	0.05	1e-7	12	1.50	2000	0.997	0.03
Kin8nm	8,192	256	0.10	0.02	40	5.00	2000	0.998	0.32
Naval	11,934	256	0.10	1e-7	200	12.50	1000	0.997	0.03
Power	9,568	256	0.20	0.05	4	1.00	1000	0.995	0.24
Protein	45,730	8192	0.10	0.5	50	5.56	3000	0.995	0.71
Wine	1,599	64	0.05	0.5	20	1.82	500	0.997	0.77
Yacht	308	64	0.05	1e-7	15	2.50	3000	0.997	0.10
Song Year	515,345	32768	0.01	0.7	2	0.02	500	0.996	0.84

## E.7 Out-of-Distribution Classification

### E.7.1 Fashion MNIST

We trained a three-layer NN on eight of ten classes of Fashion MNIST. We trained on 48,000 examples, tested on 8,000.

Experiments were repeated 5 times with a different random seed for each run.

Data categories were created as suggested by their name in table 6. Examples are shown in figure 12.

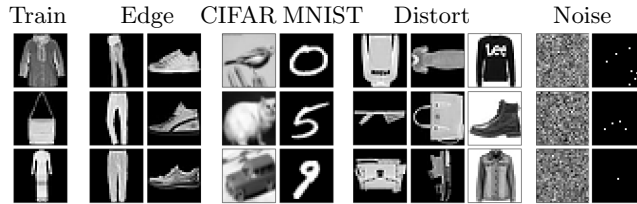


Figure 12: Fashion MNIST OOD data examples.

- **Distort** comprised of rotations, vertical flips, and pixel value inversions.
- **Noise** comprised of iid Gaussian noise, mean = 0.0, standard deviation = 2.0.
- **Sparse** comprised of iid Bernoulli noise, pixels were given a value of 50.0 with  $p = 0.005$ , else 0.0.

Hyperparameters: activation = ReLU, optimiser = adam, epochs = 30, learning rate = 0.005, batch size = 256, hidden layers = 3, hidden units = 100

### E.7.2 CIFAR-10

CIFAR-10 contains 50,000 32x32 color training images, labelled over 10 categories, and 10,000 test images.

We removed 2 categories during training (ships, dogs) so trained over 40,000 examples.

OOD data classes are as show in the images in table 2.

- **Scramble** permuted each row of pixels in a given image.
- **Invert** took the negative of the pixel values.
- **Noise** sampled pixels from bernoulli distribution ( $p=0.005$ ) of large magnitude (pixel value=50).

NN architecture: A convolutional NN was used, with the following structure, 64-64-maxpool-128-128-maxpool-256-256-256-maxpool-512-512-512-maxpool-flatten-2048fc-softmax.

All convolutional kernels were  $[3 \times 3 \times \text{number of channels in previous layer}]$ . All maxpooling kernels were  $[2 \times 2]$ . The total number of parameters was 8,689,472.

Hyperparameters: activation = ReLU, optimiser = adam, learning rate = 0.001 decreasing to 0.0005 after 10 epochs and to 0.0001 after 20 epochs, batch size = 300.

In order to bring test accuracies and confidence on the training dataset roughly in line, it was necessary to train for a different number of training epochs for each method (this effectively applies early stopping to the unconstrained case). Anchored eps = 25, Regularise eps = 30, Unconstrained eps = 15.

Experiments were repeated 3 times with a different random seed for each run.

### E.7.3 IMDb

Dataset of 25,000 movie reviews, labelled as positive or negative.

Example: *“this movie is the best horror movie bar none i love how stanley just dumps the women into the lake i have been a fan of judd nelson’s work for many years and he blew me away its a blend of horror and ... ”*

OOD data classes were generated as follows.

- **Reuters** - taken from the Reuters news dataset.

Example: *“said it has started talks on the possible ... of the company with various parties that it did not identify the company said the talks began after it ... ”*

- **Random 1** - A single integers sampled uniformly at random from {1...vocabulary size} and converted to a repeated sequence of words.

Example: *“member member member member member member member member member member member member member member ... ”*

- **Random 2** - One integer per word sampled uniformly at random from {1...vocabulary size} and converted to words.

Example: *“twists mentally superb finest will dinosaur variety models stands knew refreshing member spock might mode lose leonard resemble began happily names... ”*

- **Random 3** - As for Random 2, but now only sample from least commonly used 100 words.

Example: *“computers towers bondage braveheart threatened rear triangle refuse detectives hangs bondage firmly btw token 1990s mermaid reeves landed dylan remove hum natives insightful demonic... ”*

NN architecture: used an embedding layer (outputting 20 dimensions), followed by 1D convolutional layer using 50 filters with kernel size of 3 words. Finally a hidden layer with 200 hidden nodes.

Hyperparameters: activation = ReLU, optimiser = adam, learning rate = 0.001, batch size = 64, max sentence length = 200, vocabulary size = 6000

Experiments were repeated 5 times with a different random seed for each run.

## E.8 Reinforcement Learning

### E.8.1 Uncertainty-Aware Reinforcement Learning

The FetchPush environment from OpenAI Gym was used with the sparse rewards setting. We modified the environment slightly. The goal was positioned at a fixed radius from the block (but at varying angle). Actions were discretised and vertical movements removed so the agent had a choice of moving 0.4 units forward/backwards/left/right. Gaussian noise was added to the actions to make the problem stochastic. Inputs were preprocessed so that relative coordinates of gripper to cube and cube to goal were provided directly to the NNs.

We used fixed target NNs which were updated every 500 episodes.

The simulation was run for 40,000 episodes, with final average rewards around  $-0.4$ . Two-layer NNs of 50 nodes were used. Learning rate = 0.001, batch size = 100, episodes in between training = 100,  $\gamma = 0.98$ , buffer size = 100,000.