
Fast Predictive Uncertainty for Classification with Bayesian Deep Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In Bayesian Deep Learning, distributions over the output of classification neural
2 networks are approximated by first constructing a Gaussian distribution over the
3 weights, then sampling from it to receive a distribution over the categorical output
4 distribution. This is costly. We reconsider old work to construct a Dirichlet
5 approximation of this output distribution, which yields an analytic map between
6 Gaussian distributions in logit space and Dirichlet distributions (the conjugate
7 prior to the categorical) in the output space. We argue that the resulting Dirichlet
8 distribution has theoretical and practical advantages, in particular more efficient
9 computation of the uncertainty estimate, scaling to large datasets and networks like
10 ImageNet and DenseNet. We demonstrate the use of this Dirichlet approximation
11 by using it to construct a lightweight uncertainty-aware output ranking for the
12 ImageNet setup.

13 1 Introduction

14 Quantifying the uncertainty of Neural Networks’ (NNs) predictions is important in safety-critical
15 applications such as medical-diagnosis [1] and self-driving vehicles [2; 3]. Architectures for classifi-
16 cation tasks produce a probability distribution as their output, constructed by applying the softmax to
17 the point-estimate output of the penultimate layer. However, it has been shown that this distribution
18 is overconfident [4; 5] and thus cannot be used for predictive uncertainty quantification.

19 Approximate Bayesian methods provide quantified uncertainty over the network’s parameters and thus
20 the outputs in a tractable fashion. The commonly used Gaussian approximate posterior [6; 7; 8; 9]
21 approximately induces a Gaussian distribution over the logits of a NN [10]. However, the associated
22 predictive distribution, which is the expectation of the softmax function w.r.t. the Gaussian does not
23 have an analytic form. It is thus generally approximated by Monte Carlo (MC) integration requiring
24 multiple samples. Predictions in Bayesian Neural Networks (BNNs) are thus generally expensive
25 operations.

26 In this paper, we re-introduce an old but largely overlooked idea originally proposed by David
27 JC MacKay [11] in a different setting (arguably the inverse of the Deep Learning setting) which
28 transforms a Dirichlet distribution into a Gaussian. Dirichlet distributions are generally defined on
29 the simplex. But when its variable is defined on the inverse softmax’s domain, its shape effectively
30 approximates a Gaussian. The inverse of this approximation, which will be called the *Laplace*
31 *Bridge* here [12], analytically maps a Gaussian distribution onto a Dirichlet distribution. Given a
32 Gaussian distribution over the logits of a NN, one can thus efficiently obtain an approximate Dirichlet
33 distribution over the softmax outputs (Figure 1).

34 Our contributions in this paper are: We re-visit MacKay’s derivation with particular attention to a
35 symmetry constraint that becomes necessary in our “inverted” use of the argument from the Gaussian

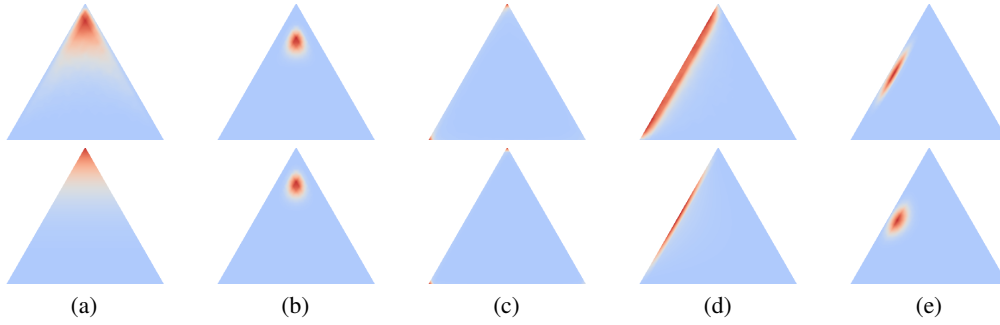


Figure 1: Densities on the simplex of the true distribution (top row, computed by MC integration) and ‘‘Laplace Bridge’’ approximation constructed in this paper (bottom row). For column (a) and (b), two different Gaussians were constructed, such that the resulting MAP estimate is the same, but the uncertainty differs. For (c), (d) and (e) the same mean with decreasing uncertainty was used. We find that in all cases the Laplace Bridge is a good approximation and captures the desired properties.

to the Dirichlet family. We then validate the quality of this approximation both by theoretical and empirical arguments and demonstrate significant speed-up over MC-integration. Finally, we show a use-case, leveraging the analytic properties of Dirichlet distributions to improve the popular top- k metric through uncertainties.

Section 2 provides the mathematical derivation. Section 3 discusses the Laplace Bridge in the context of NNs. We compare it to the recent approximations of the predictive distributions of NNs in Section 4. Experiments are presented in Section 5.

2 The Laplace Bridge

Laplace approximations¹ are a popular and light-weight method to approximate a general probability distribution $q(\mathbf{x})$ with a Gaussian $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$. It sets $\boldsymbol{\mu}$ to a mode of q , and $\boldsymbol{\Sigma} = -(\nabla^2 \log q(\mathbf{x})|_{\boldsymbol{\mu}})^{-1}$, the inverse Hessian of $\log q$ at that mode. This scheme can work well if the true distribution is unimodal and defined on the real vector space.

The Dirichlet distribution, which has the density function

$$\text{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha}) := \frac{\Gamma\left(\sum_{k=1}^K \alpha_k\right)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K \pi_k^{\alpha_k-1}, \quad (1)$$

is defined on the probability simplex and can be multimodal in the sense that the maxima of the distribution lie at the boundary of the simplex when $\alpha_k < 1$, for all $k = 1, \dots, K$. Both issues preclude a Laplace approximation, at least in the naïve form described above. However, MacKay [11] noted that both can be fixed, elegantly, by a change of variable. Details of the following argument can be found in the supplements. Consider the K -dimensional variable $\boldsymbol{\pi} \sim \text{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha})$ defined as the softmax of $\mathbf{z} \in \mathbb{R}^K$:

$$\pi_k(\mathbf{z}) := \frac{\exp(z_k)}{\sum_{l=1}^K \exp(z_l)}, \quad (2)$$

for all $k = 1, \dots, K$. We will call \mathbf{z} the logit of $\boldsymbol{\pi}$. When expressed as a function of \mathbf{z} , the density of the Dirichlet in $\boldsymbol{\pi}$ has to be multiplied by the Jacobian determinant

$$\det \frac{\partial \boldsymbol{\pi}}{\partial \mathbf{z}} = \prod_k \pi_k(z), \quad (3)$$

thus removing the -1 terms in the exponent:

$$\text{Dir}_{\mathbf{z}}(\boldsymbol{\pi}(\mathbf{z})|\boldsymbol{\alpha}) := \frac{\Gamma\left(\sum_{k=1}^K \alpha_k\right)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K \pi_k(\mathbf{z})^{\alpha_k}, \quad (4)$$

¹For clarity: Laplace approximations are *also* one out of several possible ways to construct a Gaussian approximation to the weight posterior of a NN, by constructing a second-order Taylor approximation of the empirical risk at the trained weights. This is *not* the way they are used in this section. The Laplace Bridge is agnostic to how the input Gaussian distribution is constructed. It could, e.g., also be constructed as a variational approximation, or the moments of Monte Carlo samples. See also ??.

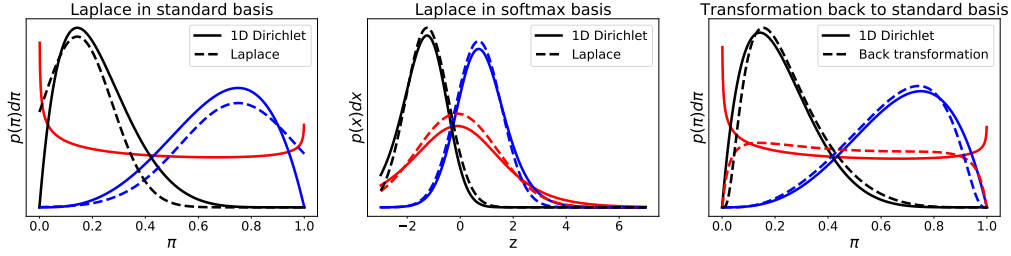


Figure 2: (Adapted from Hennig et al. [12]). Visualization of the Laplace Bridge for the Beta distribution (1D special case of the Dirichlet). **Left:** “Generic” Laplace approximations of standard Beta distributions by Gaussians. Note that the Beta Distribution (red curve) does not even have a valid approximation because the Hessian is not positive semi-definite. **Middle:** Laplace approximation to the same distributions after basis transformation through the softmax (4). The transformation makes the distributions “more Gaussian” (i.e. uni-modal, bell-shaped, with support on the real line) compared to the standard basis, thus making the Laplace approximation more accurate. **Right:** The same Beta distributions, with the back-transformation of the Laplace approximations from the middle figure to the simplex, yielding a much improved approximate distribution. In particular, in contrast to the left-most image, the dashed lines now actually are probability distributions (they integrate to 1 on the simplex).

This density of \mathbf{z} (!), the Dirichlet distribution in the *softmax basis*, can now be accurately approximated by a Gaussian through a Laplace approximation, yielding an analytic map from the parameter space $\boldsymbol{\alpha} \in \mathbb{R}_+^K$ to the parameter space of the Gaussian ($\boldsymbol{\mu} \in \mathbb{R}^K$ and symmetric positive definite $\boldsymbol{\Sigma} \in \mathbb{R}^{K \times K}$), given by

$$\mu_k = \log \alpha_k - \frac{1}{K} \sum_{l=1}^K \log \alpha_l, \quad (5)$$

$$\Sigma_{k\ell} = \delta_{k\ell} \frac{1}{\alpha_k} - \frac{1}{K} \left[\frac{1}{\alpha_k} + \frac{1}{\alpha_\ell} - \frac{1}{K} \sum_{u=1}^K \frac{1}{\alpha_u} \right]. \quad (6)$$

The corresponding derivations require care because the Gaussian parameter space is evidently larger than that of the Dirichlet and not fully identified by the transformation. A pseudo-inverse of this map was provided by Hennig et al. [12]. It maps the Gaussian parameters to those of the Dirichlet as

$$\alpha_k = \frac{1}{\Sigma_{kk}} \left(1 - \frac{2}{K} + \frac{e^{\mu_k}}{K^2} \sum_{l=1}^K e^{-\mu_l} \right) \quad (7)$$

(Note that this equation ignores off-diagonal elements of $\boldsymbol{\Sigma}$, more discussion below). Together, Eqs. 5, 6 and 7 will here be used for Bayesian Deep Learning, and jointly called the *Laplace Bridge*. Note that, even though the Laplace Bridge implies a reduction of the expressiveness of the distribution, we show in Section 3 that this map is still sufficiently accurate.

Figure 1 shows the quality of the resulting approximation. We consider multiple different $\boldsymbol{\mu}, \boldsymbol{\Sigma}$ in three dimensions. We exhaustively sample from the Gaussian and apply the softmax. The resulting histogram is compared to the PDF of the corresponding Dirichlet. The first part of the figure emphasizes that a point estimate is insufficient. Since the mean for the Dirichlet is the normalized parameter vector $\boldsymbol{\alpha}$, the parameters $(\alpha_1 = [2, 2, 6]^\top$ and $\alpha_2 = [11, 11, 51]^\top$) yield the same point estimate even though their distributions are clearly different. The second part shows how the Laplace Bridge maps w.r.t decreasing uncertainty.

3 The Laplace Bridge for BNNs

Let $f_\theta : \mathbb{R}^N \rightarrow \mathbb{R}^K$ be an L -layer Neural Network parametrized by $\theta \in \mathbb{R}^P$, with a Gaussian approximate posterior $\mathcal{N}(\theta | \boldsymbol{\mu}_\theta, \boldsymbol{\Sigma}_\theta)$. For any input $\mathbf{x} \in \mathbb{R}^N$, one way to obtain an approximate Gaussian distribution on the pre-softmax output (logit vector) $f_\theta(\mathbf{x}) =: \mathbf{z}$ is as

$$q(\mathbf{z} | \mathbf{x}) \approx \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}_\theta^\top \mathbf{x}, \mathbf{J}(\mathbf{x})^\top \boldsymbol{\Sigma}_\theta \mathbf{J}(\mathbf{x})), \quad (8)$$

where $\mathbf{J}(\mathbf{x})$ is the $P \times K$ Jacobian matrix representing the derivative $\frac{\partial \mathbf{z}}{\partial \boldsymbol{\theta}}$ [10]. Approximating the density of the softmax of this Gaussian random variable as a Dirichlet, using the Laplace Bridge, *analytically* approximates the predictive distribution in a single step, as opposed to many samples. From Eq. (7), this requires $\mathcal{O}(K)$ computations to construct the K parameters α_k of the Dirichlet. In contrast, MC-integration has computational costs of $\mathcal{O}(MJ)$, where M is the number of samples and J is the cost of sampling from $q(\mathbf{z}|\mathbf{x})$ (typically J is of order K^2 after an initial $\mathcal{O}(K^3)$ operation for a matrix decomposition of the covariance). The Monte Carlo approximation has the usual sampling error of $\mathcal{O}(1/\sqrt{M})$, while the Laplace Bridge has a fixed but small error (empirical comparison in Section 5.3).

We now discuss several qualitative properties of the Laplace Bridge relevant for the uncertainty quantification use case in Deep Learning. For output classes of “comparably high” probability (as defined below), the variance $\text{Var}(\pi_k|\boldsymbol{\alpha})$ under the Laplace Bridge increases with the variance of the underlying Gaussian. In this sense, the Laplace Bridge approximates the uncertainty information encoded in the output of a BNN.

Proposition 1 (proof in supplements). *Let $\text{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha})$ be obtained via the Laplace Bridge from a Gaussian distribution $\mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ over \mathbb{R}^K . Then, for each $k = 1, \dots, K$, letting $\alpha_{\neq k} := \sum_{l \neq k} \alpha_l$, if*

$$\alpha_k > \frac{1}{4} \left(\sqrt{9\alpha_{\neq k}^2 + 10\alpha_{\neq k} + 1} - \alpha_{\neq k} - 1 \right),$$

then the variance $\text{Var}(\pi_k|\boldsymbol{\alpha})$ of the k -th component of $\boldsymbol{\pi}$ is increasing in $\boldsymbol{\Sigma}_{kk}$.

Intuitively, this result describes the condition that needs to be fulfilled such that the variance of the resulting Dirichlet scales with the variance of the k -th component of the Gaussian. It can be seen as a proxy for a high quality approximation. An empirical evaluation showing that this condition is fulfilled in most cases can be found in the supplements.

Further benefits of this approximation arise from the convenient analytical properties of the Dirichlet exponential family. For example, a point estimate of the posterior predictive distribution is directly given by the Dirichlet’s mean,

$$\mathbb{E}\boldsymbol{\pi} = \left(\frac{\alpha_1}{\sum_{l=1}^K \alpha_l}, \dots, \frac{\alpha_K}{\sum_{l=1}^K \alpha_l} \right)^\top, \quad (9)$$

Additionally, Dirichlets have Dirichlet marginals: If $p(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha})$, then

$$p \left(\left[\pi_1, \pi_2, \dots, \pi_j, \sum_{k>j} \pi_k \right]^\top \right) = \text{Dir} \left(\alpha_1, \alpha_2, \dots, \alpha_j, \sum_{k>j} \alpha_k \right). \quad (10)$$

An additional benefit of the Laplace Bridge for BNNs is that it is more flexible than an MC-integral. If we let $p(\boldsymbol{\pi})$ be the distribution over $\boldsymbol{\pi} := \text{softmax}(\mathbf{z}) := [e^{z_1} / \sum_l e^{z_l}, \dots, e^{z_K} / \sum_l e^{z_l}]^\top$, then the MC-integral can be seen as a “point-estimate” of this distribution since it approximates $\mathbb{E}\boldsymbol{\pi}$. In contrast, the Dirichlet distribution $\text{Dir}(\boldsymbol{\pi}|\boldsymbol{\alpha})$ approximates the distribution $p(\boldsymbol{\pi})$. Thus, the Laplace Bridge enables tasks that can be done only with a distribution but not a point estimate. For instance, one could ask “what is the distribution of the first L classes?” when one is dealing with K -class ($L < K$) classification. Since the marginal distribution can be computed analytically (10), the Laplace Bridge provides a convenient yet cheap way of answering this question.

4 Related Work

In Bayesian Neural Networks, analytic approximations of posterior predictive distributions have attracted a great deal of research. In the binary classification case, for example, the probit approximation has been proposed already in the 1990s [13; 14]. However, while there exist some bounds [15] and approximations of the expected log-sum-exponent function [16; 17], in the multi-class case, obtaining a good analytic approximation of the expected softmax function under a Gaussian measure is still considered an open problem. The Laplace Bridge can be used to produce a close analytical approximation of this integral. It furthers the trend of sampling-free solutions within Bayesian Deep

Learning [18; 19, etc.]. The crucial difference is that, unlike these methods, the Laplace Bridge approximates the full distribution over the softmax outputs of a deep network.

Recently, it has been proposed to model the distribution of softmax outputs of a network directly. Similar to the Laplace Bridge, Malinin and Gales [20, 21]; Sensoy et al. [22] proposed to use the Dirichlet distribution to model the posterior predictive for non-Bayesian networks. They further proposed novel training techniques in order to directly learn the Dirichlet. In contrast, the Laplace Bridge tackles the problem of approximating the distribution over the softmax outputs of the ubiquitous Gaussian-approximated Bayesian networks [7; 8; 23; 24, etc] without any additional training procedure. This allows the Laplace Bridge to be used with pre-trained networks and emphasises its non-invasive nature.

5 Experiments

We conduct four experiments. In Section 5.1, we analyze the approximation quality of the Laplace Bridge applied to a BNN on the MNIST [25] dataset. Then, we compare the Laplace Bridge to the MC-integral in the example application of out-of-distribution (OOD) detection (Section 5.2). Their computational costs are compared in Section 5.3. Finally, in Section 5.4, we present analysis on ImageNet [26] to demonstrate the scalability of the Laplace Bridge and the advantage of having a full Dirichlet distribution over softmax outputs.

All experiments were conducted using different forms of Laplace approximations of Neural Networks. For the smaller experiments a full (all-layer) Laplace approximation with a diagonal Hessian [27] was used. For the experiments with larger networks the Laplace approximation has been applied only to the last-layer of the network. This scheme has been successfully used by Snoek et al. [28]; Wilson et al. [29]; Brosse et al. [30], etc. and it has been shown theoretically to mitigate overconfidence problems in ReLU networks [31]. For the last-layer experiments we use diagonal and Kronecker-factorized [9] (KFAC) approximations of the Hessian, since inverting the exact Hessian is too costly. A detailed mathematical explanation and setup of the experiments can be found in the supplements.

While the Laplace Bridge could also be applied to different approximations to a Gaussian posterior predictive such as Variational Inference [7; 8], we used a Laplace approximation in our experiments to construct such an approximation. This is for two reasons: (i) it is one of the fastest ways to get a Gaussian posterior predictive and (ii) it can be applied to pre-trained networks which is especially useful for ImageNet experiments. Nevertheless, we want to emphasize again that the Laplace Bridge can be applied to any Gaussian over the outputs independent of the way it was generated. Despite the overlap in nomenclature, the Laplace Bridge is *not* restricted to Gaussians arising from a Laplace approximation of the network.

5.1 Uncertainty estimates on MNIST

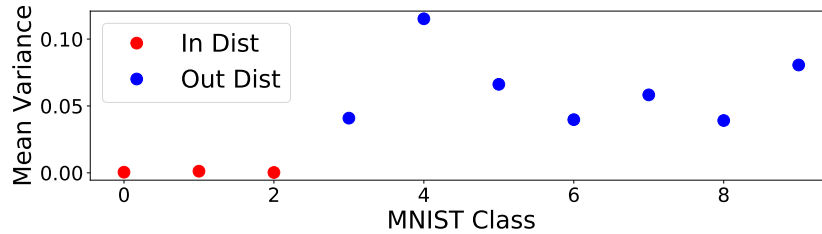


Figure 3: Average variance of the Dirichlet distributions of each MNIST class. The in-distribution uncertainty (variance) is nearly nil, while out-of-distribution variance is higher. This implies usability for OOD detection.

We empirically investigate the approximation quality of the Laplace Bridge in a “real-world” BNN on the MNIST dataset. A CNN with 2 convolutional and 2 fully-connected layers is trained on the first three digits of MNIST (the digits 0, 1, and 2). Adam optimizer with learning rate $1e-3$ and weight decay $5e-4$ is used. The batch size is 128. To obtain the posterior over the weights of this network, we perform a full (all-layer) Laplace approximation using BackPACK [32] to get the diagonal Hessian. The network is then evaluated on the full test set of MNIST (containing all ten classes). We present the results in Figure 3. We show for each $k = 1, \dots, K$, the average variance

163 $\frac{1}{D_k} \sum_{i=1}^{D_k} \text{Var}(\pi_k(f_{\theta}(\mathbf{x}_i)))$ of the resulting Dirichlet distribution over the softmax outputs, where
164 D_k is the number of test points predicted with label k . The results show that the variance of the
165 Dirichlet distribution obtained via the Laplace Bridge is useful for uncertainty quantification: OOD
166 data can be easily detected since the mean variance of the first three classes is nearly zero while that
167 of the others is higher.

Table 1: OOD detection results. The Laplace Bridge (LB) wins most comparisons with Diagonal sampling and draws even with KFAC sampling w.r.t. both metrics. However, the LB is around 400 times faster on average. 1000 samples were drawn from the Gaussian over the outputs. The (F-, K-, not-)MNIST experiments were done with a Laplace approximation of the entire network while the others only used the last layer.

Train	Test	Diag Sampling		Diag LB		KFAC Sampling		KFAC LB		Time in s \downarrow	
		MMC \downarrow	AUROC \uparrow	MMC \downarrow	AUROC \uparrow	MMC \downarrow	AUROC \uparrow	MMC \downarrow	AUROC \uparrow	Sampling	LB
MNIST	MNIST	0.932 \pm 0.007	-	0.987 \pm 0.001	-	-	-	-	-	6.6	0.016
MNIST	FMNIST	0.407 \pm 0.010	0.989 \pm 0.002	0.377 \pm 0.019	0.994 \pm 0.002	-	-	-	-	6.6	0.016
MNIST	notMNIST	0.535 \pm 0.018	0.958 \pm 0.006	0.630 \pm 0.018	0.962 \pm 0.007	-	-	-	-	12.3	0.029
MNIST	KMNIST	0.500 \pm 0.014	0.974 \pm 0.005	0.630 \pm 0.018	0.975 \pm 0.004	-	-	-	-	6.6	0.016
CIFAR-10	CIFAR-10	-	0.948	0.966	-	0.857 \pm 0.003	-	0.966	-	13.6	0.031
CIFAR-10	CIFAR-100	0.708	0.889	0.742	0.866	0.562 \pm 0.003	0.880 \pm 0.012	0.741	0.866	13.5	0.027
CIFAR-10	SVHN	0.643	0.933	0.647	0.934	0.484 \pm 0.004	0.939 \pm 0.001	0.648 \pm 0.003	0.934 \pm 0.001	35.2	0.070
SVHN	SVHN	0.986	-	0.993	-	0.947 \pm 0.002	-	0.993	-	34.5	0.073
SVHN	CIFAR-100	0.595	0.984	0.526	0.985	0.460 \pm 0.004	0.986 \pm 0.001	0.527 \pm 0.002	0.985	13.4	0.027
SVHN	CIFAR-10	0.593	0.984	0.520	0.987	0.458 \pm 0.004	0.986 \pm 0.001	0.520 \pm 0.002	0.987	13.3	0.028
CIFAR-100	CIFAR-100	0.762	-	0.590	-	0.404	-	0.593	-	24.6	0.030
CIFAR-100	CIFAR-10	0.467	0.788	0.206	0.791	0.213	0.788	0.209	0.791	24.6	0.027
CIFAR-100	SVHN	0.461	0.795	0.170	0.815	0.180 \pm 0.001	0.838 \pm 0.001	0.173	0.815	63.8	0.069

168 5.2 OOD detection

169 We compare the performance of the Laplace Bridge to the MC-integral (Diagonal and KFAC) on a
170 standard OOD detection benchmark suite, to test whether the Laplace Bridge gives similar results to
171 the MC sampling methods and compare their computational overhead. Following prior literature,
172 we use the standard mean-maximum-confidence (MMC) and area under the ROC-curve (AUROC)
173 metrics [33]. For an in-distribution dataset, a higher MMC value is desirable while for the OOD
174 dataset we want a lower MMC value (optimally, $1/K$ in K -class classification problems). For the
175 AUROC metric, the higher the better, since it represents how good a method is for distinguishing in-
176 and out-of-distribution datasets.

177 The test scenarios are as follows: (i) The same convolutional network as in Section 5.1 is trained on
178 the MNIST dataset. To approximate the posterior over the parameter of this network, a full (all-layer)
179 Laplace approximation with a diagonal Hessian is used. The OOD datasets for this case are FMNIST
180 [34], notMNIST [35], and KMNIST [36]. (ii) For larger datasets, i.e. CIFAR-10 [37], SVHN [38],
181 and CIFAR-100 [37], we use a ResNet-18 network [39]. Since this network is large, Equation (8) in
182 conjunction with a full Laplace approximation is too costly. We, therefore, use a last-layer Laplace
183 approximation to obtain the approximate diagonal and KFAC Gaussian posterior. The OOD datasets
184 for CIFAR-10, SVHN, and CIFAR-100 are SVHN and CIFAR100; CIFAR-10 and CIFAR-100; and
185 SVHN and CIFAR-10, respectively. In all scenarios, the networks are well-trained with 99% accuracy
186 on MNIST, 95.4% on CIFAR-10, 76.6% on CIFAR-100, and 100% on SVHN. For the sampling
187 baseline, we use 1000 posterior samples to compute the predictive distribution. We use the mean of
188 the Dirichlet to obtain a comparable approximation to the MC-integral. Further comparisons with
189 ensemble networks can be found in the supplements.

190 The results are presented in Table 1. The Laplace Bridge yields, on average, better results than
191 diagonal sampling and ties with KFAC sampling w.r.t both metrics. However, the Laplace Bridge is
192 around 400 times faster than both sampling based methods and is therefore preferable.

193 5.3 Time comparison

194 We compare the computational cost of the density-estimated p_{sample} distribution via sampling and the
195 Dirichlet distribution obtained from the Laplace Bridge p_{LB} for approximating the true distribution
196 p_{true} over softmax-Gaussian samples². Different amounts of samples are drawn from the Gaussian,
197 the softmax is applied and the KL divergence between the histogram of the samples with the true dis-
198 tribution is computed. We use KL-divergences $D_{\text{KL}}(p_{\text{true}} \| p_{\text{sample}})$ and $D_{\text{KL}}(p_{\text{true}} \| p_{\text{LB}})$, respectively,

²I.e. samples are obtained by first sampling from a Gaussian and transforming it via the softmax function.

199 to measure similarity between the approximations and ground truth while the number of samples for
200 p_{sample} is increased on a logarithmic scale. The true distribution p_{true} is constructed via MC with 100k
201 samples. The experiment is conducted for three different Gaussian distributions over \mathbb{R}^3 . Since the
202 softmax applied to a Gaussian does not have a closed-form analytic solution, the algebraic calculation
203 of the approximation error is not possible and an empirical evaluation via sampling is the best option.
204 The fact that there is no analytic solution is part of the justification for using the Laplace Bridge in
205 the first place.

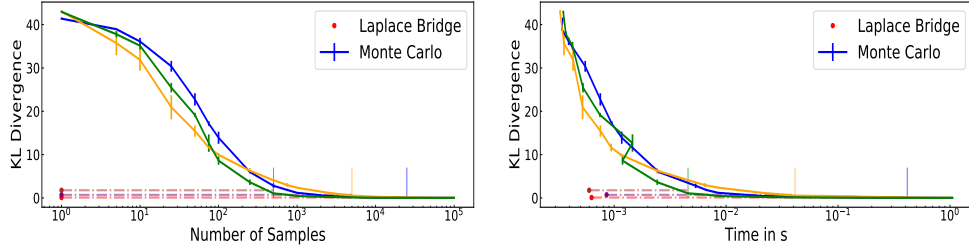


Figure 4: KL-divergence plotted against the number of samples (left) and wall-clock time (right). Monte Carlo density estimation becomes as good as the Laplace Bridge after around 750 to 10000 samples and takes at least 100 times longer. The three lines represent three different samples.

206 Figure 4 suggests that the number of samples required such that the distribution p_{sample} is approximat-
207 ing the true distribution p_{true} as good as the Dirichlet distribution obtained via the Laplace Bridge is
208 large, i.e. somewhere between 500 and 10000. This translates to a wall-clock time advantage of at
209 least a factor of 100 before sampling becomes competitive in quality with the Laplace Bridge.

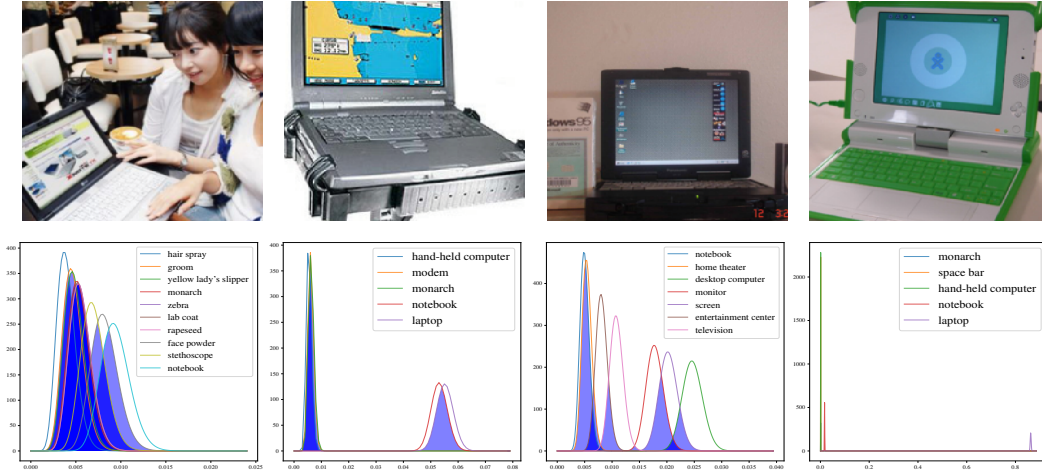


Figure 5: **Upper row:** images from the “laptop” class of ImageNet. **Bottom row:** Beta marginal distributions of the top- k predictions for the respective image. In the first column, the overlap between the marginal of all classes is large, signifying high uncertainty, i.e. the prediction is “I do not know”. In the column, “notebook” and “laptop” have confident, yet overlapping marginal densities and we, therefore, have a top-2 prediction: “either a notebook or a laptop”. In the third column “desktop computer”, “screen” and “monitor” have overlapping marginal densities, yielding a top-3 estimate. The last case shows a top-1 estimate: the network is confident that “laptop” is the only correct label.

210 5.4 Uncertainty-aware output ranking on ImageNet

211 Classification tasks on large datasets with many classes, like ImageNet, are not often done in a
212 Bayesian fashion since the posterior inference and sampling are expensive. The Laplace Bridge,
213 in conjunction with the last-layer Bayesian approximations, can be used to alleviate this problem.
214 Furthermore, having a full distribution over the softmax outputs of a BNN gives rise to new possi-
215 bilities. For example, one could subsume all classes which have sufficiently overlapping marginal
216 distributions into one if they are semantically similar as illustrated in Figure 5.

Another possibility is to improve the standard classification metrics. Large classification tasks like ImageNet are often compared along a top-5 metric, i.e. it is tested whether the correct class is within the five most probable estimates of the network. Although widely accepted, this metric has some pathologies. We can easily construct examples where the top-5 include either too many or too few classes for our purposes which a static rule (always 5) can't handle (see Figure 5).

Leveraging the probabilistic output provided by the Laplace Bridge, we propose a simple decision rule that can handle such examples and is more fine-grained due to its awareness of uncertainty. One may call such a rule *uncertainty-aware top-k*; pseudocode for the algorithm is given in the supplements. Instead of taking the top- k as a decision threshold for an arbitrary k we take the uncertainty/confidence of the model to inform the decision. This is more flexible and therefore able to handle situations in which different numbers of classes are plausible outcomes. The Dirichlet distribution obtained from the Laplace Bridge provides this capability. In particular, since the marginal distribution over each component of a Dirichlet distribution is a $\text{Beta}(\alpha_i, \sum_{j \neq i} \alpha_j)$, this can be done analytically and efficiently. The proposed decision rule uses the area of overlap between the marginal distributions of the sorted outcomes. This is similar to hypotheses testing, i.e. t -tests [40] or its Bayesian alternatives [41]. If, for example, two Beta densities overlap more than 5%, we cannot say that they are different distributions with high confidence. All distributions that have sufficient overlap should become the new top- k estimate. Figure 5 shows four examples from the “laptop” class of ImageNet.

We evaluate this decision rule on the test set of ImageNet. The overlap is calculated through the inverse CDF³ of the respective Beta marginals. The original top-1 accuracy of DenseNet on ImageNet is 0.744. Meanwhile, the uncertainty-aware top- k accuracy is 0.797, where k is on average 1.688. A more detailed analysis of the distribution of top- k estimates (see Figure 6) shows that most of the predictions given by the uncertainty-aware metric still yielded a top-1 prediction.

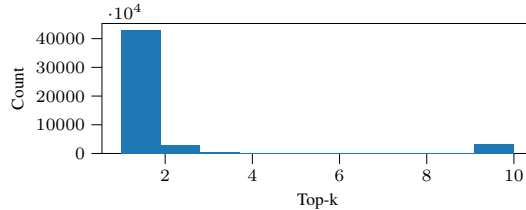


Figure 6: A histogram of ImageNet predictions’ length using the proposed uncertainty-aware top- k . Most test images are a top-1 prediction, indicating high confidence. There are some top-2, top-3, and top-10 predictions, showing an increasing uncertainty.

This means that using uncertainty does not imply adding meaningless classes to the prediction. However, there are some non-negligible cases where k equals to 2, 3, or 10. This indicates that whenever there is ambiguity in the class labels, our method is able to detect it, and thus yields a significantly higher accuracy.

6 Conclusion

We have adapted an old but overlooked approximation scheme for new use in Bayesian Deep Learning. Given a Gaussian approximation to the weight-space posterior of a neural network (which can be constructed by various means, including another Laplace approximation), and an input, the Laplace Bridge analytically maps the marginal Gaussian prediction on the logits onto a Dirichlet distribution over the softmax vectors. The associated computational cost of $\mathcal{O}(K)$ for K -class prediction compares favorably to that of Monte Carlo sampling. The proposed method both theoretically and empirically preserves predictive uncertainty, offering an attractive, low-cost, high-quality alternative to Monte Carlo sampling. In conjunction with a low-cost, last-layer Bayesian approximation, it can be useful in real-time applications wherever uncertainty is required - especially because it reduces the cost of predicting a posterior distribution at test time.

³Also known as the quantile function or percent point function

7 Broader Impact

More and more tasks are solved through Deep Learning and Neural Networks. While they often provide state-of-the-art results in terms of their accuracy, there are nearly no theoretical bounds on their behavior when confronted with new situations. It is therefore of high importance for a Neural Network to be able to provide well-calibrated uncertainty about its predictions. A network has to be able to say "I don't know" when it receives data that it can not classify sufficiently well or which are far away from the training distribution. Such uncertainty is especially important in safety-critical tasks such as self-driving vehicles or medical applications, where it can directly affect and improve decision making.

While the field of Bayesian Deep Learning (BDL) is rapidly improving, many of its applications have one of two problems: either (i) acquiring the uncertainty estimate is computationally expensive since it involves sampling or (ii) Bayesian methods yield good uncertainty estimates but don't yield the same accuracy as conventional methods.

Reducing the computational overhead of BDL is important, especially during test time, because it implies viability for applications where either (i) small differences in time can make large differences in outcome (e.g. breaking earlier to prevent an accident) or (ii) uncertainty estimates are required in rapid succession (e.g. multiple hundred frames per second). Additionally, it also implies less energy usage and, thereby, higher accessibility because of the reduced cost. However, our method mostly saves overhead during test time and not during training. Therefore, we expect the effects on the climate and access to be marginal compared to its other benefits.

While our method, the Laplace Bridge for Neural Networks, by no means offers a perfect solution to fast and precise uncertainty quantification, it amounts to a significant step in this direction. Computing a fully parameterised distribution over the outputs is faster than drawing one (!) sample from the posterior predictive Gaussian and thereby allows for the just described benefits during test time. At the same time, it can be applied to existing trained networks such that conventionally effective methods can be used for training without loss of accuracy.

References

- [1] E. Begoli, T. Bhattacharya, and D. Kusnezov. The need for uncertainty quantification in machine-assisted medical decision making. *Nat Mach Intell*, 1:20–23, 2019.
- [2] Rowan McAllister, Yarin Gal, Alex Kendall, Mark van der Wilk, Amar Shah, Roberto Cipolla, and Adrian Weller. Concrete problems for autonomous vehicle safety: Advantages of bayesian deep learning. In *IJCAI*, 2017.
- [3] Rhiannon Michelmore, Marta Kwiatkowska, and Yarin Gal. Evaluating uncertainty quantification in end-to-end autonomous driving control. *CoRR*, abs/1811.06817, 2018.
- [4] Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *CVPR*, 2015.
- [5] Matthias Hein, Maksym Andriushchenko, and Julian Bitterwolf. Why relu networks yield high-confidence predictions far away from the training data and how to mitigate the problem. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [6] David J. C. MacKay. A practical bayesian framework for backpropagation networks. *Neural Comput.*, 4(3):448–472, May 1992. ISSN 0899-7667.
- [7] Alex Graves. Practical variational inference for neural networks. In J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24*, pages 2348–2356. Curran Associates, Inc., 2011.
- [8] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. *ArXiv*, 2015.
- [9] Hippolyt Ritter, Aleksandar Botev, and David Barber. A scalable laplace approximation for neural networks. In *International Conference on Learning Representations*, 2018.
- [10] David J C Mackay. Probable networks and plausible predictions — a review of practical bayesian methods for supervised neural networks. *Network: Computation in Neural Systems*, 6(3):469–505, 1995.

- [11] David J.C. MacKay. Choice of basis for laplace approximation. *Machine Learning*, 33(1):77–86, Oct 1998. ISSN 1573-0565.
- [12] P. Hennig, D. Stern, R. Herbrich, and T. Graepel. Kernel topic models. In *Fifteenth International Conference on Artificial Intelligence and Statistics*, volume 22 of *JMLR Proceedings*, pages 511–519. JMLR.org, 2012.
- [13] David J Spiegelhalter and Steffen L Lauritzen. Sequential updating of conditional probabilities on directed graphical structures. *Networks*, 20(5):579–605, 1990.
- [14] David JC MacKay. The evidence framework applied to classification networks. *Neural computation*, 4(5): 720–736, 1992.
- [15] Michalis Titsias. One-vs-each approximation to softmax for scalable estimation of probabilities. In *NIPS*, 2016.
- [16] Amr Ahmed and Eric Xing. On tight approximate inference of the logistic-normal topic admixture model. In *Proceedings of the 11th Tenth International Workshop on Artificial Intelligence and Statistics*, 2007.
- [17] Michael Braun and Jon McAuliffe. Variational inference for large-scale models of discrete choice. *Journal of the American Statistical Association*, 105(489):324–335, 2010.
- [18] Anqi Wu, Sebastian Nowozin, Edward Meeds, Richard E. Turner, José Miguel Hernández-Lobato, and Alexander L. Gaunt. Fixing variational bayes: Deterministic variational inference for bayesian neural networks. *CoRR*, abs/1810.03958, 2018. URL <http://arxiv.org/abs/1810.03958>.
- [19] Manuel Haussmann, Sebastian Gerwinn, and Melih Kandemir. Bayesian evidential deep learning with pac regularization, 2019.
- [20] Andrey Malinin and Mark Gales. Predictive uncertainty estimation via prior networks. In *Advances in Neural Information Processing Systems*, pages 7047–7058, 2018.
- [21] Andrey Malinin and Mark Gales. Reverse kl-divergence training of prior networks: Improved uncertainty and adversarial robustness. In *Advances in Neural Information Processing Systems*, pages 14520–14531, 2019.
- [22] Murat Sensoy, Lance Kaplan, and Melih Kandemir. Evidential deep learning to quantify classification uncertainty. In *Advances in Neural Information Processing Systems*, pages 3179–3189, 2018.
- [23] Christos Louizos and Max Welling. Structured and efficient variational deep learning with matrix gaussian posteriors. In *ICML*, 2016.
- [24] Shengyang Sun, Changyou Chen, and Lawrence Carin. Learning structured weight uncertainty in Bayesian neural networks. In *Artificial Intelligence and Statistics*, pages 1283–1292, 2017.
- [25] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010.
- [26] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Fei-Fei Li. Imagenet large scale visual recognition challenge. *CoRR*, abs/1409.0575, 2014.
- [27] Yann LeCun, John S. Denker, and Sara A. Solla. Optimal brain damage. In D. S. Touretzky, editor, *Advances in Neural Information Processing Systems 2*, pages 598–605. Morgan-Kaufmann, 1990. URL <http://papers.nips.cc/paper/250-optimal-brain-damage.pdf>.
- [28] Jasper Snoek, Oren Rippel, Kevin Swersky, Ryan Kiros, Nadathur Satish, Narayanan Sundaram, Mostofa Patwary, Mr Prabhat, and Ryan Adams. Scalable bayesian optimization using deep neural networks. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 2171–2180, Lille, France, 07–09 Jul 2015. PMLR.
- [29] Andrew Gordon Wilson, Zhiting Hu, Ruslan Salakhutdinov, and Eric P. Xing. Deep kernel learning. In Arthur Gretton and Christian C. Robert, editors, *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, volume 51 of *Proceedings of Machine Learning Research*, pages 370–378, Cadiz, Spain, 09–11 May 2016. PMLR.
- [30] Nicolas Brosse, Carlos Riquelme, Alice Martin, Sylvain Gelly, and Éric Moulines. On last-layer algorithms for classification: Decoupling representation from uncertainty estimation. *arXiv preprint arXiv:2001.08049*, 2020.

- 355 [31] Agustinus Kristiadi, Matthias Hein, and Philipp Hennig. Being bayesian, even just a bit, fixes overconfi-
356 dence in relu networks. *arXiv preprint arXiv:2002.10118*, 2020.
- 357 [32] Felix Dangel, Frederik Kunstner, and Philipp Hennig. Backpack: Packing more into backprop. *arXiv*
358 *preprint arXiv:1912.10985*, 2019.
- 359 [33] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples
360 in neural networks. *CoRR*, abs/1610.02136, 2016.
- 361 [34] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking
362 machine learning algorithms. *CoRR*, abs/1708.07747, 2017.
- 363 [35] Yaroslav Bulatov. notmnist dataset. 2011. URL [http://yaroslavvb.blogspot.com/2011/09/](http://yaroslavvb.blogspot.com/2011/09/notmnist-dataset.html)
364 [notmnist-dataset.html](http://yaroslavvb.blogspot.com/2011/09/notmnist-dataset.html).
- 365 [36] Tarin Clanuwat, Mikel Bober-Irizar, Asanobu Kitamoto, Alex Lamb, Kazuaki Yamamoto, and David Ha.
366 Deep learning for classical japanese literature. *CoRR*, abs/1812.01718, 2018.
- 367 [37] Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- 368 [38] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. Reading digits in
369 natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised*
370 *Feature Learning 2011*, 2011.
- 371 [39] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
372 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- 373 [40] Raymond S Nickerson. Null hypothesis significance testing: a review of an old and continuing controversy.
374 *Psychological methods*, 5(2):241, 2000.
- 375 [41] Michael E. J. Masson. A tutorial on a practical bayesian alternative to null-hypothesis significance testing.
376 *Behavior Research Methods*, 43(3):679–690, Sep 2011. ISSN 1554-3528.