
Fast Predictive Uncertainty for Classification with Bayesian Deep Networks

Anonymous Author(s)

Affiliation

Address

email

1 Appendix

In Section 1 we give some background knowledge and a proof for the proposition of the main experiment. The Laplace approximation of the Dirichlet has already been done by David JC MacKay in [1]. Subsection 2 shows a concise summary of said approximation with some additional explanations. The Inversion of this Laplace approximation has already been done by Philipp Hennig in their PhD Thesis [2]. Section 3 shows a brief overview of the main steps. Note that we are not claiming any of their original ideas but rather show their work here as a look up for interested readers. In Section 4 we give a more detailed overview over the setup of some experiments that are beyond what is necessary to understand them but might be interesting for some readers.

10 1 Appendix A: Background and Proofs

11 Change of Variable for pdf

Let \mathbf{x} be an n -dimensional continuous random variable with joint density function $p_{\mathbf{x}}$. If $\mathbf{y} = G(\mathbf{x})$, where G is a differentiable function, then \mathbf{y} has density $p_{\mathbf{y}}$:

$$g(\mathbf{y}) = f\left(G^{-1}(\mathbf{y})\right) \left| \det \left[\frac{dG^{-1}(\mathbf{z})}{d\mathbf{z}} \Big|_{\mathbf{z}=\mathbf{y}} \right] \right| \quad (1)$$

where the differential is the Jacobian of the inverse of G evaluated at \mathbf{y} . This procedure, also known as ‘change of basis’, is at the core of the Laplace bridge since it is used to transform the Dirichlet into the softmax basis.

17 Proof for Proposition

Proof. Considering that α_k is a decreasing function of Σ_{kk} by definition (19), it is sufficient to show that under the hypothesis, the derivative of $\frac{\partial}{\partial \alpha_k} \text{Var}(\pi_k | \boldsymbol{\alpha})$ is negative.

By definition, the variance $\text{Var}(\pi_k | \boldsymbol{\alpha})$ is

$$\text{Var}(\pi_k | \boldsymbol{\alpha}) = \frac{\frac{\alpha_k}{\alpha_k + \alpha_{\neq k}} - \frac{\alpha_k^2}{(\alpha_k + \alpha_{\neq k})^2}}{\alpha_k + \alpha_{\neq k} + 1}.$$

The derivative is therefore

$$\begin{aligned} \frac{\partial}{\partial \alpha_k} \text{Var}(\pi_k | \boldsymbol{\alpha}) = \\ \frac{\alpha_{\neq k}(\alpha_{\neq k}^2 - \alpha_{\neq k}\alpha_k + \alpha_{\neq k} - \alpha_k(2\alpha_k + 1))}{(\alpha_k + \alpha_{\neq k})^3(\alpha_k + \alpha_{\neq k} + 1)^2}. \end{aligned}$$

22 Solving $\frac{\partial}{\partial \alpha_k} \text{Var}(\pi_k | \alpha) < 0$ for α_k yields

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

23 Therefore, under this hypothesis, $\text{Var}(\pi_k | \alpha)$ is a decreasing function of α_k . \square

24 Experimental Evaluation of the Proposition

25 [To test how often the condition is fulfilled we count its frequency. The fact that the condition
26 is fulfilled implies a good approximation. The fact that the condition is not fulfilled does not
27 automatically imply a bad approximation.]

		frequency
MNIST	MNIST	-
MNIST	FMNIST	-
MNIST	notMNIST	-
MNIST	KMNIST	-
CIFAR-10	CIFAR-10	0.998
CIFAR-10	CIFAR-100	0.925
CIFAR-10	SVHN	0.832
SVHN	SVHN	0.999
SVHN	CIFAR-100	0.668
SVHN	CIFAR-10	0.653
CIFAR-100	CIFAR-100	0.662
CIFAR-100	CIFAR-10	0.214
CIFAR-100	SVHN	0.166

Table 1

28 2 Appendix B: Laplace Approximation of the Dirichlet

29 Assume we have a Dirichlet in the standard basis with parameter vector α and probability density
30 function:

$$\text{Dir}(\pi | \alpha) := \frac{\Gamma\left(\sum_{k=1}^K \alpha_k\right)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K \pi_k^{\alpha_k - 1}, \quad (2)$$

31 We aim to transform the basis of this distribution via the softmax transform to be in the new base π :

$$\pi_k(\mathbf{z}) := \frac{\exp(z_k)}{\sum_{l=1}^K \exp(z_l)}, \quad (3)$$

32 Usually, to transform the basis we would need the inverse transformation $H^{-1}(\mathbf{z})$ as described in the
33 main paper. However, the softmax does not have an analytic inverse. Therefore David JC MacKay
34 uses the following trick. Assume we know that the distribution in the transformed basis is:

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

35 then we can show that the original distribution is the result of the basis transform by the softmax.

36 **The Dirichlet in the softmax basis:** We show that the density over π shown in Equation 4
37 transforms into the Dirichlet over \mathbf{z} . First, we consider the special case where π is confined to a
38 $I - 1$ dimensional subspace satisfying $\sum_i \pi_i = c$. In this subspace we can represent π by an $I - 1$
39 dimensional vector \mathbf{a} such that

$$\pi_i = \mathbf{a}_i \quad i, \dots, I-1 \quad (5)$$

$$\pi_I = c - \sum_i^{I-1} \mathbf{a}_i \quad (6)$$

and similarly we can represent \mathbf{z} by an $I-1$ dimensional vector φ :

$$z_i = \varphi_i \quad i, \dots, I-1 \quad (7)$$

$$z_I = 1 - \sum_i^{I-1} \varphi_i \quad (8)$$

then we can find the density over φ (which is proportional to the required density over \mathbf{z}) from the density over $\boldsymbol{\pi}$ (which is proportional to the given density over $\boldsymbol{\pi}$) by finding the determinant of the $(I-1) \times (I-1)$ Jacobian \mathbf{J} given by

$$\begin{aligned} J_{ik} &= \frac{\partial \xi_i}{\partial \tau_i} = \sum_j^I \frac{\partial x_i}{\partial \pi_j} \frac{\partial \pi_j}{\partial \tau_k} \\ &= \delta_{ik} \mathbf{x}_i - \mathbf{x}_i \mathbf{x}_k + \mathbf{x}_i \mathbf{x}_I = \mathbf{x}_i (\delta_{ik} - (\mathbf{x}_k - \mathbf{x}_I)) \end{aligned} \quad (9)$$

We define two additional $I-1$ dimensional helper vectors $\mathbf{x}_k^+ := \mathbf{x}_k - \mathbf{x}_I$ and $n_k := 1$, and use $\det(I - xy^T) = 1 - x \cdot y$ from linear algebra. It follows that

$$\begin{aligned} \det J &= \prod_{i=1}^{I-1} \mathbf{x}_i \times \det[I - n \mathbf{x}^{+T}] \\ &= \prod_{i=1}^{I-1} \mathbf{x}_i \times (1 - n \cdot \mathbf{x}^+) \\ &= \prod_{i=1}^{I-1} \mathbf{x}_i \times \left(1 - \sum_k \mathbf{x}_k^+\right) = I \prod_{i=1}^I \mathbf{x}_i \end{aligned} \quad (10)$$

Therefore, using Equation 4 we find that

$$P(\mathbf{z}) = \frac{P(\boldsymbol{\pi})}{|\det \mathbf{J}|} \propto \prod_{i=1}^I \mathbf{z}_i^{\alpha_i - 1} \quad (11)$$

This result is true for any constant c since it can be put into the normalizing constant. Thereby we make sure that the integral of the distribution is 1 and we have a valid probability distribution.

3 Appendix C: Inverting the Laplace Approximation of the Dirichlet

Through the figures of the 1D Dirichlet approximation in the main paper we have already established that the mode of the Dirichlet lies at the mean of the Gaussian distribution and therefore $\boldsymbol{\pi}(\mathbf{y}) = \frac{\alpha}{\sum_i \alpha_i}$. Additionally, the elements of \mathbf{y} must sum to zero. These two constraints combined yield only one possible solution for $\boldsymbol{\mu}$.

$$\mu_k = \log \alpha_k - \frac{1}{K} \sum_{l=1}^K \log \alpha_l \quad (12)$$

54 Calculating the covariance matrix Σ is more complicated but layed out in the following. The
 55 logarithm of the Dirichlet is, up to additive constants

$$\log p_y(y|\alpha) = \sum_k \alpha_k \pi_k \quad (13)$$

56 Using π_k as the softmax of \mathbf{y} as shown in Equation 3 we can find the elements of the Hessian \mathbf{L}

$$L_{kl} = \hat{\alpha}(\delta_{kl}\hat{\pi}_k - \hat{\pi}_k\hat{\pi}_l) \quad (14)$$

57 where $\hat{\alpha} := \sum_k \alpha_k$ and $\hat{\pi} = \frac{\alpha_k}{\hat{\alpha}}$ for the value of π at the mode. Analytically inverting \mathbf{L} is done
 58 via a lengthy derivation using the fact that we can write $\mathbf{L} = \mathbf{A} + \mathbf{XBX}^\top$ and inverting it with the
 59 Schur-complement. This process results in the inverse of the Hessian

$$L_{kl}^{-1} = \delta_{kl} \frac{1}{\alpha_k} - \frac{1}{K} \left[\frac{1}{\alpha_k} + \frac{1}{\alpha_l} - \frac{1}{K} \left(\sum_u \frac{1}{\alpha_u} \right) \right] \quad (15)$$

60 We are mostly interested in the diagonal elements, since we desire a sparse encoding for computational
 61 reasons and we otherwise needed to map a $K \times K$ covariance matrix to a $K \times 1$ Dirichlet parameter
 62 vector which would be a very overdetermined mapping. Note that K is a scalar not a matrix. The
 63 diagonal elements of $\Sigma = \mathbf{L}^{-1}$ can be calculated as

$$\Sigma_{kk} = \frac{1}{\alpha_k} \left(1 - \frac{2}{K} \right) + \frac{1}{K^2} \sum_l \frac{1}{\alpha_l}. \quad (16)$$

64 To invert this mapping we transform Equation 12 to

$$\alpha_k = e^{\mu_k} \prod_l \alpha_l^{1/K} \quad (17)$$

65 by applying the logarithm and re-ordering some parts. Inserting this into Equation 16 and re-arranging
 66 yields

$$\prod_l \alpha_l^{1/K} = \frac{1}{\Sigma_{kk}} \left[e^{-\mu} \left(1 - \frac{2}{K} \right) + \frac{1}{K^2} \sum_u e^{-\mu_u} \right] \quad (18)$$

67 which can be re-inserted into Equation 17 to give

$$\alpha_k = \frac{1}{\Sigma_{kk}} \left(1 - \frac{2}{K} + \frac{e^{-\mu_k}}{K^2} \sum_l e^{-\mu_l} \right) \quad (19)$$

68 which is the final mapping. With Equations 12 and 16 we are able to map from Dirichlet to Gaussian
 69 and with Equation 19 we are able to map the inverse direction.

70 **4 Appendix D: Experiments Details**

71 The exact experimental setups, i.e. network architectures, learning rates, random seeds, etc. can be
 72 found in the accompanying github repository (not linked for anonymity. Code is attached as zip file).
 73 This section is mostly used to justify some of the decisions we made during the process in more detail
 74 and highlight some miscellaneous interesting things.

75 Mathematical description of the setup

76 In principle, the Gaussian over the weights required by the Laplace Bridge for BNNs can be
 77 constructed by any Gaussian approximate Bayesian method such as variational Bayes [3; 4] and
 78 Laplace approximations for NNs [5; 6]. We will focus on the Laplace approximation, which uses the
 79 same principle as the Laplace Bridge. However, in the Laplace approximation for neural networks,
 80 the posterior distribution over the weights of a network is the one that is approximated as a Gaussian,
 81 instead of a Dirichlet distribution over the outputs as in the Laplace Bridge.

82 Given a dataset $\mathcal{D} := \{(\mathbf{x}_i, t_i)\}_{i=1}^D$ and a prior $p(\boldsymbol{\theta})$, let

$$p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\boldsymbol{\theta})p(\mathcal{D}|\boldsymbol{\theta}) = p(\boldsymbol{\theta}) \prod_{(\mathbf{x}, t) \in \mathcal{D}} p(y = t|\boldsymbol{\theta}, \mathbf{x}), \quad (20)$$

83 be the posterior over the parameter $\boldsymbol{\theta}$ of an L -layer network $f_{\boldsymbol{\theta}}$. Then we can get an approximation of
 84 the posterior $p(\boldsymbol{\theta}|\mathcal{D})$ by fitting a Gaussian $\mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\mu}_{\boldsymbol{\theta}}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}})$ where

$$\begin{aligned} \boldsymbol{\mu}_{\boldsymbol{\theta}} &= \boldsymbol{\theta}_{\text{MAP}}, \\ \boldsymbol{\Sigma}_{\boldsymbol{\theta}} &= (-\nabla^2|_{\boldsymbol{\theta}_{\text{MAP}}} \log p(\boldsymbol{\theta}|\mathcal{D}))^{-1} =: \mathbf{H}_{\boldsymbol{\theta}}^{-1}. \end{aligned}$$

85 That is, we fit a Gaussian centered at the mode $\boldsymbol{\theta}_{\text{MAP}}$ of $p(\boldsymbol{\theta}|\mathcal{D})$ with the covariance determined by the
 86 curvature at that point. We assume that the prior $p(\boldsymbol{\theta})$ is a zero-mean isotropic Gaussian $\mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, \sigma^2 \mathbf{I})$
 87 and the likelihood function is the Categorical density

$$p(\mathcal{D}|\boldsymbol{\theta}) = \prod_{(\mathbf{x}, t) \in \mathcal{D}} \text{Cat}(y = t | \text{softmax}(f_{\boldsymbol{\theta}}(\mathbf{x}))).$$

88 For various applications in Deep Learning, an approximation with full Hessian is often computationally
 89 too expensive. Indeed, for each input $\mathbf{x} \in \mathbb{R}^N$, one has to do K backward passes to compute
 90 the Jacobian $\mathbf{J}(\mathbf{x})$. Moreover, it requires an $\mathcal{O}(PK)$ storage which is also expensive since P is
 91 often in the order of millions. A cheaper alternative is to fix all but the last layer of $f_{\boldsymbol{\theta}}$ and only
 92 apply the Laplace approximation on \mathbf{W}_L , the last layer's weight matrix. This scheme has been
 93 used successfully by Snoek et al. [7]; Wilson et al. [8]; Brosse et al. [9], etc. and has been shown
 94 theoretically that it can mitigate overconfidence problems in ReLU networks [10]. In this case, given
 95 the approximate last-layer posterior

$$p(\mathbf{W}^L|\mathcal{D}) \approx \mathcal{N}(\text{vec}(\mathbf{W}^L) | \text{vec}(\mathbf{W}_{\text{MAP}}^L), \mathbf{H}_{\mathbf{W}^L}^{-1}), \quad (21)$$

96 one can efficiently compute the distribution over the logits. That is, let $\phi : \mathbb{R}^N \rightarrow \mathbb{R}^Q$ be the first
 97 $L - 1$ layers of $f_{\boldsymbol{\theta}}$, seen as a feature map. Then, for each $\mathbf{x} \in \mathbb{R}^N$, the induced distribution over the
 98 logit $\mathbf{W}^L \phi(\mathbf{x}) =: \mathbf{z}$ is given by

$$p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z} | \mathbf{W}_{\text{MAP}}^L \phi(\mathbf{x}), (\phi(\mathbf{x})^\top \otimes \mathbf{I}) \mathbf{H}_{\mathbf{W}^L}^{-1} (\phi(\mathbf{x}) \otimes \mathbf{I})), \quad (22)$$

99 where \otimes denotes the Kronecker product.

100 An even more efficient last-layer approximation can be obtained using a Kronecker-factored matrix
 101 normal distribution [11; 12; 6]. That is, we assume the posterior distribution to be

$$p(\mathbf{W}^L|\mathcal{D}) \approx \mathcal{MN}(\mathbf{W}^L | \mathbf{W}_{\text{MAP}}^L, \mathbf{U}, \mathbf{V}), \quad (23)$$

102 where $\mathbf{U} \in \mathbb{R}^{K \times K}$ and $\mathbf{V} \in \mathbb{R}^{Q \times Q}$ are the Kronecker factorization of the inverse Hessian matrix
 103 $\mathbf{H}_{\mathbf{W}^L}^{-1}$ [13]. In this case, for any $\mathbf{x} \in \mathbb{R}^N$, one can easily show that the distribution over logits is
 104 given by

$$p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z} | \mathbf{W}_{\text{MAP}}^L \phi(\mathbf{x}), (\phi(\mathbf{x})^\top \mathbf{V} \phi(\mathbf{x})) \mathbf{U}), \quad (24)$$

105 which is easy to implement and computationally cheap. Finally, and even more efficient, is a last-layer
 106 approximation scheme with a diagonal Gaussian approximate posterior, i.e. the so-called mean-field
 107 approximation. In this case, we assume the posterior distribution to be

$$p(\mathbf{W}^L|\mathcal{D}) \approx \mathcal{N}(\text{vec}(\mathbf{W}^L) | \text{vec}(\mathbf{W}_{\text{MAP}}^L), \text{diag}(\boldsymbol{\sigma}^2)), \quad (25)$$

108 where $\boldsymbol{\sigma}^2$ is obtained via the diagonal of the Hessian of the log-posterior w.r.t. $\text{vec}(\mathbf{W}^L)$ at
 109 $\text{vec}(\mathbf{W}_{\text{MAP}}^L)$.

Uncertainty estimates on MNIST

Most of the experimental setup is already explained in the main paper. The exact details can be found in the accompanying code. Every experiment has been conducted with 5 different seeds.

OOD detection

Every experiment has been conducted with 5 different seeds. In the tables the mean and standard deviations are presented. The reason why the sampling procedure for the CIFAR-10 and CIFAR-100 case are similarly fast even though we draw from a 10- vs 100-dimensional Gaussian is because the sampling procedures were parallelized on a GPU. All prior uncertainties over the weights were chosen such that the MMC of the sampling averages was around 5% lower than the MAP estimate. In the following we show the results including a KFAC approximation of the last layer.

Table 2: OOD detection results. Same table as in the main text but a different GPU (GTX 1080 Ti) was used. Results for the all-layer Laplace approximation have been changed since a bug was fixed.

Train	Test	Diag Sampling		Diag LB		KFAC Sampling		KFAC LB		Time in s	
		MMC ↓	AUROC ↑	MMC ↓	AUROC ↑	MMC ↓	AUROC ↑	MMC ↓	AUROC ↑	Sampling	LB
MNIST	MNIST	0.942 ± 0.007	-	0.987 ± 0.000	-	-	-	-	-	26.8	0.062
MNIST	FMNIST	0.397 ± 0.001	0.992 ± 0.000	0.363 ± 0.000	0.996 ± 0.000	-	-	-	-	26.8	0.062
MNIST	notMNIST	0.543 ± 0.000	0.960 ± 0.000	0.649 ± 0.000	0.961 ± 0.000	-	-	-	-	50.3	0.117
MNIST	KMNIST	0.513 ± 0.001	0.974 ± 0.000	0.637 ± 0.000	0.973 ± 0.000	-	-	-	-	26.9	0.062
CIFAR-10	CIFAR-10	0.948 ± 0.000	-	0.966 ± 0.000	-	0.857 ± 0.003	-	0.966 ± 0.000	-	6.58	0.017
CIFAR-10	CIFAR-100	0.708 ± 0.000	0.889 ± 0.000	0.742 ± 0.000	0.866 ± 0.000	0.562 ± 0.003	0.880 ± 0.012	0.741 ± 0.000	0.866 ± 0.000	6.59	0.016
CIFAR-10	SVHN	0.643 ± 0.000	0.933 ± 0.000	0.647 ± 0.000	0.934 ± 0.000	0.484 ± 0.004	0.939 ± 0.001	0.648 ± 0.003	0.934 ± 0.001	17.0	0.040
SVHN	SVHN	0.986 ± 0.000	-	0.993 ± 0.000	-	0.947 ± 0.002	-	0.993 ± 0.000	-	17.1	0.042
SVHN	CIFAR-100	0.595 ± 0.000	0.984 ± 0.000	0.526 ± 0.000	0.985 ± 0.000	0.460 ± 0.004	0.986 ± 0.001	0.527 ± 0.002	0.985 ± 0.000	6.62	0.017
SVHN	CIFAR-10	0.593 ± 0.000	0.984 ± 0.000	0.520 ± 0.000	0.987 ± 0.000	0.458 ± 0.004	0.986 ± 0.001	0.520 ± 0.002	0.987 ± 0.000	6.62	0.017
CIFAR-100	CIFAR-100	0.762 ± 0.000	-	0.590 ± 0.000	-	0.404 ± 0.000	-	0.593 ± 0.000	-	6.76	0.016
CIFAR-100	CIFAR-10	0.467 ± 0.000	0.788 ± 0.000	0.206 ± 0.000	0.791 ± 0.000	0.213 ± 0.000	0.788 ± 0.000	0.209 ± 0.000	0.791 ± 0.000	6.71	0.017
CIFAR-100	SVHN	0.461 ± 0.000	0.795 ± 0.000	0.170 ± 0.000	0.815 ± 0.000	0.180 ± 0.001	0.838 ± 0.001	0.173 ± 0.000	0.815 ± 0.000	17.3	0.041

Train	Test	Sampling (100)			Laplace Bridge		
		MMC ↓	AUROC ↑	Time in s ↓	MMC ↓	AUROC ↑	Time in s ↓
MNIST	MNIST	0.989 ± 0.000	-	76.3	0.987 ± 0.000	-	0.062
MNIST	FMNIST	0.465 ± 0.004	0.992 ± 0.000	76.3	0.363 ± 0.000	0.996 ± 0.000	0.062
MNIST	notMNIST	0.645 ± 0.001	0.955 ± 0.001	42.6	0.649 ± 0.000	0.961 ± 0.000	0.117
MNIST	KMNIST	0.632 ± 0.001	0.967 ± 0.001	76.0	0.637 ± 0.000	0.973 ± 0.000	0.062

Table 3: Results for sampling from all weights instead of the last layer. Number of samples was 100. Sampling all weights and doing a forward pass seems to be worse than using the Laplace Bridge even though it takes much longer.

Time comparison

Every experiment has been conducted with 5 different seeds. The presented curves are the averages over these 5 experiments with errorbars. The reason why taking one sample is slower than two is because of the way random numbers are generated for the normal distribution. For further information read up on the Box-Mueller Transform.

Uncertainty-aware output ranking on ImageNet

The prior covariances for the Laplace approximation of the Hessian over the weights were chosen such that uncertainty estimate of the Laplace bridge MMC over the outputs was not more than 5% lower than the MAP estimate. The length of list generated by our uncertainty aware method was chosen such that it contained at least one and maximally ten samples. Originally we wanted to choose the maximal length according to the size of the largest category (e.g. fishes or dogs) but the class tree hierarchy of ImageNet does not answer this question meaningfully. We chose ten because there are no reasonable bins larger than ten when looking at a histogram. Pseudo code for the uncertainty aware top-k ranking can be found in Algorithm 1.

Algorithm 1 Uncertainty-aware top- k

Input: A Dirichlet parameter $\alpha \in \mathbb{R}^K$ obtained by applying the Laplace Bridge to the Gaussian over the logit of an input, a percentile threshold T e.g. 0.05, a function `class_of` that returns the underlying class of a sorted index.

```
 $\tilde{\alpha} = \text{sort\_descending}(\alpha)$  // start with the highest confidence
 $\alpha_0 = \sum_i \alpha_i$ 
 $\mathcal{C} = \{\text{class\_of}(1)\}$  // initialize top-k, must include at least one class
for  $i = 2, \dots, K$  do
     $F_{i-1} = \text{Beta}(\tilde{\alpha}_{i-1}, \alpha_0 - \tilde{\alpha}_{i-1})$  // the previous marginal CDF
     $F_i = \text{Beta}(\tilde{\alpha}_i, \alpha_0 - \tilde{\alpha}_i)$  // the current marginal CDF
     $l_{i-1} = F_{i-1}^{-1}(T/2)$  // left  $\frac{T}{2}$  percentile of the previous marginal
     $r_i = F_i^{-1}(1 - T/2)$  // right  $\frac{T}{2}$  percentile of the current marginal
    if  $r_i > l_{i-1}$  then
         $\mathcal{C} = \mathcal{C} \cup \{\text{class\_of}(i)\}$  // overlap detected, add the current class
    else
        break // No more overlap, end the algorithm
    end if
end for
```

Output: \mathcal{C} // return the resulting top-k prediction

References

- [1] David J.C. MacKay. Choice of basis for laplace approximation. *Machine Learning*, 33(1): 77–86, Oct 1998. ISSN 1573-0565.
- [2] P. Hennig. *Approximate Inference in Graphical Models*. PhD thesis, November 2010.
- [3] 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.
- [4] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. *ArXiv*, 2015.
- [5] David J. C. MacKay. A practical bayesian framework for backpropagation networks. *Neural Comput.*, 4(3):448–472, May 1992. ISSN 0899-7667.
- [6] Hippolyt Ritter, Aleksandar Botev, and David Barber. A scalable laplace approximation for neural networks. In *International Conference on Learning Representations*, 2018.
- [7] 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.
- [8] 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.
- [9] 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.
- [10] Agustinus Kristiadi, Matthias Hein, and Philipp Hennig. Being bayesian, even just a bit, fixes overconfidence in relu networks. *arXiv preprint arXiv:2002.10118*, 2020.

- 161 [11] Christos Louizos and Max Welling. Structured and efficient variational deep learning with
162 matrix gaussian posteriors. In *ICML*, 2016.
- 163 [12] Shengyang Sun, Changyou Chen, and Lawrence Carin. Learning structured weight uncertainty
164 in Bayesian neural networks. In *Artificial Intelligence and Statistics*, pages 1283–1292, 2017.
- 165 [13] James Martens and Roger Grosse. Optimizing neural networks with Kronecker-factored
166 approximate curvature. In *ICML*, 2015.