

A APPENDIX

A.1 IMPLEMENTATION DETAILS

To assess our proposed model’s predictive performance and uncertainty estimation capabilities, we conducted experiments on synthetic two moons data (Pedregosa et al., 2011), a mixture of Gaussians, the CIFAR-100 dataset (Krizhevsky et al., 2009), and the Imagenet dataset (Deng et al., 2009). We compare against a standard deterministic ResNet model as a baseline (He et al., 2016), against the heteroscedastic method (Collier et al., 2020; 2021) and the SNGP (Liu et al., 2020) (which form the basis for our combined model) and against the recently proposed Posterior Network model (Charpentier et al., 2020), which also offers distance-aware uncertainties, similarly to the SNGP. We used the same backbone neural network architecture for all models, which was a fully-connected ResNet for the synthetic data, a WideResNet18 on CIFAR and a ResNet50 in Imagenet.

For most baselines, we used the hyperparameters from the `uncertainty_baselines` library (Nado et al., 2021). On CIFAR, we trained our HetSNGP with a learning rate of 0.1 for 300 epochs and used $R = 6$ factors for the heteroscedastic covariance, a softmax temperature of $\tau = 0.5$ and $S = 5000$ Monte Carlo samples. On Imagenet, we trained with a learning rate of 0.07 for 270 epochs and used $R = 15$ factors, a softmax temperature of $\tau = 1.25$ and $S = 5000$ Monte Carlo samples. We implemented all models in TensorFlow in Python and trained on Tensor Processing Units (TPUs) in the Google Cloud.

We train all Imagenet-21k models for 90 epochs with batch size 1024 on 8×8 TPU slices. We train using the Adam optimizer with initial learning rate of 0.001 using a linear learning rate decay schedule with termination point 0.00001 and a warm-up period of 10,000 steps. We train using the sigmoid cross-entropy loss function and L2 weight decay with multiplier 0.03. The heteroscedastic method uses a temperature of 0.4, 1,000 Monte Carlo samples and $R = 50$ for the low rank approximation. HetSNGP has the same heteroscedastic hyperparameters except the optimal temperature is 1.5. For SNGP and HetSNGP the GP covariance is approximated using the momentum scheme presented in Liu et al. (2020) with momentum parameter 0.999.

A.2 LAPLACE APPROXIMATION

In this section, we will derive the Laplace posterior in Eq. (5). The derivation follows mostly from the sections 3.4 and 3.5 in Rasmussen & Williams (2006).

First note that the log posterior of β_c given the data is

$$\log p(\beta_c | \mathbf{x}, y) = \log p(y | \beta_c) + \log p(\beta_c) - Z \quad (7)$$

where Z is a normalization constant that does not depend on β_c . Following Rasmussen & Williams (2006), we will denote the unnormalized log posterior as

$$\Psi(\beta_c) = \log p(y | \beta_c) + \log p(\beta_c) \quad (8)$$

Recall that the first term is the likelihood and the second term is our prior from Eq. (4).

The Laplace approximation now approximates the posterior with a local second-order expansion around the MAP solution, that is

$$p(\beta_c | \mathbf{x}, y) \approx \mathcal{N}(\hat{\beta}_c, \mathbf{\Lambda}^{-1}) \quad (9)$$

with the MAP solution $\hat{\beta}_c = \arg \max_{\beta_c} \Psi(\beta_c)$ and the Hessian $\mathbf{\Lambda} = -\nabla^2 \Psi(\beta_c)|_{\beta_c=\hat{\beta}_c}$.

The MAP solution can be found using standard (stochastic) gradient descent, while the Hessian is given by

$$\begin{aligned} \nabla^2 \Psi(\beta_c) &= \nabla^2 \log p(y | \beta_c) + \nabla^2 \log p(\beta_c) \\ &= \nabla_{\beta} (\nabla_{\mathbf{u}} \log p(y | \mathbf{u}) \nabla_{\beta} \mathbf{u}) - \mathbf{I}_m \\ &= \nabla_{\beta} (\nabla_{\mathbf{u}} \log p(y | \mathbf{u}) \Phi) - \mathbf{I}_m \\ &= \Phi^{\top} \nabla_{\mathbf{u}}^2 \log p(y | \mathbf{u}) \Phi - \mathbf{I}_m \\ &= -W \Phi^{\top} \Phi - \mathbf{I}_m \end{aligned}$$

where we used the chain rule and the fact that $\mathbf{u} = \Phi\beta$ and W is a diagonal matrix of point-wise second derivatives of the likelihood, that is, $W_{ii} = -\nabla^2 \log p(y_i | \mathbf{u}_i)$ (Rasmussen & Williams, 2006). For instance, in the case of the logistic likelihood, $W_{ii} = \mathbf{p}_i (1 - \mathbf{p}_i)$, where \mathbf{p}_i is a vector of output probabilities for logits \mathbf{u}_i . To get the Hessian at the MAP, we then just need to compute this quantity for $\hat{\mathbf{u}} = \Phi\hat{\beta}$.

The approximate posterior is therefore

$$p(\beta_c | \mathbf{x}, y) \approx \mathcal{N}(\hat{\beta}_c, (W\Phi^T\Phi + \mathbf{I}_m)^{-1}) \quad (10)$$

where the precision matrix can be computed over data points (recovering Eq. (5)) as

$$\Lambda = \mathbf{I}_m + \sum_{i=1}^N \mathbf{p}_i (1 - \mathbf{p}_i) \Phi_i \Phi_i^\top \quad (11)$$