

# RELIABLE UNCERTAINTY ESTIMATES IN DEEP NEURAL NETWORKS USING NOISE CONTRASTIVE PRIORS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Obtaining reliable uncertainty estimates of neural network predictions is a long standing challenge. Bayesian neural networks have been proposed as a solution, but it remains open how to specify the prior. In particular, the common practice of a standard normal prior in weight space imposes only weak regularities, causing the function posterior to possibly generalize in unforeseen ways on out-of-distribution inputs. We propose noise contrastive priors (NCPs). The key idea is to train the model to output high uncertainty for data points outside of the training distribution. NCPs do so using an input prior, which adds noise to the inputs of the current mini batch, and an output prior, which is a wide distribution given these inputs. NCPs are compatible with any model that represents predictive uncertainty, are easy to scale, and yield reliable uncertainty estimates throughout training. Empirically, we show that NCPs prevent overfitting outside of the training distribution and result in uncertainty estimates that are useful for active learning. We demonstrate the scalability of our method on the flight delays data set, where we significantly improve upon previously published results.

## 1 INTRODUCTION

Many successful applications of neural networks (Krizhevsky et al., 2012; Sutskever et al., 2014; van den Oord et al., 2016) are in restricted settings where predictions are only made for inputs similar to the training distribution. In real-world scenarios, neural networks can face truly novel data points during inference, and in these settings it can be valuable to have good estimates of the model’s uncertainty. For example, in healthcare, reliable uncertainty estimates can prevent overconfident decisions for rare or novel patient conditions (Schulam and Saria, 2015). Similarly, autonomous agents that actively explore their environment can use uncertainty estimates to decide what data points will be most informative (MacKay, 1992a).

Epistemic uncertainty describes the degree of missing knowledge about the data generating function. Uncertainty can in principle be completely reduced by observing more data points at the right locations and training on them. In contrast, the data generating function may also have inherent randomness, which we call aleatoric noise. This noise can be captured by models outputting a distribution rather than a point prediction. Obtaining more data points will move the noise estimate closer to its true value, which is usually different from zero. For active learning, it is crucial to separate the two types of randomness: we want to acquire labels in regions of high uncertainty but low noise.

Bayesian analysis provides a principled approach to modeling uncertainty in neural networks (Denker et al., 1987; MacKay, 1992b). Namely, one places a prior over the network’s weights and biases. This effectively places a distribution over the functions that the network can represent, capturing uncertainty in which function best fits the data. Specifying the prior remains an open challenge. Common practice is to use a standard normal prior in weight space, which imposes weak shrinkage regularities analogous to weight decay. It is

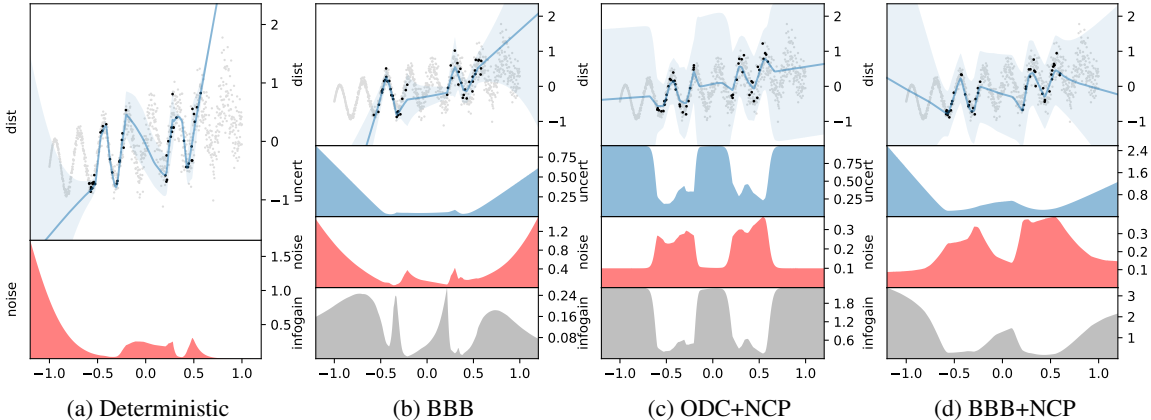


Figure 1: Predictive distributions on a low-dimensional active learning task. The predictive distributions are visualized as mean and two standard deviations shaded. They decompose into epistemic uncertainty ■ and aleatoric noise ■. Data points are only available within two bands, and are selected using the expected information gain ■. **(a)** A deterministic network conflates uncertainty as part of the noise and is overconfident outside of the data distribution. **(b)** A variational Bayesian neural network with standard normal prior represents uncertainty and noise separately but is overconfident outside of the training distribution. **(c)** On the OOD classifier model, NCP prevents overconfidence. **(d)** On the Bayesian neural network, NCP produces smooth uncertainty estimates that generalize well to unseen data points. Models trained with NCP also separate uncertainty and noise well. The experimental setup is described in Section 5.1.

neither informative about the induced function class nor the data (e.g., it is sensitive to parameterization). This can cause the induced function posterior to generalize in unforeseen ways on out-of-distribution (OOD) inputs, which are inputs outside of the distribution that generated the training data.

Motivated by these challenges, we introduce noise contrastive priors (NCPs), which encourage uncertainty outside of the training distribution using a loss in data space. NCPs are compatible with any model that represents functional uncertainty as a random variable, are easy to scale, and yield reliable uncertainty estimates that show significantly improved active learning performance.

## 2 RELATED WORK

**Priors for neural networks** Classic work has investigated entropic priors (Buntine and Weigend, 1991) and hierarchical priors (MacKay, 1992b; Neal, 2012; Lampinen and Vehtari, 2001). More recently, Depeweg et al. (2018) introduce networks with latent variables in order to disentangle forms of uncertainty, and Flam-Shepherd et al. (2017) propose general-purpose weight priors based on approximating Gaussian processes. Other works have analyzed priors for compression and model selection (Ghosh and Doshi-Velez, 2017; Louizos et al., 2017). Instead of a prior in weight space (or latent inputs as in Depeweg et al. (2018)), NCPs take the functional view by imposing explicit regularities in terms of the network’s inputs and outputs.

**Input and output regularization** There is classic work on adding noise to inputs for improved generalization (Matsuoka, 1992; An, 1996; Bishop, 1995). For example, denoising autoencoders (Vincent et al., 2008) encourage reconstructions given noisy encodings. Output regularization is also a classic idea from the maximum entropy principle (Jaynes, 1957), where it has motivated label smoothing (Szegedy et al., 2016) and entropy penalties (Pereyra et al., 2017). Also related is virtual adversarial training (Miyato et al., 2015), which includes examples that are close to the current input but cause a maximal change in the model output, and

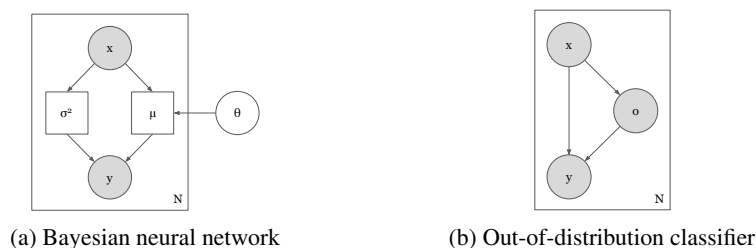


Figure 2: Graphical representations of the two uncertainty-aware models we consider. Circles denote random variables, squares denote deterministic variables, shading denotes observations during training. **(a)** The Bayesian neural network captures a belief over parameters for the predictive mean, while the predictive variance is a deterministic function of the input. In practice, we only use weight uncertainty for the mean’s output layer and share earlier layers between the mean and variance. **(b)** The out-of-distribution classifier model uses a binary auxiliary variable  $o$  to determine if a given input is out-of-distribution; given its value, the output is drawn from either a neural network prediction or a wide output prior.

mixup (Zhang et al., 2018), which includes examples under the vicinity of training data. These methods are all orthogonal to NCPs: they aim to improve generalization from finite data (but under the same distribution); our approach aims to improve uncertainty estimates outside of the training distribution.

**Classifying out-of-distribution inputs** A simple approach for neural network uncertainty is to classify whether data points belong to the data distribution, or are OOD (Hendrycks and Gimpel, 2017). This is core to noise contrastive estimation (Gutmann and Hyvärinen, 2010b), a training method for intractable probabilistic models. More recently, Lee et al. (2017) introduce a GAN to generate OOD samples, and Liang et al. (2018) add perturbations to the input, applying an “OOD detector” to improve softmax scores on OOD samples according to a scaled temperature. Extending these directions of research, we connect to Bayesian principles and focus on uncertainty estimates that are useful for active data acquisition.

### 3 NOISE CONTRASTIVE PRIORS

Specifying priors is intuitive for small probabilistic models, as each variable typically has a clear interpretation (Blei, 2014). It is less intuitive for neural networks: the parameters serve more as adaptive basis coefficients in a nonparametric function. For example, neural network models are non-identifiable due to weight symmetries that can produce the same function output (Müller and Insua, 1998). This makes it difficult to express informative priors on the weights, such as how to express high uncertainty on unfamiliar examples.

**Data priors** Unlike a prior in weight space, a *data prior* lets one easily express informative assumptions about input-output relationships. Here, we use the example of a prior over a labeled data set  $(x, y)$ , although the prior can also be on  $x$  and another variable in the model that represents uncertainty and has a clear interpretation. The prior takes the form  $p_{\text{prior}}(x, y) = p_{\text{prior}}(x) p_{\text{prior}}(y | x)$ , where  $p_{\text{prior}}(x)$  represents the *input prior* and  $p_{\text{prior}}(y | x)$  denotes the *output prior*.

To prevent overconfident predictions, a good input prior  $p_{\text{prior}}(x)$  should include OOD examples so that it is enforced beyond the training distribution. A good output prior  $p_{\text{prior}}(y | x)$  should be a high-variance distribution, representing high uncertainty about the model output given OOD inputs.

**Generating OOD inputs** Exactly generating OOD data is difficult. A priori, we must uniformly represent the input domain. A posteriori, we must represent the complement of the training distribution. Both distributions are uniform over infinite support, making them ill-defined. To estimate the OOD inputs, we

develop an algorithm inspired by noise contrastive estimation (Gutmann and Hyvärinen, 2010a; Mnih and Kavukcuoglu, 2013), where a complement distribution is approximated using random noise.

A hypothesis of our work is that in practice it is enough to encourage high uncertainty output near the *boundary* of the training distribution, and that this effect will propagate to the entire OOD space. This hypothesis is backed up by previous work (Lee et al., 2017) as well as our experiments (see Figure 1). This means we no longer need to sample arbitrary OOD inputs. It is enough to sample OOD points that lie close to the boundary of the training distribution, and to apply our desired prior at those points.

**Loss function** Noise contrastive priors are data priors where we approximate OOD inputs by perturbing training inputs  $x$  with noise. For example, in binary and categorical input domains, we approximate OOD inputs by randomly flipping the features to different classes with a certain probability. For continuous valued inputs  $x$ , we can apply additive Gaussian noise to obtain noised up inputs  $\tilde{x} = x + \epsilon$ . This expresses the noise contrastive prior where inputs are distributed according to the convolved distribution,

$$p_{\text{prior}}(\tilde{x}) = \int_x p_{\text{train}}(x) \text{Normal}(\tilde{x} - x \mid \mu_x, \sigma_x^2) dx \quad p_{\text{prior}}(\tilde{y} \mid \tilde{x}) = \text{Normal}(\mu_y, \sigma_y^2). \quad (1)$$

The variances  $\sigma_x^2$  and  $\sigma_y^2$  are hyperparameters that tune how far from the boundary we sample, and how large we want the output uncertainty to be. We choose  $\mu_x = 0$  to apply the prior equally in all directions from the data manifold. The output mean  $\mu_y$  determines the default prediction of the model outside of the training distribution, for example  $\mu_y = 0$ . We set  $\mu_y = y$  which corresponds to data augmentation (Matsuoka, 1992; An, 1996), where a model is trained to recover the true labels from perturbed inputs. This way, NCP makes the model uncertain while still trying to be accurate on OOD inputs.

NCPs add noise to all data inputs rather than manually selecting the subset on the boundary. In our experiments, we found that this does not affect performance: noised-up inputs that remain in the training distribution’s support can be seen as a form of label smoothing, and it avoids a potentially sensitive preprocessing step.

For training, we minimize the loss function

$$\mathcal{L}(\theta) = \mathbb{E}_{p_{\text{train}}(x)} [D_{\text{KL}}[p_{\text{train}}(y \mid x) \parallel p_{\text{model}}(y \mid x, \theta)]] + \gamma \mathbb{E}_{p_{\text{prior}}(\tilde{x})} [D_{\text{KL}}[p_{\text{prior}}(\tilde{y} \mid \tilde{x}) \parallel p_{\text{model}}(\tilde{y} \mid \tilde{x}, \theta)]] \quad (2)$$

The first term represents typical maximum likelihood, in which one minimizes the KL divergence to the empirical training distribution  $p_{\text{train}}(y \mid x)$  over training inputs. The second term is added by our method: it represents the analogous term on a data prior. The hyperparameter  $\gamma$  sets the relative trade-off between them.

**Interpretation as function prior** The noise contrastive prior can be interpreted as inducing a function prior. This is formalized through the predictive distribution, which takes the form

$$p(y \mid x) = \int p_{\text{model}}(y \mid x, \theta) p_{\text{model}}(\theta \mid \tilde{x}, \tilde{y}) p_{\text{prior}}(\tilde{x}, \tilde{y}) d\theta d\tilde{x} d\tilde{y}. \quad (3)$$

The distribution marginalizes over network parameters  $\theta$  as well as data fantasized from the data prior. The distribution  $p(\theta \mid \tilde{x}, \tilde{y})$  represents the distribution of model parameters after fitting the prior data. That is, the belief over weights is shaped to make  $p(y \mid x)$  highly variable. This parameter belief makes uncertain predictions outside of the training distribution, which we could not specify in weight space directly.

Because network weights are constrained to fit the data prior, the prior acts as “pseudo-data.” This is similar to classical work on conjugate priors: a Beta( $\alpha, \beta$ ) prior on the probability of a Bernoulli likelihood implies a Beta posterior, and if the posterior mode is chosen as an optimal parameter setting, then the prior translates to  $\alpha - 1$  successes and  $\beta - 1$  failures. It is also similar to pseudo-data in sparse Gaussian processes (Quiñero-Candela and Rasmussen, 2005).

Data priors encourage learning parameters that not only capture the training data well but also the prior data. In practice, we may combine NCP with other priors, for example the typical standard normal prior in weight space for Bayesian neural networks, although we did not find this necessary in our experiments.

## 4 BAYESIAN NEURAL NETWORKS WITH NCP

Noise contrastive priors are applicable to any model that represents uncertainty in a random variable. The NCP can then be added to that random variable to make the model uncertain on OOD inputs. In this section, we apply NCP to a Bayesian neural network (BNN) trained via variational inference. Blundell et al. (2015) introduce such a model under the name Bayes by Backprop (BBB) that uses a standard normal prior in weight space. We extend this model with a NCP on the mean predicted by the neural network.

Consider a regression task with data  $\{x, y\}$  that we model as  $p(y | x, \theta) = \text{Normal}(\mu(x), \sigma^2(x))$  with mean and variance predicted by a neural network from the inputs. This model is heteroskedastic, meaning that it can predict a different aleatoric noise for every point in the input space. We use a weight prior for only the output layer (Lázaro-Gredilla and Figueiras-Vidal, 2010; Calandra et al., 2014) that predicts the mean, resulting in the model

$$\theta \sim \text{Normal}(0, 0.1) \quad y \sim \text{Normal}(\mu(x, \theta), \sigma^2(x)). \quad (4)$$

In this model, the distribution of the mean induced by the weight prior represents epistemic uncertainty, i.e.,  $q(\mu(x)) = \int \mu(x, \theta) q_\phi(\theta) d\theta$ . Note that this is different from the predictive distribution, which combines both uncertainty and noise. We place an NCP on the distribution of the mean, resulting in the loss function

$$\mathcal{L}(\phi) = -\mathbb{E}_{q_\phi(\theta)}[\ln p(y | x, \theta)] + \beta D_{\text{KL}}[q_\phi(\theta) \| p(\theta)] + \underbrace{\gamma D_{\text{KL}}[\text{Normal}(\mu_\mu, \sigma_\mu^2) \| q(\mu(\tilde{x}))]}_{\text{NCP loss}}. \quad (5)$$

Here,  $\tilde{x}$  are the perturbed inputs and  $q_\phi(\theta)$  forms an approximate posterior over weights.<sup>1</sup> Because we only use the weight belief for the linear output layer, we can compute the KL-divergence of the NCP loss analytically. In other models, it may be approximated using samples.

The loss function applies weight regularization in order for network weights to regress to standard normal prior; like other regularization techniques, this assists in improving the network’s generalization in-distribution. The NCP loss encourages the network’s generalization OOD by matching the mean distribution to the output prior. Minimizing the KL divergence to a wide output prior results in high uncertainty on OOD inputs, so the model will explore these data points during active learning.

In practice, we find that NCP is sufficient as a prior for the BNN and set  $\beta = 0$ . The appendix (Appendix B) includes an alternative interpretation explaining why NCP might be sufficient, which represents the weight space KL-divergence after a change of variables.

## 5 EXPERIMENTS

To demonstrate their usefulness, we evaluate NCPs on various tasks where uncertainty estimates are desired. Our focus is on active learning for regression tasks, where only few targets are visible in the beginning, and additional targets are selected regularly based on an acquisition function. We use two data sets—a toy example and a large flights data set—and we also evaluate how sensitive our method is to the choice of input noise. Finally, we show that NCP scales to large data sets by training on the full flights data set in a passive learning setting. Our implementation uses TensorFlow Probability (Dillon et al., 2017; Tran et al., 2016) and is open-sourced at <https://<hidden-for-review>>.

<sup>1</sup>To derive the loss, set  $p(y | x, \theta) = \mathbb{E}_{q_\phi(\theta)}[p(y | x, \theta)]$  in Equation 2 and apply Jensen’s inequality (Blundell et al., 2015; Higgins et al., 2016).

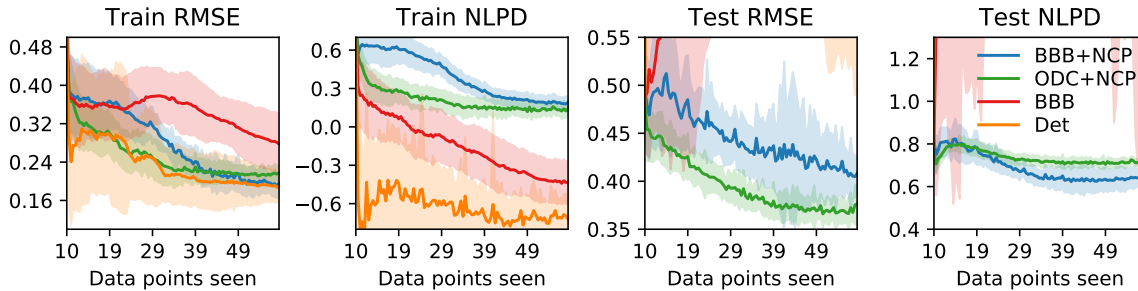


Figure 3: Active learning on the 1-dimensional regression problem, mean and standard deviation over 20 seeds. The test root mean squared error (RMSE) and negative log predictive density (NLPD) of the models trained with NCP decreases during the active learning run, while the baseline models select less informative data and overfit. The deterministic network is barely visible in the plots as it overfits quickly. Figure 1 shows the predictive distributions of the models.

We compare four neural network models, all using leaky ReLU activations (Maas et al., 2013) and trained using Adam (Kingma and Ba, 2014). The four models are:

- **Deterministic neural network (Det)** A neural network that predicts the mean and variance of a normal distribution. The name standard for deterministic, as there is no weight uncertainty.
- **Bayes by Backprop (BBB)** A Bayesian neural network trained via gradient-based variational inference with a standard normal prior in weight space (Blundell et al., 2015; Kucukelbir et al., 2017). We use the same model as in Section 4 but without the NCP loss term.
- **Bayes by Backprop with noise contrastive prior (BBB+NCP)** Bayes by Backprop with NCP on the predicted mean distribution as described in Section 4.
- **Out-of-distribution classifier with noise contrastive prior (OCD+NCP)** An uncertainty classifier model described in Appendix A. It is a deterministic neural network combined with NCP which we use as a baseline alternative to Bayes by Backprop with NCP.

For active learning, we select new data points that maximize the expected information gain. We use the approximation for Gaussian posterior predictive distributions described from MacKay (1992a). We place a softmax distribution on the information gain for all available data points and acquire labels by sampling with a temperature of  $\tau = 0.5$  to get diversity when selecting batches of labels,

$$\{x_{\text{new}}\} \sim p_{\text{new}}(x) = \frac{1}{Z} \exp\left(\frac{1}{2\tau} \ln(1 + \text{Var}[q(\mu(x))]/\sigma^2(x))\right) = \frac{1}{Z}(1 + \text{Var}[q(\mu(x))]/\sigma^2(x)). \quad (6)$$

Since we only have a weight belief for the output weights,  $\text{Var}[q(\mu(x))]$  is Gaussian and can be computed in closed form. Modeling uncertainty only for the predicted mean fits well with the information gain approximation that also only considers uncertainty around the mean. In the classifier model, we use the OOD probability  $p(o = 1|x)$  as a proxy. For the deterministic neural network, we use  $\text{Var}[p(y|x)]$  as proxy since it has no uncertainty estimate.

## 5.1 LOW-DIMENSIONAL ACTIVE LEARNING

For visualization purposes, we start with experiments on a 1-dimensional regression task that consists of a sine function with a small slope and increasing variance for higher inputs. Training data can be acquired within two bands, and the model is evaluated on all data points that are not visible to the model. This structured split

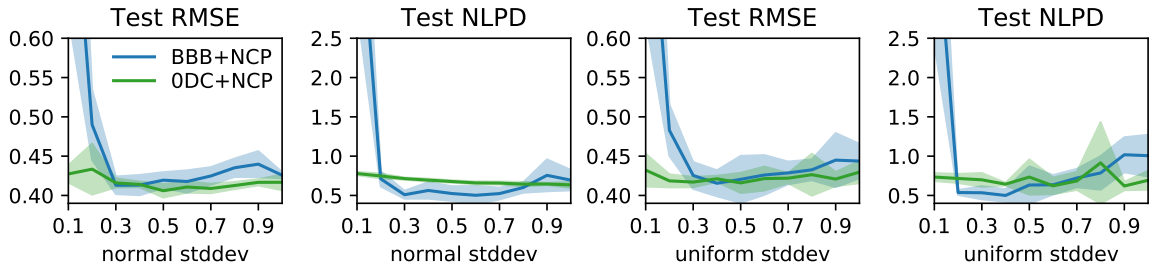


Figure 4: Robustness to different noise patterns. Plots show the final test performance on the low-dimensional active learning task explained in the paper (mean and stddev over 5 seeds). Lower is better. The baseline performances are RMSE: BBB ( $0.75 \pm 0.31$ ), Det ( $1.46 \pm 0.64$ ) and NLPD: BBB ( $10.29 \pm 8.05$ ), Det ( $1.3 \times 10^8 \pm 1.7 \times 10^8$ ). NCP works with both Gaussian and uniform input noise  $\epsilon$  and is robust to  $\sigma_x^2$ .

between training and testing data causes a strong distributional shift at test time, forcing successful models to have reliable uncertainty estimates to avoid mispredictions for OOD inputs.

We use two layers of 200 hidden units, a batch size of 10, and a learning rate of  $3 \times 10^{-4}$  for all models. NCP models use noise  $\epsilon \sim \text{Normal}(0, 0.5)$ . We start with 10 randomly selected initial target, and select 1 additional target every 1000 epochs. Figure 3 shows the root mean squared error (RMSE) and negative log predictive density (NLPD) throughout learning and Figure 1 visualizes the predictive distributions at the end of training. The two baseline models severely overfit to the training distribution early on when only few data points are visible. Models with NCP outperform BBB, which in turn outperforms Det. Figure 1 visualizes the models’ predictive distributions at the end of learning, showing that NCP prevents overconfident generalization.

## 5.2 ROBUSTNESS TO NOISE PATTERNS

The choice of input noise might seem like a critical hyperparameter for NCP. In this experiment, we find that our method is actually quite robust to the choice of input noise. We use the same setup as in Section 5.1 but with uniform or normal input noise with different variance ( $\sigma_x^2 \in \{0.1, 0.2, \dots, 1.0\}$ ). For uniform input noise, this means noise is drawn from the interval  $[-2\sigma_x, 2\sigma_x]$ . Figure 4 shows the final root mean squared error and negative log predictive density at the end of the experiment. Above a small threshold (BBB+NCP  $\sigma_x^2 \geq 0.3$  and ODC+NCP  $\sigma_x^2 \geq 0.1$ ), NCP consistently consistently performs well.

## 5.3 ACTIVE LEARNING ON FLIGHT DELAYS

We consider the flight delay data set (Hensman et al., 2013; Deisenroth and Ng, 2015; Lakshminarayanan et al., 2016), a large scale regression benchmark with several published results. The data set has 8 input variables describing a flight, and the target is the delay of the flight in minutes. There are 700K training examples and 100K test examples. The test set has a subtle distributional shift, since the 100K data points temporally follow after the training data.

We use two layers with 50 units each, a batch size of 10, and a learning rate of  $10^{-4}$  (hyperparameters were tuned across a grid for all models). For NCP models,  $\epsilon \sim \text{Normal}(0, 0.1)$ . Starting from 10 labels, the models select a batch of 10 additional labels every 50 epochs. The 700K data points of the training data set are available for acquisition, and we evaluate performance on the typical test split. Figure 5 shows the performance for the visible data points and the test set respectively. BBB overfits on this task, while the error continues decreasing for the two models with NCP. We note that BBB and BBB+NCP show similar error on the visible data points, but the NCP models generalize better to unseen data.

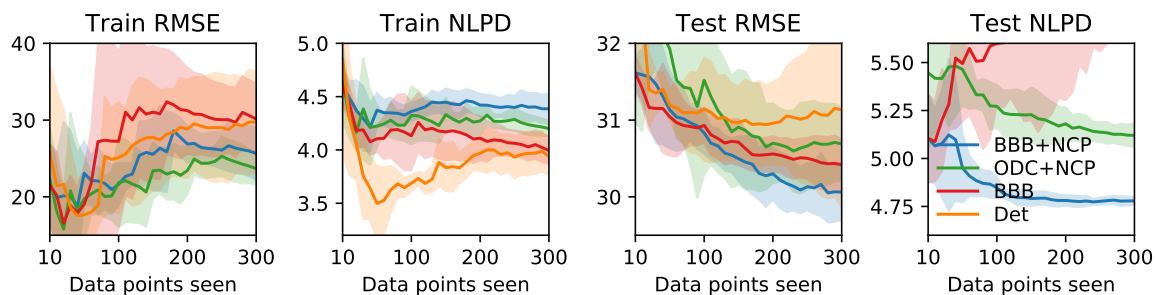


Figure 5: Active learning on the flights data set. The models trained with NCP achieve significantly lower negative log predictive density (NLPD) on the test set, and Bayes by Backprop with NCP achieves the lowest root mean squared error (RMSE). The test NLPD for the baseline models diverges as they overfit to the visible data points. Plots show mean and std over 10 runs.

#### 5.4 LARGE SCALE REGRESSION OF FLIGHT DELAYS

In addition to the active learning experiments, we perform a passive learning run on all 700K data points of the flights data set to explore the scalability of NCP. We use networks of 3 layers with 1000 units and a learning rate of  $10^{-4}$ . Table 1 compares the performance of our models to previously published results. We significantly improve state of the art performance on this data set.

## 6 DISCUSSION

We develop *noise contrastive priors* (NCPs), a prior for neural networks in data space. NCPs encourage network weights that not only explain the training data but also capture high uncertainty on OOD inputs. We show that NCPs offer strong improvements over baselines and scale to large regression tasks.

We focused on active learning for regression tasks, where uncertainty is crucial for determining which data points to select next. In future work it would be interesting to apply NCPs to alternative settings where uncertainty is important, such as image classification and learning with sparse or missing data. In addition, NCPs are only one form of a data prior, designed to encourage uncertainty on OOD inputs. Priors in data space can easily capture other properties such as periodicity or spatial invariance, and they may provide a scalable alternative to Gaussian process priors.

Table 1: Performance on all 700K data points of the flights data set. While uncertainty estimates are not necessary when a large data set that is similar to the test data set is available, it shows that our method scales easily to large data sets.

| Model                             | NLPD        | RMSE         |
|-----------------------------------|-------------|--------------|
| gPoE (Deisenroth & Ng 2015)       | 8.1         | —            |
| SAVIGP (Bonilla et al. 2016)      | 5.02        | —            |
| SVI GP (Hensman et al. 2013)      | —           | 32.60        |
| HGP (Ng & Deisenroth 2014)        | —           | 27.45        |
| MF (Lakshminarayanan et al. 2016) | 4.89        | 26.57        |
| BBB                               | <b>4.38</b> | <b>24.59</b> |
| BBB+NCP                           | <b>4.38</b> | <b>24.71</b> |
| ODC+NCP                           | <b>4.38</b> | <b>24.68</b> |

**Acknowledgements** Hidden for review.

## REFERENCES

- G. An. The effects of adding noise during backpropagation training on a generalization performance. *Neural computation*, 8(3):643–674, 1996.
- C. M. Bishop. Training with noise is equivalent to tikhonov regularization. *Neural computation*, 7(1):108–116, 1995.
- D. M. Blei. Build, compute, critique, repeat: Data analysis with latent variable models. *Annual Review of Statistics and Its Application*, 1:203–232, 2014.
- C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra. Weight uncertainty in neural networks. *arXiv preprint arXiv:1505.05424*, 2015.
- W. L. Buntine and A. S. Weigend. Bayesian back-propagation. *Complex Systems*, 5(6):603–643, 1991.
- R. Calandra, J. Peters, C. E. Rasmussen, and M. P. Deisenroth. Manifold gaussian processes for regression. *arXiv preprint arXiv:1402.5876*, 2014.
- S. Dasgupta. Analysis of a greedy active learning strategy. In *Neural Information Processing Systems (NIPS)*, 2004.
- M. P. Deisenroth and J. W. Ng. Distributed gaussian processes. *arXiv preprint arXiv:1502.02843*, 2015.
- J. Denker, D. Schwartz, B. Wittner, S. Solla, R. Howard, L. Jackel, and J. Hopfield. Large automatic learning, rule extraction, and generalization. *Complex Systems*, 1(5):877–922, 1987.
- S. Depeweg, J. M. Hernández-Lobato, F. Doshi-Velez, and S. Udluft. Uncertainty decomposition in bayesian neural networks with latent variables. In *International Conference on Machine Learning*, 2018.
- J. V. Dillon, I. Langmore, D. Tran, E. Brevdo, S. Vasudevan, D. Moore, B. Patton, A. Alemi, M. Hoffman, and R. A. Saurous. Tensorflow distributions. *arXiv preprint arXiv:1711.10604*, 2017.
- D. Flam-Shepherd, J. Requeima, and D. Duvenaud. Mapping gaussian process priors to bayesian neural networks. In *NIPS Workshop*, 2017.
- Y. Freund, H. S. Seung, E. Shamir, and N. Tishby. Selective sampling using the query by committee algorithm. *Machine learning*, 28(2-3):133–168, 1997.
- Y. Gal, R. Islam, and Z. Ghahramani. Deep bayesian active learning with image data. *arXiv preprint arXiv:1703.02910*, 2017.
- S. Ghosh and F. Doshi-Velez. Model selection in bayesian neural networks via horseshoe priors. *arXiv preprint arXiv:1705.10388*, 2017.
- M. Gutmann and A. Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 297–304, 2010a.
- M. Gutmann and A. Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 297–304, 2010b.
- D. Hendrycks and K. Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *International Conference on Learning Representations*, 2017.

- J. Hensman, N. Fusi, and N. D. Lawrence. Gaussian processes for big data. *arXiv preprint arXiv:1309.6835*, 2013.
- I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. *International Conference on Learning Representations*, 2016.
- N. Houlsby, F. Huszár, Z. Ghahramani, and M. Lengyel. Bayesian active learning for classification and preference learning. *arXiv preprint arXiv:1112.5745*, 2011.
- E. T. Jaynes. Information theory and statistical mechanics. *Physical review*, 106(4):620, 1957.
- A. J. Joshi, F. Porikli, and N. Papanikolopoulos. Multi-class active learning for image classification. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 2372–2379. IEEE, 2009.
- A. Kapoor, K. Grauman, R. Urtasun, and T. Darrell. Active learning with gaussian processes for object categorization. In *Computer Vision, 2007. ICCV 2007. IEEE 11th International Conference on*, pages 1–8. IEEE, 2007.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Neural Information Processing Systems*, pages 1097–1105, 2012.
- A. Kucukelbir, D. Tran, R. Ranganath, A. Gelman, and D. M. Blei. Automatic differentiation variational inference. *The Journal of Machine Learning Research*, 18(1):430–474, 2017.
- B. Lakshminarayanan, D. M. Roy, and Y. W. Teh. Mondrian forests for large-scale regression when uncertainty matters. In *Artificial Intelligence and Statistics*, pages 1478–1487, 2016.
- J. Lampinen and A. Vehtari. Bayesian approach for neural networks—review and case studies. *Neural networks*, 14(3):257–274, 2001.
- M. Lázaro-Gredilla and A. R. Figueiras-Vidal. Marginalized neural network mixtures for large-scale regression. *IEEE transactions on neural networks*, 21(8):1345–1351, 2010.
- K. Lee, H. Lee, K. Lee, and J. Shin. Training confidence-calibrated classifiers for detecting out-of-distribution samples. *arXiv preprint arXiv:1711.09325*, 2017.
- D. D. Lewis and W. A. Gale. A sequential algorithm for training text classifiers. In *Proceedings of the 17th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 3–12. Springer-Verlag New York, Inc., 1994.
- X. Li and Y. Guo. Adaptive active learning for image classification. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 859–866. IEEE, 2013.
- S. Liang, Y. Li, and R. Srikant. Principled detection of out-of-distribution examples in neural networks. In *International Conference on Learning Representations*, 2018.
- C. Louizos, K. Ullrich, and M. Welling. Bayesian compression for deep learning. In *Neural Information Processing Systems*, pages 3290–3300, 2017.
- A. L. Maas, A. Y. Hannun, and A. Y. Ng. Rectifier nonlinearities improve neural network acoustic models. In *International Conference on Machine Learning (ICML)*, 2013.

- D. J. MacKay. Information-based objective functions for active data selection. *Neural computation*, 4(4): 590–604, 1992a.
- D. J. MacKay. A practical bayesian framework for backpropagation networks. *Neural computation*, 4(3): 448–472, 1992b.
- K. Matsuoka. Noise injection into inputs in back-propagation learning. *IEEE Transactions on Systems, Man, and Cybernetics*, 22(3):436–440, 1992.
- A. K. McCallum and K. Nigam. Employing em and pool-based active learning for text classification. In *Proc. International Conference on Machine Learning (ICML)*, pages 359–367. Citeseer, 1998.
- T. Miyato, S. Maeda, M. Koyama, K. Nakae, and S. Ishii. Distributional smoothing with virtual adversarial training. *arXiv preprint arXiv:1507.00677*, 2015.
- A. Mnih and K. Kavukcuoglu. Learning word embeddings efficiently with noise-contrastive estimation. In *Advances in neural information processing systems*, pages 2265–2273, 2013.
- P. Müller and D. R. Insua. Issues in bayesian analysis of neural network models. *Neural Computation*, 10(3): 749–770, 1998.
- R. M. Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 2012.
- G. Pereyra, G. Tucker, J. Chorowski, Ł. Kaiser, and G. Hinton. Regularizing neural networks by penalizing confident output distributions. *arXiv preprint arXiv:1701.06548*, 2017.
- J. Quiñero-Candela and C. E. Rasmussen. A unifying view of sparse approximate gaussian process regression. *Journal of Machine Learning Research*, 6(Dec):1939–1959, 2005.
- H. A. Rowley, S. Baluja, and T. Kanade. Neural network-based face detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 1998.
- N. Roy and A. McCallum. Toward optimal active learning through monte carlo estimation of error reduction. *ICML, Williamstown*, pages 441–448, 2001.
- P. Schulam and S. Saria. A framework for individualizing predictions of disease trajectories by exploiting multi-resolution structure. In *Advances in Neural Information Processing Systems*, pages 748–756, 2015.
- B. Settles. Active learning literature survey. Computer Sciences Technical Report 1648, University of Wisconsin–Madison, 2009.
- H. S. Seung, M. Opper, and H. Sompolinsky. Query by committee. In *Proceedings of the fifth annual workshop on Computational learning theory*, pages 287–294. ACM, 1992.
- K.-K. Sung. Learning and example selection for object and pattern detection. *MIT A.I. Memo No. 1521*, 1994.
- I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *Neural Information Processing Systems*, 2014.
- C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2818–2826, 2016.
- S. Tong and D. Koller. Support vector machine active learning with applications to text classification. *Journal of machine learning research*, 2(Nov):45–66, 2001.

- D. Tran, A. Kucukelbir, A. B. Dieng, M. Rudolph, D. Liang, and D. M. Blei. Edward: A library for probabilistic modeling, inference, and criticism. *arXiv preprint arXiv:1610.09787*, 2016.
- A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
- P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pages 1096–1103. ACM, 2008.
- K. Wang, D. Zhang, Y. Li, R. Zhang, and L. Lin. Cost-effective active learning for deep image classification. *IEEE Transactions on Circuits and Systems for Video Technology*, 27(12):2591–2600, 2017.
- H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018.

## A OOD CLASSIFIER MODEL WITH NCP

In Section 4, we showed how to apply NCP to a Bayesian neural network model that captures function uncertainty in a belief over parameters. An alternative approach to capture uncertainty is to make explicit predictions about whether an input is OOD. Figure 2b shows such a mixture model via a binary variable  $o$ ,

$$o \sim \text{Bernoulli}(\pi(x, \theta))$$

$$y \sim \begin{cases} \text{Normal}(\mu(x, \theta), \sigma^2(x, \theta)) & \text{if } o = 0 \\ \text{Normal}(\mu_y, \sigma_y^2) & \text{if } o = 1, \end{cases} \quad (7)$$

where  $p(o = 1 | x)$  is the OOD probability of  $x$ . If  $o = 0$  (“in distribution”), the model outputs the neural network prediction. Otherwise, if  $o = 1$  (“out of distribution”), the model uses a fixed output prior. The neural network weights  $\theta$  are estimated using a point estimate, so we do not maintain a belief distribution over them.

The classifier prediction  $p(o | x, \theta)$  captures uncertainty in this model. We apply the NCP  $p(o | \tilde{x}, \theta) = \delta(o = 1 | \tilde{x}, \theta)$  to this variable, which assumes noised-up inputs to be OOD. During training on the data set,  $\{x, y\}$  and  $o = 0$  are observed, as training data are in-distribution by definition. Following Equation 2, the loss function is

$$\begin{aligned} \mathcal{L}(\theta) &= D_{\text{KL}}[p_{\text{train}}(y | x) \| p_{\text{model}}(y | x, o = 0, \theta)] + D_{\text{KL}}[p_{\text{prior}}(\tilde{o} | \tilde{x}) \| p_{\text{model}}(\tilde{o} | \tilde{x}, \theta)] \\ &= -\ln p(y, o = 0 | x, \theta) - \ln p(y, o = 1 | \tilde{x}, \theta) \\ &= -\ln \text{Normal}(y | \mu(x, \theta), \sigma^2(x, \theta)) - \ln \text{Bernoulli}(0 | \pi(x, \theta)) - \underbrace{\ln \text{Bernoulli}(1 | \pi(\tilde{x}, \theta))}_{\text{NCP loss}}. \end{aligned} \quad (8)$$

Analogously to the Bayesian neural network model in Section 4, we can either set  $\mu_y, \sigma_y^2$  manually or use the neural network prediction. In our experiments, we implement the OOD classifier model using a single neural network with two output layers that parameterize the Gaussian distribution and the binary distribution.

## B BNN WITH NCP USING REVERSE KL

In Section 4, we derived the Bayes by Backprop model with NCP by adding a forward KL-divergence from the mean prior to the model mean to the loss. An alternative derivation uses the fact that the KL-divergence is invariant to parameterization to replace the reverse KL-divergence in weight space by a KL-divergence in output space.

$$\begin{aligned} & \mathbb{E}_{p(x,y)} [\ln p(y | x)] \\ &= \mathbb{E}_{p(x,y)} \left[ \ln \int p(y | x, \theta) p(\theta) \frac{q(\theta)}{q(\theta)} d\theta \right] \\ &\geq \mathbb{E}_{p(x,y)} \left[ \int q(\theta) \ln p(y | x, \theta) \frac{p(\theta)}{q(\theta)} d\theta \right] \\ &= \mathbb{E}_{p(x,y)} [\mathbb{E}_{q(\theta)} [\ln p(y | x, \theta)] - D_{\text{KL}}[q(\theta) \| p(\theta)]] \\ &= \mathbb{E}_{p(x,y)} [\mathbb{E}_{q(\theta)} [\ln p(y | x, \theta)] - \mathbb{E}_{p(\tilde{x}|x)} [D_{\text{KL}}[q(\theta) \| p(\theta)]]] \\ &\approx \mathbb{E}_{p(x,y)} [\mathbb{E}_{q(\theta)} [\ln p(y | x, \theta)] - \mathbb{E}_{p(\tilde{x}|x)} [D_{\text{KL}}[q(\mu(\tilde{x})) \| p(\mu(\tilde{x}) | x)]]] \end{aligned} \quad (9)$$

where  $p(\mu(\tilde{x})) = \int \mu(\tilde{x}, \theta) p(\theta) d\theta$  and  $q(\mu(\tilde{x})) = \int \mu(\tilde{x}, \theta) q(\theta) d\theta$  are the distributions of the predicted mean induces by the weight beliefs. As a result, instead of specifying a prior in weight space, we can specify a prior in output space.

Above, we reparameterized the KL in weight space as a KL in output space; by the change of variables, this is equivalent if the mapping  $\mu(\cdot, \theta)$  is continuous and 1-1 with respect to  $\theta$ . This assumption does not hold for neural nets as multiple parameter vectors can lead to the same predictive distribution, thus the approximation above. A compact reparameterization of the neural network (equivalence class of parameters) would make this an equality.

## C OTHER RELATED WORK

**Active learning** Active learning is often employed in domains where data is cheap but labeling is expensive, and is motivated by the idea that not all data points are equally valuable when it comes to learning (Settles, 2009; Dasgupta, 2004). Active learning techniques can be coarsely grouped into three categories. Ensemble methods (Seung et al., 1992; McCallum and Nigam, 1998; Freund et al., 1997) generate queries that have the greatest disagreement between a set of classifiers. Error reduction approaches incorporate the select data based on the predicted reduction in classifier error based on information (MacKay, 1992a), Monte Carlo estimation (Roy and McCallum, 2001), or hard-negative example mining (Sung, 1994; Rowley et al., 1998).

Uncertainty-based techniques select samples for which the classifier is most uncertain. Approaches include maximum entropy (Joshi et al., 2009), distance from the decision boundary (Tong and Koller, 2001), pseudo labelling high confidence examples (Wang et al., 2017), and mixtures of information density and uncertainty measures (Li and Guo, 2013). Within this category, the area most related to our work are Bayesian methods. Kapoor et al. (2007) estimate expected improvement using a Gaussian process. Other approaches use classifier confidence (Lewis and Gale, 1994), predicted expected error (Roy and McCallum, 2001), or model disagreement (Houlsby et al., 2011). Recently, Gal et al. (2017) applied a convolutional neural network with dropout uncertainty to images.