

---

# Predictive Uncertainty Quantification with Compound Density Networks

---

Agustinus Kristiadi<sup>1</sup> Sina Däubener<sup>2</sup> Asja Fischer<sup>2</sup>

## Abstract

Despite the huge success of deep neural networks (NNs), finding good mechanisms for quantifying their prediction uncertainty is still an open problem. Bayesian neural networks are one of the most popular approaches to uncertainty quantification. On the other hand, it was recently shown that ensembles of NNs, which belong to the class of mixture models, can be used to quantify prediction uncertainty. In this paper, we build upon these two approaches. First, we increase the mixture model’s flexibility by replacing the fixed mixing weights by an adaptive, input-dependent distribution (specifying the probability of each component) represented by NNs, and by considering uncountably many mixture components. The resulting class of models can be seen as the continuous counterpart to mixture density networks and is therefore referred to as *compound density networks* (CDNs). We employ both maximum likelihood and variational Bayesian inference to train CDNs, and empirically show that they yield better uncertainty estimates on out-of-distribution data and are more robust to adversarial examples than the previous approaches.

## 1. Introduction

Deep neural networks (NNs) have achieved state-of-the-art performance in many application areas, such as computer vision (Krizhevsky et al., 2012) and natural language processing (Collobert et al., 2011). However, despite achieving impressive prediction accuracy on these supervised machine learning tasks, NNs do not provide good ways of quantifying predictive uncertainty. This is undesirable for many mission-critical applications, where taking wrong predictions with high confidence could have fatal consequences (e.g. in medical diagnostics or autonomous driving).

<sup>1</sup>Department of Computer Science, University of Tübingen, Germany <sup>2</sup>Department of Mathematics, Ruhr University Bochum, Germany. Correspondence to: Agustinus Kristiadi <agustinus.kristiadi@uni-tuebingen.de>.

A principled and the most explored way to quantify the uncertainty in NNs is through Bayesian inference. In the so-called Bayesian neural networks (BNNs) (Neal, 1995), the NN parameters are treated as random variables and the goal of learning is to infer the posterior probability distribution of the parameters given the training data. Since exact Bayesian inference in NNs is computationally intractable, different approximation techniques have been proposed (Neal, 1995; Blundell et al., 2015; Hernández-Lobato & Adams, 2015; Ritter et al., 2018, etc.). Given the (approximate) posterior, the final predictive distribution is obtained as the expected predictive distribution under the posterior. This expectation can be seen as an ensemble of an uncountably infinite number of predictors, where the prediction of each model is weighted by the posterior probability of the corresponding parameters.

Based on a Bayesian interpretation of dropout (Srivastava et al., 2014), Gal & Ghahramani (2016) proposed to apply it not only during training but also when making predictions to estimate predictive uncertainty. Interestingly, dropout has been also interpreted as ensemble model (Srivastava et al., 2014) where the predictions are averaged over the different NNs resulting from different dropout-masks. Inspired by this, Lakshminarayanan et al. (2017) proposed to use a simple NN ensemble to quantify the prediction uncertainty, i.e. to train a set of independent NNs using a proper scoring rule and defining the final prediction as the arithmetic mean of the outputs of the individual models, which corresponds to defining a uniformly-weighted mixture model. It is argued, that the model is able to encode two sources of uncertainty by calibrating the target uncertainty (i.e. uncertainty in target  $y$  given input  $x$ ) in each component and capturing the model uncertainty by averaging over the components.

In this paper, we therefore aim at further investigating the potential that lies in employing mixture distributions for uncertainty quantification. The flexibility of the mixture model can be increased by learning input-conditioned mixture weights like it is done by mixture density networks (MDNs) (Bishop, 1994). Furthermore, one can consider uncountably many component distributions instead of a finite set, which turns the mixture distribution into a compound distribution. We combine both by deriving the continuous counterpart of MDNs, which we call *compound density*

networks (CDNs). These networks can be trained by likelihood maximization. Moreover, variational Bayes can be employed to infer the posterior distribution over the CDN parameters, leading to a combination of the mixture model and the Bayesian approach to uncertainty modeling. We experimentally show that CDNs allow for better uncertainty quantification and are more robust to adversarial examples than previous approaches.

This paper is organized as follows. In Section 2 we give a brief introduction to MDNs. We then formally define CDNs in Section 3. We review related work in Section 4 and present a detailed experimental analysis in Section 5. Finally, we conclude our paper in Section 6.

## 2. Mixture Density Networks

Let  $\mathcal{D} = \{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N$  be an i.i.d dataset. Let us define the following conditional mixture model

$$p(\mathbf{y}|\mathbf{x}) = \sum_{k=1}^K p(\mathbf{y}; \phi_k(\mathbf{x})) p(\phi_k(\mathbf{x}); \boldsymbol{\pi}(\mathbf{x})), \quad (1)$$

and an NN that maps  $\mathbf{x}$  onto both the parameters  $\boldsymbol{\pi}(\mathbf{x})$  of the mixing distribution and the parameters  $\{\phi_k(\mathbf{x})\}_{k=1}^K$  of the  $K$  mixture components. The complete system is called mixture density network (MDN) and was proposed by Bishop (1994). That is, an MDN is an NN parametrizing a conditional mixture distribution, where both the mixture components and the mixture coefficients depend on input  $\mathbf{x}$ .<sup>1</sup> MDNs can be trained by maximizing the log-likelihood of the parameters of the NN given the training set  $\mathcal{D}$  using gradient-based optimizers such as stochastic gradient descent (SGD) and its variants.

MDNs belong to a broader class of models called mixture of experts (MoE) (Jacobs et al., 1991) which differ from standard mixture models by assuming that the mixture coefficients depend on the input.<sup>2</sup> Because of its formulation as a mixture distribution, the predictive distribution of an MDN can handle multimodality better than a standard discriminative neural network.

## 3. Compound density networks

We aim at generalizing the MDN from a finite mixture distribution to a mixture of an uncountable set of components. The continuous counterpart of a conditional mixture distribution in eq. (1) is given by the conditional *compound*

<sup>1</sup>For instance, as in the original publication, the mixture components could be  $K$  Gaussians, with  $\phi_k(\mathbf{x})$  being input-specific means and variances, and the mixture probabilities could be given by (applying the softmax function to the unnormalized) mixing weights  $\boldsymbol{\pi}(\mathbf{x})$ , both computed by one NN.

<sup>2</sup>See Bishop (2006, ch. 5.6 and ch. 14.5.3) and Murphy (2012, ch. 11.2.4) for a detailed discussion of MDNs.

*probability distribution*

$$p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}; \phi(\mathbf{x})) p(\phi(\mathbf{x}); \boldsymbol{\pi}(\mathbf{x})) d\phi(\mathbf{x}), \quad (2)$$

where  $\phi(\mathbf{x})$  turns from a discrete into a continuous random variable.

We now want to follow the approach of MDNs to model the parameters of the components and the mixing distribution by NNs. To handle the continuous state space of  $\phi(\mathbf{x})$  in the case of a compound distribution, the key idea is now to let  $\phi(\mathbf{x})$  be given by a stochastic NN  $f(\mathbf{x}; \boldsymbol{\theta}) = \phi(\mathbf{x})$  with stochastic parameters  $\boldsymbol{\theta}$ . Since given  $\mathbf{x}$ ,  $f$  is a deterministic map from  $\boldsymbol{\theta}$  to  $\phi(\mathbf{x})$ , it is possible to replace the mixing distribution  $p(\phi(\mathbf{x}); \boldsymbol{\pi}(\mathbf{x})) = p(f(\mathbf{x}; \boldsymbol{\theta}); \boldsymbol{\pi}(\mathbf{x}))$  by a distribution  $p(\boldsymbol{\theta}; \boldsymbol{\pi}(\mathbf{x}))$  over  $\boldsymbol{\theta}$ . We further assume, that the parameters  $\boldsymbol{\pi}(\mathbf{x})$  of the mixing distribution are given by some parametrized function  $g(\mathbf{x}; \boldsymbol{\psi}) = \boldsymbol{\pi}(\mathbf{x})$  which can also be modeled based on NNs. In correspondence to MDNs, the complete system is called *compound density network* (CDN) and it is summarized by the following equation

$$\begin{aligned} p(\mathbf{y}|\mathbf{x}; \boldsymbol{\psi}) &:= \int p(\mathbf{y}; f(\mathbf{x}; \boldsymbol{\theta})) p(\boldsymbol{\theta}; g(\mathbf{x}; \boldsymbol{\psi})) d\boldsymbol{\theta} \\ &= \mathbb{E}_{p(\boldsymbol{\theta}; g(\mathbf{x}; \boldsymbol{\psi}))} [p(\mathbf{y}; f(\mathbf{x}; \boldsymbol{\theta}))]. \end{aligned} \quad (3)$$

As MDNs, CDNs can be trained by maximizing the log-likelihood of  $\boldsymbol{\psi}$  given the data set  $\mathcal{D}$ . Moreover, one can add a regularization term encouraging the mixing distribution to stay close to some distribution  $p(\boldsymbol{\theta})$ , which leads to the objective

$$\begin{aligned} \mathcal{L}_{\text{ML}}(\boldsymbol{\psi}) &:= \sum_{n=1}^N \log \mathbb{E}_{p(\boldsymbol{\theta}; g(\mathbf{x}_n; \boldsymbol{\psi}))} [p(\mathbf{y}_n; f(\mathbf{x}_n; \boldsymbol{\theta}))] \\ &\quad - \lambda \sum_{n=1}^N D_{\text{KL}}[p(\boldsymbol{\theta}; g(\mathbf{x}_n; \boldsymbol{\psi})) || p(\boldsymbol{\theta})], \end{aligned} \quad (4)$$

where  $\lambda$  is a hyperparameter controlling the strength of the regularization.<sup>3</sup>

Alternatively, we can turn CDNs into Bayesian models by defining a prior  $p(\boldsymbol{\psi})$  over  $\boldsymbol{\psi}$  and employing variational Bayes (VB) (Hinton & Van Camp, 1993) to infer an approximate posterior  $q(\boldsymbol{\psi}; \boldsymbol{\omega}) \approx p(\boldsymbol{\psi}|\mathcal{D})$  by maximizing the evidence lower bound (ELBO):

$$\begin{aligned} \mathcal{L}_{\text{VB}}(\boldsymbol{\omega}) &:= \sum_{n=1}^N \mathbb{E}_{q(\boldsymbol{\psi}; \boldsymbol{\omega})} [\log \mathbb{E}_{p(\boldsymbol{\theta}; g(\mathbf{x}_n; \boldsymbol{\psi}))} [p(\mathbf{y}_n; f(\mathbf{x}_n; \boldsymbol{\theta}))]] \\ &\quad - D_{\text{KL}}[q(\boldsymbol{\psi}; \boldsymbol{\omega}) || p(\boldsymbol{\psi})]. \end{aligned} \quad (5)$$

<sup>3</sup>Note, that the objective gets equivalent to the one proposed by Alemi et al. (2017) when it is approximated based on a single sample of  $\boldsymbol{\theta}$ . We experimentally show that this objective leads to better results than theirs (when approximated with more samples) in the supplement.

**Algorithm 1** The training procedure of CDNs with  $\mathcal{L}_{\text{ML}}$ .

**Require:**

Mini-batch size  $M$ , number of samples  $S$  of  $\theta$ , regularization strength  $\lambda$ , and learning rate  $\alpha$ .

- 1: **while** the stopping criterion is not satisfied **do**
- 2:    $\{\mathbf{x}_m, \mathbf{y}_m\}_{m=1}^M \sim \mathcal{D}$
- 3:   **for**  $m = 1, \dots, M; s = 1, \dots, S$  **do**
- 4:      $\theta_{ms} \sim p(\theta; g(\mathbf{x}_m; \psi))$
- 5:      $\phi_s(\mathbf{x}_m) = f(\mathbf{x}_m; \theta_{ms})$
- 6:   **end for**
- 7:    $\mathcal{L}(\psi) = \sum_{m=1}^M \log \frac{1}{S} \sum_{s=1}^S p(\mathbf{y}_m; \phi_s(\mathbf{x}_m)) - \lambda \sum_{m=1}^M D_{\text{KL}}[p(\theta; g(\mathbf{x}_m; \psi)) || p(\theta)]$
- 8:    $\psi \leftarrow \psi + \alpha \nabla \mathcal{L}(\psi)$
- 9: **end while**

Given the approximate posterior the predictive distribution of the Bayesian CDN is defined as

$$p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}|\mathbf{x}; \psi) q(\psi; \omega) d\psi, \quad (6)$$

which is approximated based on samples of  $\psi$  and  $\theta$ .

We present pseudocode for training CDNs with  $\mathcal{L}_{\text{ML}}$  in Algorithm 1.<sup>4</sup> Note that CDNs correspond to an abstract framework for modeling compound distributions with NNs, i.e. we still need to concretely define the stochastic NN  $f(\mathbf{x}; \theta)$  and choose the statistical models for the mixture components and the mixing distribution.

### 3.1. Probabilistic hypernetworks

Ha et al. (2017) proposed to (deterministically) generate the parameters of an NN by another NN, which they call the *hypernetwork*<sup>5</sup>. We would like to follow this approach for modeling a CDN, that is, we aim at modeling the mixing distribution  $p(\theta; g(\mathbf{x}; \psi))$  over network parameters by NNs. Since now the hypernetworks map  $\mathbf{x}$  to a distribution over parameters instead of a specific value  $\theta$ , we refer to them as *probabilistic hypernetworks*. In the following, we will describe this idea in more detail.

Let  $f$  be a multi-layer perceptron (MLP) with  $L$ -layers, parametrized by a set of layers' weight matrices<sup>6</sup>  $\theta = \{\mathbf{W}_l\}_{l=1}^L$ , that computes the parameters  $\phi(\mathbf{x}) = f(\mathbf{x}; \theta)$  of the CDNs component distribution in eq. (3). Let  $\mathbf{h}_1, \dots, \mathbf{h}_{L-1}$  denote the states of the hidden layers, and let us define  $\mathbf{h}_0 = \mathbf{x}$ , and  $\mathbf{h}_L = f(\mathbf{x}; \theta)$ . We now assume

<sup>4</sup>Pseudocode for training CDNs with  $\mathcal{L}_{\text{VB}}$  is presented in the supplement.

<sup>5</sup>Specifically, they propose to apply a hypernetwork to compute the weight matrix of a recurrent NN at each time-step, given the current input and the previous hidden state.

<sup>6</sup>We assume that the bias parameters are absorbed into the weight matrix.

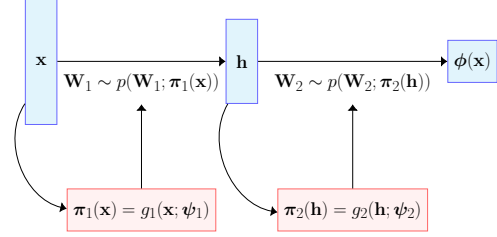


Figure 1: An example of a probabilistic hypernetwork applied to a two-layer MLP.

the weight matrices  $\{\mathbf{W}_l\}_{l=1}^L$  to be random variables and to be independent of each other given the state of the previous hidden layer. We define a series of probabilistic hypernetworks  $g = \{g_l\}_{l=1}^L$  (parametrized by  $\psi = \{\psi_l\}_{l=1}^L$ ), where  $g_l$  maps  $\mathbf{h}_{l-1}$  to the parameters of the distribution of  $\mathbf{W}_l$ , and let the joint distribution over  $\theta$  be given by

$$p(\theta; g(\mathbf{x}; \psi)) := \prod_{l=1}^L p(\mathbf{W}_l; g_l(\mathbf{h}_{l-1}; \psi_l)). \quad (7)$$

An illustration of a stochastic two-layer network  $f(\mathbf{x}; \theta)$  computing  $\phi(\mathbf{x})$ , where the distribution of the parameters is given by probabilistic hypernetworks as defined in eq. (7), is given in Figure 1.

### 3.2. Probabilistic hypernetworks with matrix variate normals

A statistical model that was recently applied as the posterior over weight matrices in BNNs (Louizos & Welling, 2016; Sun et al., 2017; Zhang et al., 2018; Ritter et al., 2018) is the matrix variate normal (MVN) distribution (Gupta & Nagar, 1999). An MVN is parametrized by three matrices: a mean matrix  $\mathbf{M}$  and two covariance factor matrices  $\mathbf{A}$  and  $\mathbf{B}$ . It is connected to the multivariate Gaussian by the following equivalence

$$\begin{aligned} \mathbf{X} &\sim \mathcal{MN}(\mathbf{X}; \mathbf{M}, \mathbf{A}, \mathbf{B}) \\ \iff \text{vec}(\mathbf{X}) &\sim \mathcal{N}(\text{vec}(\mathbf{X}); \text{vec}(\mathbf{M}), \mathbf{B} \otimes \mathbf{A}), \end{aligned} \quad (8)$$

where  $\text{vec}(\mathbf{X})$  denotes the vectorization of matrix  $\mathbf{X}$ . Due to the Kronecker factorization of the covariance, an MVN requires fewer parameters compared to a multivariate Gaussian, which motivates us to use it as the distribution over weight matrices in this work. Furthermore, we assume that the covariance factor matrices are diagonal matrices, following Louizos & Welling (2016). That is, we choose the mixture distribution of the CDN to be

$$\begin{aligned} p(\theta; g(\mathbf{x}; \psi)) &= \prod_{l=1}^L \mathcal{MN}(\mathbf{W}_l; g_l(\mathbf{h}_{l-1}; \psi_l)) \\ &= \prod_{l=1}^L \mathcal{MN}(\mathbf{W}_l; \mathbf{M}_l, \text{diag}(\mathbf{a}_l), \text{diag}(\mathbf{b}_l)), \end{aligned} \quad (9)$$

where  $g_l$  maps the state  $\mathbf{h}_{l-1}$  of the previous hidden layer onto the  $l$ -th MVN’s parameters  $\{\mathbf{M}_l, \mathbf{a}_l, \mathbf{b}_l\}$  defining the distribution over  $\mathbf{W}_l$ . Suppose  $\mathbf{W}_l \in \mathbb{R}^{r \times c}$ , then the corresponding MVN distribution has  $rc + r + c$  parameters, which is more efficient compared to  $rc + rc$  parameters when modeling  $\mathbf{W}_l$  as fully-factorized Gaussian random variable. To further reduce the model complexity we use a vector-scaling parametrization similar to the one used by Ha et al. (2017) and Krueger et al. (2017) for the mean matrices  $\{\mathbf{M}_l\}_{l=1}^L$ . We detail this parametrization in the supplement.

For the regularization during ML training (eq. (4)), we define the prior over  $\theta$  to be  $p(\theta) := \prod_{l=1}^L \mathcal{MN}(\mathbf{W}_l; \mathbf{0}, \mathbf{I}, \mathbf{I})$ . Meanwhile, to perform variational Bayes (eq. (5)), the prior over  $\psi$  is assumed to be the product of standard MVNs (similar to  $p(\theta)$ ), while the variational posterior is defined as product of diagonal MVNs with variational parameters  $\omega$ , following Louizos & Welling (2016). That is, for CDNs trained with VB, we assume each probabilistic hypernetwork  $g_l$  to be a variational matrix Gaussian (VMG), the BNN proposed by Louizos & Welling (2016).

Note that using the mixing distribution and approximate posterior as defined above allows us to apply the reparametrization trick (Louizos & Welling, 2016). Furthermore, the KL-divergence term in eq. (5) can be computed in closed form, as also noted by Louizos & Welling (2016).

#### 4. Related work

Various approaches for quantifying predictive uncertainty in NNs have been proposed. Applying Bayesian inference to NNs, i.e. treating the network parameters as random variables and estimating the posterior distribution given the training data based on Bayes’ theorem, results in BNNs (MacKay, 1992; Neal, 1995; Graves, 2011; Blundell et al., 2015; Louizos & Welling, 2016; Sun et al., 2017; Louizos & Welling, 2017; Ritter et al., 2018; Zhang et al., 2018, etc). Since the true posterior distribution is intractable, BNNs are trained based on approximate inference methods such as variational inference (VI) (Peterson, 1987; Hinton & Van Camp, 1993; Graves, 2011; Blundell et al., 2015), Markov Chain Monte Carlo (Neal, 1995), or Laplace approximation (MacKay, 1992; Ritter et al., 2018). The final prediction is then given by the expectation of the network prediction (given the parameters) w.r.t. the approximate posterior distribution. In VI, many modeling choices for BNNs’ approximate posterior have been proposed. Louizos & Welling (2016) proposed to train BNNs with an MVN as the approximate posterior of each weight matrix (leading to a model they refer to as VMG). Multiplicative normalizing flow (MNF) (Louizos & Welling, 2017) models the approximate posterior as a compound distribution, where the mixing density is given by a normalizing flow. Zhang et al. (2018) also use an MVN approximate posterior and apply approx-

imate natural gradient (Amari, 1998) based maximization on the VI objective, which results in an algorithm called noisy K-FAC. Meanwhile, the Kronecker-factored Laplace approximation (KFLA) (Ritter et al., 2018) extends the classical Laplace approximation by using an MVN approximate posterior with tractable and efficiently computed covariance factors, based on the Fisher information matrix.

Models with similar functional form as CDNs, i.e.  $p(\mathbf{y}|\mathbf{x}; \psi) = \int p(\mathbf{y}|\mathbf{x}; \theta)p(\theta|\mathbf{x}; \psi) d\theta$ , have been previously studied in various settings. The deep variational information bottleneck (VIB) (Alemlı et al., 2017) assumes  $\theta$  to be the hidden units of a certain layer instead of the parameters of an NN and trains the model with an objective derived from the information bottleneck method (Tishby et al., 2001). Interestingly, this objective gets equivalent to the ML objective for CDNs when approximated based on a single sample of  $\theta$ . Malinin & Gales (2018) proposed Prior Networks (PNs) which can be described by  $p(\mathbf{y}|\mathbf{x}) = \int p(\mathbf{y}; \theta)p(\theta|\mathbf{x}; \psi) d\theta$  and aim at modeling data uncertainty with the component distribution and what they call “distributional uncertainty” (i.e. uncertainty due to mismatch between the distributions of test and training data) with the mixture distribution. Specifically, they propose Dirichlet Prior Networks (DPNs) for uncertainty quantification in classification tasks, where  $p(\theta|\mathbf{x}; \psi)$  is assumed to be Dirichlet distribution. In contrast to CDNs, DPNs use only a single NN to parametrize the model and an objective that augments likelihood maximization/KL-divergence minimization by a term explicitly making use of out-of-distribution samples. Depeweg et al. (2017; 2018) investigated BNNs with latent variables  $\theta$  (referred to as BNN+LV), which can be described by  $p(\mathbf{y}|\mathbf{x}) = \iint p(\mathbf{y}|\mathbf{x}, \theta; \psi)p(\theta|\mathbf{x})p(\psi) d\theta d\psi$ . This model is similar to VB-CDNs, but in contrast assumes that  $\theta$  is independent of  $\psi$ , which serves as additional input to the component distribution instead. Furthermore, BNN+LV employ an  $\alpha$ -divergence-based objective, instead of the ELBO.

There have been several concurrent works (Krueger et al., 2017; Louizos & Welling, 2017; Pawłowski et al., 2017; Sheikh et al., 2017) applying hypernetworks (Jia et al., 2016; Ha et al., 2017) to model the posterior distribution over network parameters in BNNs. Krueger et al. (2017) and Louizos & Welling (2017) use normalizing flows, while Pawłowski et al. (2017) and Sheikh et al. (2017) use arbitrary NNs as their hypernetworks. Note, that in Bayesian CDNs the hypernetworks themselves become BNNs.

Gal & Ghahramani (2016) developed a theoretical framework that relates dropout training in NNs to approximate Bayesian inference and, as a result, proposed to approximate the predictive distribution by an average over the different networks resulting from independently sampled dropout-

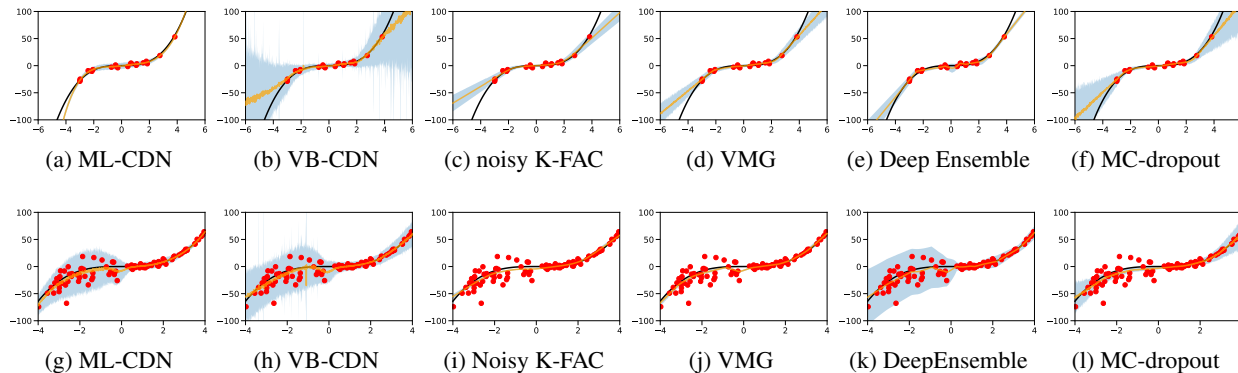


Figure 2: Comparison of the predictive distributions given by the CDNs and the baselines on toy datasets with homoscedastic noise and few samples (first row) and heteroscedastic noise and more samples (second row). Black lines correspond to the true noiseless function, red dots correspond to samples, orange lines and shaded regions correspond to the empirical mean and the  $\pm 3$  standard deviation of the predictive distribution, respectively. The VB-CDN is the only model able to capture both the epistemic uncertainty on the first and the aleatoric on the second dataset.

masks, a technique which they referred to as MC-dropout (MCD) and which they applied to estimate the prediction uncertainty in NNs. Recently, Lakshminarayanan et al. (2017), proposed to use an ensemble of NNs in conjunction with a proper scoring rule and adversarial training to quantify the prediction uncertainty of deep NNs, leading to a model referred to as Deep Ensemble (DE). The DE provides a non-Bayesian way to quantify prediction uncertainty, and is in this sense related to the approaches of Guo et al. (2017) and Hendrycks & Gimpel (2017).

## 5. Experiments

We consider several standard tasks in our experimental analysis: 1D toy regression problems inspired by Hernández-Lobato & Adams (2015) (Section 5.1), classification under out-of-distribution data (Section 5.2), and detection of and defense against adversarial examples (Szegedy et al., 2014) (Section 5.3). We refer to the CDNs that are trained via  $\mathcal{L}_{\text{ML}}$  (eq. (4)) and  $\mathcal{L}_{\text{VB}}$  (eq. (5)) as ML-CDNs and VB-CDNs, respectively. The following recent models (described in Section 4) are considered as the baselines: VMG, MNF, DPN, noisy K-FAC, MC-dropout, and Deep Ensemble.<sup>7</sup>

We estimate the predictive distribution  $p(\mathbf{y}|\mathbf{x})$  of the CDNs, based on 100 joint samples of  $\boldsymbol{\psi} \sim q(\boldsymbol{\psi}; \boldsymbol{\omega})$ ,  $\boldsymbol{\theta} \sim p(\boldsymbol{\theta}; g(\mathbf{x}; \boldsymbol{\psi}))$  for VB-CDNs and 100 samples of  $\boldsymbol{\theta} \sim p(\boldsymbol{\theta}; g(\mathbf{x}; \boldsymbol{\psi}))$  for ML-CDNs. We also draw 100 samples from the posterior to approximate the predictive distribution of BNN baselines. If not stated otherwise, we use a single sample to perform Monte Carlo integration during training. We pick the regularization hyperparameter  $\lambda$  for ML-CDNs (eq. (4)) out of the set

<sup>7</sup>Additional baselines are investigated in the supplement.

$\{10^{-4}, 10^{-5}, 10^{-6}, 10^{-7}, 10^{-8}\}$  which maximizes the validation accuracy. We use Adam (Kingma & Ba, 2015) with default hyperparameters for optimization in all experiments and the implementations provided by Louizos & Welling (2017)<sup>8</sup> and Zhang et al. (2018)<sup>9</sup> for MNF and noisy K-FAC, respectively. Where mini-batching is necessary, e.g. on MNIST and CIFAR-10, we use mini-batches of size 200. All models are optimized over 10000 iterations in the toy regression experiments, 20000 iterations ( $\approx 67$  epochs) in experiments on MNIST and Fashion-MNIST, and 100 epochs in experiments on CIFAR-10. We chose ReLU and hyperbolic tangent as the nonlinearity of the ML-CDNs’ and VB-CDNs’ hypernetworks, respectively. For the details on the selection of the model-specific hyperparameters of the baselines, we refer the reader to the supplementary material. The source code for all our experiments is available at <https://anonymous.com>.

### 5.1. Toy regression

Following Hernández-Lobato & Adams (2015), we generate the first toy regression dataset as follows: We sample 20 input points  $x \sim \mathcal{U}[-4, 4]$  and their target values  $y = x^3 + \epsilon$ , where  $\epsilon \sim \mathcal{N}(0, 3^2)$ , i.e. the data noise is homoscedastic. We aim at analyzing how well the target function is modeled over the larger interval  $[-6, 6]$ . Having only few data points, it is a desirable property of a model to express high (epistemic) uncertainty in regions with no or only few samples, e.g. between  $-6$  and  $-4$  or  $4$  and  $6$ . The second toy regression dataset is constructed by sampling 100 data points as above, this time with different scale of noise in different

<sup>8</sup>[https://github.com/AMLab-Amsterdam/MNF\\_VBNN](https://github.com/AMLab-Amsterdam/MNF_VBNN)

<sup>9</sup><https://github.com/gd-zhang/noisy-K-FAC>

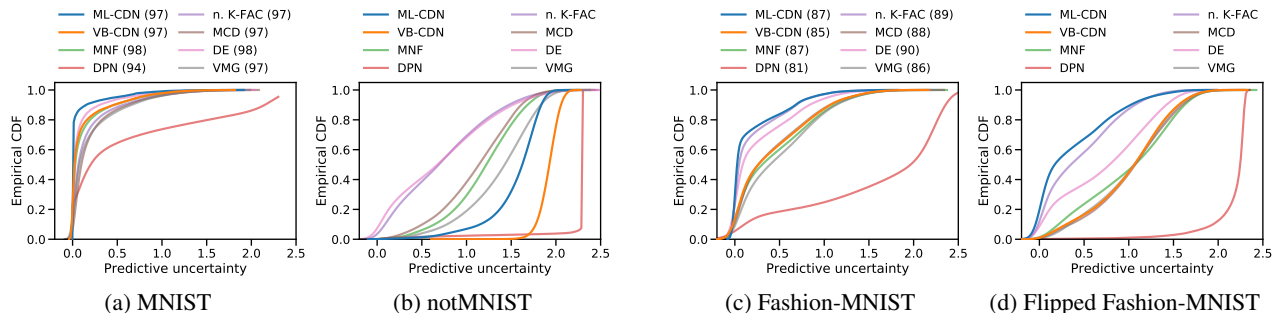


Figure 3: CDFs of the empirical entropy of the predictive distribution of the models, trained on MNIST (the first two figures) and Fashion-MNIST (the last two figures). The caption of each figure indicates the test set used, the y-axis denotes the fraction of predictions having entropy less than the corresponding value on the x-axis, and the number next to each model name indicates its test accuracy, in percent.

intervals:  $\epsilon \sim \mathcal{N}(0, 3^2)$ , if  $x \geq 0$ . and  $\epsilon \sim \mathcal{N}(0, 15^2)$ , otherwise. This dataset is designated for testing whether a model can capture heteroscedastic aleatoric uncertainty.

In these experiments, we use a two-layer MLP with 100 hidden units as the predictive network, while the hypernetworks of the CDNs ( $g_1$  and  $g_2$ ) are modeled with two-layer MLPs with 10 hidden units each. Three samples of  $\theta$  (along with a single sample of  $\psi$  in the case of the VB-CDN) are used to approximate the objectives during training of both CDNs and BNNs.<sup>10</sup> A regularization hyperparameter of  $\lambda = 10^{-3}$  is used for training the ML-CDNs.

The results for the first data set (shown in the first row of Figure 2) demonstrate that the VB-CDN is capable of capturing the epistemic uncertainty like other Bayesian models. This is not the case for the ML-CDN (which displays high confidence everywhere) and the DE (which captures only the uncertainty on the left side). This demonstrates the benefits of using a Bayesian approach for capturing parameter uncertainty. On the other hand, the mixture models, i.e. the CDNs and the DE, are the only ones able to capture the aleatoric uncertainty on the second dataset, as shown in the second row of Figure 2. This can be explained by the ability of CDNs and DEs to model input-dependent variance.

To further investigate the different roles in uncertainty modeling of the mixing distribution and the approximate posterior of VB-CDNs, we compare their average variance (over parameters and samples).<sup>11</sup> On the first data set, the aver-

<sup>10</sup>On these toy datasets, we found that using more than one sample is crucial for the results of the CDNs (i.e. results for using just one sample look similar to that of the VMG and Noisy K-FAC as can be seen in the supplement), while it does not significantly change the behaviour of the BNNs.

<sup>11</sup>We picked 1000 evenly spaced points from  $[-6, 6]$  and  $[-4, 4]$  respectively and approximated the means over the posterior with 100 samples.

age variance of the mixing distribution is 0.356 and that of the posterior distribution is 0.916. On the second data set the average variance of the posterior distribution is 0.429 and that of the mixing distribution is 0.618 for  $x < 0$  and 0.031 for  $x \geq 0$ . Therefore, the variance of the posterior is reduced on the second data set (as desired for more training data) while the mixing distribution successfully captures the higher data uncertainty for  $x < 0$ , indicating that the approximate posterior successfully models epistemic and the mixing distribution aleatoric uncertainty.

## 5.2. Out-of-distribution classification

Following Lakshminarayanan et al. (2017), we train all models on the MNIST training set<sup>12</sup> and investigate their performance on the MNIST test set and the notMNIST dataset<sup>13</sup>, which contains images (of the same size and format as MNIST) of letters from the alphabet instead of handwritten digits. On such an out-of-distribution (OOD) test set, the predictive distribution of an ideal model should have maximum entropy, i.e. it should have a value of  $\ln 10 \approx 2.303$  which would be achieved if all ten classes are equally probable. The predictive NN used for this experiment is an MLP with a 784-100-10 architecture.

We present the results in Figure 3a and b, where we plotted the cumulative distribution function (CDF) of the empirical entropy of the predictive distribution, following Louizos & Welling (2017). A CDF curve close to the top-left corner of the figure implies that the model yields mostly low entropy predictions, indicating that the model is very confident. While one wishes to observe high confidence on data points similar to those seen during training, the model should express uncertainty when exposed to OOD data. That is, we

<sup>12</sup>We use Fashion-MNIST as OOD data for training the DPN.

<sup>13</sup><http://yaroslavvb.blogspot.com/2011/09/notmnist-dataset.html>.

prefer a model to have a CDF curve closer to the bottom-right corner on notMNIST, as this implies it makes mostly uncertain (high entropy) predictions, and a curve closer to the upper-left corner for MNIST. As the results show, the VB-CDN yields high confidence on the test set of MNIST while having significantly lower confidence on notMNIST compared to all baseline models, except the DPN. Note however, that training DPNs requires additional data (which makes the comparison unfair) and that the DPN’s prediction accuracy and confidence on the MNIST test set are low compared to all other models. For the ML-CDN, we observe that it is more confident than all other models on within-distribution data, at the expense of showing lower uncertainty on OOD data than the VB-CDN.

Quantitatively, we calculated the mean maximal confidence for in-distribution (MMC-in) and OOD data (MMC-out) as well as the area under receiver operating characteristic (AUROC). The results can be found in Table 1.

Table 1: Mean maximal confidence (MMC) for in distribution (MNIST) and OOD data (notMNIST) and area under receiver operating characteristic (AUROC).

Algorithm	MMC-in	MMC-out	AUROC
CDN	<b>0.978</b>	<b>0.430</b>	<b>0.993</b>
VMG	0.938	0.507	0.964
MNF	0.959	0.504	0.977
MCD	0.950	0.665	0.928
DE	0.970	0.740	0.862
noisy-KFAC	0.949	0.744	0.848

Our model clearly has the highest MMC value for in-distribution data and the highest AUROC, while having the lowest MMC value for OOD data.

On the more challenging OOD task introduced by Alemi et al. (2018) where Fashion-MNIST (Xiao et al., 2017) is used as training set<sup>14</sup>, while the vanilla and the up-down flipped test set of Fashion-MNIST are used for evaluation (Figure 3c and d), the results are less pronounced, but the CDNs still show a performance competitive to that of the baseline models.

### 5.3. Adversarial examples

To investigate the robustness and detection performance of CDNs w.r.t. adversarial examples (Szegedy et al., 2014), we apply the Fast Gradient Sign Method (FGSM) (Goodfellow et al., 2015) to a 10% fraction (i.e. 1000 samples) of the MNIST, Fashion-MNIST (Xiao et al., 2017), and CIFAR-10

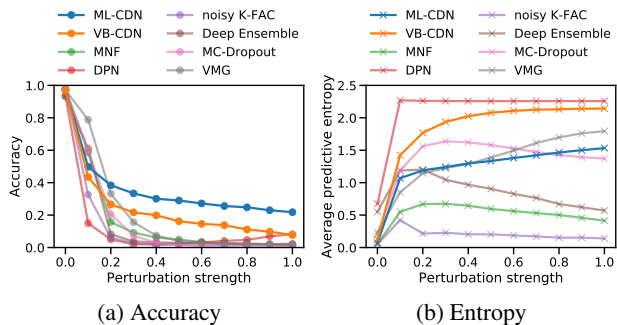


Figure 4: Prediction accuracy and average entropy of models trained on MNIST when attacked by FGSM-based adversarial examples (Goodfellow et al., 2015) with varying perturbation strength.

test set.<sup>15</sup> We do so, by making use of the implementation provided by Cleverhans (Papernot et al., 2018). We employ a transfer learning scheme by using DenseNet-121 (Huang et al., 2017) trained on ImageNet, as a fixed feature extractor for CIFAR-10. The predictive network for both MNIST and CIFAR-10 is a two-layer MLP with 100 hidden units. The probabilistic hypernetworks are two-layer MLPs with 50 hidden units. Note, that we do not use adversarial training when training the Deep Ensemble in this experiment to allow for a fair comparison. Furthermore, we use SVHN (Netzer et al., 2011) as the OOD training set for DPN baseline.

**MNIST:** Figure 4 presents the accuracy and the average empirical entropy of the predictive distribution w.r.t. adversarial examples for MNIST with varying levels of perturbation strength (between 0 and 1). We observe that the CDNs are more robust to adversarial examples than all baseline models. More specifically, the ML-CDN is significantly more robust in terms of accuracy to adversarial examples than all other models, while showing a competitive and nicely increasing entropy. The VB-CDN has only slightly better prediction accuracy but attains higher uncertainty than all the baselines except the DPN. Moreover, it shows uncertainties close to that of the DPN, while having higher accuracy and without needing additional data during training. Furthermore, we found that using more samples of  $\theta$  during training is beneficial for the robustness of both ML-CDNs and VB-CDNs, as shown in Figure 5. This behavior is significantly more pronounced for CDNs than for BNNs (as exemplarily shown for Noisy K-FAC and VMG). When using 10 samples per iteration during training the accuracy stays over 0.7 and 0.5 for ML-CDNs and VB-CDNs respectively, even for strong perturbations. As shown in Figure 6, even when the adversarial examples are stronger, i.e. estimated

<sup>14</sup>We use MNIST as OOD data for training the DPN.

<sup>15</sup>We generate the adversarial examples based on a single forward-backward pass.

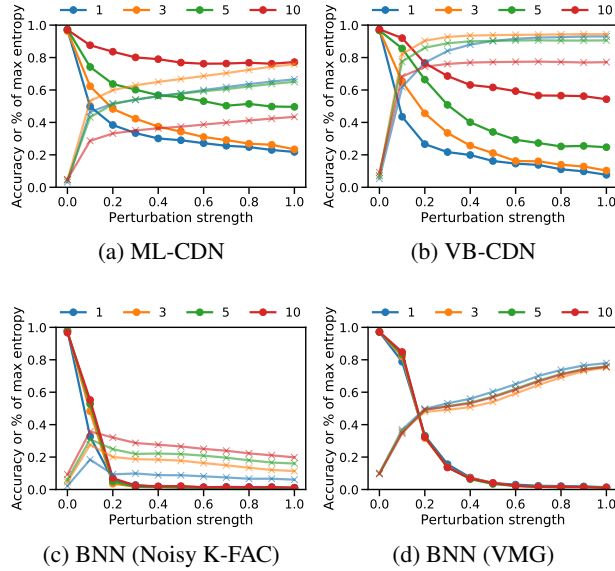


Figure 5: Accuracy and average entropy of CDNs and a BNN (Noisy K-FAC) under FGSM attack, with varying number of samples of  $\theta$  used during training. Circles indicate accuracy, while crosses indicate entropy. The y-axis represents both the accuracy and the entropy relative to the maximum entropy (i.e.  $\ln 10$ ). While using more samples during training significantly improves the overall performance of CDNs, it only has little impact for BNNs (Noisy K-FAC and VMG).

by averaging over multiple forward-backward passes, the performance of both CDNs is only marginally decreased (for VB-CDNs, it stays almost unchanged).

**Fashion-MNIST:** The results on Fashion-MNIST are shown in Figure 7: Overall the same observations and conclusions can be made as for MNIST. We note that strangely, the DPN’s uncertainty estimate is decreasing with increasing perturbation strength.

**CIFAR-10:** The results shown in Figure 8 demonstrate that the VB-CDN is competitive to other state-of-the-art models on CIFAR-10. The ML-CDN does not reflect uncertainty very well but has slightly higher accuracy than other models.

## 6. Conclusion

We introduce compound density networks (CDNs), a new class of models that allows for better uncertainty quantification in neural networks (NNs) and corresponds to a compound distribution (i.e. a mixture with uncountable components) in which both the component distribution and the input-dependent mixing distribution are parametrized by NNs. CDNs are inspired by the success of recently pro-

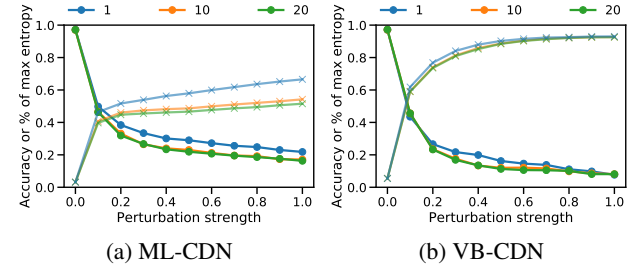


Figure 6: Prediction accuracy and average entropy of CDNs for stronger adversarial examples, constructed by averaging over multiple forward-backward passes.

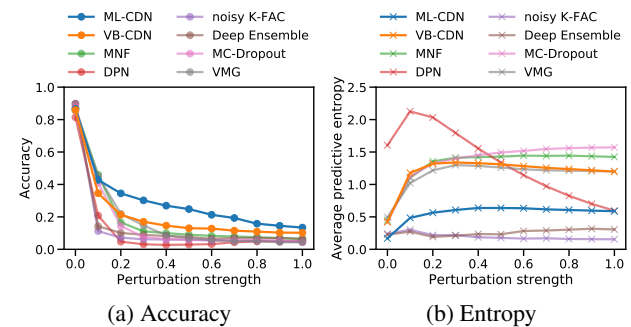


Figure 7: Prediction accuracy and average entropy of models trained on Fashion-MNIST when attacked by FGSM-based adversarial examples with varying perturbation strength.

posed ensemble methods and represent a continuous generalization of mixture density networks (MDNs) (Bishop, 1994). They can be implemented by using hypernetworks to map the input to a distribution over the parameters of the target NN, that models a predictive distribution. For training CDNs, regularized maximum likelihood or variational Bayes can be employed. Extensive experimental analyses showed that CDNs are able to produce promising results in terms of uncertainty quantification. Specifically, Bayesian CDNs are able to capture epistemic as well as aleatoric uncertainty, and yield very high uncertainty on out-of-distribution samples while still making high confidence predictions on within-distribution samples. Furthermore, when facing FGSM-based adversarial attacks, the predictions of CDNs are significantly more robust in terms of accuracy than those of previous models. This robustness under adversarial attacks is especially pronounced for CDNs trained with a maximum likelihood objective, but also clearly visible for Bayesian CDNs, which also provide a better chance of detecting the attack by showing increased uncertainty compared to the baselines. These promising experimental results indicate the benefits of applying a mixture model approach in conjunction with Bayesian inference



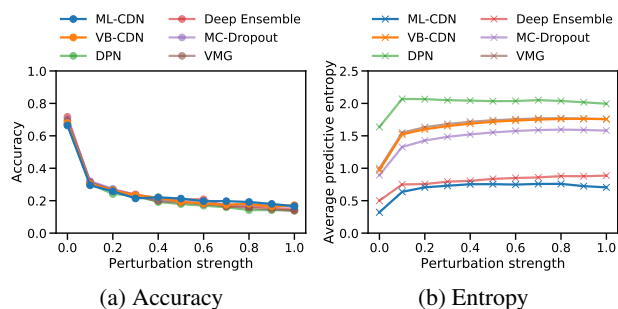


Figure 8: Prediction accuracy and average entropy of models trained on CIFAR-10 when attacked by FGSM-based adversarial examples with varying perturbation strength.

for uncertainty quantification in NNs. We will investigate other implementations of CDNs and adaptations to recurrent and convolutional NNs in future.

## Acknowledgements

SD and AF are funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s Excellence Strategy - EXC 2092 CASA - 390781972.

## References

- Alemi, A., Fischer, I., Dillon, J., and Murphy, K. Deep variational information bottleneck. In *ICLR*, 2017. URL <https://arxiv.org/abs/1612.00410>.
- Alemi, A. A., Fischer, I., and Dillon, J. V. Uncertainty in the variational information bottleneck. *arXiv preprint arXiv:1807.00906*, 2018.
- Amari, S.-I. Natural gradient works efficiently in learning. *Neural computation*, 10(2):251–276, 1998.
- Bishop, C. M. Mixture density networks. 1994.
- Bishop, C. M. *Pattern Recognition and Machine Learning*. Springer, 2006. ISBN 978-0387-31073-2.
- Blundell, C., Cornebise, J., Kavukcuoglu, K., and Wierstra, D. Weight uncertainty in neural networks. pp. 1613–1622, 2015.
- Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., and Kuksa, P. Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12(Aug):2493–2537, 2011.
- Depeweg, S., Hernández-Lobato, J. M., Doshi-Velez, F., and Udluft, S. Learning and policy search in stochastic dynamical systems with bayesian neural networks. In *Proceedings of the Second International Conference on Learning Representations (ICLR 2017)*, 2017.
- Depeweg, S., Hernandez-Lobato, J.-M., Doshi-Velez, F., and Udluft, S. Decomposition of uncertainty in bayesian deep learning for efficient and risk-sensitive learning. In *International Conference on Machine Learning*, pp. 1192–1201, 2018.
- Gal, Y. and Ghahramani, Z. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pp. 1050–1059, 2016.
- Goodfellow, I., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015.
- Graves, A. Practical Variational Inference for Neural Networks. In *Advances in Neural Information Processing Systems 24*, pp. 2348–2356. 2011.
- Guo, C., Pleiss, G., Sun, Y., and Weinberger, K. Q. On calibration of modern neural networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pp. 1321–1330, 06–11 Aug 2017.
- Gupta, A. K. and Nagar, D. K. *Matrix variate distributions*. Chapman and Hall/CRC, 1999.
- Ha, D., Dai, A., and Le, Q. V. HyperNetworks. In *Proceedings of the Second International Conference on Learning Representations (ICLR 2017)*, 2017.
- Hendrycks, D. and Gimpel, K. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *Proceedings of International Conference on Learning Representations*, 2017.
- Hernández-Lobato, J. M. and Adams, R. Probabilistic back-propagation for scalable learning of bayesian neural networks. In *International Conference on Machine Learning*, pp. 1861–1869, 2015.
- Hinton, G. E. and Van Camp, D. Keeping the neural networks simple by minimizing the description length of the weights. In *Proceedings of the sixth annual conference on Computational learning theory*, pp. 5–13. ACM, 1993.
- Huang, G., Liu, Z., Van Der Maaten, L., and Weinberger, K. Q. Densely connected convolutional networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2261–2269. IEEE, 2017.
- Jacobs, R. A., Jordan, M. I., Nowlan, S. J., and Hinton, G. E. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87, 1991.

- Jia, X., De Brabandere, B., Tuytelaars, T., and Gool, L. V. Dynamic filter networks. In *Advances in Neural Information Processing Systems*, pp. 667–675, 2016.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference for Learning Representations*, 2015.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- Krueger, D., Huang, C.-W., Islam, R., Turner, R., Lacoste, A., and Courville, A. Bayesian Hypernetworks. *arXiv:1710.04759 [cs, stat]*, October 2017. arXiv: 1710.04759.
- Lakshminarayanan, B., Pritzel, A., and Blundell, C. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems*, pp. 6402–6413, 2017.
- Louizos, C. and Welling, M. Structured and efficient variational deep learning with matrix gaussian posteriors. In *International Conference on Machine Learning*, pp. 1708–1716, 2016.
- Louizos, C. and Welling, M. Multiplicative normalizing flows for variational Bayesian neural networks. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 2218–2227, 2017.
- MacKay, D. J. A practical bayesian framework for backpropagation networks. *Neural computation*, 4(3):448–472, 1992.
- Malinin, A. and Gales, M. Predictive uncertainty estimation via prior networks. *arXiv preprint arXiv:1802.10501*, 2018.
- Murphy, K. P. *Machine Learning: A Probabilistic Perspective*. The MIT Press, 2012. ISBN 0262018020, 9780262018029.
- Neal, R. M. *BAYESIAN LEARNING FOR NEURAL NETWORKS*. PhD thesis, University of Toronto, 1995.
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, volume 2011, pp. 5, 2011.
- Papernot, N., Faghri, F., Carlini, N., Goodfellow, I., Feinman, R., Kurakin, A., Xie, C., Sharma, Y., Brown, T., Roy, A., Matyasko, A., Behzadan, V., Hambardzumyan, K., Zhang, Z., Juang, Y.-L., Li, Z., Sheatsley, R., Garg, A., Uesato, J., Gierke, W., Dong, Y., Berthelot, D., Hendricks, P., Rauber, J., and Long, R. Technical report on the cleverhans v2.1.0 adversarial examples library. *arXiv preprint arXiv:1610.00768*, 2018.
- Pawlowski, N., Brock, A., Lee, M. C. H., Rajchl, M., and Glocker, B. Implicit Weight Uncertainty in Neural Networks. *arXiv:1711.01297 [cs, stat]*, November 2017. arXiv: 1711.01297.
- Peterson, C. A mean field theory learning algorithm for neural networks. *Complex systems*, 1:995–1019, 1987.
- Ritter, H., Botev, A., and Barber, D. A scalable laplace approximation for neural networks. In *International Conference on Learning Representations*, 2018.
- Sheikh, A.-S., Rasul, K., Merentitis, A., and Bergmann, U. Stochastic maximum likelihood optimization via hypernetworks. In *Advances in Neural Information Processing Systems*, 2017.
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- Sun, S., Chen, C., and Carin, L. Learning structured weight uncertainty in bayesian neural networks. In *Artificial Intelligence and Statistics*, pp. 1283–1292, 2017.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. Intriguing properties of neural networks. 2014.
- Tishby, N., C. Pereira, F., and Bialek, W. The information bottleneck method. *Proceedings of the 37th Allerton Conference on Communication, Control and Computation*, 49, 07 2001.
- Xiao, H., Rasul, K., and Vollgraf, R. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017.
- Zhang, G., Sun, S., Duvenaud, D., and Grosse, R. Noisy natural gradient as variational inference. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 5852–5861, 2018.