# Pairing an arbitrary regressor with an artificial neural network estimating aleatoric uncertainty

Pavel Gurevich,[*] Hannes Stuke[‡]

September 5, 2018

### Abstract

We suggest a general approach to quantification of different forms of aleatoric uncertainty in regression tasks performed by artificial neural networks. It is based on the simultaneous training of two neural networks with a joint loss function and a specific hyperparameter $\lambda > 0$ that allows for automatically detecting noisy and clean regions in the input space and controlling their *relative contribution* to the loss and its gradients. After the model has been trained, one of the networks performs predictions and the other quantifies the uncertainty of these predictions by estimating the locally averaged loss of the first one. Unlike in many classical uncertainty quantification methods, we do not assume any a priori knowledge of the ground truth probability distribution, neither do we, in general, maximize the likelihood of a chosen parametric family of distributions. We analyze the learning process and the influence of clean and noisy regions of the input space on the loss surface, depending on $\lambda$. In particular, we show that small values of $\lambda$ increase the relative contribution of clean regions to the loss and its gradients. This explains why choosing small $\lambda$ allows for better predictions compared with neural networks without uncertainty counterparts and those based on classical likelihood maximization. Finally, we demonstrate that one can naturally form ensembles of pairs of our networks and thus capture both aleatoric and epistemic uncertainty and avoid overfitting.

***Keywords.*** Uncertainty quantification, aleatoric noise, regression, artificial neural networks, ensembles, learning speed

# Contents

[*]Free University of Berlin, Berlin, Germany; RUDN University, Moscow, Russia; email: gurevich@math.fu-berlin.de

[†]Free University of Berlin, Berlin, Germany; email: h.stuke@fu-berlin.de

[‡]Equal contribution.

# 1 Introduction

Neural networks (NNs) are among the main tools that are used nowadays for solving regression and forecasting problems [7, 8, 16, 23]. One theoretical limitation of standard NNs with regression is that they generate averaged predictions of the target variables, but do not provide information on how certain the predictions are. Obviously, quantification of this uncertainty is crucial from the viewpoint of the real world applications [4, 22]. One of the first approaches to learning uncertainty by NNs was the delta method [17, 42, 44] originating from nonlinear regression theory. It assumes that the noise is input-independent (homoscedastic), while the prediction error is proportional to the noise and to the gradient of the output of a NN with respect to the NN's weights. Nowadays, there are two major approaches to learning uncertainty by NNs, both having a probabilistic background. The first exploits Bayesian NNs allowing one to capture epistemic uncertainty (related to

lack of data). The second, sometimes called the frequentist approach, treats the weights as deterministic parameters and rather captures aleatoric uncertainty by directly reconstructing probability distributions of the observed data.

In the Bayesian approach to NNs [14, 31, 34], the weights are treated as random variables. Based on the likelihood of the target variables, a prior distribution of the weights and Bayes' theorem, one obtains a posterior distribution of the weights, and hence a predictive distribution for the data. However, in practice, the latter two are computationally intractable, especially for large data sets and large network architectures. Different ways to approximate them form a field of ongoing research [3–5, 12–14, 18, 21, 26, 27, 29, 30, 34, 43]. In [5], a connection between NNs trained with dropout [15, 39] and Gaussian processes [38] was established. The latter are Bayesian, but generally "non-NN" regression methods. The typical choice of the likelihood in Bayesian NNs is Gaussian with an input-independent variance. Although this variance is often estimated during fitting, the input-dependence of the variance of the predictive distribution occurs only due to the variance of the weights.

In the frequentist approach, one directly approximates the ground truth input-dependent distribution of the target variables by a parametrically given predictive distribution. The parameters of the latter are approximated by NNs with deterministic weights. The typical choice of the approximating distribution is Gaussian [35], explicitly yielding the mean and the variance of the data, but other options were also considered, e.g., Laplacian [20], Students' t [10], and mixture of Gaussians [1, 2, 46].

To explain both epistemic and aleatoric uncertainty at the same time, a combination of the Bayesian approach (using the dropout variational inference) with the mean-variance estimate was used in [20]. A combination of the mean-variance estimate and ensembles of NNs was suggested in [24, 25].

We refer the reader to [4, 10, 24, 25, 33] for further references and experimental comparison of the above methods.

In this paper, we propose an approach that generalizes the frequentist methods (in which one maximizes the likelihood of the data such as Gaussian or Laplace). However, unlike in the above methods, we do not require a priori knowledge about the probabilistic structure of the noise and do not necessarily maximize a likelihood of the data. Our approach quantifies aleatoric uncertainty by automatically estimating a locally averaged loss of the regression network (called the *regressor*) with the help of the second network (called the *uncertainty quantifier*). This quantification is applicable to any loss of the regressor and thus allows for estimating uncertainty exactly in terms of the objective one wants to minimize. We will see that it not only quantifies how certain the predictions are, but also allows for better predictions (compared with standard NNs and NNs based on classical likelihood maximization), especially in regions of the input space with relatively small noise in the target variables (clean regions). This is achieved due to simultaneously training the regressor and the uncertainty quantifier by minimizing a joint loss.

As we said, our approach does not use any explicit form of probability distribution of the data. However, once this assumption is done, one can give a natural interpretation for ensembles of pairs of our networks in terms of mixture distributions, similarly to [24, 25]. This allows one to capture both aleatoric and epistemic uncertainty and avoid overfitting. Again, we will see that our ensembles typically outperform those in [24, 25].

The paper is organized as follows. In Sec. 2, we informally explain our main idea. In Sec. 3, we introduce our model and describe the joint loss function. In Sec. 4, we give an explicit formula

for estimating the expected regressor's loss, introduce ensembles, and discuss similarities and dissimilarities of our method to the classical likelihood maximization approach. In Sec. 5, we consider different functional forms of the joint loss in the cases where the output of the uncertainty quantifier is implemented as the sigmoid and softplus activation functions, respectively. In Sec. 6, we analyze how the hyperparameters of the joint loss affect the overall learning speed and how the loss surface is influenced by clean and noisy regions. In particular, we will see that small values of the hyperparameter reduce the relative contribution of samples from noisy regions to the loss and explain why the usage of the joint loss may improve regressor's predictions. In Sec. 7, we illustrate our results with synthetic one-dimensional data. In Sec. 8, we apply our approach to publicly available data sets and compare it with other NN methods. Section 9 contains a conclusion. In Appendix A, we present the values of hyperparameters of different methods that are compared in Sec. 8.

## 2  Main idea

Suppose we have a standard neural network $\mathcal{N}_r = \mathcal{N}_r(x)$ for regression (the *regressor*) with a loss function $\mathcal{L}_r$. We complement it by another neural network $\mathcal{N}_q = \mathcal{N}_q(x)$ (the *uncertainty quantifier*), and train both networks by minimizing a joint loss of the form

$$\mathcal{L}_{\text{joint}} = \mathcal{L}_r(\mathcal{N}_r)f(\mathcal{N}_q) + \lambda g(\mathcal{N}_q), \tag{2.1}$$

where $f(z)$ and $g(z)$ are some fixed functions and $\lambda > 0$ is a hyperparameter. The main assumption concerning the functions $f$ and $g$ is that the former is positive ($f(z) > 0$) and increasing ($f'(z) > 0$) and the latter is decreasing ($g'(z) < 0$). We will see below that *small* values of $\mathcal{N}_q$ correspond to uncertain predictions and *large* values to certain predictions. Intuitively, the more certain the quantifier $\mathcal{N}_q$ is, the larger $f(\mathcal{N}_q)$ is and hence the smaller the loss $\mathcal{L}_r$ tends to be. On the other hand, the less certain the quantifier $\mathcal{N}_q$ is, the smaller $f(\mathcal{N}_q)$ is and hence the larger the loss $\mathcal{L}_r$ can be. The second term $\lambda g(\mathcal{N}_q)$ penalizes uncertain predictions. Thus, the regressor $\mathcal{N}_r$ can "afford" to fit worse in uncertain noisy regions and use this freedom to fit better in certain clean regions. Moreover, while the regressor $\mathcal{N}_r$ predicts the target value, it turns out (Interpretation 4.1) that the quantifier $\mathcal{N}_q$ allows for estimating the expected regressor's loss via the formula

$$\text{expected regressor's loss} = -\frac{\lambda g'(\mathcal{N}_q)}{f'(\mathcal{N}_q)}. \tag{2.2}$$

We emphasize that relation (2.2) need not involve any likelihood maximization and directly estimates any regressor's loss (see Sec. 4.3 for more details). However, e.g., in the case where the loss $\mathcal{L}_r$ is chosen as the mean square error (MSE), this yields the empirical variance of the data for *any* choice of $f$, $g$, and $\lambda$ and *any* (a priori unknown) ground truth distribution (Example 4.1).

In Sec. 6, we explain that choosing small $\lambda$ may significantly facilitate regressor's ability to learn in clean regions (compared with standard NNs and NNs based on classical likelihood maximization). We perform a rigorous analysis that indicates how clean and noisy regions contribute to $\mathcal{L}_{\text{joint}}$ and its gradients, depending on $\lambda$, on the functions $f$ and $g$, and on the activation function of the output node of the quantifier $\mathcal{N}_q$. In particular, we show that choosing small $\lambda$ decreases the relative contribution of samples from noisy regions to the loss. As a result, the minima of the original loss $\mathcal{L}_{\text{joint}}$ get close to the minima of the loss containing samples from clean regions *only*. Table 2.1 informally summarizes this influence for the sigmoid and softplus activation functions, while Fig. 2.1 and 2.2 illustrate the qualitative dependence of the overall learning speed and of the relative contribution of clean and noisy regions, depending on $\lambda$.

4

| $\lambda$ | sigmoid | | softplus | |
|---|---|---|---|---|
| | regressor | quantifier | regressor | quantifier |
| small | clean > noisy | clean = noisy | clean > noisy | clean = noisy |
| large | clean = noisy | clean < noisy | clean > noisy | clean ≪ noisy |

Table 2.1: Relative contribution of clean and noisy regions to fitting the parameters of the regressor and uncertainty quantifier in case of sigmoid and softplus activation functions in the output of the quantifier $\mathcal{N}_q$. Notation "=" stands for "comparable" contribution, "<" and ">" for a "lower" and "higher" contribution, and "≪" for a "much lower" contribution.



Figure 2.1: The overall learning speed and the relative contribution of clean and noisy regions to fitting the parameters of the regressor $\mathcal{N}_r$ and the uncertainty quantifier $\mathcal{N}_q$, depending on $\lambda$ for the *sigmoid output* of $\mathcal{N}_q$. Left: overall learning speed. Right: Relative contribution of clean and noisy regions.
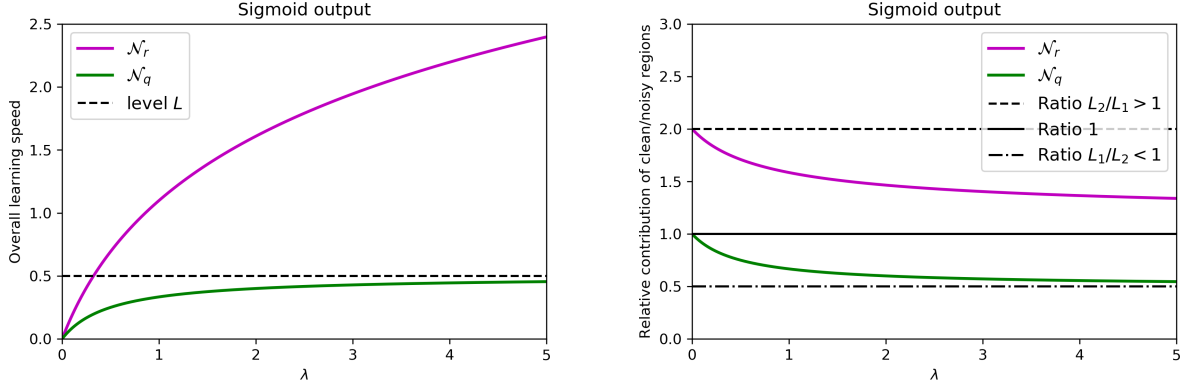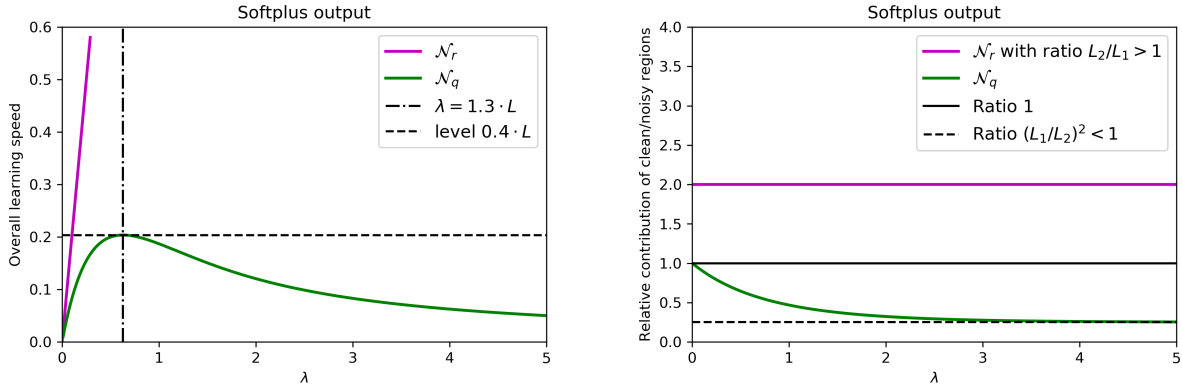


Figure 2.2: The overall learning speed and the relative contribution of clean and noisy regions to fitting the parameters of the regressor $\mathcal{N}_r$ and the uncertainty quantifier $\mathcal{N}_q$, depending on $\lambda$ for the *softplus output* of $\mathcal{N}_q$. Left: overall learning speed. Right: Relative contribution of clean and noisy regions.

# 3   General model

Let $X = \{x^1, \dots, x^N\} \subset \mathbb{R}^n$, $n \in \mathbb{N}$, be a (training) data set consisting of $N \in \mathbb{N}$ samples, and $Y = \{y^1, \dots, y^N\} \subset \mathbb{R}^m$, $m \in \mathbb{N}$, the corresponding target set of observations. The model we propose consists of two NNs that are trained simultaneously. The first network $\mathcal{N}_r : \mathbb{R}^n \to \mathbb{R}^m$ is supposed to perform the regression task, while the second network $\mathcal{N}_q : \mathbb{R}^n \to \mathbb{R}$ is supposed to quantify the uncertainty associated with the predictions of the regression network. Thus, to each sample $x \in \mathbb{R}^n$, we assign the pair $\big(\mathcal{N}_r(x), \mathcal{N}_q(x)\big)$, where $\mathcal{N}_r(x)$ is the prediction and $\mathcal{N}_q(x)$ is its certainty. The networks $\mathcal{N}_r$ and $\mathcal{N}_q$ are parametrized by learnable weights $\theta_r$ and $\theta_q$, respectively, that may be shared but need not be. We omit the dependence of the networks on these weights whenever it does not lead to confusion. In what follows, we call $\mathcal{N}_r$ a *regressor* and $\mathcal{N}_q$ an *uncertainty quantifier*. We will see that the smaller $\mathcal{N}_q(x)$ is, the more uncertain the prediction of $\mathcal{N}_r(x)$ is.

Let
$$\mathcal{L}_r(y^i, y_r^i), \quad y_r^i = \mathcal{N}_r(x^i),$$

be a loss of the regressor $\mathcal{N}_r$. This can be any loss function used in artificial NNs for regression. We only assume that $\mathcal{L}_r$ takes positive values. Now we define a joint loss as follows (cf. (2.1)):

$$\mathcal{L}_{\text{joint}} = \frac{1}{N} \sum_{i=1}^{N} \big(\mathcal{L}_r(y^i, y_r^i) f(z^i) + \lambda g(z^i)\big), \quad y_r^i = \mathcal{N}_r(x^i), \ z^i = \mathcal{N}_q(x^i), \tag{3.1}$$

where $f, g : I \mapsto \mathbb{R}$ are some fixed functions defined on an open interval $I \subset \mathbb{R}$ and $\lambda > 0$ is a hyperparameter. We discuss their role in Secs. 5 and 6. We will see that the choice of the functions $f(z)$ and $g(z)$ depends on concrete implementations of the uncertainty quantifier $\mathcal{N}_q$, while $\lambda$ affects the overall learning speed and the ratio between the learning speeds in clean and noisy regions. To keep notations uncluttered, we do not indicate the arguments of $\mathcal{L}_{\text{joint}}$. Depending on the context, we will treat it as a function of $y_r^i, z^i$ $(i = 1, \dots, N)$ or $\theta_r, \theta_q$.

We assume throughout the following.

**Condition 3.1.**    *1.  $f(z) > 0$, $f'(z) > 0$, $g'(z) < 0$ for all $z \in I$.*

*2.  For any $L > 0$, the function $Lf(z) + \lambda g(z)$ of variable $z$ achieves its finite minimum on $I$.*

Due to item 1, the more certain the quantifier $z^i = \mathcal{N}_q(x^i)$ is, the larger $f(z^i)$ is and hence the smaller the loss $\mathcal{L}_r$ in (3.1) tends to be. The term $\lambda g(z^i)$ in (3.1) penalizes uncertain predictions.

# 4   Probabilistic aspects and generalizations to ensembles

In this section, we give a probabilistic interpretation of our approach. In Sec. 6, we analyze the influence of clean and noisy regions on the loss function and its gradients. In both cases, it is convenient to use an equivalent representation of the loss function (3.1). We will group the points in $X$ in pairwise disjoint sets $X^1, \dots, X^J$ in such a way that $X = \bigcup_{j=1}^{J} X^j$ and both $\mathcal{N}_r(x)$ and $\mathcal{N}_q(x)$ are almost constant on each $X^j$. We denote these constants by $y_r^j$ and $z^j$, respectively (which have slightly different meaning compared with (3.1)). Denote by $Y^1, \dots, Y^J$ the corresponding

subdivision of $Y$. Let $M_j$ be the number of points in $X^j$. Then

$$\begin{aligned}
\mathcal{L}_{\text{joint}} &\approx \sum_{j=1}^{J} \frac{M_j}{N} \left( \frac{1}{M_j} \sum_{y \in Y^j} \mathcal{L}_r(y, y_r^j) f(z^j) + \lambda g(z^j) \right) \\
&= \sum_{j=1}^{J} \frac{M_j}{N} \left( \mathbb{E}_j[\mathcal{L}_r(\cdot, y_r^j)] f(z^j) + \lambda g(z^j) \right),
\end{aligned} \tag{4.1}$$

where $\mathbb{E}_j[\mathcal{L}_r(y, y_r^j)]$ stands for the empirical mean of regressor's loss $\mathcal{L}_r(y, y_r^j)$, $y \in Y^j$, and we used the relation $M_1 + \cdots + M_J = N$. For the further theoretical analysis, we will use the right-hand side of (4.1) for $\mathcal{L}_{\text{joint}}$.

## 4.1 Similarity to probabilistic models

We fix the region $X^j$ for some $j$. The first simple observation is as follows.

**Observation 4.1.** *1. Let a pair*

$$(\overline{y_r}, \overline{z}) \in \mathbb{R}^m \times \mathbb{R} \tag{4.2}$$

*be such that $\overline{z}$ is a critical point of the joint loss function*

$$\mathbb{E}_j[\mathcal{L}_r(\cdot, \overline{y_r})] f(z) + \lambda g(z)$$

*with respect to $z$. Then*

$$\mathbb{E}_j[\mathcal{L}_r(\cdot, \overline{y_r})] = -\frac{\lambda g'(\overline{z})}{f'(\overline{z})}.$$

*2. If, additionally, the pair in (4.2) is a critical point of the joint loss function*

$$\mathbb{E}_j[\mathcal{L}_r(\cdot, y_r)] f(z) + \lambda g(z) \tag{4.3}$$

*with respect to $(y_r, z)$, then $\overline{y_r}$ is a critical point of $\mathbb{E}_j[\mathcal{L}_r(\cdot, y_r)]$. Moreover, if $\mathbb{E}_j[\mathcal{L}_r(\cdot, y_r)]$ is convex with respect to $y_r$, then $\overline{y_r}$ is its global minimum.*

Note that, due to Condition 3.1, a critical point $\overline{z} \in \mathbb{R}$ always exists and $f'(\overline{z}) \neq 0$. We emphasize that $(\overline{y_r}, \overline{z})$ in Observation 4.1 must be a *critical point*, but need not be a local minimum. However, in our concrete implementations in Secs. 5 and 6, it is a global minimum.

Observation 4.1 together with representation of the joint loss function in (4.1) implies the following interpretation in terms of neural networks.

**Interpretation 4.1.** *Let $\mathcal{N}_r(x)$ be a prediction of the regressor and $y$ the ground truth value. Then the uncertainty quantifier $z = \mathcal{N}_q(x)$ provides*

$$\text{expected loss } \mathcal{L}_r(y, \mathcal{N}_r(x)) = -\frac{\lambda g'(z)}{f'(z)}.$$

Now we illustrate the relation between our approach and learning probability distributions.

**Example 4.1.**  1. Assume the observations $y$ are scalar sampled from a ground truth probability distribution $\mathbf{P}(y|x)$. We define regressor's loss as

$$\mathcal{L}_r(y, y_r) = |y - y_r|^2. \tag{4.4}$$

Assume that $(\overline{y_r}, \overline{z})$ in (4.2) is a critical point of the joint loss (4.3). Then, due to (4.4) and Observation 4.1,

$$\overline{y_r} = \arg\min_{\mu \in \mathbb{R}} \mathbb{E}_{y \sim \mathbf{P}(y|x)} \left[ |y - \mu|^2 \right] = \mathbb{E}_{y \sim \mathbf{P}(y|x)}[y], \tag{4.5}$$

$$-\frac{\lambda g'(\overline{z})}{f'(\overline{z})} = \mathbb{E}_{y \sim \mathbf{P}(y|x)} \left[ |y - \overline{y_r}|^2 \right] = \mathrm{Var}_{y \sim \mathbf{P}(y|x)}[y]. \tag{4.6}$$

Thus, the regressor $\mathcal{N}_r(x)$ learns the expectation of $y$ due to (4.5), and the uncertainty quantifier $\mathcal{N}_q(x)$ yields the variance of $y$ according to (4.6). We emphasize that we do *not* need to know the exact form of the ground truth distribution $\mathbf{P}(y|x)$ in order to reconstruct its mean and variance.

2. Additionally to the special choice (4.4) of regressor's loss, we consider a particular choice of the functions $f, g$ and the hyperparameter $\lambda$:

$$f(z) = z, \quad g(z) = -\ln z, \quad \lambda = 1. \tag{4.7}$$

Then our method becomes equivalent to maximization of the log-likelihood of an approximating Gaussian distribution

$$\mathbf{P}_{\mathrm{approx}}(y|x, \mu, \tau) = \sqrt{\frac{\tau}{2\pi}} e^{-|y-\mu|\tau/2}. \tag{4.8}$$

In particular, (4.5) and (4.6) take the form $\overline{y_r} = \mu$ and $-\dfrac{\lambda g'(\overline{z})}{f'(\overline{z})} = \overline{z} = \tau$. However, we will see in Secs. 7 and 8 that the choice $\lambda = 1$ is not optimal from the practical viewpoint.

**Example 4.2.**  1. Similarly to Example 4.1 (item 1) it is easy to see that if $\mathcal{L}_r(y, y_r) = |y - y_r|$, then

$$\overline{y_r} = \mathbb{E}_{y \sim \mathbf{P}(y|x)}[y], \qquad -\frac{\lambda g'(\overline{z})}{f'(\overline{z})} = \mathbb{E}_{y \sim \mathbf{P}(y|x)} \left[ |y - \overline{y_r}| \right] \tag{4.9}$$

for *any* ground truth distribution $\mathbf{P}(y|x)$.

2. For the particular choice (4.7) of the functions $f, g$ and the hyperparameter $\lambda$, our method becomes equivalent to maximization of the log-likelihood of the approximating Laplace distribution

$$\mathbf{P}_{\mathrm{approx}}(y|x, \mu, \tau) = \frac{\tau}{2} e^{-|y-\mu|\tau} \tag{4.10}$$

and yields $\overline{y_r} = \mu$ and $-\dfrac{\lambda g'(\overline{z})}{f'(\overline{z})} = \overline{z} = \tau$.

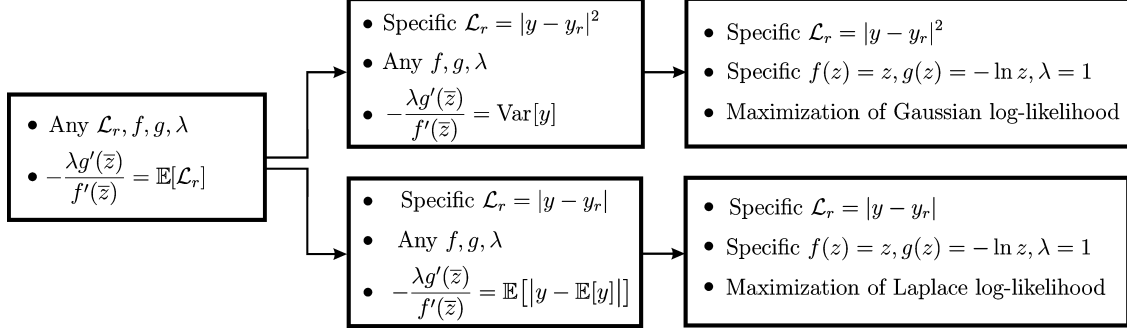Examples 4.1 and 4.2 are illustrated with a flowchart in Fig. 4.1.

Figure 4.1: A flowchart illustrating particular cases of our method. Expectation and variance are taken with respect to (a priori unknown) ground truth distribution $\mathbf{P}(y|x)$.

## 4.2   Ensembles

Similarly to [24, 25], one can fit an ensemble of $K$ pairs of our regressor-quantifier networks and treat their predictions as components of a mixture model. Suppose we are interested in predicting means and variance of observed data. We assume that the data is sampled from an unknown ground truth distribution $\mathbf{P}(y|x, \mu, V)$ with unknown mean $\mu$ and variance $V$. Each individual pair of networks estimates the values $(\mu_j, V_j)$, $j = 1, \ldots, K$. If we interpret them as admissible values for the mean and variance, respectively, we can marginalize over them and, similarly to [24, 25], obtain the predictive (approximating) distribution

$$\mathbf{P}_{\text{approx}}(y|x) = \frac{1}{K} \sum_{j=1}^{K} \mathbf{P}(y|x, \mu_j, V_j)$$

of the ensemble. One can easily check that the mean and the variance of $\mathbf{P}_{\text{approx}}(y|x)$ are

$$\mu = \frac{1}{K} \sum_{j=1}^{K} \mu_j, \qquad V = \frac{1}{K} \sum_{j=1}^{K} \left( V_j + (\mu_j - \mu)^2 \right), \tag{4.11}$$

*independently* of the concrete form of the unknown ground truth distribution $\mathbf{P}(y|x, \mu, V)$.

As another example, suppose we want to predict means and expected absolute errors (EAE) of the observed data. To give a probabilistic interpretation to the ensemble in this case, we assume that the data is sampled from a ground truth Laplace distribution $\mathbf{P}(y|x, \mu, \tau)$ given by the right-hand side in (4.10) with unknown $\mu$ and $\tau$. Due to Example 4.2, each individual pair of networks (for any choice of $f, g, \lambda$) estimates the values $(\mu_j, \tau_j)$, which can be interpreted as admissible values for the parameters of the Laplace distribution (4.10). Marginalizing over them, we obtain the predictive (approximating) distribution

$$\mathbf{P}_{\text{approx}}(y|x) = \frac{1}{K} \sum_{j=1}^{K} \mathbf{P}(y|x, \mu_j, \tau_j).$$

One can easily check that, for this distribution,

$$\mu = \mathbb{E}[y] = \frac{1}{K} \sum_{j=1}^{K} \mu_j, \qquad \mathbb{E}[|y - \mu|] = \frac{1}{K} \sum_{j=1}^{K} \left( |\mu - \mu_j| + \frac{1}{\tau_j} e^{-|\mu - \mu_j|\tau_j} \right). \tag{4.12}$$

9

## 4.3 Dissimilarity to the classical likelihood maximization approach

Let us explain in detail why our approach is different from the classical likelihood maximization. For each fixed $z$, one can consider the (approximating) probability distribution

$$\mathbf{P}_{\text{approx}}(y|x, y_r, z) = \frac{e^{-[\mathcal{L}_r(y, y_r)f(z) + \lambda g(z)]}}{Z(z)} \qquad (4.13)$$

associated with regressor's loss $\mathcal{L}_r(y, y_r)$, cf. [40]. Assuming that regressor's loss depends only on $y - y_r$, the normalization constant

$$Z(z) = e^{-\lambda g(z)} \int_{-\infty}^{\infty} e^{-\mathcal{L}_r(y, y_r)f(z)} dy \qquad (4.14)$$

does not depend on $y_r$, but depends on $z$. For a fixed $z$, the minimization of the joint loss $\mathcal{L}_r(y, y_r)f(z) + \lambda g(z)$ with respect to $y_r$ is obviously equivalent to the maximization of the log-likelihood of $\mathbf{P}_{\text{approx}}(y|x, y_r, z)$ with respect to $y_r$, in which case neither the functions $f(z), g(z)$, nor the normalization constant $Z(z)$ play any role. However, we minimize the joint loss with respect to *both* $y_r$ and $z$ in our approach. One could try to maximize the log-likelihood of $\mathbf{P}_{\text{approx}}(y|x, y_r, z)$ with respect to both $y_r$ and $z$ as well. This is also a valid approach, but then two issues would arise.

First, the normalization constant $Z(z)$ is in general not explicitly given or the integral in (4.14) may even diverge, which makes it impossible in practice to maximize the log-likelihood, nor to find expected regressor's loss[1]

$$\mathbb{E}_{y \sim \mathbf{P}_{\text{approx}}(y|x, y_r, z)}[\mathcal{L}_r(y, y_r)] = \int_{-\infty}^{\infty} \mathcal{L}_r(y, y_r) \mathbf{P}_{\text{approx}}(y|x, y_r, z) \, dy. \qquad (4.15)$$

On the other hand, the minimization of the joint loss $\mathcal{L}_r(y, y_r)f(z) + \lambda g(z)$ is still straightforward.

Second, even if $Z(z)$ can be explicitly calculated, it follows from (4.13) and (4.14) that $\mathbf{P}_{\text{approx}}(y|x, y_r, z)$ is actually independent of $\lambda$ and $g(z)$. Hence, the value $z_{\max}$ for which the maximum of log-likelihood is achieved is also independent of $\lambda$ and of the choice of the function $g$. This is obviously not the case for the minimization of $\mathcal{L}_r(y, y_r)f(z) + \lambda g(z)$. The value $z_{\min}(\lambda)$ that minimizes the joint loss depends on $\lambda$ and $f, g$. This value $z_{\min}(\lambda)$ need *not* maximize the log-likelihood of $\mathbf{P}_{\text{approx}}(y|x, y_r, z)$, but still allows one to estimate the expected loss according to Interpretation 4.1. Thus, the learning process in our case is completely different. Moreover, we will see in Sec. 6 that the freedom to choose $\lambda$ together with the fact that $z_{\min}(\lambda)$ and hence $f(z_{\min}(\lambda))$ depend on $\lambda$ allow one to vary the relative contribution of clean and noisy regions to the joint loss and to effectively downweight samples from noisy regions.

## 5 Sigmoid and softplus activation functions in the output of $\mathcal{N}_q$

In this section, we discuss different choices of the functions $f$ and $g$ for concrete activation functions and rewrite the general Interpretation 4.1 accordingly.

---

[1] In this formula, we assume that we managed to maximize the log-likelihood and that the found approximating distribution $\mathbf{P}_{\text{approx}}(y|x, y_r, z)$ coincides with the ground truth distribution. In general, the family of distributions $\mathbf{P}_{\text{approx}}(y|x, y_r, z)$ may even be not rich enough to contain $\mathbf{P}(y|x)$. Then formula (4.15) would provide only an approximation for $\mathbb{E}_{y \sim \mathbf{P}(y|x)}[\mathcal{L}_r(y, y_r)]$.

## 5.1 Sigmoid output of $\mathcal{N}_q$

Assume the output $z$ of the uncertainty quantifier $\mathcal{N}_q$ is implemented as the sigmoid activation function

$$z = Z(\xi) = \frac{1}{1 + e^{-\xi}}, \quad \xi = (\text{output of the last hidden layer}) \cdot w, \tag{5.1}$$

where $w$ is a column of weights connecting the last hidden layer with the sigmoid activation function. We define the functions $f$ and $g$ as follows:

$$f(z) = -\ln(1 - z), \quad g(z) = -\ln z, \quad z \in I = (0, 1). \tag{5.2}$$

It is easy to see that Condition 3.1 is fulfilled. The general Interpretation 4.1 takes the following form.

**Interpretation 5.1.** *Let $\mathcal{N}_r(x)$ be a prediction of the regressor and $y$ the ground truth value. Then the uncertainty quantifier $z = \mathcal{N}_q(x)$ provides*

$$\text{expected loss } \mathcal{L}_r(y, \mathcal{N}_r(x)) = \lambda \left( \frac{1}{z} - 1 \right).$$

## 5.2 Softplus output of $\mathcal{N}_q$

Assume the output $z$ of the uncertainty quantifier $\mathcal{N}_q$ is implemented as the softplus nonlinearity

$$z = Z(\xi) = \ln(1 + e^{\xi}), \quad \xi = (\text{output of the last hidden layer}) \cdot w, \tag{5.3}$$

where $w$ has the same meaning as in (5.1). We define the functions $f$ and $g$ as follows:

$$f(z) = z, \quad g(z) = -\ln z, \quad z \in I = (0, \infty), \tag{5.4}$$

Condition 3.1 is again fulfilled. The general Interpretation 4.1 takes the following form.

**Interpretation 5.2.** *Let $\mathcal{N}_r(x)$ be a prediction of the regressor and $y$ the ground truth value. Then the uncertainty quantifier $z = \mathcal{N}_q(x)$ provides*

$$\text{expected loss } \mathcal{L}_r(y, \mathcal{N}_r(x)) = \frac{\lambda}{z}.$$

# 6 Contribution of clean and noisy regions to the joint loss, depending on $\lambda$

In this section, we clarify the role of the hyperparameter $\lambda$. We study in detail the case of the sigmoid activation function in the output of $\mathcal{N}_q$. The analysis of the softplus output is analogous, and hence we will formulate only final conclusions.

We distinguish *clean* and *noisy* regions in the data set $X$. At each learning stage, we say that a region is *clean* if the loss of the regressor $\mathcal{N}_r(x)$ is small for $x$ in this region. Otherwise, we call a region *noisy*. The notions clean and noisy are understood relative to each other. It turns out that $\lambda$ affects to what extent the samples from clean and noisy regions contribute to the loss. The influence is illustrated in Fig. 2.1 and 2.2 and in Table 2.1, which we justify below.

We discuss two different (but closely related) aspects. First, we analyze the contribution of different regions to the gradients of the loss (see Sec. 6.1.1 and 6.1.2 for the sigmoid output and Sec. 6.2.1 and 6.2.2 for the softplus output). Second, we analyze the influence of different regions on the loss surface itself and its minima (see Sec. 6.1.3 for the sigmoid output and Sec. 6.2.3 for the softplus output).

In both cases, we are interested in the structure of the joint loss $\mathcal{L}_{\text{joint}} = \mathcal{L}_{\text{joint}}(\theta_r, \theta_q)$ in the $\theta_r$-space of the regressor's weights and in the $\theta_q$-space of the uncertainty quantifier's weights. Due to (4.1), we can write the joint loss as

$$\mathcal{L}_{\text{joint}} \approx \sum_{j=1}^{J} C_j \mathcal{M}(L_j, \xi^j), \tag{6.1}$$

where

$$\mathcal{M}(L, \xi) = L f(Z(\xi)) + \lambda g(Z(\xi)), \qquad C_j = \frac{M_j}{N}, \tag{6.2}$$

$L_j = \mathbb{E}_j[\mathcal{L}_r(\cdot, y_r^j)]$, the function $Z(\xi)$ is given either by (5.1) or (5.3), $\xi^j = $ (output of the last hidden layer)$\cdot$ $w$, the "output of the last hidden layer" is determined by the inputs from $X^j$ and all the weights of $\mathcal{N}_q$ except for the last layer, and $w$ is a column of weights connecting the last hidden layer with the corresponding activation function of the output of $\mathcal{N}_q$; finally, $f$ and $g$ are given by (5.2).

In this section, we assume that the weights $\theta_q$ of the uncertainty quantifier satisfy the following.

**Condition 6.1.** *The weights $\theta_q$ of the uncertainty quantifier are in a neighborhood of the global minimum $\overline{\theta}_q$ of $\mathcal{L}_{\text{joint}}$ as a function of $\theta_q$. Moreover, this global minimum corresponds to the values $\overline{\xi^j} = \xi^j(\overline{\theta}_q)$ that minimize $\mathcal{M}(L_j, \xi^j)$ for all $j = 1, \ldots, J$.*

Since the gradient of $\mathcal{L}_{\text{joint}}$ equals the sum of the gradients of the individual terms in (6.1), we will concentrate on these individual terms, corresponding to different regions $X_j$.

## 6.1 Sigmoid output of $\mathcal{N}_q$

In this subsection, we assume that the output $z$ of the uncertainty quantifier $\mathcal{N}_q$ is implemented as the sigmoid activation function, see Sec. 5.1. We will justify the column "sigmoid" in Table 2.1.

We begin with the following simple lemma, which shows in particular that, for any $L > 0$, the function $\mathcal{M}(L, \xi)$ has a unique global minimum with respect to $\xi$.

**Lemma 6.1.**　　*1. The function $\mathcal{M}(L, \xi)$ in (6.2) is convex with respect to $\xi$. For each $L > 0$, it achieves a global minimum with respect to $\xi$ at the point $\overline{\xi} = \overline{\xi}(L) = \ln(\lambda/L)$.*

*2. $\dfrac{\partial^2 \mathcal{M}(L, \xi)}{\partial \xi^2}\bigg|_{\xi = \overline{\xi}} = \dfrac{L\lambda}{L + \lambda}$.*

*3. $f(Z(\overline{\xi})) = \ln(1 + \lambda/L)$.*

*Proof.* Due to (6.2), (5.1), and (5.2), $\mathcal{M}(L, \xi) = -L \ln\left(\dfrac{e^{-\xi}}{1 + e^{-\xi}}\right) - \lambda \ln\left(\dfrac{1}{1 + e^{-\xi}}\right)$. Hence, the result follows from the formulas

$$\frac{\partial \mathcal{M}(L, \xi)}{\partial \xi} = \frac{L - \lambda e^{-\xi}}{1 + e^{-\xi}}, \qquad \frac{\partial^2 \mathcal{M}(L, \xi)}{\partial \xi^2} = (L + \lambda)\frac{e^{-\xi}}{(1 + e^{-\xi})^2}.$$

$\square$

### 6.1.1 Impact of $\lambda$ on the overall learning speed

To analyze the overall learning speed of $\mathcal{N}_r$, we check how the gradients of $\mathcal{M}(L, \xi)$ (where $L = L(\theta_r)$ and $\xi = \xi(\theta_q)$) with respect to $\theta_r$ and $\theta_q$ depend on $\lambda$.

**Regressor.** Due to Condition 6.1,

$$\nabla_{\theta_r} \mathcal{M}(L, \xi) = \nabla_{\theta_r} L(\theta_r) \cdot f(Z(\xi)) \approx \nabla_{\theta_r} L(\theta_r) \cdot f(Z(\bar{\xi})). \tag{6.3}$$

Hence, by Lemma 6.1, item 3,

$$|\nabla_{\theta_r} \mathcal{M}(L, \xi)| \propto \ln\left(1 + \frac{\lambda}{L}\right) \quad \text{(the magenta line in Fig. 2.1, left)}, \tag{6.4}$$

where $\propto$ stands for "approximately proportional". Thus, for a fixed $L$, the smaller $\lambda$ is, the closer to 0 the learning speed of $\mathcal{N}_r$ is (asymptotically proportionally to $\lambda/L$). On the other hand, the larger $\lambda$ is, the larger the learning speed of $\mathcal{N}_r$ is.

**Uncertainty quantifier.** By Condition 6.1,

$$\nabla_{\theta_q} \mathcal{M}(L, \xi(\theta_q)) \approx \mathbf{H}(L, \bar{\theta}_q) \cdot (\theta_q - \bar{\theta}_q) = \frac{\partial^2 \mathcal{M}(L, \bar{\xi})}{\partial \xi^2} (\nabla_{\theta_q} \xi) \cdot (\nabla_{\theta_q} \xi)^T \cdot (\theta_q - \bar{\theta}_q), \tag{6.5}$$

where $\mathbf{H}(L, \bar{\theta}_q)$ is the Hessian of $\mathcal{M}(L, \xi(\theta_q))$ with respect to $\theta_q$ evaluated at $\bar{\theta}_q$, and $\bar{\xi}$ is defined in Lemma 6.1. Therefore, by Lemma 6.1, item 2,

$$\left|\nabla_{\theta_q} \mathcal{M}(L, \xi(\theta_q))\right| \propto \frac{L\lambda}{L + \lambda} \quad \text{(the green line in Fig. 2.1, left)}. \tag{6.6}$$

Hence, for a fixed $L$, the smaller $\lambda$ is, the closer to 0 the learning speed of $\mathcal{N}_q$ is (asymptotically proportionally to $\lambda$, independently of $L$). On the other hand, the larger $\lambda$ is, the closer to $L$ the learning speed of $\mathcal{N}_q$ is.

### 6.1.2 Relative contribution of clean and noisy regions to the gradients of the loss, depending on $\lambda$

Now we analyze the *relative* amplitudes of the gradients in clean and noisy regions. Their ratio determines which regions are downweighted or upweighted, respectively, during training. Consider two regions $X^{j_1}$ and $X^{j_2}$ with the corresponding values of regressor's loss $L_1$ and $L_2$. Assume that $X^{j_1}$ is a clean region and $X^{j_2}$ is a noisy region in the sense that

$$L_2 \gg L_1.$$

**Regressor.** Due to (6.1) and (6.4), the relative contribution of the clean and noisy regions $X^{j_1}$ and $X^{j_2}$ to $\nabla_{\theta_r} \mathcal{L}_{\text{joint}}$ is determined by the value

$$R_{\text{sigmoid}}(\lambda) = \frac{\ln\left(1 + \dfrac{\lambda}{L_1}\right)}{\ln\left(1 + \dfrac{\lambda}{L_2}\right)} \quad \text{(the magenta line in Fig. 2.1, right)}. \tag{6.7}$$

The smaller $\lambda$ is, the closer to $L_2/L_1$ ($\gg 1$) the ratio $R_{\text{sigmoid}}(\lambda)$ is and the lower the relative contribution of the noisy region to $\nabla_{\theta_r} \mathcal{L}_{\text{joint}}$ is. *This is the regime in which $\mathcal{N}_r$ can learn more efficient*

*in clean regions compared with the usual regression network without the uncertainty quantifier.* On the other hand, the larger $\lambda$ is, the closer to 1 the ratio $R_{\text{sigmoid}}(\lambda)$ is and the higher the relative contribution of the noisy region to $\nabla_{\theta_r}\mathcal{L}_{\text{joint}}$ is.

**Uncertainty quantifier.** Due to (6.1) and (6.6), the relative contribution of the clean and noisy regions $X^{j_1}$ and $X^{j_2}$ to $\nabla_{\theta_q}\mathcal{L}_{\text{joint}}$ is determined by the value

$$Q_{\text{sigmoid}}(\lambda) = \frac{L_1\lambda}{L_1 + \lambda}\left(\frac{L_2\lambda}{L_2 + \lambda}\right)^{-1} = \frac{L_1(L_2 + \lambda)}{L_2(L_1 + \lambda)} \quad \text{(the green line in Fig. 2.1, right).} \qquad (6.8)$$

The larger $\lambda$ is, the closer to $L_1/L_2$ ($\ll 1$) the ratio $Q_{\text{sigmoid}}(\lambda)$ is and the higher the relative contribution of the noisy region to $\nabla_{\theta_q}\mathcal{L}_{\text{joint}}$ is. On the other hand, the smaller $\lambda$ is, the closer to 1 the ratio $Q_{\text{sigmoid}}(\lambda)$ is and the more balanced the relative contributions of clean and noisy regions to $\nabla_{\theta_r}\mathcal{L}_{\text{joint}}$ is.

### 6.1.3 Relative contribution of clean and noisy regions to the loss surface and its minima, depending on $\lambda$

For clarity, assume we have only two regions with regressor's losses $L_1 = L_1(\theta_r)$ and $L_2 = L_2(\theta_r)$, respectively. As before, let $L_1 \ll L_2$.

**Regressor.** Due to (6.1), (6.2) and Lemma 6.1 (item 3),

$$\mathcal{L}_{\text{joint}} \approx \ln\left(1 + \frac{\lambda}{L_1}\right)\left(C_1 L_1(\theta_r) + \frac{C_2}{R_{\text{sigmoid}}(\lambda)}L_2(\theta_r) + C\right), \qquad (6.9)$$

where $C$ does not depend on $\theta_r$ and $R_{\text{sigmoid}}(\lambda)$ is given by (6.7). As $\lambda \to 0$, we have $1/R_{\text{sigmoid}}(\lambda) \to L_1/L_2 \ll 1$, see the magenta line in Fig. 2.1 (right). Thus, for small $\lambda$, the minimum $\theta_r^*$ of $\mathcal{L}_{\text{joint}}$ is generically close to the minimum $\theta_r^1$ of the loss $L_1(\theta_r)$ defined *only* with samples from the clean region. Figure 6.1 schematically shows the loss surfaces of the loss $\mathcal{L}_{\text{joint}}(\theta_r)$ defined with all samples and the loss surfaces defined with samples only from clean regions $L_1(\theta_r)$ and only from noisy regions $L_2(\theta_r)$. We see that, for small values of $\lambda$, the loss surface $\mathcal{L}_{\text{joint}}$ is close[2] to that given by $L_1(\theta_r)$.

**Uncertainty quantifier.** Now we analyze the loss surface $\mathcal{L}_{\text{joint}}$ in the $\theta_q$-space for fixed values $L_1$ and $L_2$ of regressor's loss. Analogously to (6.5), we have

$$\mathcal{L}_{\text{joint}} \approx C_1 \tilde{\mathcal{M}}_1(\theta_q) + C_2 \tilde{\mathcal{M}}_2(\theta_q) + C,$$

where $C$ does not depend on $\theta_q$,

$$\begin{aligned}
\tilde{\mathcal{M}}_j(\theta_q) &= \frac{1}{2}\frac{\partial^2 \mathcal{M}(L, \xi^j(\theta_q^j))}{\partial \xi^j(\theta_q^j)}\nabla_{\theta_q}\xi^j(\theta_q^j) \cdot \nabla_{\theta_q}\xi^j(\theta_q^j)^T \cdot (\theta_q - \theta_q^j)^2 \\
&= \frac{1}{2}\frac{L_j\lambda}{L_j + \lambda}\nabla_{\theta_q}\xi^j(\theta_q^j) \cdot \nabla_{\theta_q}\xi^j(\theta_q^j)^T \cdot (\theta_q - \theta_q^j)^2,
\end{aligned}$$

and $\theta_q^j$ is the minimum of $\mathcal{M}(L_j, \xi^j)$ (which is the loss taking into account only the samples from the $j$th region, see (6.1) and (6.2)). Now we see that the distance between the minimum $\theta_q^*$ of the

---

[2]After an appropriate vertical shift and rescaling, which do not affect its minima.
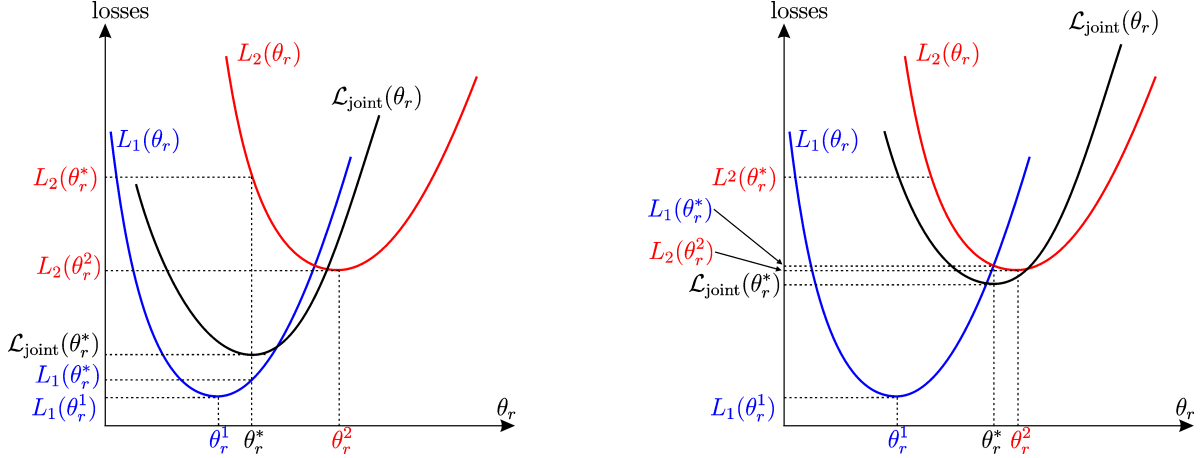
Figure 6.1: Regressor's losses $L_1(\theta_r)$ (blue) and $L_2(\theta_r)$ (red) defined with samples from clean and noisy regions, respectively, and the joint loss $\mathcal{L}_{\text{joint}}(\theta_r)$ defined with all samples as functions of regressor's weights $\theta_r$. The values $\theta_r^1$ and $\theta_r^2$ are the minima of $L_1(\theta_r)$ and $L_2(\theta_r)$, while $\theta_r^*$ is the minimum of $\mathcal{L}_{\text{joint}}(\theta_r)$. Left: Small values of $\lambda$ imply that $\theta_r^*$ is close to $\theta_r^1$; hence the loss $L_1(\theta_r^*)$ in the clean region is close to its optimal value $L_1(\theta_r^1)$. Right: Large values of $\lambda$ imply that $\theta_r^*$ is far away from $\theta_r^1$; hence the loss $L_1(\theta_r^*)$ in the clean region is far away from its optimal value $L_1(\theta_r^1)$.

joint loss $\mathcal{L}_{\text{joint}}$ and the minima $\theta_q^j$ depends on the ratio $Q_{\text{sigmoid}}(\lambda)$ in (6.8). In particular, the larger $\lambda$ is the closer $\theta_q^*$ to $\theta_q^2$ (the minimum corresponding to the noisy region) is.

Note that, due to the assumption in Condition 6.1, the above argument is valid only if the uncertainty quantifier correctly estimates both values $L_1$ and $L_2$ of regressor's losses in clean and noisy regions. These estimates dynamically change during training. Again, due to (6.8), if $\lambda$ is large, the uncertainty quantifier "ignores" clean regions and rather learns from samples in noisy regions. As a result, it may wrongly quantify clean regions as noisy. Hence, it erroneously suppresses their contribution to the gradient $\nabla_{\theta_r}\mathcal{L}_{\text{joint}}$ with respect to the regressor's weights $\theta_r$. On the other hand, if $\lambda$ is small, clean and noisy regions equally contribute to the learning of the quantifier's weights $\theta_q$, which facilitates the fit of the regressor in clean regions.

## 6.2 Softplus output of $\mathcal{N}_q$

Assume the output $z$ of the uncertainty quantifier $\mathcal{N}_q$ is implemented as the softplus nonlinearity, see Sec. 5.2. We will justify the column "softplus" in Table 2.1. Similarly to Lemma 6.1, one can show that the function $\mathcal{M}(L,\xi)$ is convex with respect to $\xi$, and, for each $L > 0$, it achieves a global minimum with respect to $\xi$ at the point $\overline{\xi} = \overline{\xi}(L) = \ln\left(e^{\lambda/L} - 1\right)$ and

$$f(Z(\overline{\xi})) = \frac{\lambda}{L}. \tag{6.10}$$

### 6.2.1 Impact of $\lambda$ on the overall learning speed

**Regressor.** Similarly to Sec. 6.1.1, due to (6.10) the overall learning speed of $\mathcal{N}_r$ is determined by

$$|\nabla_{\theta_r}\mathcal{M}(L,\xi)| \propto \frac{\lambda}{L} \quad \text{(the magenta line in Fig. 2.2, left).} \tag{6.11}$$

**Uncertainty quantifier.** The overall learning speed of $\mathcal{N}_q$ is determined by

$$\left| \nabla_{\theta_q} \mathcal{M}(L, \xi(\theta_q)) \right| \propto \frac{L^2}{\lambda} \left( 1 - e^{-\lambda/L} \right)^2 \quad \text{(the green line in Fig. 2.2, left).} \tag{6.12}$$

We note that the function in (6.12) asymptotically equals $\lambda$ as $\lambda \to 0$ and $L^2/\lambda$ as $\lambda \to \infty$. In particular, it tends to zero both as $\lambda \to 0$ and $\lambda \to \infty$. It is easy to calculate that it achieves its maximum

$$\frac{4\mu_0}{(1 + 2\mu_0)^2} L \approx 0.4 \cdot L \quad \text{(the dashed line in Fig. 2.2, left).}$$

at $\lambda = \mu_0 L$, where $\mu_0 \approx 1.3$ is a positive root of the equation $e^{\mu_0} = 1 + 2\mu_0$.

### 6.2.2 Relative contribution of clean and noisy regions to the gradients of the loss, depending on $\lambda$

As in Sec. 6.1.2, assume that $X^{j_1}$ is a clean region and $X^{j_2}$ is a noisy region in the sense that $L_2 \gg L_1$.

**Regressor.** Due to (6.1) and (6.11), the relative contribution of the clean and noisy regions to $\nabla_{\theta_r} \mathcal{L}_{\text{joint}}$ is determined by the value

$$R_{\text{soft}} = \frac{L_2}{L_1} \quad \text{(the magenta line in Fig. 2.2, right).} \tag{6.13}$$

In particular, it does not depend on $\lambda$, and clean regions always dominate (provided $L_1$ and $L_2$ are correctly estimated by the uncertainty quantifier).

**Uncertainty quantifier.** However, the relative contribution of the clean and noisy regions to $\nabla_{\theta_q} \mathcal{L}_{\text{joint}}$ does depend on $\lambda$. Due to (6.1) and (6.12), it is determined by the value

$$Q_{\text{soft}}(\lambda) = \left( \frac{L_1 \left( 1 - e^{-\lambda/L_1} \right)}{L_2 \left( 1 - e^{-\lambda/L_2} \right)} \right)^2 \quad \text{(the green line in Fig. 2.2, right).} \tag{6.14}$$

As $\lambda$ increases, $Q_{\text{soft}}(\lambda)$ decays to $(L_1/L_2)^2 \ll 1$. Note that this limit is even smaller than the limit $L_1/L_2$ of $Q_{\text{sigmoid}}(\lambda)$ in the case of the sigmoid output. Hence, in this case the main contribution to the gradient $\nabla_{\theta_q} \mathcal{L}_{\text{joint}}$ is due to noisy regions. On the other hand, for small $\lambda$, $Q_{\text{soft}}(\lambda)$ is close to 1, i.e., the contributions of clean and noisy regions get balanced.

### 6.2.3 Relative contribution of clean and noisy regions to the loss surface and its minima, depending on $\lambda$

**Regressor.** Repeating the argument in Sec. 6.1.3 and using (6.10), we see that the distance between the minimum $\theta_r^*$ of $\mathcal{L}_{\text{joint}}(\theta_r)$ and the minimum $\theta_r^1$ of $L_1(\theta_r)$ depends on the value $R_{\text{soft}}$ in (6.13). In particular, it is small if $R_{\text{soft}} \ll 1$, independently of whether $\lambda$ is small or large.

**Uncertainty quantifier.** However, the contribution of the clean region to the loss surface $\mathcal{L}_{\text{joint}}$ in the $\theta_q$-space is governed by the ratio $Q_{\text{soft}}(\lambda)$ in (6.14), which does depend on $\lambda$. As $\lambda$ increases, $Q_{\text{soft}}(\lambda)$ decays to $(L_1/L_2)^2 \ll 1$, see the green line in Fig. 2.2 (right). Saying differently, the minimum $\theta_q^*$ of the joint loss $\mathcal{L}_{\text{joint}}$ with respect to $\theta_q$ gets close to the minimum $\theta_q^2$ of the loss defined with the samples from the noisy region only. As a result, clean regions can be misidentified as noisy. On the other hand, as in Sec. 6.1.3, we see that this does not happen for small $\lambda$, and, therefore, the fit of the regressor in clean regions gets facilitated.
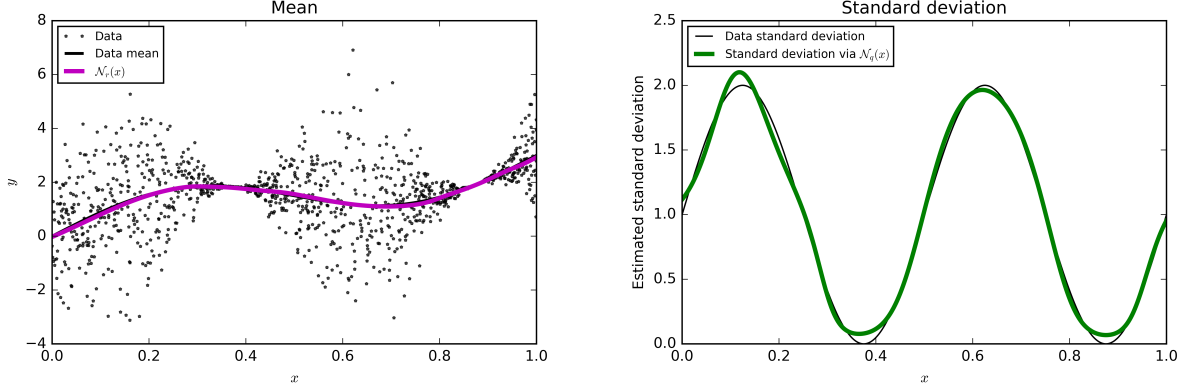
Figure 7.1: Fit of the networks $(\mathcal{N}_r, \mathcal{N}_q)$ in (7.1) with $\lambda = 0.1$, MSE regressor's loss for the data $y$ sampled from the normal distribution with mean $3x + \sin(2\pi x)$ and standard deviation $1 + \sin(4\pi x)$. Left: $\mathcal{N}_r(x)$ (the curve is almost indistinguishable from the black line indicating the mean of the data). Right: standard deviation via $\mathcal{N}_q(x)$ according to (2.2).

# 7 Synthetic data

In this section, we generate data with $X \subset [0, 1]$ and $Y \subset \mathbb{R}$ and implement the uncertainty quantifiers with the sigmoid output, see Sec. 5.1. We choose regressor's loss to be the MSE and generate Gaussian noise. We predict the uncertainty in terms of the standard deviation, using (4.6).

## 7.1 Smooth data

First, we consider smoothly varying mean and variance. Namely, we sample $y$ from the normal distribution with mean $3x + \sin(2\pi x)$ and standard deviation $1 + \sin(4\pi x)$, $x \in [0, 1]$. We implement the network pair $(\mathcal{N}_r, \mathcal{N}_q)$ as follows:

$$\begin{cases} \mathcal{N}_r : \text{ input(1), 2 x hidden(10, tanh), output(1, linear),} \\ \mathcal{N}_q : \text{ input(1), 2 x hidden(10, tanh), output(1, sigmoid).} \end{cases} \tag{7.1}$$

We take $\lambda = 0.1$. The fit of the regressor $\mathcal{N}_r(x)$ is illustrated by Fig. 7.1 (left) and the predictions of the standard deviation via (2.2) are shown in Fig. 7.1 (right). The corresponding choice of $\lambda$, $f$, and $g$ is discussed in Sec. 7.1. Figure 7.1 represents the situation where both the mean of the data and its variance vary smoothly. In this case, the choice of $\lambda$ is not too important from the practical point of view. However, it becomes crucial once the data exhibits rather sharp interfaces. Figure 7.2 illustrates "clean" data given by $y = 3x + \sin(2\pi x)$, $x \in [0, 1]$, complemented by two vertical strips of width 0.1 with the Gaussian noise (mean is $-2$ in each strip and standard deviation is 1 in the first strip and 5 in the second). One can see that the standard NN (red line) estimates the mean in the clean regions outside the two noisy strips much worse than the regressor $\mathcal{N}_r$ (blue line). This is especially evident in the lower picture, where the noisy regions contain 80% of the data. These and other simulations are discussed in more detail in Sec. 7. We refer to Sec. 8 for the analysis of real world data sets.
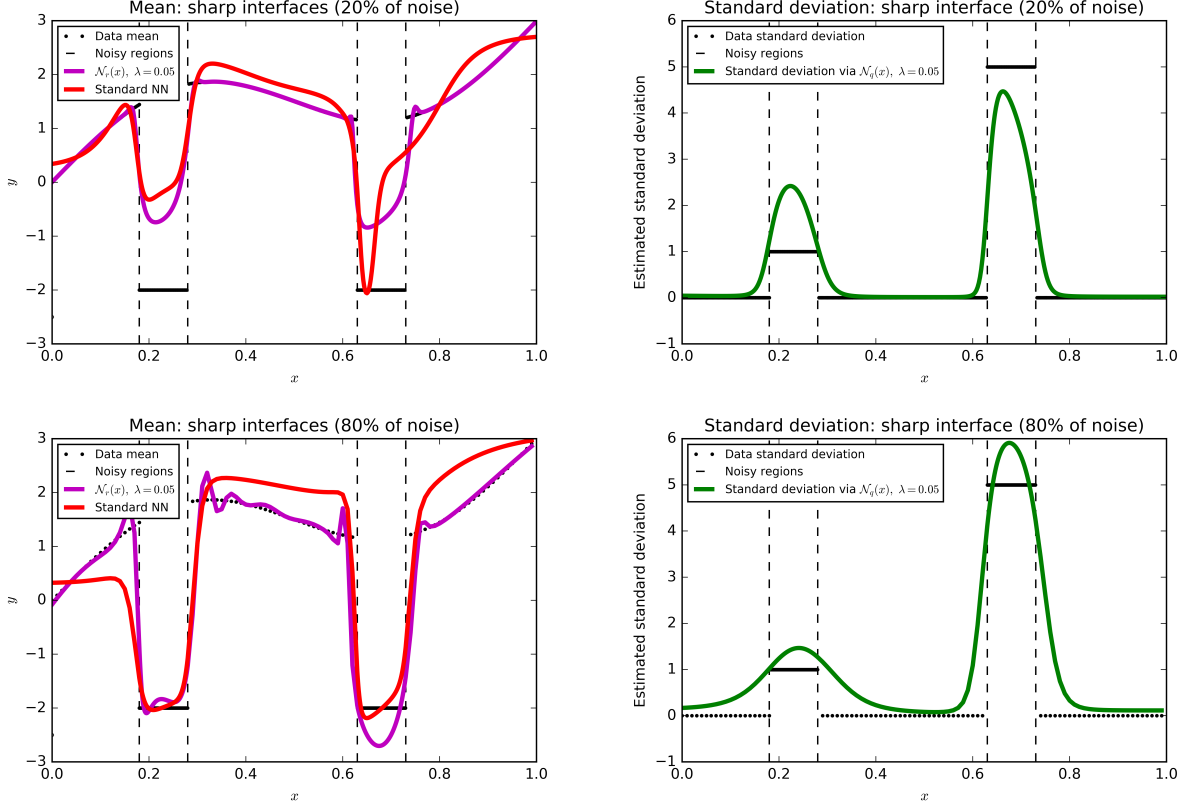
Figure 7.2: Fit of the networks $(\mathcal{N}_r, \mathcal{N}_q)$ in (7.1) with MSE regressor's loss in comparison with a standard NN (in red). The dashed vertical lines indicate the two regions of width 0.1 each with Gaussian noise with mean $-2$ and standard deviations 1 and 5, respectively. Top: noisy regions contain 20% of the data. Bottom: noisy regions contain 80% of the data. In both cases, $\lambda = 0.05$.

## 7.2 Data with sharp interfaces

In our second example, we generate "clean" data given by $y = 3x + \sin(2\pi x)$, $x \in [0, 1]$, and complement them by two vertical strips of width 0.1 with Gaussian noise (mean is $-2$ in each strip and standard deviation is 1 in the first strip and 5 in the second). The fit of the network (7.1) is illustrated in Fig. 7.2. One can see that even if the data contains 80% of noise, the uncertainty quantifier allows the regressor to fit well enough for the remaining 20%. This can be explained by formula (4.1) for the loss and by the results in Sec. 6. Indeed, small coefficients $M_j/N$ in (4.1) corresponding to regions $X_j$ with low density would also be present in the loss function of a standard NN for regression. As a consequence, the samples from these regions would contribute little to the gradient of the loss, and the gradient descent would be mostly governed by the samples from the noisy regions. On the other hand, as we saw in Sec. 6 (cf. also Table 2.1 and Fig. 2.1 and 2.2 (right)), in case of the joint loss (4.1), small values of $\lambda$ yield larger gradients of the terms $\mathbb{E}_j[\mathcal{L}_r(\cdot, y_r^j)]f(z^j) + \lambda g(z^j)$ corresponding to clean regions $X_j$, which compensate the small values of $M_j/N$.

Next, we compare the case $\lambda = 0.05$ (Fig. 7.3, the solid lines) with the case $\lambda = 2$ (Fig. 7.3, the dashed lines) and illustrate how $\lambda$ affects the relative learning speeds in clean and noisy regions (cf.

Sec. 6.1.2 and Fig. 2.1, right). We fill in the two noisy strips with 80% of the data and plot the fit of the network $\mathcal{N}_r$ and the standard deviation via $\mathcal{N}_q$ according to (4.6) on three different learning stages. The four vertical dashed lines in Fig. 7.3 divide the interval $[0, 1]$ in five subintervals. We refer to them as regions 1–5 from left to right (with regions 2 and 4 filled with the Gaussian noise).

Consider $\lambda = 0.05$ (bold lines in Fig. 7.3). In this case, the regressor $\mathcal{N}_r$ learns faster in the regions that have smaller variance according to $\mathcal{N}_q$. Indeed, we see that $\mathcal{N}_r$ first starts to learn in regions 1 and 2, where $\mathcal{N}_q$ is smaller (top figure), and then additionally in region 3, where $\mathcal{N}_q$ was moderately larger (middle figure). As $\mathcal{N}_q$ learns that the variance in region 5 is less than in region 4, the regressor $\mathcal{N}_r$ accelerates its learning in region 5 (bottom picture).

Now consider $\lambda = 2$ (dashed lines in Fig. 7.3). In this case, $\mathcal{N}_q$ learns faster in the noisy regions. In particular, this results (bottom figure) in the domination of the noisy region 2 over the clean region 1. This prevents $\mathcal{N}_q$ from learning that region 1 has a smaller variance compared with region 2. On the other hand, the learning speeds of the regressor $\mathcal{N}_r$ in clean and noisy regions are closer to each other. As a result, it is not able to estimate regions 1, 3, and 5 (outside of the noisy strips) as well as it does for $\lambda = 0.05$.

Finally, we illustrate the implementation of the output of $\mathcal{N}_q$ via the softplus, see Sec. 5.2. The architecture of $(\mathcal{N}_r, \mathcal{N}_q)$ is identical to (7.1), except for the sigmoid replaced by the softplus. In Fig. 7.4, we compare the cases $\lambda = 0.03$ (the solid lines) and $\lambda = 1$ (the dashed lines). Recall that the minimization of the joint loss (3.1) with $\lambda = 1$ is equivalent to maximization of the log-likelihood of the Gaussian distribution, see Example 4.1. In Fig. 7.4, we see that the choice $\lambda = 0.03$ is more optimal. For larger $\lambda$, $\mathcal{N}_q$ learns much faster in the noisy regions compared with the clean regions (cf. Sec. 6.2.2 and Fig. 2.2, right). As a result, it cannot learn properly in the clean regions 1, 3, and 5, especially in region 1.

# 8 Real world data

**Data sets.** We analyzed the following publicly available data sets: Boston House Prices [11] (506 samples, 13 features), Concrete Compressive Strength [45] (1030 samples, 8 features), Combined Cycle Power Plant [19, 41] (9568 samples, 4 features), Yacht Hydrodynamics [6, 37] (308 samples, 6 features), Kinematics of an 8 Link Robot Arm Kin8Nm[3] (8192 samples, 8 feature), and Year Prediction MSD [28] (515345 samples, 90 features). For each data set, a one-dimensional target variable is predicted. Each data set, except for the year prediction MSD, is randomly split into 50 train-test folds with 95% of samples in each train subset. All the measure values reported below are the averages of the respective measure values over 50 folds. For the year prediction MSD, we used a single split recommended in [28]. The data are normalized so that the input features and the targets have zero mean and unit variance in the training set.

**Methods.** We compare the following methods:

1. the maximum likelihood method (ML), in which one maximizes the likelihood of the normal distribution; like in our method, two networks (one predicting the mean and another predicting the variance) are trained simultaneously,

2. Stein variational gradient descent (SVGD) [29]; Bayesian method, in which one uses a particle approximation of the posterior distribution of the weights,

---

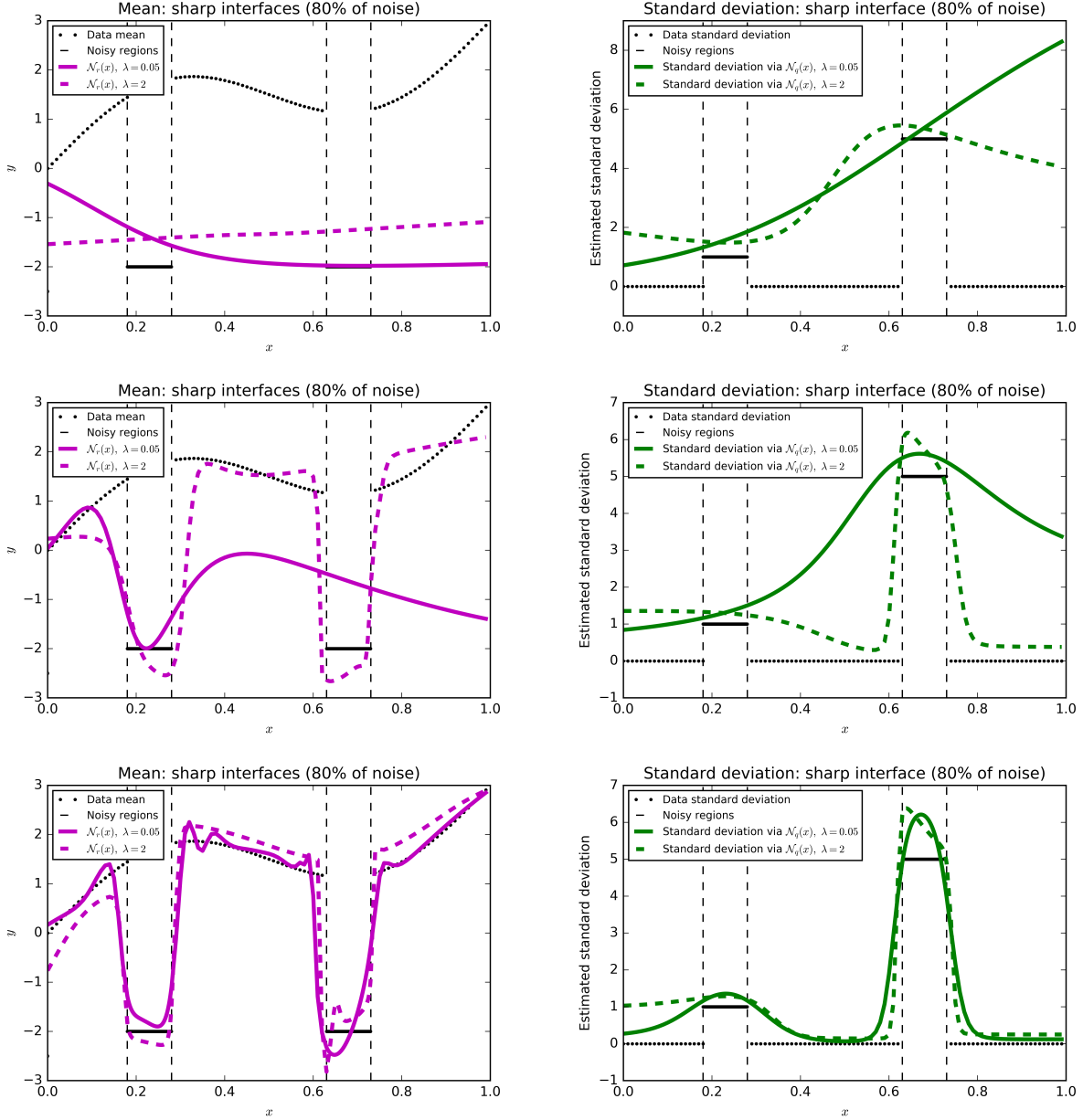[3]http://mldata.org/repository/data/viewslug/regression-datasets-kin8nm/

Figure 7.3: Fit of the networks $\mathcal{N}_r$ (left) and $\mathcal{N}_q$ with the sigmoid output (right). Solid lines correspond to $\lambda = 0.05$ and dashed lines to $\lambda = 2$. Top: early stage. Middle: middle stage. Bottom: final stage.

3. the probabilistic back propagation (PBP) [12]; Bayesian method, in which one minimizes the KL divergence from the exact posterior to the approximating one, using assumed density filtering [36] and expectation-propagation [32] methods,

4. our method (with $\lambda = 0.1, 0.2, 0.5, 1, 2, 5$),

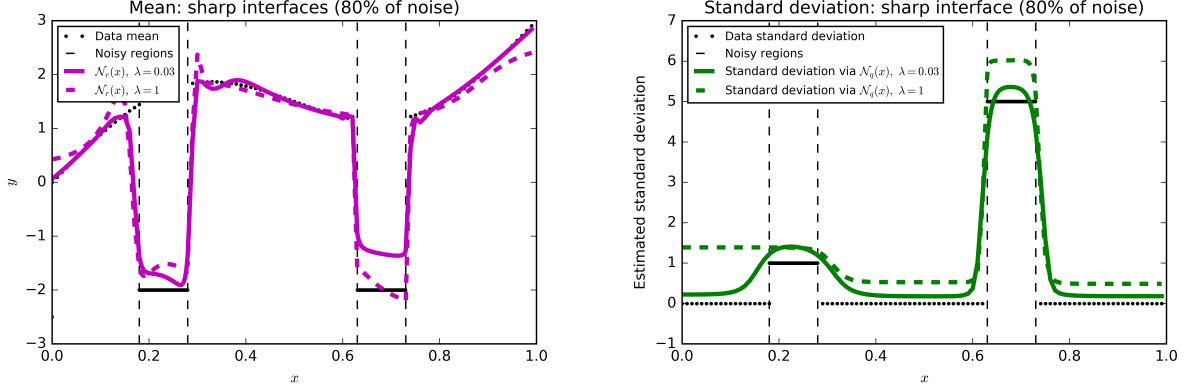5. ensemble of 5 MLs (EnsML) [24, 25],

20

Figure 7.4: Fit of the networks $\mathcal{N}_r$ (left) and $\mathcal{N}_q$ with the softplus output (right). Solid lines correspond to $\lambda = 0.03$ and dashed lines to $\lambda = 1$.

6. ensemble of 5 pairs of networks according to our method.

To show that our approach works for different types of regressor's loss, we optimize the above methods for the root mean squared error (RMSE) and the mean absolute error (MAE). In the case of RMSE/MAE, we use the Gaussian (4.8) / Laplacian (4.10) likelihood in the ML, SVGD, and EnsML, and we use the MSE/MAE regressor's loss in our method. In the case of MAE for the PBP, we still use the Gaussian likelihood because changing the likelihood would require a new approximation of the normalization constant in the assumed density filtering method applied in [12].

In the case of RMSE, our method quantifies uncertainty in terms of the expected squared error (ESA) due to (4.6). The other non-ensemble methods quantify ESA as the variance of the predictive distribution. For the ensemble methods 5 and 6, we use (4.11) for the predictive mean and ESA.

In the case of MAE, our method quantifies uncertainty in terms of the expected absolute error (EAE) due to (4.9). The other non-ensemble methods (except the PBP) quantify EAE as $1/\tau$, where $\tau$ is the parameter in (4.10). For the PBP, we use its predictive Gaussian distribution $\mathcal{N}(y|\mu, V)$ to calculate the expectation of $|y - \mu|$, which yields $\sqrt{2V\pi}$ as EAE. For the ensemble methods, we use (4.12) for the predictive mean and EAE.

Note that the ML and our method estimate aleatoric uncertainty, the Bayesian methods SVGD and PBP estimate epistemic uncertainty, and ensembles take into account both types of uncertainty.

**Architectures.** We use architectures proposed in [12, 24, 25], namely, 1-hidden layer regressors and 1-hidden layer uncertainty quantifiers with ReLU nonlinearities. Each NN contains 50 hidden units for all the data sets, except for MSD, where we use 100 hidden units. For our method, we use the sigmoid output of the uncertainty quantifier (see Secs. 5.1 and 6.1). We refer to Appendix A for the values of hyperparameters that we used for the different methods.

**Measures.** We use two measures to estimate the quality of the fit.

1. The overall error: *RMSE* and *MAE*.

2. The area under the following curve (*AUC*), measuring the trade-off between properly learning the mean and estimating uncertainty. Assume the test set contains $N$ samples. If we are interested in the RMSE, then we order the samples with respect to their predicted ESA. For
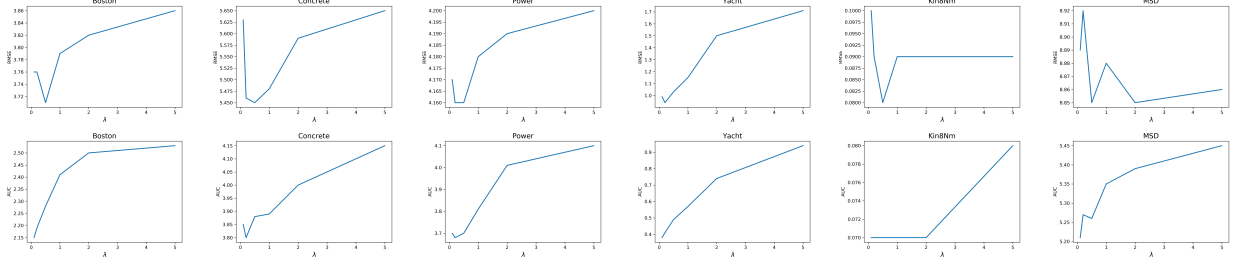
21

Figure 8.1: RMSE (upper row) and AUC (lower row) of our method with MSE regressor's loss versus $\lambda$ for different data sets.

| | **Boston** | | **Concrete** | | **Power** | |
| | RMSE | AUC | RMSE | AUC | RMSE | AUC |
|---|---|---|---|---|---|---|
| ML | 3.79±1.44 | 2.40±0.59 | 5.60±0.63 | 4.00±0.59 | 4.09±0.31 | 3.75±0.24 |
| SVGD | **3.10±1.14** | 2.34±0.54 | 6.11±0.59 | 4.93±0.86 | 4.22±0.30 | 4.26±0.28 |
| PBP | **3.54±1.29** | 2.28±0.47 | 5.58±0.60 | 4.66±0.64 | **4.09±0.26** | 3.86±0.26 |
| Our method | 3.76±1.50 | 2.15±0.51 | 5.46±0.63 | 3.80±0.60 | 4.13±0.30 | **3.63±0.24** |
| EnsML | **3.60±1.36** | 2.35±0.49 | **4.76±0.64** | **3.34±0.44** | **4.05±0.31** | **3.71±0.24** |
| Our ensemble | **3.14±1.31** | **1.86±0.46** | **4.71±0.65** | **3.17±0.38** | **4.13±0.30** | **3.64±0.26** |

Table 8.1: Values of RMSE and AUC for the different methods on the Boston, Concrete, and Power data sets.

each $n = 0, \ldots, N - 1$, we remove $n$ samples with the highest ESA and calculate the RMSE for the remaining $N - n$ samples. We denote the result by $\mathrm{RMSE}(n)$ and plot it versus $n$ as a continuous piecewise linear curve. The AUC is the area under this curve normalized by $N - 1$:

$$\mathrm{AUC} = \frac{1}{N - 1} \sum_{n=0}^{N-2} \frac{\mathrm{RMSE}(n) + \mathrm{RMSE}(n + 1)}{2}.$$

When we use MAE instead of RMSE, the AUC is calculated similarly with $\mathrm{MAE}(n)$ instead of $\mathrm{RMSE}(n)$.

**RMSE results.** Figure 8.1 shows the dependence of RMSE and AUC in our method on $\lambda$ for different data sets. We see that typically our method performs better when $\lambda < 1$, which is due to a better fit on clean regions (see Sec. 6). Tables 8.1 and 8.2 show RMSE and AUC for the different methods[4]. The values with the best mean and the values that are not significantly different from those with the best mean (due to the two-tailed paired difference test with $p = 0.05$) are marked in bold. Our ensemble achieves the best or not significantly different from the best RMSE and AUC on all the data sets (except Yacht, where it is second best after the EnsML). Further, note that our non-ensemble networks are the best among the other non-ensemble methods, again except for the Yacht data set. Figure 8.2 shows the curves $\mathrm{RMSE}(n)$ for the different methods and data sets. The curve corresponding to our ensemble is typically located below the others.

**MAE results.**

---

[4]We were not able to fit the PBP for the largest MSD data set.

|  | Yacht | | Kin8nm | | MSD | |
|---|---|---|---|---|---|---|
|  | RMSE | AUC | RMSE | AUC | RMSE | AUC |
| ML | 0.76±0.38 | 0.25±0.08 | 0.086±0.006 | 0.074±0.004 | 9.09±NA | 5.40±NA |
| SVGD | 2.00±0.72 | 1.20±0.49 | 0.152±0.007 | 0.129±0.008 | 9.09±NA | 8.38±NA |
| PBP | 1.09±0.35 | 0.64±0.19 | 0.098±0.005 | 0.081±0.005 | NA | NA |
| Our method | 0.94±0.42 | 0.41±0.13 | 0.085±0.005 | 0.067±0.003 | 8.87±NA | 5.23±NA |
| EnsML | **0.48±0.24** | **0.21±0.07** | 0.078±0.004 | 0.062±0.003 | 8.99±NA | 5.32±NA |
| Our ensemble | 0.73±0.33 | 0.34±0.15 | **0.077±0.004** | **0.058±0.003** | **8.86±NA** | **5.13±NA** |

Table 8.2: Values of RMSE and AUC for the different methods on the Yacht, Kin8nm, and MSD data sets.



Figure 8.2: The curves RMSE($n$) for the different methods.

Tables 8.3 and 8.4 show MAE and AUC for the different methods. Our ensemble achieves the best or not significantly different from the best MAE and AUC on all the data sets (except for MAE and AUC on Yacht and MAE on Kin8nm, where it is second best after the EnsML). Our non-ensemble networks are the best among the other non-ensemble methods, with the same exceptions as for the ensembles. Figure 8.3 shows the curves MAE($n$) for the different methods.

# 9 Conclusion

We introduced a general approach to uncertainty quantification in artificial NNs, based on a specific joint loss (3.1) for two NNs: one for the regression and another for the uncertainty quantification. We analyzed in detail how the functions $f$ and $g$ and the hyperparameter $\lambda$ in the loss affect the learning process. We showed that the uncertainty quantifier provides an estimate of how certain

|  | **Boston** | | **Concrete** | | **Power** | |
|---|---|---|---|---|---|---|
|  | MAE | AUC | MAE | AUC | MAE | AUC |
| ML | 2.34±0.57 | 1.68±0.46 | 4.09±0.49 | 2.69±0.38 | 3.08±0.12 | 2.72±0.15 |
| SVGD | 2.14±0.53 | 1.61±0.33 | 4.82±0.70 | 3.29±0.59 | 3.23±0.12 | 3.22±0.18 |
| PBP | 2.21±0.44 | 1.65±0.32 | 4.21±0.64 | 3.59±0.60 | 3.19±0.15 | 3.11±0.19 |
| Our method | **2.20±0.55** | **1.39±0.39** | 4.05±0.48 | 2.57±0.38 | **3.09±0.13** | **2.66±0.15** |
| EnsML | 2.17±0.55 | 1.67±0.35 | 3.63±0.61 | 2.32±0.39 | **3.07±0.13** | 2.71±0.15 |
| Our ensemble | **2.06±0.57** | **1.46±0.38** | **3.56±0.63** | **2.19±0.40** | **3.07±0.13** | **2.63±0.16** |

Table 8.3: Values of MAE and AUC for the different methods on the Boston, Concrete, and Power data sets.

|  | **Yacht** | | **Kin8nm** | | **MSD** | |
|---|---|---|---|---|---|---|
|  | MAE | AUC | MAE | AUC | MAE | AUC |
| ML | 0.50±0.25 | 0.18±0.06 | 0.069±0.005 | 0.053±0.003 | 6.11±NA | 3.78±NA |
| SVGD | 1.35±0.54 | 0.46±0.14 | 0.116±0.006 | 0.089±0.005 | 6.12±NA | 5.66±NA |
| PBP | 0.73±0.19 | 0.50±0.15 | 0.076±0.004 | 0.064±0.004 | NA | NA |
| Our method | 0.57±0.20 | 0.22±0.07 | 0.072±0.00 | 0.052±0.003 | 5.89±NA | 3.65±NA |
| EnsML | **0.33±0.14** | **0.13±0.05** | **0.063±0.003** | **0.048±0.003** | 6.04±NA | 3.69±NA |
| Our ensemble | 0.42±0.17 | 0.15±0.04 | 0.069±0.004 | **0.048±0.002** | **5.84±NA** | **3.55±NA** |

Table 8.4: Values of MAE and AUC for the different methods on the Yacht, Kin8nm, and MSD data sets.



Figure 8.3: The curves MAE($n$) for the different methods.

the predictions are in terms of *any* regressor's loss and without the knowledge of the underlying distribution (whose form may even vary in different regions of the input space). Moreover, we explained how the presence of the uncertainty quantifier improves the predictions of the regressor and of NNs based on the classical likelihood maximization. We showed that the crucial role here is

played by the hyperparameter $\lambda$, which allows for better fits on clean regions. Finally, we compared our NNs and their ensemble counterparts with the other NN methods, using two measures: overall error (RMSE and MAE) and AUC. We showed that our approach typically yields the best or not significantly different from the best results.

It is also worth mentioning that one could fit a regressor first and then quantify the uncertainty of its predictions by training only the neural network $\mathcal{N}_q$ with the loss (3.1). Now the function $\mathcal{L}_r(y^i, y_r^i)$ need not coincide with the loss function that was used for training the regressor, but can represent *any* error whose local average we want to estimate by $\mathcal{N}_q$. With this modification, one loses the benefit of the joined training of $\mathcal{N}_r$ and $\mathcal{N}_q$ that may improve regressor's predictions in the clean regions, but, on the other hand, one can choose any type of regressor, not only a neural network. For the uncertainty quantification, one still has the full freedom in the choice of the functions $f$ and $g$, while the parameter $\lambda$ still has the same influence on the learning dynamics of $\mathcal{N}_q$ as in Sec. 6. In both settings, it would be interesting to analyze other choices of the functions $f$ and $g$ and to develop an automatic procedure that could choose an optimal $\lambda$ and properly adjust it during the learning process.

# References

[1] C. Bishop. Mixture Density Networks. Neural Computing Research Group Report: NCRG/94/004. 1994.

[2] C. Bishop. Pattern Recognition and Machine Learning. Springer, 2006.

[3] C. Blundell, J. Cornebise, K. Kavukcuoglu and D. Wierstra, Weight uncertainty in neural networks, In Proceedings of the 32nd International Conference on Machine Learning, 7-9 July 2015, Lille, France, JMLR: W&CP **37**, 1613 (2016).

[4] Y. Gal. Uncertainty in Deep Learning. PhD thesis, University of Cambridge, 2016.

[5] Y. Gal and Z. Ghahramani, Dropout as a bayesian approximation: Representing model uncertainty in deep learning, In Proceedings of the 33rd International Conference on Machine Learning, 20-22 June 2016, New York, New York, USA, JMLR: W&CP **48**, 1050 (2016).

[6] J. Gerritsma, R. Onnink, and A. Versluis. Geometry, resistance and stability of the delft systematic yacht hull series. In International Shipbuilding Progress, **28** (1981), 276–297.

[7] I. Goodfellow, Y. Bengio, A. Courville, Deep Learning, MIT Press, 2016.

[8] J. G. De Gooijer, R. J. Hyndman, 25 years of time series forecasting, Int. J. Forecast., **22** (2006) 443–473.

[9] A. Graves. Generating sequences with recurrent neural networks. arXiv:1308.0850 [cs.NE] (2014).

[10] P. Gurevich and H. Stuke: Gradient conjugate priors and deep neural networks, arXiv:1802.02643 [math.ST] (2018).

[11] D. Harrison, D. L. Rubinfeld, Hedonic prices and the demand for clean air, J. Environ. Economics and Management, **5** (1978), 81–102.

[12] J. M. Hernández-Lobato, R. P. Adams. Probabilistic backpropagation for scalable learning of bayesian neural networks. arXiv:1502.05336 [stat.ML] (2015).

[13] J. M. Hernández-Lobato, Y. Li, M. Rowland, T. Bui, D. Hernández-Lobato and R. Turner, Black-box alpha divergence minimization, In Proceedings of the 33rd International Conference on Machine Learning, 20-22 June 2016, New York, New York, USA, JMLR: W&CP **48**, 1511 (2016).

[14] G. Hinton, D. V. Camp. Keeping neural networks simple by minimizing the description length of the weights, In: Proceedings of the Sixth Annual Conference on Computational Learning Theory (1993), 5–13.

[15] G. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, R. Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. arXiv:1207.0580 (2012).

[16] M. A. Hussain, Review of the applications of neural networks in chemical process control – simulation and online implementation, Artif. Intell. Eng., **13**, no 1 (1999), 55–68.

[17] J. T. G. Hwang, A. A. Ding, Prediction intervals for artificial neural networks, J. Amer. Stat. Assoc, **92**, no 438 (1997), 748–757.

[18] P. Jylänki, A. Nummenmaa, A. Vehtari. Expectation propagation for neural networks with sparsity-promoting priors. The Journal of Machine Learning Research, **15**(1) (2014), 1849–1901.

[19] H. Kaya, P. Tüfekci , S. F. Gürgen: Local and global learning methods for predicting power of a combined gas and steam turbine, Proceedings of the International Conference on Emerging Trends in Computer and Electronics Engineering ICETCEE (2012), 13–18.

[20] A. Kendall and Y. Gal. What uncertainties do we need in Bayesian deep learning for computer Vision?, 31st Conference on Neural Information Processing Systems (NIPS 2017).

[21] D. P. Kingma, T. Salimans and M. Welling, Variational Dropout and the Local Reparameterization Trick, 29th Conference on Neural Information Processing Systems (NIPS 2015), 7-12 December 2015, Palais des Congrès de Montréal, Montréal, Canada, Advances in Neural Information Processing Systems **28**, 2575 (2015).

[22] M. Krzywinski, N. Altman. Points of significance: importance of being uncertain. Nature methods, **10**(9) (2013).

[23] M. Kuhn, J. Johnson, Applied Predictive Modeling, Springer, New York, 2013.

[24] B. Lakshminarayanan, A. Pritzel, C. Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. Workshop on Bayesian Deep Learning, NIPS 2016, Barcelona, Spain.

[25] B. Lakshminarayanan, A. Pritzel, C. Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. 31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

[26] Y. Li and Y. Gal, Dropout inference in Bayesian neural networks with alpha-divergences, In Proceedings of the 34th International Conference on Machine Learning, 6-11 August 2017, International Convention Centre, Sydney, Australia, PMLR **70**, 2052 (2017)

[27] Y. Li and R. Turner, Rényi divergence variational inference, 30th Conference on Neural Information Processing Systems (NIPS 2016), 5-10 December 2016, Centre Convencions Internacional Barcelona, Barcelona, Spain, Advances in Neural Information Processing Systems **29** (2016).

[28] M. Lichman. UCI Machine Learning Repository [http://archive.ics.uci.edu/ml]. Irvine, CA: University of California, School of Information and Computer Science (2013).

[29] Q. Liu, D. Wang, Stein variational gradient descent: a general purpose Bayesian inference algorithm. 30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain (2016).

[30] C. Louizos and M. Welling, Multiplicative normalizing flows for variational Bayesian neural networks, In Proceedings of the 34th International Conference on Machine Learning, 6-11 August 2017, International Convention Centre, Sydney, Australia, PMLR **70**, 2218 (2017)

[31] D. J. C. MacKay, Practical Bayesian framework for backpropagation networks, Neural Computation, **4** (1992), 448–472.

[32] T. Minka. A family of algorithms for approximate Bayesian inference. PhD thesis, Massachusetts Institute of Technology, 2001.

[33] P. Myshkov, S. Julier. Posterior distribution analysis for Bayesian inference in neural networks. Workshop on Bayesian Deep Learning, NIPS 2016, Barcelona, Spain.

[34] R. M. Neal, Bayesian Learning for Neural Networks, Ph.D. Thesis, Dept. of Computer Science, University of Toronto, 1995.

[35] D. A. Nix and A. S. Weigend, Estimating the mean and variance of the target probability distribution, in Proc. IEEE Int. Conf. Neural Netw., **1**. Orlando, FL, Jun.Jul. (1994), 55–60.

[36] M. Opper, O. Winther. A Bayesian approach to online learning. On-line Learning in Neural Networks, ed. D. Saad (1998), 363–378.

[37] I. Ortigosa, R. Lopez and J. Garcia. A neural networks approach to residuary resistance of sailing yachts prediction. In Proceedings of the International Conference on Marine Engineering MARINE 2007, 2007.

[38] C. E. Rasmussen and C. K. I. Williams: Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning), The MIT Press Cambridge (2006).

[39] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever and R. Salakhutdinov: Dropout: A simple way to prevent neural networks from overfitting, The Journal of Machine Learning Research **15**(1), 1929 (2014)

| | Boston | | | | Concrete | | | | Power | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\lambda$ | Learning rate | Dropout | Number of epochs | $\lambda$ | Learning rate | Dropout | Number of epochs | $\lambda$ | Learning rate | Dropout | Number of epochs |
| **ML** | | 0.00004 | 0.4 | 500 | | 0.0001 | 0.1 | 700 | | 0.00005 | 0 | 150 |
| **Our method** | 0.1 | 0.0008 | 0.4 | 500 | 0.2 | 0.0002 | 0.15 | 700 | 0.2 | 0.0003 | 0 | 80 |

Table A.1: Hyperparameters for the ML, our method, and the respective ensembles optimising RMSE on the Boston, Concrete, and Power data sets.

[40] N. Tishby, E. Levin, S. Solla. Consistent inference of probabilities in layered networks: Predictions and generalizations. In: International Joint Conference on Neural Networks (1989), 403–409.

[41] P. Tüfekci, Prediction of full load electrical power output of a base load operated combined cycle power plant using machine learning methods, International Journal of Electrical Power and Energy Systems, **60** (2014), 126–140.

[42] R. D. De Veaux, J. Schumi, J. Schweinsberg, L. H. Ungar, Prediction intervals for neural networks via nonlinear regression, Technometrics, **40**, no. 4 (1998), 273–282.

[43] M. Welling, Y. W. Teh, Bayesian learning via stochastic gradient Langevin dynamics, In Proceedings of the 28th International Conference on Machine Learning, 28 June - 2 July 2011, Bellevue, Washington, USA, 681 (2011).

[44] C. J. Wild and G. A. F. Seber, Nonlinear Regression. New York, Wiley, 1989.

[45] I-C. Yeh. Modeling of strength of high performance concrete using artificial neural networks, Cement and Concrete Research, **28**, No. 12 (1998), 1797–1808.

[46] H. Zen, A. Senior. Deep mixture density networks for acoustic modeling in statistical parametric speech synthesis. Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP), IEEE (2014), 3872–3876.

# A Hyperparameters

When we fit different methods on the real world data sets, we normalize them so that the input features and the targets have zero mean and unit variance in the training set. We used minibatch 5 on Boston, Concrete, and Yacht, minibatch 10 on Power and Kin8nm, and minibatch 5000 on MSD. For the ML and for our method, we used the dropout regularization [15] during training, but not for the prediction, and we did not regularize the ensembles. We used Nesterov momentum (with momentum 0.9) optimizer for fitting the ML, our method and the ensembles. We performed a grid search for the learning rate and the dropout rate for the non-ensemble methods. For the ensemble methods, we used the same learning rate as for the corresponding individual predictors. The parameters yielding the best AUC in case of RMSE are presented in Tables A.1 and A.2, and the parameters yielding the best AUC in case of MAE are presented in Tables A.3 and A.4. For the SVGD and PBP, we used the hyperparameters default settings in the authors' code[5].

---

[5]See        https://github.com/DartML/Stein-Variational-Gradient-Descent        for        the        SVGD        and https://github.com/HIPS/Probabilistic-Backpropagation for the PBP.

| | Yacht | | | | Kin8nm | | | | MSD | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\lambda$ | Learning rate | Dropout | Number of epochs | $\lambda$ | Learning rate | Dropout | Number of epochs | $\lambda$ | Learning rate | Dropout | Number of epochs |
| **ML** | | 0.0001 | 0.1 | 2000 | | 0.00005 | 0 | 200 | | 0.005 | 0.1 | 50 |
| **Our method** | 0.2 | 0.0004 | 0 | 500 | 0.5 | 0.0002 | 0 | 300 | 0.1 | 0.2 | 0.1 | 50 |

Table A.2: Hyperparameters for the ML, our method, and the respective ensembles optimising RMSE on the Yacht, Kin8nm, and MSD data sets.

| | Boston | | | | Concrete | | | | Power | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\lambda$ | Learning rate | Dropout | Number of epochs | $\lambda$ | Learning rate | Dropout | Number of epochs | $\lambda$ | Learning rate | Dropout | Number of epochs |
| **ML** | | 0.00003 | 0.4 | 600 | | 0.00003 | 0.1 | 800 | | 0.00005 | 0 | 150 |
| **Our method** | 0.2 | 0.0002 | 0.4 | 2000 | 0.2 | 0.0003 | 0.1 | 600 | 1 | 0.0001 | 0 | 200 |

Table A.3: Hyperparameters for the ML, our method, and the respective ensembles optimising MAE on the Boston, Concrete, and Power data sets.

| | Yacht | | | | Kin8nm | | | | MSD | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\lambda$ | Learning rate | Dropout | Number of epochs | $\lambda$ | Learning rate | Dropout | Number of epochs | $\lambda$ | Learning rate | Dropout | Number of epochs |
| **ML** | | 0.0001 | 0.1 | 1000 | | 0.00005 | 0 | 400 | | 0.01 | 0.1 | 50 |
| **Our method** | 0.2 | 0.0004 | 0 | 500 | 0.2 | 0.0002 | 0 | 400 | 0.5 | 0.1 | 0.1 | 70 |

Table A.4: Hyperparameters for the ML, our method, and the respective ensembles optimising MAE on the Yacht, Kin8nm, and MSD data sets.