
Uncertainty Decomposition in Bayesian Neural Networks with Latent Variables

Stefan Depeweg^{1,2} José Miguel Hernández-Lobato³ Finale Doshi-Velez⁴ Steffen Udluft¹

Abstract

Bayesian neural networks (BNNs) with latent variables are probabilistic models which can automatically identify complex stochastic patterns in the data. We describe and study in these models a decomposition of predictive uncertainty into its epistemic and aleatoric components. First, we show how such a decomposition arises naturally in a Bayesian active learning scenario by following an information theoretic approach. Second, we use a similar decomposition to develop a novel risk sensitive objective for safe reinforcement learning (RL). This objective minimizes the effect of model bias in environments whose stochastic dynamics are described by BNNs with latent variables. Our experiments illustrate the usefulness of the resulting decomposition in active learning and safe RL settings.

1. Introduction

Recently, there has been an increased interest in Bayesian neural networks (BNNs) and their possible use in reinforcement learning (RL) problems (Gal et al., 2016; Blundell et al., 2015; Houthoofd et al., 2016). In particular, recent work has extended BNNs with a latent variable model to describe complex stochastic functions (Depeweg et al., 2016; Moerland et al., 2017). The proposed approach enables the automatic identification of arbitrary stochastic patterns such as multimodality and heteroskedasticity, without having to manually incorporate these into the model.

In model-based RL, the aforementioned BNNs with latent variables can be used to describe complex stochastic dynamics. The BNNs encode a probability distribution over stochastic functions, with each function serving as an estimate of the ground truth continuous Markov Decision Process (MDP). Such probability distribution can then be

used for policy search, by finding the optimal policy with respect to state trajectories simulated from the model. The BNNs with latent variables produce improved probabilistic predictions and these result in better performing policies (Depeweg et al., 2016; Moerland et al., 2017).

We can identify two distinct forms of uncertainties in the class of models given by BNNs with latent variables. As described by Kendall & Gal (2017), "Aleatoric uncertainty captures noise inherent in the observations. On the other hand, epistemic uncertainty accounts for uncertainty in the model." In particular, epistemic uncertainty arises from our lack of knowledge of the values of the synaptic weights in the network, whereas aleatoric uncertainty originates from our lack of knowledge of the value of the latent variables.

In the domain of model-based RL the epistemic uncertainty is the source of model bias (or representational bias, see e.g. Joseph et al. (2013)). When there is high discrepancy between model and real-world dynamics, policy behavior may deteriorate. In analogy to the principle that "a chain is only as strong as its weakest link" a drastic error in estimating the ground truth MDP at a single transition step can render the complete policy useless (see e.g. Schneegass et al. (2008)).

In this work we address the decomposition of the uncertainty present in the predictions of BNNs with latent variables into its epistemic and aleatoric components. We show the usefulness of such decomposition in two different domains: active learning and risk-sensitive RL.

First we consider an active learning scenario with stochastic functions. We derive an information-theoretic objective that decomposes the entropy of the predictive distribution of BNNs with latent variables into its epistemic and aleatoric components. By building on that decomposition, we then investigate safe RL using a risk-sensitive criterion (García & Fernández, 2015) which focuses only on risk related to model bias, that is, the risk of the policy performing at test time significantly different from at training time. The proposed criterion quantifies the amount of epistemic uncertainty (model bias risk) in the model's predictive distribution and ignores any risk stemming from the aleatoric uncertainty. Our experiments show that, by using this risk-sensitive criterion, we are able to find policies that, when evaluated on the

¹Siemens AG ²Technical University of Munich ³University of Cambridge ⁴Harvard University. Correspondence to: Stefan Depeweg <stefan.depeweg@siemens.com>.

ground truth MDP, are safe in the sense that on average they do not deviate significantly from the performance predicted by the model at training time.

We focus on the off-policy batch RL scenario, in which we are given an initial batch of data from an already-running system and are asked to find a better policy. Such scenarios are common in real-world industry settings such as turbine control, where exploration is restricted to avoid possible damage to the system.

2. Bayesian Neural Networks with Latent Variables

Given data $\mathcal{D} = \{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N$, formed by feature vectors $\mathbf{x}_n \in \mathbb{R}^D$ and targets $\mathbf{y}_n \in \mathbb{R}^K$, we assume that $\mathbf{y}_n = f(\mathbf{x}_n, z_n; \mathcal{W}) + \epsilon_n$, where $f(\cdot, \cdot; \mathcal{W})$ is the output of a neural network with weights \mathcal{W} and K output units. The network receives as input the feature vector \mathbf{x}_n and the latent variable $z_n \sim \mathcal{N}(0, \gamma)$. The activation functions for the hidden layers are rectifiers: $\varphi(x) = \max(x, 0)$. The activation functions for the output layers are the identity function: $\varphi(x) = x$. The network output is corrupted by the additive noise variable $\epsilon_n \sim \mathcal{N}(\mathbf{0}, \Sigma)$ with diagonal covariance matrix Σ . The role of the latent variable z_n is to capture unobserved stochastic features that can affect the network’s output in complex ways. Without z_n , randomness is only given by the additive Gaussian observation noise ϵ_n , which can only describe limited stochastic patterns. The network has L layers, with V_l hidden units in layer l , and $\mathcal{W} = \{\mathbf{W}_l\}_{l=1}^L$ is the collection of $V_l \times (V_{l-1} + 1)$ weight matrices. The $+1$ is introduced here to account for the additional per-layer biases. We approximate the exact posterior distribution $p(\mathcal{W}, \mathbf{z} | \mathcal{D})$ with the factorized Gaussian distribution

$$q(\mathcal{W}, \mathbf{z}) = \left[\prod_{l=1}^L \prod_{i=1}^{V_l} \prod_{j=1}^{V_{l-1}+1} \mathcal{N}(w_{ij,l} | m_{ij,l}^w, v_{ij,l}^w) \right] \quad (1)$$

$$\left[\prod_{n=1}^N \mathcal{N}(z_n | m_n^z, v_n^z) \right]. \quad (2)$$

The parameters $m_{ij,l}^w$, $v_{ij,l}^w$ and m_n^z , v_n^z are determined by minimizing a divergence between $p(\mathcal{W}, \mathbf{z} | \mathcal{D})$ and the approximation q . For more detail the reader is referred to the work of [Hernández-Lobato et al. \(2016\)](#); [Depeweg et al. \(2016\)](#). In all our experiments we use black-box α -divergence minimization with $\alpha = 1.0$, as it seems to produce a better decomposition of uncertainty into its epistemic and aleatoric components, although further studies are needed to strengthen this claim.

The described BNNs with latent variables can describe complex stochastic patterns while at the same time account for model uncertainty. They achieve this by jointly learning

$q(\mathbf{z})$, which captures the specific values of the latent variables in the training data, and $q(\mathcal{W})$, which represents any uncertainty about the model parameters. The result is a principled Bayesian approach for inference of stochastic functions.

3. Active Learning of Stochastic Functions

Active learning is the problem of choosing which data points to incorporate next into the training data so that the resulting gains in predictive performance are as high as possible. In this section, we derive a Bayesian active learning procedure for stochastic functions. This procedure illustrates how to separate two sources of uncertainty, that is, aleatoric and epistemic, in the predictive distribution of BNNs with latent variables.

Within a Bayesian setting, active learning can be formulated as choosing data based on the expected reduction in entropy of the posterior distribution ([MacKay, 1992](#)). [Hernández-Lobato & Adams \(2015\)](#) apply this entropy-based approach to scalable BNNs. In ([Houthoofd et al., 2016](#)), the authors use a similar approach as an exploration scheme in RL problems in which a BNN is used to represent current knowledge about the transition dynamics. These previous works only assume additive Gaussian noise and, unlike the BNNs with latent variables from Section 2, they cannot capture complex stochastic patterns.

We start by deriving the expected reduction in entropy in BNNs with latent variables. We assume a scenario in which a BNN with latent variables has been fitted to a batch of data $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$ to produce a posterior approximation $q(\mathcal{W}, \mathbf{z})$. We now want to estimate the expected reduction in posterior entropy for \mathcal{W} when a particular data point \mathbf{x} is incorporated in the training set. The expected reduction in entropy is

$$\mathbf{H}(\mathcal{W} | \mathcal{D}) - \mathbf{E}_{y | \mathbf{x}, \mathcal{D}} \left[\mathbf{H}(\mathcal{W} | \mathcal{D} \cup \{\mathbf{x}, y\}) \right] \quad (3)$$

$$= \mathbf{H}(\mathcal{W}) - \mathbf{H}(\mathcal{W} | y) \quad (4)$$

$$= \mathbf{I}(\mathcal{W}; y) \quad (5)$$

$$= \mathbf{H}(y) - \mathbf{H}(y | \mathcal{W}) \quad (6)$$

$$= \mathbf{H} \left[\int_{\mathcal{W}, z} p(y | \mathcal{W}, \mathbf{x}, z) p(z) p(\mathcal{W} | \mathcal{D}) dz d\mathcal{W} \right] \quad (7)$$

$$- \mathbf{E}_{\mathcal{W} | \mathcal{D}} \left[\mathbf{H} \left(\int_z p(y | \mathcal{W} = \mathcal{W}_i, \mathbf{x}, z) p(z) dz \right) \right]$$

where $\mathbf{H}(\cdot)$ denotes the entropy of a random variable and $\mathbf{H}(\cdot | \cdot)$ and $\mathbf{I}(\cdot; \cdot)$ denote the conditional entropy and the mutual information between two random variables. In (6) and (7) we see that the active learning objective is given by the difference of two terms. The first term is the entropy of the predictive distribution, that is, $\mathbf{H}(y)$. The sec-

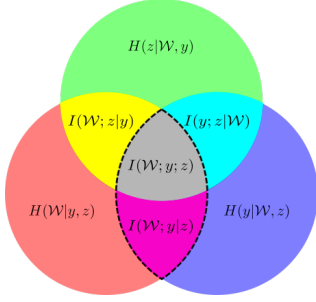


Figure 1. Information Diagram illustrating quantities of entropy with three variables. The area surrounded by a dashed line indicates reduction in entropy given by equation (7). The conditioning to \mathbf{x} is omitted for readability.

ond term, that is, $H(y|\mathcal{W})$, is a conditional entropy. To compute this term, we average across $\mathcal{W}_i \sim q(\mathcal{W})$ the entropy $H[\int_z p(y|\mathcal{W} = \mathcal{W}_i, \mathbf{x}, z)p(z)dz]$. As shown in this expression, the randomness in $p(y|\mathcal{W} = \mathcal{W}_i, \mathbf{x})$ has its origin in the latent variable z (and the constant output noise $\epsilon \sim \mathcal{N}(0, \Sigma)$ which is not shown here). Therefore, this second term can be interpreted as the ‘aleatoric uncertainty’ present in the predictive distribution, that is, the average entropy of y that originates from the latent variable z and not from the uncertainty about \mathcal{W} . We can refer to the whole objective function in (7) as an estimate of the epistemic uncertainty: the full predictive uncertainty about y given \mathbf{x} minus the corresponding aleatoric uncertainty.

The previous decomposition is also illustrated by the information diagram from Figure 1. The entropy of the predictive distribution is composed of the blue, cyan, grey and pink areas. The blue area is constant: when both \mathcal{W} and z are determined, the entropy of y is constant and given by the entropy of the additive Gaussian noise $\epsilon \sim \mathcal{N}(0, \Gamma)$. $H(y|\mathcal{W})$ is given by the light and dark blue areas. The reduction in entropy is therefore obtained by the grey and pink areas.

The quantity in equation (7) can be approximating using standard entropy estimators, e.g. nearest-neighbor methods (Kozachenko & Leonenko, 1987; Kraskov et al., 2004; Gao et al., 2016). For that, we repeatedly sample \mathcal{W} and z and do forward passes through the neural network to sample y . The resulting samples of y can then be used to approximate the respective entropies for each \mathbf{x} using the nearest-neighbor approach:

$$H(y|\mathbf{x}, \mathcal{D}) - \mathbf{E}_{\mathcal{W}|\mathcal{D}} \left[H\left(\int_z p(y|\mathcal{W}, \mathbf{x}, z)p(z)dz\right) \right] \\ \approx \hat{H}(y_1, \dots, y_L) - \frac{1}{M} \sum_{i=1}^M \left[\hat{H}(y_1^{\mathcal{W}_i}, \dots, y_L^{\mathcal{W}_i}) \right]. \quad (8)$$

where $\hat{H}(\cdot)$ computes the nearest-neighbor estimate of the entropy given an empirical sample of points, $y_1, \dots, y_L \sim$

$p(y|\mathbf{x}, \mathcal{D})$, $\mathcal{W}_1, \dots, \mathcal{W}_M \sim q(\mathcal{W})$ and $y_1^{\mathcal{W}_i}, \dots, y_L^{\mathcal{W}_i} \sim p(y|\mathbf{x}, \mathcal{D}, \mathcal{W} = \mathcal{W}_i)$ for $i = 1, \dots, M$.

3.1. Toy Problems

We now will illustrate the active learning procedure described in the previous section on two toy examples. In each problem we will first train a BNN with 2 layers and 20 units in each layer on the available data. Afterwards, we approximate the information-theoretic measures as outlined in the previous section

We first consider a toy problem given by a regression task with heteroskedastic noise. For this, we define the stochastic function $y = 7 \sin(x) + 3|\cos(x/2)|\epsilon$ with $\epsilon \sim \mathcal{N}(0, 1)$. The data availability is limited to specific regions of x . In particular, we sample 750 values of x from a mixture of three Gaussians with mean parameters $\{\mu_1 = -4, \mu_2 = 0, \mu_3 = 4\}$, variance parameters $\{\sigma_1 = \frac{2}{5}, \sigma_2 = 0.9, \sigma_3 = \frac{2}{5}\}$ and with each Gaussian component having weight equal to $1/3$ in the mixture. Figure 2a shows the raw data. We have lots of points at both borders of the x axis and in the center, but little data available in between.

Figure 2 visualizes the respective quantities. We see that the BNN with latent variables does an accurate decomposition of its predictive uncertainty between epistemic uncertainty and aleatoric uncertainty: the reduction in entropy approximation, as shown in Figure 2f, seems to be inversely proportional to the density used to sample the data (shown in Figure 2b). This makes sense, since in this toy problem the most informative data points are expected to be located in regions where data is scarce. Note that, in more complicated settings, the most informative data points may not satisfy this property.

Next we consider a toy problem given by a regression task with bimodal data. We define $x \in [-0.5, 2]$ and $y = 10 \sin(x) + \epsilon$ with probability 0.5 and $y = 10 \cos(x) + \epsilon$, otherwise, where $\epsilon \sim \mathcal{N}(0, 1)$ and ϵ is independent of x . The data availability is not uniform in x . In particular we sample 750 values of x from an exponential distribution with $\lambda = 0.5$

Figure 3 visualizes the respective quantities. The predictive distribution shown in Figure 3c suggests that the BNN has learned the bimodal structure in the data. The predictive distribution appears to get increasingly ‘washed out’ as we increase x . This increase in entropy as a function of x is shown in Figure 3d. The conditional entropy $H(y|\mathcal{W})$ of the predictive distribution shown in Figure 3e appears to be symmetric around $x = 0.75$. This suggest that the BNN has correctly learned to separate the aleatoric component from the full uncertainty for the problem at hand: the ground truth function is symmetric around $x = 0.75$ at which point it changes from a bimodal to a unimodal stochastic function.

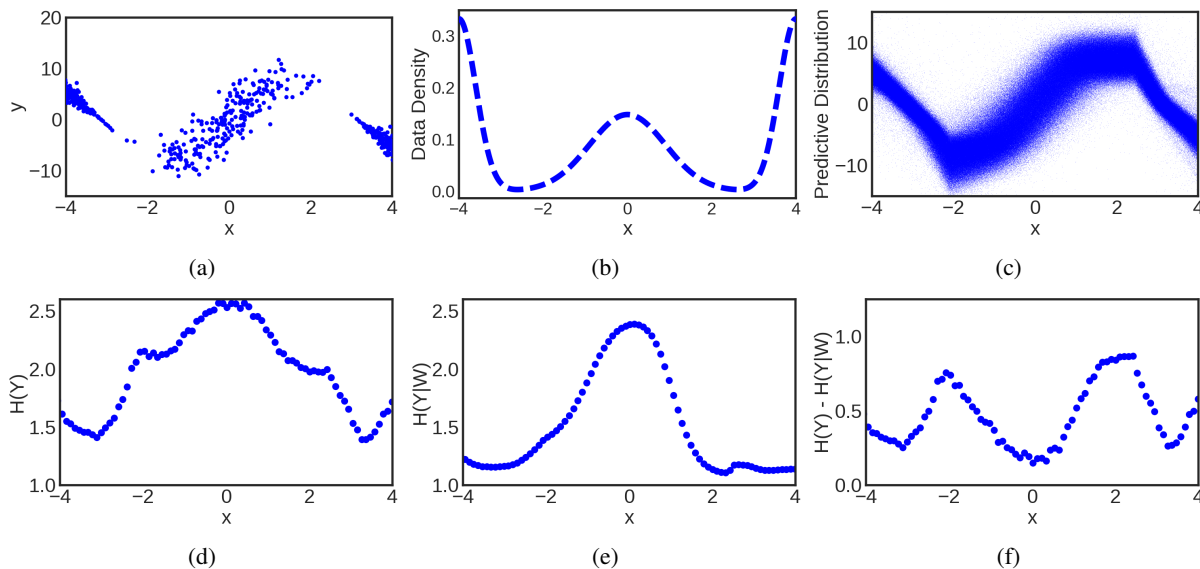


Figure 2. Active learning example using heteroskedastic data. (a): Raw data. (b): Density of x in raw data. (c): Predictive distribution: $p(y|x)$ of BNN. (d): Entropy estimate $H(y|x)$ of predictive distribution for each x . (e): Conditional Entropy estimate $\mathbb{E}_{\mathcal{W}} H(y|x, \mathcal{W})$ of predictive distribution for each x . (f): Estimate of reduction in entropy for each x .

Figure 3f shows the estimate of reduction in entropy for each x . Here we can observe two effects: First, as expected, the expected entropy reduction will increase with higher x . Second, we see a slight decrease from $x = -0.5$ to $x = 0.75$. We believe the reason for this is twofold: because the data is limited to $[-0.5, 2]$ we expect a higher level of uncertainty in the vicinity of both borders. Furthermore we expect that learning a bimodal function requires more data to reach the same level of confidence than a unimodal function.

4. Risk-Sensitive Reinforcement Learning

In the previous section we studied how BNNs with latent variables can be used for active learning of stochastic functions. The resulting algorithm is based on a decomposition of predictive uncertainty into its aleatoric and epistemic components. In this section we build up on this result to derive a new risk-sensitive objective in model-based RL with the aim to minimize the effect of model bias. Our new risk criterion enforces that the learned policies, when evaluated on the ground truth system, are safe in the sense that on average they do not deviate significantly from the performance predicted by the model at training time.

Similar to (Depeweg et al., 2016), we consider the domain of batch reinforcement learning. In this setting we are given a batch of state transitions $\mathcal{D} = \{(s_t, \mathbf{a}_t, s_{t+1})\}$ formed by triples containing the current state s_t , the action applied \mathbf{a}_t and the next state s_{t+1} . For example, \mathcal{D} may be formed by measurements taken from an already-running system.

In addition to \mathcal{D} , we are also given a cost function c . The goal is to obtain from \mathcal{D} a policy in parametric form that minimizes c on average under the system dynamics.

The aforementioned problem can be solved using model-based policy search methods. These methods include two key parts (Deisenroth et al., 2013). The first part consists in learning a dynamics model from \mathcal{D} . We assume that the true dynamical system can be expressed by an unknown neural network with stochastic inputs:

$$\mathbf{s}_t = f_{\text{true}}(\mathbf{s}_{t-1}, \mathbf{a}_{t-1}, z_t; \mathcal{W}_{\text{true}}), \quad z_t \sim \mathcal{N}(0, \gamma), \quad (9)$$

where $\mathcal{W}_{\text{true}}$ denotes the synaptic weights of the network and \mathbf{s}_{t-1} , \mathbf{a}_{t-1} and z_t are the inputs to the network. In the second part of our model-based policy search algorithm, we optimize a parametric policy given by a deterministic neural network with synaptic weights \mathcal{W}_{π} . This parametric policy computes the action \mathbf{a}_t as a function of \mathbf{s}_t , that is, $\mathbf{a}_t = \pi(\mathbf{s}_t; \mathcal{W}_{\pi})$. We optimize \mathcal{W}_{π} to minimize the expected cost $C = \sum_{t=1}^T c_t$ over a finite horizon T with respect to our belief $q(\mathcal{W})$, where $c_t = c(\mathbf{s}_t)$. This expected cost is obtained by averaging over multiple virtual roll-outs. For each roll-out we choose \mathbf{s}_0 randomly from the states in \mathcal{D} , sample $\mathcal{W}_i \sim q$ and then simulate state trajectories using the model $\mathbf{s}_{t+1} = f(\mathbf{s}_t, \mathbf{a}_t, z_t; \mathcal{W}_i) + \epsilon_{t+1}$ with policy $\mathbf{a}_t = \pi(\mathbf{s}_t; \mathcal{W}_{\pi})$, input noise $z_t \sim \mathcal{N}(0, \gamma)$ and additive noise $\epsilon_{t+1} \sim \mathcal{N}(\mathbf{0}, \Sigma)$. This procedure allows us to obtain estimates of the policy’s expected cost for any particular cost function. If model, policy and cost function are differentiable, we are then able to tune \mathcal{W}_{π} by stochastic gradient descent over the roll-out average.

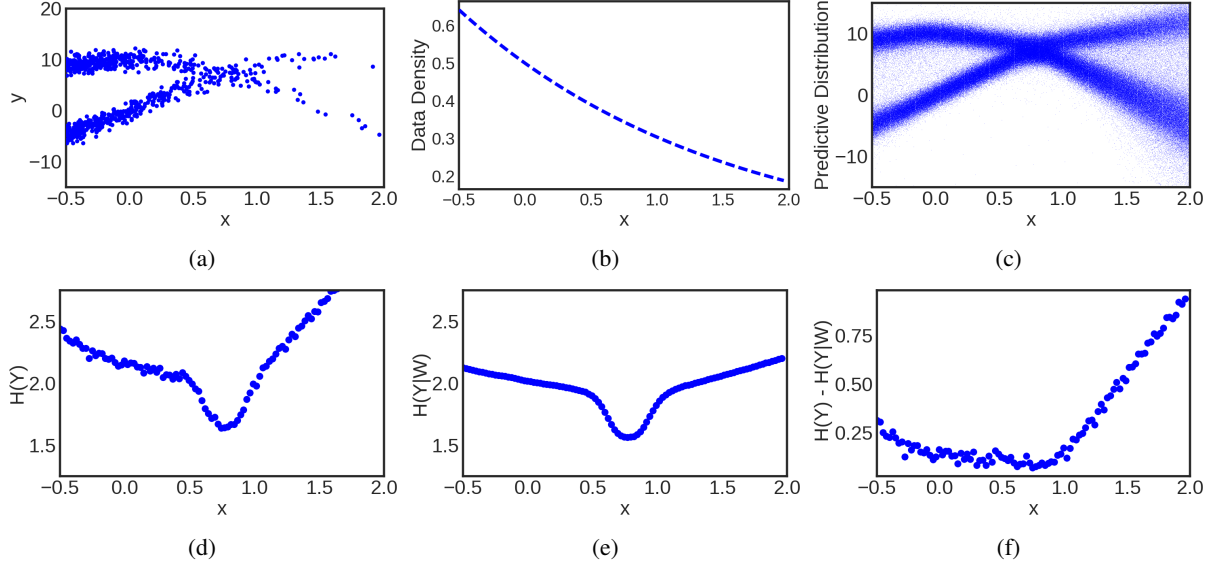


Figure 3. Active learning example using bimodal data. (a): Raw data. (b): Density of x in raw data. (c): Predictive distribution: $p(y|x)$ of BNN. (d): Entropy estimate $H(y|x)$ of predictive distribution for each x . (e): Conditional Entropy estimate $\mathbf{E}_{\mathcal{W}} H(y|x, \mathcal{W})$ of predictive distribution for each x . (f): Estimate of reduction in entropy for each x .

Given the cost function c , the objective to be optimized by the policy search algorithm is

$$J(\mathcal{W}_\pi) = \mathbf{E}_{q(\mathcal{W})} [C] = \mathbf{E}_{q(\mathcal{W})} \left[\sum_{t=1}^T c_t \right]. \quad (10)$$

In practice, s_0 is sampled uniformly at random from the training data \mathcal{D} . Standard approaches in risk-sensitive RL (García & Fernández, 2015; Maddison et al., 2017; Mihatsch & Neuneier, 2002) use the standard deviation of the cost C as risk measure. High risk is associated to high variability in the cost C . To penalize risk, the new objective to be optimized is given by

$$J(\mathcal{W}_\pi) = \mathbf{E}[C] + \beta \sigma(C), \quad (11)$$

where $\sigma(C)$ denotes the standard deviation of the cost and the free parameter β determines the amount of risk-avoidance ($\beta \geq 0$) or risk-seeking behavior ($\beta < 0$). In this standard setting, the variability of the cost $\sigma(C)$ originates from two different sources. First, from the existing uncertainty over the model parameters and secondly, from the intrinsic stochasticity of the dynamics.

One of the main dangers of model-based RL is model bias: the discrepancy of policy behavior under a) the assumed model and b) the ground truth MDP. While we cannot avoid the existence of such discrepancy when data is limited, we wish to guide the policy search towards policies that stay in state spaces where the risk for model-bias is low. For this, we can define the model bias b as follows:

$$b(\mathcal{W}_\pi) = \sum_{t=1}^T |\mathbf{E}_{\text{true}}[c_t] - \mathbf{E}_{q(\mathcal{W})}[c_t]|, \quad (12)$$

where $\mathbf{E}_{\text{true}}[c_t]$ is the expected cost obtained at time t when starting at the initial state s_0 and the ground truth dynamics are evolved according to policy $\pi(s_t; \mathcal{W}_\pi)$. Note that we focus on having similar expectations of the individual state costs c_t instead of having similar expectations of the final episode cost C . The former is a more strict criterion since it may occur that model and ground truth diverge, but both give roughly the same cost C on average.

As indicated in (9), we assume that the true dynamics are determined by a neural network with latent variables and weights given by $\mathcal{W}_{\text{true}}$. By using the approximate posterior $q(\mathcal{W})$, and assuming that $\mathcal{W}_{\text{true}} \sim q(\mathcal{W})$, we can obtain an upper bound on the expected model bias as follows:

$$\begin{aligned} \mathbf{E}_{q(\mathcal{W})}[b(\mathcal{W}_\pi)] &= \mathbf{E}_{\mathcal{W}_{\text{true}} \sim q(\mathcal{W})} \sum_{t=1}^T |\mathbf{E}[c_t | \mathcal{W}_{\text{true}}] - \mathbf{E}_{q(\mathcal{W})}[c_t]| \\ &= \sum_{t=1}^T \mathbf{E}_{\mathcal{W}_{\text{true}} \sim q(\mathcal{W})} \sqrt{(\mathbf{E}[c_t | \mathcal{W}_{\text{true}}] - \mathbf{E}_{q(\mathcal{W})}[c_t])^2} \\ &\leq \sum_{t=1}^T \sqrt{\mathbf{E}_{\mathcal{W}_{\text{true}} \sim q(\mathcal{W})} (\mathbf{E}[c_t | \mathcal{W}_{\text{true}}] - \mathbf{E}_{q(\mathcal{W})}[c_t])^2} \\ &= \sum_{t=1}^T \sqrt{\sigma_{\mathcal{W}_{\text{true}} \sim q(\mathcal{W})}^2 (\mathbf{E}[c_t | \mathcal{W}])} \\ &= \sum_{t=1}^T \sigma_{\mathcal{W}_{\text{true}} \sim q(\mathcal{W})} (\mathbf{E}[c_t | \mathcal{W}_{\text{true}}]). \end{aligned} \quad (13)$$

We note that $\mathbf{E}[c_t | \mathcal{W}]$ is the expected reward of a policy \mathcal{W}_π under the dynamics given by \mathcal{W} . The expectation integrates out the influence of the latent variables z_1, \dots, z_t and the

output noise $\epsilon_1, \dots, \epsilon_t$. The last equation in (13) can thereby be interpreted as the variability of the reward, that originates from our uncertainty over the dynamics given by distribution $q(\mathcal{W})$.

In Section 3 we showed how (7) encodes decomposition of the entropy of the predictive distribution into its aleatoric and epistemic components. The resulting decomposition naturally arises from an information-theoretic approach for active learning. We can express $\sigma_{\mathcal{W}_{\text{true}} \sim q(\mathcal{W})}^2(\mathbf{E}[c_t | \mathcal{W}_{\text{true}}])$ in a similar way using the law of total variance:

$$\sigma_{\mathcal{W}_{\text{true}} \sim q(\mathcal{W})}^2(\mathbf{E}[c_t | \mathcal{W}_{\text{true}}]) = \sigma^2(c_t) - \mathbf{E}_{\mathcal{W} \sim q(\mathcal{W})}[\sigma^2(c_t | \mathcal{W})].$$

We extend the policy search objective of (10) with a risk component given by an approximation to the model bias. Similar to Depeweg et al. (2016), we derive a Monte Carlo approximation that enables optimization by gradient descent. For this, we perform $M \times N$ roll-outs by first sampling $\mathcal{W} \sim q(\mathcal{W})$ a total of M times and then, for each of these samples of \mathcal{W} , performing N roll-outs in which \mathcal{W} is fixed and we only sample the latent variables and the additive Gaussian noise. In particular,

$$J(\mathcal{W}_\pi) = \sum_{t=1}^T \{ \mathbf{E}_{q(\mathcal{W})} [c_t] + \beta \sigma_{\mathcal{W}_{\text{true}} \sim q(\mathcal{W})}(\mathbf{E}[c_t | \mathcal{W}_{\text{true}}]) \} \approx \sum_{t=1}^T \left\{ \frac{1}{MN} \left[\sum_{m=1}^M \sum_{n=1}^N c_{m,n}(t) \right] + \beta \hat{\sigma}_M \left(\frac{1}{N} \sum_{n=1}^N c_{m,n}(t) \right) \right\}, \quad (14)$$

where $c_{m,n}(t) = c(\mathbf{s}_t^{\mathcal{W}^m, \{z_1^{m,n}, \dots, z_t^{m,n}\}, \{\epsilon_1^{m,n}, \dots, \epsilon_t^{m,n}\}}, \mathcal{W}_\pi)$ is the cost that is obtained at time t in a roll-out generated by using a policy with parameters \mathcal{W}_π , a transition function parameterized by \mathcal{W}^m and latent variable values $z_1^{m,n}, \dots, z_t^{m,n}$, with additive noise values $\epsilon_1^{m,n}, \dots, \epsilon_t^{m,n}$. $\hat{\sigma}_M$ is an empirical estimate of the standard deviation calculated over M draws of \mathcal{W} .

The free parameter β determines the importance of the risk criterion. As described above, the proposed approximation generates $M \times N$ roll-out trajectories for each starting state s_0 . For this, we sample $\mathcal{W}^m \sim q(\mathcal{W})$ for $m = 1, \dots, M$ and for each m we then do N roll-outs with different draws of the latent variables $z_t^{m,n}$ and the additive Gaussian noise $\epsilon_t^{m,n}$. We average across the $M \times N$ roll-outs to estimate $\mathbf{E}_{\mathcal{W} \sim q(\mathcal{W})}[c_t]$. Similarly, for each m , we average across the corresponding N roll-outs to estimate $\mathbf{E}[c_t | \mathcal{W}^m]$. Finally, we compute the empirical standard deviation of the resulting estimates to approximate $\sigma_{\mathcal{W}_{\text{true}} \sim q(\mathcal{W})}(\mathbf{E}[c_t | \mathcal{W}_{\text{true}}])$.

4.1. Application: Industrial Benchmark

We show now the effectiveness of the proposed method on a stochastic dynamical system. For this we use the industrial benchmark, a high-dimensional stochastic model inspired by

properties of real industrial systems. A detailed description and example experiments can be found in (Hein et al., 2016; Depeweg et al., 2016), with python source code available¹².

In our experiments, we first define a behavior policy that is used to collect data by interacting with the system. This policy is used to perform three roll-outs of length 1000 for each setpoint value in $\{0, 10, 20, \dots, 100\}$. The setpoint is a hyper-parameter of the industrial benchmark that indicates the complexity of its dynamics. The setpoint is included in the state vector \mathbf{s}_t as a non-controllable variable which is constant throughout the roll-outs. Policies in the industrial benchmark specify changes Δ_v , Δ_g and Δ_s in three steering variables v (velocity), g (gain) and s (shift) as a function of \mathbf{s}_t . In the behavior policy these changes are stochastic and sampled according to

$$\Delta_v \sim \begin{cases} \mathcal{N}(0.5, \frac{1}{\sqrt{3}}), & \text{if } v(t) < 40 \\ \mathcal{N}(-0.5, \frac{1}{\sqrt{3}}), & \text{if } v(t) > 60 \\ \mathcal{U}(-1, 1), & \text{otherwise} \end{cases} \quad (15)$$

$$\Delta_g \sim \begin{cases} \mathcal{N}(0.5, \frac{1}{\sqrt{3}}), & \text{if } g(t) < 40 \\ \mathcal{N}(-0.5, \frac{1}{\sqrt{3}}), & \text{if } g(t) > 60 \\ \mathcal{U}(-1, 1), & \text{otherwise} \end{cases} \quad (16)$$

$$\Delta_s \sim \mathcal{U}(-1, 1). \quad (17)$$

The velocity $v(t)$ and gain $g(t)$ can take values in $[0, 100]$. Therefore, the data collection policy will try to keep these values only in the medium range given by the interval $[40, 60]$. Because of this, large parts of the state space will be unobserved. After collecting the data, the 30,000 state transitions are used to train a BNN with latent variables with the same hyperparameters as in (Depeweg et al., 2016).

After this, we train different policies using the Monte Carlo approximation described in equation (14). We consider different choices of $\beta \in [0, 5]$ and use a horizon of $T = 100$ steps, with $M = 50$ and $N = 25$ and a minibatch size of 1.

Performance is measured using two different objectives. The first one is the expected cost obtained under the ground truth dynamics of the system, that is $\sum_{t=1}^T \mathbf{E}_{\text{true}}[c_t]$. The second objective is the model bias as defined in equation (12). We compare with two baselines. The first one ignores any risk and, therefore, is obtained by just optimizing equation (10). The second baseline uses the standard deviation $\sigma(c_t)$ as risk criterion and, therefore, is similar to equation (11), which is the standard approach in risk-sensitive RL.

In Figure 4 we show the results obtained by our method and by the second baseline when performance is evaluated under the model (Figure 4a) or under the ground truth (Figure

¹<https://github.com/siemens/industrialbenchmark>

²https://github.com/siemens/policy_search_bb-alpha

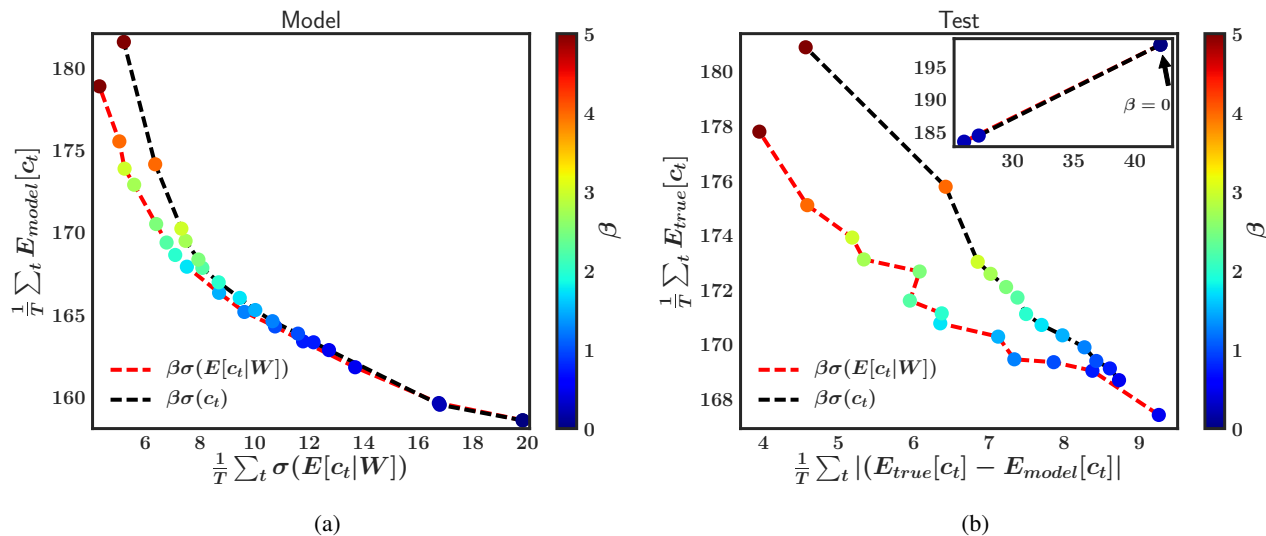


Figure 4. Results on Industrial Benchmark. Performances of policies trained using equation (14) (red curve) and baseline that minimizes (11) (black curve) for different values of β . Figure (a) shows results under the model and Figure (a) shows results under the ground truth.

4b). Each plot shows empirical estimates of the model bias vs. the expected cost, for various choices of β . We also highlight the result obtained with $\beta = 0$, the first baseline.

Our novel approach for risk-sensitive reinforcement learning produces policies that attain at test time better trade-offs between expected cost and model bias. As β increases, the policies gradually put more emphasis on the expected model bias. This leads to higher costs but lower discrepancy between model and real-world performance.

5. Conclusion

We have studied a decomposition of predictive uncertainty into its epistemic and aleatoric components when working with Bayesian neural networks with latent variables. This decomposition naturally arises in an information-theoretic active learning setting. The decomposition also inspired us to derive a novel risk objective for safe reinforcement learning that minimizes the effect of model bias in stochastic dynamical systems.

References

- Blundell, Charles, Cornebise, Julien, Kavukcuoglu, Koray, and Wierstra, Daan. Weight uncertainty in neural networks. *arXiv preprint arXiv:1505.05424*, 2015.
- Deisenroth, Marc Peter, Neumann, Gerhard, Peters, Jan, et al. A survey on policy search for robotics. *Foundations and Trends® in Robotics*, 2(1–2):1–142, 2013.
- Depeweg, Stefan, Hernández-Lobato, José Miguel, Doshi-Velez, Finale, and Udluft, Steffen. Learning and policy search in stochastic dynamical systems with bayesian neural networks. *arXiv preprint arXiv:1605.07127*, 2016.
- Gal, Yarin, McAllister, Rowan Thomas, and Rasmussen, Carl Edward. Improving pilco with bayesian neural network dynamics models. In *Data-Efficient Machine Learning workshop*, volume 951, pp. 2016, 2016.
- Gao, Weihao, Oh, Sewoong, and Viswanath, Pramod. Breaking the bandwidth barrier: Geometrical adaptive entropy estimation. In *Advances in Neural Information Processing Systems*, pp. 2460–2468, 2016.
- García, Javier and Fernández, Fernando. A comprehensive survey on safe reinforcement learning. *The Journal of Machine Learning Research*, 16(1):1437–1480, 2015.
- Hein, Daniel, Hentschel, Alexander, Sterzing, Volkmar, Tokic, Michel, and Udluft, Steffen. Introduction to the “industrial benchmark”. *arXiv preprint arXiv:1610.03793*, 2016.
- Hernández-Lobato, José Miguel and Adams, Ryan P. Probabilistic backpropagation for scalable learning of bayesian neural networks. *arXiv preprint arXiv:1502.05336*, 2015.
- Hernández-Lobato, José Miguel, Li, Yingzhen, Rowland, Mark, Hernández-Lobato, Daniel, Bui, Thang, and Turner, Richard E. Black-box α -divergence minimization. In *Proceedings of The 33rd International Conference on Machine Learning (ICML)*, 2016.
- Houthoofd, Rein, Chen, Xi, Duan, Yan, Schulman, John, De Turck, Filip, and Abbeel, Pieter. VIME: Variational

- information maximizing exploration. In *Advances in Neural Information Processing Systems*, pp. 1109–1117, 2016.
- Joseph, Joshua, Geramifard, Alborz, Roberts, John W, How, Jonathan P, and Roy, Nicholas. Reinforcement learning with misspecified model classes. In *Robotics and Automation (ICRA), 2013 IEEE International Conference on*, pp. 939–946. IEEE, 2013.
- Kendall, Alex and Gal, Yarin. What uncertainties do we need in bayesian deep learning for computer vision? *arXiv preprint arXiv:1703.04977*, 2017.
- Kozachenko, LF and Leonenko, Nikolai N. Sample estimate of the entropy of a random vector. *Problemy Peredachi Informatsii*, 23(2):9–16, 1987.
- Kraskov, Alexander, Stögbauer, Harald, and Grassberger, Peter. Estimating mutual information. *Physical review E*, 69(6):066138, 2004.
- MacKay, David JC. Information-based objective functions for active data selection. *Neural computation*, 4(4):590–604, 1992.
- Maddison, Chris J, Lawson, Dieterich, Tucker, George, Heess, Nicolas, Doucet, Arnaud, Mnih, Andriy, and Teh, Yee Whye. Particle value functions. *arXiv preprint arXiv:1703.05820*, 2017.
- Mihatsch, Oliver and Neuneier, Ralph. Risk-sensitive reinforcement learning. *Machine learning*, 49(2-3):267–290, 2002.
- Moerland, Thomas M, Broekens, Joost, and Jonker, Catholijn M. Learning multimodal transition dynamics for model-based reinforcement learning. *arXiv preprint arXiv:1705.00470*, 2017.
- Schneegass, Daniel, Udluft, Steffen, and Martinetz, Thomas. Uncertainty propagation for quality assurance in reinforcement learning. In *Neural Networks, 2008. IJCNN 2008.(IEEE World Congress on Computational Intelligence). IEEE International Joint Conference on*, pp. 2588–2595. IEEE, 2008.