
Robust uncertainty estimates with out-of-distribution pseudo-inputs training

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Probabilistic models often use neural networks to control their predictive uncer-
2 tainty. However, when making *out-of-distribution (OOD)* predictions, the often-
3 uncontrollable extrapolation properties of neural networks yield poor uncertainty
4 predictions. Such models then don't *know what they don't know*, which directly lim-
5 its their robustness w.r.t unexpected inputs. To counter this, we propose to explicitly
6 train the uncertainty predictor where we are not given data to make it reliable. As
7 one cannot train without data, we provide mechanisms for generating *pseudo-inputs*
8 in informative low-density regions of the input space, and show how to leverage
9 these in a practical Bayesian framework that casts a prior distribution over the
10 model uncertainty. With a holistic evaluation, we demonstrate that this yields
11 robust and interpretable predictions of uncertainty while retaining state-of-the-art
12 performance on diverse tasks such as regression and generative modeling.

13 1 Introduction

14 Neural networks generally extrapolate arbitrarily [Xu et al., 2020], and high quality predictions are
15 limited to regions of the input space where the networks have been trained. This is to be expected and
16 is only problematic if the associated predictions are not accompanied with a well-calibrated measure
17 of uncertainty. If a neural network is used for estimating such a measure of uncertainty, we, however,
18 quickly run into trouble, as the reported uncertainty then exhibits arbitrary behaviour in regions with
19 no training data. Alarming, these are exactly the regions where evaluating the uncertainty is most
20 important to the safe deployment of machine learning models in real world applications [Amodè
21 et al., 2016]. One potential solution is to avoid using directly the output of neural networks for
22 predicting uncertainty, and let it emerge from another mechanism, e.g. an *ensemble* [Hansen and
23 Salamon, 1990, Lakshminarayanan et al., 2017] or some notion of *Monte Carlo* [MacKay, 1992, Gal
24 and Ghahramani, 2016]. Here we explore the alternative view that the networks should simply be
25 trained where there is no data.

26 But can we train without data? The Bayesian formalism often
27 does so implicitly: most *conjugate priors* can be seen as addi-
28 tional training data [Bishop, 2006], e.g. in Gaussian models,
29 a mean prior $\mathcal{N}(\mu_0, \sigma_0^2)$ can be realised by additional training
30 data of μ_0 with σ_0^2 setting the amount of observations. Placing
31 a prior over the output of a neural network can, thus, be inter-
32 preted as additional training data. Unfortunately, this view is
33 not practical as it implies additional data *for all* possible inputs
34 to a neural network, resulting in infinite data. Our approach
35 is simple: we locate regions of low data density in *input space*
36 and implicitly place observations here in *output space* by min-

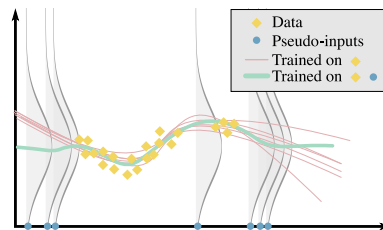


Figure 1: Pseudo-inputs are generated out of distribution, and there we train towards a prior (grey density).

37 imizing an appropriate KL divergence towards a prior (see Fig. 1). The result is a remarkably simple
38 algorithm that drastically improves uncertainty estimates in both regression and generative modeling.

39 1.1 Background and related work

40 The predictive performance of machine learning models has drastically increased in the past decade,
41 but the quality of the accompanying uncertainties have not followed. Uncertainties are reported as be-
42 ing miscalibrated [Guo et al., 2017] and overconfident [Lakshminarayanan et al., 2017, Hendrycks and
43 Gimpel, 2016]. Some models even see higher likelihoods of out-of-distribution than in-distribution
44 data [Nalisnick et al., 2019, Nguyen et al., 2015, Louizos and Welling, 2017].

45 **Neural networks** commonly output distributions which gives a notion of predictive uncertainty. Clas-
46 sifiers trained with *soft-max* is an ever-present example of such. These predictions are generally ob-
47 served to be *overconfident* [Lakshminarayanan et al., 2017, Hendrycks and Gimpel, 2016] and to carry
48 little meaning outside the support of the training data [Skafte et al., 2019, Lee et al., 2017]. The latter
49 is an artifact of the hard-to-control extrapolation that comes with neural networks [Xu et al., 2021].
50 In general, since extrapolation is difficult to control, uncertainties predicted by neural networks will
51 exhibit seemingly arbitrary behavior outside the support of the data, yielding untrustworthy results.

52 **Mean-variance networks** for regression [Nix and Weigend, 1994] model the conditional target
53 density as a normal $p(y|x) = \mathcal{N}(y|\mu(x), \sigma(x)^2)$ with mean and variance predicted by neural
54 networks. The reported predictive uncertainty is generally accurate in regions near training data,
55 but otherwise unreliable [Hauberg, 2019]. To counter this, Arvanitidis et al. [2017] and Skafte et al.
56 [2019] proposed variance network architectures to enforce a specified extrapolation value, but these
57 heuristics tend to be difficult to tune, and lack principle. Mean-variance networks have seen a recent
58 uptake within generative modeling, where they are applied as an *encoder* distribution in *variational*
59 *autoencoders* (VAEs) [Kingma and Welling, 2013, Rezende et al., 2014].

60 **Which uncertainty?** A commonly called-upon dichotomy [Der Kiureghian and Ditlevsen, 2009]
61 is that the uncertainty of a model’s *prediction* can be decomposed into the uncertainty of the *model*
62 (*epistemic*) and of the *data* (*aleatoric*). The epistemic uncertainty can be lowered by increasing
63 the amount of data, simplifying the model or otherwise reducing the complexity of the learning
64 problem. The aleatoric uncertainty, on the other hand, is a property of the world, and cannot be
65 changed; no prediction should ever be more certain than the uncertainty displayed by the associated
66 data. Depending on the task at hand, we may be interested in different types of uncertainty: In *active*
67 *learning* [Settles, 2012] and *Bayesian optimization* [Moćkus, 1975] we request data for which we
68 have high epistemic, but low aleatoric uncertainty to ensure maximal information gain; while for
69 classification and regression we often just want to minimize the overall predictive uncertainty.

70 **Bayesian methods** are often used to quantify uncertainty due to their explicit formulation of
71 uncertainty. *Gaussian processes* (GPs) [Rasmussen and Williams, 2005] provide an elegant
72 framework that provide state-of-the-art uncertainty estimates, but, alas, the corresponding mean
73 predictions are often not up to the standards of neural networks. GPs are tightly linked to *Bayesian*
74 *neural networks* (BNNs) [MacKay, 1992] that place a prior over the network weights and seek the
75 corresponding posterior. Despite advances in *variational approximations* [Graves, 2011, Kingma
76 and Welling, 2013, Blundell et al., 2015], *expectation propagation* [Hernández-Lobato and Adams,
77 2015, Hasenclever et al., 2017], or *Monte Carlo* methods [Welling and Teh, 2011, Springenberg et al.,
78 2016], training BNNs remains difficult. Furthermore, the predictive uncertainty seems dependent
79 on the degree of approximation and is thus controlled by the available compute power.

80 **Ensemble methods** have long been used to produce aggregated predictions with uncertainty estimates
81 [Hansen and Salamon, 1990, Breiman, 1996]. *Deep ensembles* [Lakshminarayanan et al., 2017],
82 a collection of differently initialized networks trained on the same data, are generally reported as
83 state-of-the-art for uncertainty quantification in deep models [Thagaard et al., 2020, Ovidia et al.,
84 2019]. As the models in the ensemble are trained on overlapping data, they are correlated, which
85 influence the ensemble uncertainty in ways that remains unclear [Breiman, 2001]. *Monte-Carlo*
86 *dropout* [Gal and Ghahramani, 2016] casts dropout training [Srivastava et al., 2014] as an ensemble
87 model. It is computationally cheap, but experiments [Ovidia et al., 2019, Skafte et al., 2019] show
88 that the increased correlation of ensemble elements amplifies the method’s overconfidence.

89 **Robustness to distribution shift** is paramount to a well-behaved uncertainty predictor [Ovidia
90 et al., 2019] and must be evaluated accordingly. For out-of-distribution detection, Liang et al.

91 [2017] proposes a pre-processing perturbation step inspired by adversarial attacks [Goodfellow et al.,
 92 2014a] that helps the model distinguish in-distribution and out-of-distribution inputs. Hendrycks
 93 et al. [2018] used a *Generative Adversarial Network (GAN)* [Goodfellow et al., 2014b] to generate
 94 out-of-distribution pseudo-inputs whose inclusion in the training under an additional regularizing
 95 term in the loss function, called *outlier exposure*, enhances the predictor’s ability to discriminate
 96 out-of-distribution inputs [Lee et al., 2017, Dai et al., 2017].

97 1.2 Robust uncertainty estimates

98 Our work is strongly inspired by the critical assessment of the issues that undermine variance estima-
 99 tion ran by [Skafte et al., 2019] and by the proposal of [Stirn and Knowles 2020] which we detail here.

100 **Notation.** Let the observed variable $x \in \mathcal{X}$ follow the data generating distribution $p_{\text{data}}(x)$, only
 101 known through the training dataset of N i.i.d samples $\mathcal{D}_{\text{train}} = \{x_n\}_{n=1}^N$. In the case of supervised
 102 learning, the observed variables $x = (x, y)$, with $x \in \mathbb{R}^d$ being the input and $y \in \mathbb{R}^{d'}$ the target
 103 for the model, follow the joint decomposition $p_{\text{data}}(x, y) = p_{\text{data}}(y|x)p_{\text{data}}(x)$. The proposed
 104 probabilistic model $p_{\theta}(x)$, whose weights are indicated by θ , aims to accurately emulate $p_{\text{data}}(x)$.

105 **Practical problems in variance estimation.** Gaussian likelihoods in the form of $p_{\theta}(x) =$
 106 $\mathcal{N}(x|\mu_{\theta}(x), \sigma_{\theta}(x)^2)$ are widely adopted to model continuous covariates. Real world data cannot be
 107 expected to be *homoscedastic*, i.e constant throughout input space, and thus the predictive uncertainty,
 108 $\sigma_{\theta}(x)$, most often uses neural networks to map continuously the observed x onto the parameter space.
 109 Beyond the well-known unreliable extrapolation properties of neural networks, this parametrisation
 110 of predictive uncertainty is hamstrung by serious defects. Firstly, the predictive variance scales the
 111 learning rates of the mean and variance updates by $1/2\sigma_{\theta}(x)^2$, resulting in a bias for data regions with
 112 low uncertainty [Nix and Weigend, 1994]. Secondly, the maximisation of the modeled likelihood is
 113 particularly sensitive to scarce data, as local gradient updates for the variance point towards the then
 114 undefined *maximum likelihood estimate (MLE)* [Skafte et al., 2019]. Lastly, and more worryingly, such
 115 model’s likelihood is ill-defined [Mattei and Frelsen, 2018a], as it can arbitrarily and without bound
 116 increase when the variance estimates collapse towards a detrimental 0. Overall, the naive maximisa-
 117 tion of model likelihood seems insufficient to generate robust and well-behaved uncertainty estimates.

118 **Student-t likelihood.** The Bayesian formalism, by imposing to learn a parametrised distribution
 119 over the predictive uncertainty, offers an attractive view to approaching the problem of uncertainty
 120 estimation. [Skafte et al., 2019] notably adopts a Gamma distributed precision, $1/\sigma^2 = \lambda \sim \Gamma(\alpha, \beta)$,
 121 as the conjugate of an unknown precision for a Gaussian likelihood, to yield a non-standard Student-t
 122 distributed marginal likelihood¹. It is known to offer a more robust likelihood, especially in the scarce
 123 data regime [Gelman et al., 2013],

$$p_{\theta}(x) = \int \mathcal{N}(x|\mu, \lambda)\Gamma(\lambda|\alpha, \beta)d\lambda = T\left(x|\nu = 2\alpha, \hat{\mu} = \mu, \hat{\sigma} = \sqrt{\beta/\alpha}\right). \quad (1)$$

124 Interestingly, its variance $\text{Var}[x] = \beta/(\alpha - 1) = (\beta/\alpha) \cdot (\alpha/(\alpha - 1))$ can be explicitly decomposed
 125 to an aleatoric term β/α and an epistemic term¹ $\alpha/(\alpha - 1)$ [Jørgensen, 2020, p16], and offers a
 126 direct verification of whether a model knows what it knows.

127 **Variational variance.** [Stirn and Knowles 2020] assumes a latent model precision λ . This is
 128 generated by a prior $p(\lambda)$ and its posterior is approximated variationally by the family of Gamma
 129 distributions, conditioned on the inputs to reflect heteroscedasticity. Through *amortized variational*
 130 *inference (AVI)* [Kingma and Welling, 2013] neural networks f_{ϕ} map to the posterior parameters from
 131 data, $q(z|f_{\phi}(x))$. As such, variational variance preserves the modelling capacity and robustness of
 132 the non-standard Student-t marginal likelihood, without modifying its parameter architecture, while
 133 the definition of a prior over the latent precision induces a more robust training objective. Assuming
 134 the likelihood precision is the unique latent code, the *evidence lower bound (ELBO)*,

$$\begin{aligned} \mathcal{L}(q; x) &= \mathbb{E}_{q(\lambda)} [\log p(x|\lambda)] - D_{\text{KL}}(q(\lambda|x) || p(\lambda)) \\ &= \frac{1}{2} \left(\psi(\alpha) - \log \beta - \log(2\pi) - \frac{\alpha}{\beta}(x - \mu)^2 \right) - D_{\text{KL}}(q(\lambda|x) || p(\lambda)), \end{aligned} \quad (2)$$

135 takes the form of a regularised log-likelihood, exposing the benefits of the prior regularisation. It
 136 penalises predicted variances that would unrealistically get arbitrarily close to either the detrimental

¹See Section [1](#) of the supplementary materials.

137 limits of 0 or ∞ , reducing the concerns regarding the ill-definition of the objective. Additionally,
 138 the scaling effect of the learning rates of the likelihood parameters is reduced. Naturally, the effect
 139 of the regularisation will be highly dependent on the prior selected. Here, because we are mostly
 140 interested in enforcing a constant desired uncertainty extrapolation, we adopt an homoscedastic
 141 Gamma distributed prior, $p(\lambda) = \Gamma(\lambda|a, b)$, that matches the level of uncertainty observed in data,
 142 and leave it for future practitioners to adopt the most adequate prior for the task at hand.

143 2 Out-of-distribution pseudo-inputs training

144 2.1 Dissipative loss

145 In the variational variance formalism, due to AVI, the predictive uncertainty is controlled by α and β ,
 146 the independent neural networks paramtrising the posterior distribution, $\text{Var}[x] = \beta(x)/(\alpha(x) - 1)$.
 147 The unreliable extrapolation properties of neural networks therefore directly challenge the robustness
 148 of the method’s uncertainty estimates outside of its training support, limiting the applicability of the
 149 method. We consider that this flawed extrapolation is not inevitable.

150 Inspired by outlier exposure [Hendrycks et al., 2018], we propose to include deliberately generated
 151 out-of-distribution *pseudo-inputs*, $\{\hat{x}_k\}_{k=1}^K$ where $\hat{x}_k \sim p_{\text{out}}(x)$, in the training of our variational
 152 objective to constrain the extrapolation of the posterior paramtrisation. The optimal variational
 153 objective q^* is chosen such that it minimises our proposed *dissipative loss* over the consolidated
 154 dataset $\mathcal{D} = \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{out}}$, where $\mathcal{D}_{\text{out}} = \{\hat{x}_k\}_{k=1}^K$,

$$\text{Loss}(q; \mathcal{D}) = -\left[\mathcal{L}_{\text{in}}(q; \mathcal{D}_{\text{train}}) + \mathcal{L}_{\text{out}}(q; \mathcal{D}_{\text{out}})\right]. \quad (3)$$

155 The in-distribution component of the loss function $\mathcal{L}_{\text{in}}(q; \mathcal{D})$ naturally arises as the standard ELBO
 156 over the training set. The out-of-distribution component $\mathcal{L}_{\text{out}}(q; \mathcal{D})$ operates on a fundamentally
 157 different source of data. As the only information available regarding the pseudo-inputs is that they
 158 are out-of-distribution, we assert for them a constant, non-informative likelihood $p(\hat{x}|\lambda) = c$, that
 159 has thus no influence on optimisation. This is similar to the strategy of *censoring* [Lee and Wang,
 160 2003] where different likelihoods are used for observations with different properties. As a result, the
 161 dissipative loss becomes,

$$\text{Loss}(q; \mathcal{D}) = -\left[\sum_{x \in \mathcal{D}_{\text{train}}} \mathbb{E}_{q(\lambda|x)} [p_{\theta}(x|\lambda)] - D_{\text{KL}}(q(\lambda|x) || p(\lambda)) - \sum_{\hat{x} \in \mathcal{D}_{\text{out}}} D_{\text{KL}}(q(\lambda|\hat{x}) || p(\lambda))\right]. \quad (4)$$

162 It share the same motivating intuition as the *confidence loss* of [Lee et al., 2017] and completes the
 163 variational variance formalism with a principled mechanism to learn robust variance estimates with
 164 the desired extrapolation properties. It indeed explicitly forces the predictor to match our high-entropy
 165 prior expectations on out-of-distribution samples while learning the low-entropy covariate dependent
 166 distribution, hence the name of dissipative. The reliance of the model’s predictive uncertainty
 167 on its mean predictions implies that it is primordial here to safeguard its generative performance.
 168 We guarantee it with the implementation of a split training procedure [Skafte et al., 2019]; the
 169 out-of-distribution regularisation is only applied after the model’s mean has been trained.

170 2.2 Pseudo-input generators (PIGs)

171 Minimising the posterior KL divergence out-of-distribution requires an efficient sampling procedure
 172 of pseudo-inputs. As exposed in Fig. 2, their generation should leverage a-priori knowledge about
 173 $p_{\text{data}}(x)$ to resolve the undefined nature of $p_{\text{out}}(x)$. In this simple regression case, we show the
 174 predictive uncertainty of variational variance models trained on artificial heteroscedastic data. We
 175 use a prior uncertainty level that matches the maximum of the data uncertainty. As anticipated,
 176 without pseudo-inputs, the model extrapolates uncertainty to a constant, arbitrary level, and only the
 177 introduction of pseudo-inputs near the training data results in the desired uncertainty extrapolation.
 178 Reassuringly, this suggests that we do not need to regularise our model’s extrapolation in the entire
 179 out-of-distribution space. Instead, we can focus on the simpler task of generating pseudo-inputs
 180 in low-density regions of the input space that neighbours training data, as they can enforce correct
 181 extrapolation in the rest of the out-of-distribution space. [Lee et al., 2017] gives supporting arguments.

182 Recent contributions have relied on GANs for generating a useful representation of $p_{\text{out}}(x)$ [Lee
 183 et al., 2017, Dai et al., 2017]. Although conceptually intuitive, GANs incur a heavy computational

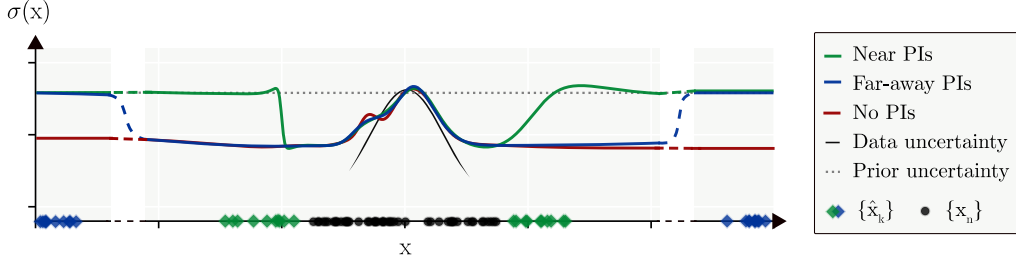


Figure 2: Effect of different pseudo-input distributions on the predictive uncertainty of variational variance models. Training data (black points) is generated uniformly on $[-5, 5]$, with a variance that scales as $\exp(-0.5(\|x\|/s)^2)$. The *near* pseudo-inputs (green diamonds) are generated uniformly in $[-10, -5] \cup [5, 10]$, while the *far-away* (blue diamonds) are on $[-200, -190] \cup [190, 200]$. Dashes amount for the empty space that separates far away pseudo-inputs.

burden and most likely induce serious practical challenges as a result of the instability of their training (Shrivastava et al., 2017). Furthermore, as one need to understand what is in-distribution to model what it is not, we instead propose to directly leverage the information at hand about the data.

Algorithm 1 gives a simple procedure for generating pseudo-inputs using the data density. Pseudo-inputs are originally sampled from $p_{\text{data}}(x)$, and their positions iteratively updated with gradient descent, with step size δ , by following the directions that minimise their likelihood under $p_{\text{data}}(x)$, similarly to reversed adversarial steps (Goodfellow et al., 2014a).

Algorithm 1: Pseudo-Input Generator (PIG)

```

 $\forall k \in [1, K], \hat{x}_k \sim p_{\text{data}}(x). \text{ iterations} = 0. \epsilon = \infty;$ 
while ( $\text{iterations} < \text{max\_iterations}$ ) & ( $\epsilon > \text{tolerance}$ ) do
  compute  $\forall k \in [1, K], \nabla_x p(x)(\hat{x}_k);$ 
   $\epsilon = \max_{k \in [1, K]} (\delta \nabla_x p(x)(\hat{x}_k));$ 
   $\forall k \in [1, K], \hat{x}_k = \hat{x}_k - \delta \nabla_x p(x)(\hat{x}_k);$ 
   $\text{iterations} = \text{iterations} + 1;$ 
end

```

The procedure can run prior to training, in parallel for all \hat{x}_k with automatic differentiation, and thus results in limited additional complexity for the optimisation². It relies on the availability of a differentiable density estimate of the data, which is, depending on the use case, either directly available (see Sec. 3.2), or can be approximated through a variety of methods such as *Bayesian Gaussian mixture models* (Bishop, 2006), or various *normalising flows* (Rezende and Mohamed, 2015) based methods such as *masked autoregressive flows* (Papamakarios et al., 2017) (see Sec. 3.1). A caveat here is that depending on the PIG’s parameters, and on the quality of the density estimate available, pseudo-inputs might be generated in undesired regions of the input space, e.g uninformative density minima. In practice, we adopted conservative density estimates and parameters and did not observe any significant degradation of the predictive uncertainty due to the addition of pseudo-inputs.

3 Experiments

Holistic evaluation of uncertainty estimates. The ground truth for uncertainty estimates is usually unknown, making their evaluation non-trivial. Similarly as in (Stirn and Knowles, 2020), we propose to assess them using a collection of metrics. Calibration, which evaluates probabilistic predictions w.r.t the long-run frequencies that actually occur (Dawid, 1982) can be measured by *proper scoring rules* (Lakshminarayanan et al., 2017) such as the model log-likelihood $\log p_\theta(x|\lambda)$. Additionally, the *root mean squared error (RMSE)* between the predictive and empirical variance, $\text{Var}[x] - (\mathbb{E}_{q(z|x)} [p_\theta(x|\lambda)] - x)^2$, offers a quantification of the model’s awareness of its own uncertainty. It nevertheless requires an understanding of the model’s mean predictive performance, as commonly measured by the RMSE of the mean residuals, $\mathbb{E}_{q(\lambda|x)} [p_\theta(x|\lambda)] - x$. We further propose to evaluate the cooperation of mean and uncertainty estimates for the generation of credible samples, which constitutes a consistency check for the learned precision distribution (Gelman et al., 2013), by measuring the RMSE of sample residuals $x^* - x$, with $x^* \sim p_\theta(x)$. Finally, The ELBO, despite the absence of theoretical grounding for it (Blei et al., 2017), is commonly reported as an approximation of the marginal likelihood, and thus of the overall model’s predictive performance.

²Running times are reported in Sec. IV of the supplementary materials.

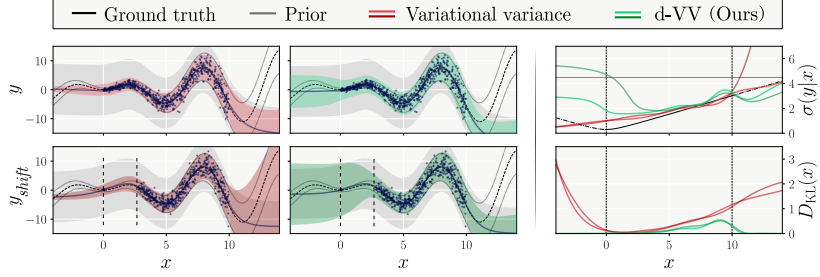


Figure 4: Toy regression results. On the left, the mean predictions are surrounded by ± 2 standard deviations, with the training data of the bottom row presenting a shift. On the right are displayed the predictive uncertainty fit and the prior KL divergence.

222 A complete assessment of a model’s uncertainty estimates further requires their evaluation under
 223 distributional shift [Ovadia et al., 2019], which we either introduce voluntarily through deliberate
 224 splitting of the training and test sets, as in Sec. 3.1 or by using test data from a different dataset
 225 altogether, as in Sec. 3.2.

226 3.1 Regression

227 In a regression setting where the proposed model must capture the conditioning between targets and
 228 inputs $y|x$, the precision λ of a Gaussian likelihood is the only assumed latent code.

229 Faithfully to variational variance [Stirn and Knowles, 2020] we
 230 adopt a Gamma heteroscedastic variational posterior $q_\phi(\lambda|x) =$
 231 $\Gamma(\lambda|\alpha_\phi(x), \beta_\phi(x))$ parametrised by the independent α_ϕ and β_ϕ
 232 networks, with weights ϕ , uniquely conditioned on the inputs (see
 233 Fig. 3). This approximate posterior, independent of the targets, gives
 234 up on the dependency of the true posterior on both covariates to
 235 guarantee heteroscedasticity³.

236 For strictly more than 2 degrees of freedom, or equivalently,
 237 $\alpha_\phi(x) > 1$, the marginal predictive probability $p_{\theta,\phi}(y|x) =$
 238 $T(y|2\alpha_\phi(x), \mu_\theta(x), \sqrt{\beta_\phi(x)/\alpha_\phi(x)})$, has its first two moments
 239 defined, $\mathbb{E}[y|x] = \mu_\theta(x)$ and $\text{Var}[y|x] = \beta_\phi(x)/(\alpha_\phi(x) - 1)$, pro-

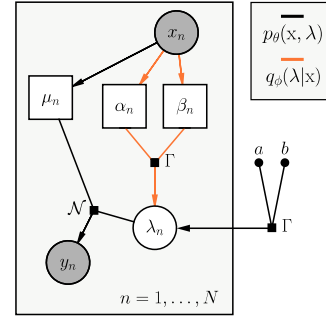


Figure 3: PGM for regression

240 viding explicit mean and uncertainty estimates with a single forward
 241 pass in the single layered, fully connected, α_ϕ , β_ϕ and μ_θ networks used here. To ensure definition of
 242 both the posterior distribution and of the marginal distribution’s variance, the parameter maps use a
 243 soft-plus activation on their last layer to ensure positivity, and the α_ϕ network is further shifted by 1.

244 The unique dependence of the posterior on the inputs implies that the generation of pseudo-inputs
 245 should only rely on the input density. In a general regression setting, it is unknown, and we estimate
 246 it here prior to training with a Bayesian Gaussian mixture model [Bishop, 2006]. We refer to it
 247 henceforth as *dissipative variational variance* (*d-VV*). The specific implementation details are listed
 248 in Section II. of the supplementary materials.

249 3.1.1 Toy regression

250 The desiderata for our method are clear: capture of the data heteroscedasticity, extrapolation to a
 251 higher uncertainty level, no underestimation of the predictive uncertainty, and posterior extrapolation
 252 to the prior out-of-distribution. Skafte et al. [2019] first showed on the toy regression task, $y =$
 253 $x \sin(x) + 0.3 \epsilon_1 + 0.3 x \epsilon_2$, where $\epsilon_1, \epsilon_2 \sim \mathcal{N}(0, 1)$, that amongst a collection of methods, only their
 254 proposed variance network architecture could realise our first three expectations. Fig. 4 demonstrates
 255 that our more principled approach also fulfills all of our requirements, without the need for arbitrarily
 256 enforcing the desired extrapolation in our architecture. The importance of out-of-distribution training
 257 is also revealed as the standard variational variance approach fails to produce uncertainty estimates
 258 that extrapolate correctly and are robust to distributional shift (bottom row of Fig. 4).

³See Section II. of the supplementary materials for the expression of the true posterior.

Table 1: UCI benchmarks. Each square shows the performance of a given model (rows) on a given dataset (columns). The intensity of the colouring represents the certitude that the associated model performed best on the given dataset. Grey rows mean impossible evaluation for a metric.

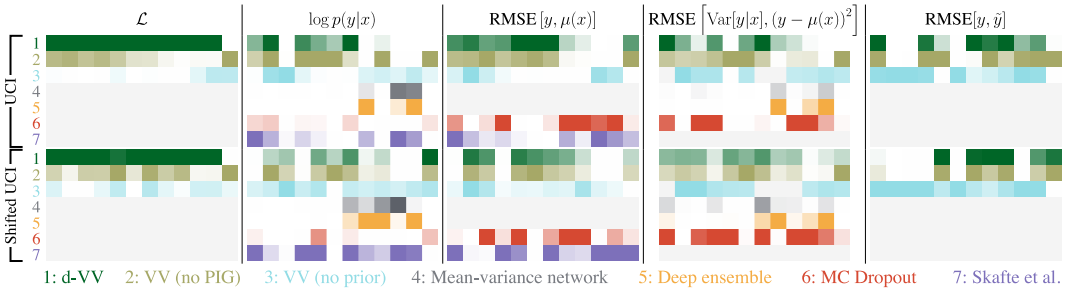


Table 2: Evaluation of the generative modeling. For each dataset, we report mean \pm std over 5 trials.

| | | FashionMNIST | SVHN | CIFAR |
|-----------------------|--------|---------------------------------------|---------------------------------------|---------------------------------------|
| $\log p(x)$ | VAE | 2215.54 \pm 68.81 | 4304.90 \pm 58.45 | 2930.64 \pm 14.82 |
| | d-V3AE | 2349.71 \pm 11.80 | 4133.41 \pm 64.28 | 2668.85 \pm 13.23 |
| RMSE(x, \tilde{x}) | VAE | 0.171 \pm 0.003 | 0.097 \pm 7e-4 | 0.154 \pm 5e-4 |
| | d-V3AE | 0.158 \pm 0.003 | 0.087 \pm 0.002 | 0.129 \pm 7e-4 |

259 **Decomposition of the model and data uncertainty.** Fig. 5
 260 presents the decomposition of the predictive uncertainty. The
 261 aleatoric component captures the heteroscedastic increase of
 262 uncertainty in the training data while the epistemic uncertainty,
 263 constant in distribution, extrapolates to higher values. The
 264 proposed method therefore demonstrates, to the best of our
 265 knowledge, a principled decomposition of uncertainty factors.

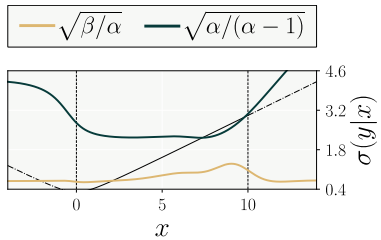


Figure 5: Aleatoric (yellow) and epistemic (dark) uncertainties.

266 3.1.2 UCI Benchmarks

267 Real world regression datasets from the UCI repository⁴ are used to evaluate our model against
 268 curated baselines, analogically to the setup from Hernández-Lobato and Adams [2015] and Skafte
 269 et al. [2019]⁵. As revealed by the summarising Tab. 1⁶, our method retains the mean predictive power
 270 of the variational variance method. The log-likelihood and RMSE of variance and sample residuals
 271 further show the improved calibration resulting from the imposition of a prior on the variance, as
 272 the VV methods generally outperform the MLE Student-t (VV (no prior)) that shares the same
 273 architecture. The table thus proves that holistically, the dissipative loss strengthens the variational
 274 variance model’s performance, which itself generally surpasses the chosen baselines.

275 The robustness of the methods to distributional shift is further evaluated as in Foong et al. [2019]. For
 276 each input feature, a hole is created in the training data by assigning the middle third of observations
 277 to the test set, when sorted w.r.t that feature. Interestingly, we see that our method’s calibration
 278 slightly improves under the shift, highlighting the robustness benefits of the OOD prior regularisation.

279 We note that both MC Dropout and the combined method of Skafte et al. [2019] generally perform
 280 well, confirming their interest for regression tasks requiring uncertainty quantification, but the
 281 former’s calibration is not robust to data shifts, as is also reported in Ovadia et al. [2019], and the
 282 latter is in practice difficult to tune and lacks principle.

⁴<https://archive.ics.uci.edu/ml/index.php>

⁵See Sec. II for details about the chosen baselines, datasets and implementations specifics.

⁶The full numbers are included in Sec. II of the supplements.

283 3.2 Generative models

284 We extend the evaluation of our proposal to the case of generative models through the lens of VAEs
 285 [Kingma and Welling, 2013, Rezende et al., 2014]. Variational auto-encoders infer a low dimensional
 286 latent encoding of the data $z \in \mathbb{R}^D$, on which is conditioned the generative process $p_\theta(x|z)$. Its
 287 predictive uncertainty, which evaluates the confidence of the model in its ability to adequately
 288 reconstruct inputs is known to be untrustworthy.

289 In the case of continuous or seemingly continuous inputs,
 290 the adoption of a Gaussian decoder $p_\theta(x|z) =$
 291 $\mathcal{N}(x|\mu_\theta(z), \sigma_\theta(z)^2)$ results in an ill-defined model likelihood
 292 [Mattei and Frellsen, 2018a] that encourages decoder variance
 293 collapse, making the training of the model notoriously harder
 294 [Skafte et al., 2019]. Most implementations therefore choose
 295 to fix the variance to a set level e.g $\sigma_\theta(z) = 0.1$, or elude the
 296 challenge by adopting a Bernoulli likelihood.

297 Motivated by our previous results, we now aim to demonstrate
 298 that VAEs, whose decoder is fitted with our method, are
 299 able to provide robust uncertainty estimates. Assuming a
 300 latent generative precision, the latent variables of the model
 301 are decomposed into $z = \{z, \lambda\}$, with z the latent input
 302 representations. The marginalisation of the Gamma distributed
 303 latent variance results in a Student-T decoder, as detailed in
 304 Eq. [1]. The overall architecture of the *variational variance variational auto encoder (V3AE)* [Stirn
 305 and Knowles, 2020] is shown in Fig. [6] and yields, with the addition of our out-of-distribution
 306 pseudo-inputs training, the dissipative loss function[7]

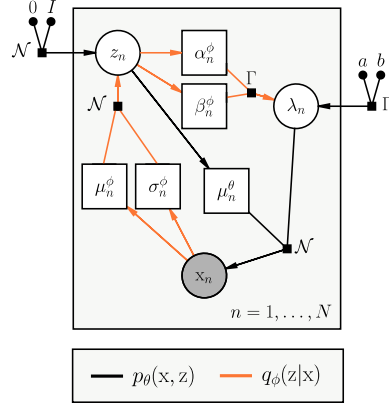


Figure 6: PGM for V3AE

$$\text{Loss}(q_\phi, \theta; \mathcal{D}_{\text{train}}) = - \left[\sum_{x \in \mathcal{D}_{\text{train}}} \mathcal{L}(q_\phi, \theta; x) - \mathbb{E}_{q_{\text{out}}(z)} [D_{\text{KL}}(q_\phi(\lambda|z) || p(\lambda))] \right]. \quad (5)$$

307 Because only the decoder is regularised, the pseudo-inputs lie in the space of latent representations,
 308 $\mathcal{D}_{\text{out}} = \{\hat{z}_k\}_{k=1}^K \in \mathbb{R}^D$. The distribution of training inputs is therefore readily accessible as the
 309 aggregate variational posterior $q_\phi(z|\mathcal{D}_{\text{train}}) = q_\phi(z|x_1) \cdot \dots \cdot q_\phi(z|x_N)$. Here again, we rely on a split
 310 training procedure to leverage this perk; the encoder parameter maps μ_θ and σ_θ , as well as the
 311 decoder mean μ_ϕ are first trained until convergence, allowing the generation of the out-of-distribution
 312 pseudo-inputs and subsequently, the training of the decoder variance.

313 **Image data.** We evaluate the performance of our proposed *dissipative*
 314 *V3AE (d-V3AE)* against a fully Gaussian VAE on image data, coming from
 315 FashionMNIST, SVHN and CIFAR10. For both models, all parameter
 316 maps share the same underlying architecture, with the addition of either
 317 a softplus and/or a shifting last layer to ensure definition of both the
 318 variational and the generative distribution’s moments[8]

319 Tab. [2] compares model performance on two metrics, the log-likelihood
 320 and the RMSE between the original inputs x and reconstructed samples \tilde{x} ,
 321 where $\tilde{x} \sim p_\theta(x|\lambda, z)$, $(\lambda, z) \sim q_\phi(\lambda, z|x)$. Unlike most previous imple-
 322 mentations, we focus on actual samples, and not the mean, of the generative
 323 distributions. This comparison emphasize the cooperation between the
 324 decoder’s mean and variance, allowing evaluation of the models’ uncer-
 325 tainty estimates. Our method both qualitatively (Fig. [7]), and quantitatively
 326 improves on a Gaussian VAE’s sampling ability. The prior smoothens the
 327 uncertainty estimates, resulting in more realistic and less crisp samples.
 328 The log-likelihoods, evaluated at test time using truncation, i.e. $p_{\text{trunc}}(x) =$
 329 $p_\theta(x)/(F_x(1) - F_x(0))$, to account for the finite support of data, reveal
 330 that our model can achieve a better fit, if the prior is selected correctly. In
 331 SVHN and CIFAR10, the presence of color channels complicates the selection process and challenges
 332 our choice of a single homoscedastic prior for all pixels and channels. We note that the dissipative loss
 333 also applies to classic VAEs with Bernoulli-only decoders; see Sec. [III] of the supplements for details.

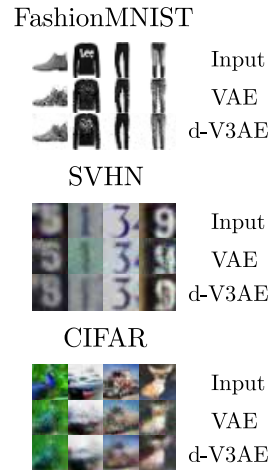


Figure 7: Generated samples

⁷The derivation of the dissipative loss function is provided in Sec. [III] of the supplementary materials.

⁸Implementation details are provided in Sec. [III] of the supplementary materials.

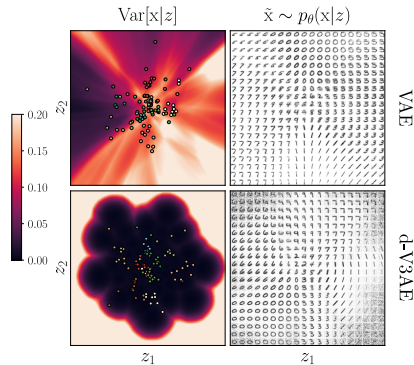


Figure 8: Decoder’s aggregated variance (left) and generated samples (right) from the latent space. Coloured points correspond to latent representations of test data, with per-class colours.

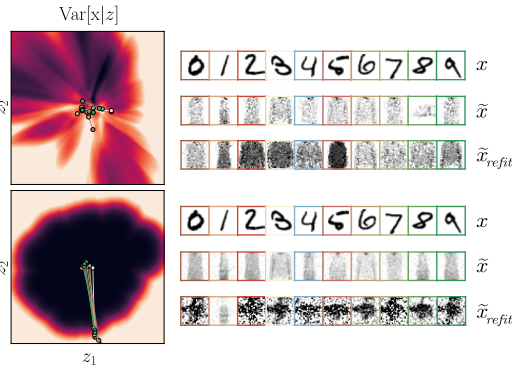


Figure 9: Effect of encoder refitting on the latent representations (left) and corresponding samples (right). OOD inputs (first rows, x) initially result in in-distribution samples (second rows, \tilde{x}). The refitted encoder displaces the encodings (coloured trajectories), modifying the generated samples (third rows, \tilde{x}_{refit}).

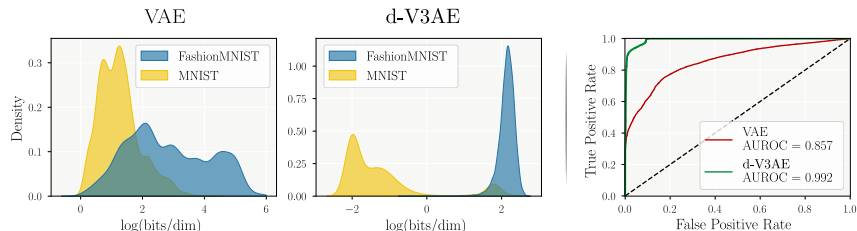


Figure 10: Empirical densities of likelihoods for FashionMNIST (ID) and MNIST (OOD). The clear separation of distributions offered by our method is reflected in the high AUROC shown on the right.

334 **Applications of robust generative uncertainty.** In Figs. 8 & 9 the colouring of the 2D latent space
 335 represent the aggregated decoder variance $\sum_{i=1}^d (\sigma_{\theta}(z)^2)_i$. It is clear that our method displays more
 336 regular uncertainty estimates, and provides the extrapolation guarantees we strove for. Beyond in-
 337 creased robustness and better generative power, this unlocks meaningful out-of-distribution detection,
 338 beating previous state-of-the-art [Havtorn et al., 2021]. For Figs. 9 & 10, as argued in [Mattei and
 339 Frelsen [2018b], we refit at test time the encoder of models trained on FashionMNIST on MNIST.
 340 The regularity and structure of the decoder variance rewards the encoder for learning to place represen-
 341 tations of OOD data outside of the region of in-distribution latent encodings, resulting in a model that
 342 is aware of its own inability to reconstruct plausible data, as displayed by the row \tilde{x}_{refit} of d-V3AE.

343 4 Conclusion

344 We have introduced a novel loss, the dissipative loss, that leverages artificial out-of-distribution
 345 pseudo-inputs for learning robust uncertainty estimates. We demonstrate through a Bayesian approach
 346 that casts a prior distribution over the model’s variance a principled mechanism for controlling the
 347 extrapolation properties of neural networks governing the predictive uncertainty. Our experimental
 348 results reflect the benefits of our principled and scalable approach, displaying better calibrated and
 349 more robust uncertainty estimates, while matching the predictive power of known baselines. Finally,
 350 and most interestingly, our approach can instill into probabilistic models a notion of their own
 351 ignorance, increasing their ability to *know what they don’t know*.

352 **Limitations.** The largest limitation of our approach is that it depends on an estimate of the density of
 353 the input data. In our experience, even coarse-grained densities are sufficient to significantly improve
 354 upon current approaches. However, as one rarely has guaranteed good estimates of the input density,
 355 our method cannot be approached as a black-box. One exception seems to be the application to VAEs,
 356 where the aggregated posterior, in our experience, always provide a suitable density estimate.

357 **Negative societal impact.** Improving the ability of predictive models to assess their own confidence
 358 is solely a positive contribution as it can help alleviate potential consequences of incorrect predictions.
 359 We are therefore not aware of any potential negative impacts of our work.

360 References

- 361 Keyulu Xu, Mozhi Zhang, Jingling Li, Simon S Du, Ken-ichi Kawarabayashi, and Stefanie Jegelka.
362 How neural networks extrapolate: From feedforward to graph neural networks. *arXiv preprint*
363 *arXiv:2009.11848*, 2020.
- 364 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané.
365 Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- 366 Lars Kai Hansen and Peter Salamon. Neural network ensembles. *IEEE transactions on pattern*
367 *analysis and machine intelligence*, 12(10):993–1001, 1990.
- 368 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive
369 uncertainty estimation using deep ensembles. In *Advances in neural information processing*
370 *systems*, pages 6402–6413, 2017.
- 371 David JC MacKay. A practical bayesian framework for backpropagation networks. *Neural computa-*
372 *tion*, 4(3):448–472, 1992.
- 373 Yarín Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model
374 uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059,
375 2016.
- 376 Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer, August 2006. ISBN
377 0387310738.
- 378 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural
379 networks. *arXiv preprint arXiv:1706.04599*, 2017.
- 380 Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution
381 examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.
- 382 Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, and Balaji Lakshminarayanan. Detecting out-of-
383 distribution inputs to deep generative models using typicality. *arXiv preprint arXiv:1906.02994*,
384 2019.
- 385 Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence
386 predictions for unrecognizable images. In *Proceedings of the IEEE conference on computer vision*
387 *and pattern recognition*, pages 427–436, 2015.
- 388 Christos Louizos and Max Welling. Multiplicative normalizing flows for variational bayesian neural
389 networks. *arXiv preprint arXiv:1703.01961*, 2017.
- 390 Nicki Skafté, Martin Jørgensen, and Søren Hauberg. Reliable training and estimation of variance
391 networks. In *Advances in Neural Information Processing Systems*, pages 6326–6336, 2019.
- 392 Kimin Lee, Honglak Lee, Kibok Lee, and Jinwoo Shin. Training confidence-calibrated classifiers for
393 detecting out-of-distribution samples. *arXiv preprint arXiv:1711.09325*, 2017.
- 394 Keyulu Xu, Mozhi Zhang, Jingling Li, Simon S. Du, Ken ichi Kawarabayashi, and Stefanie Jegelka.
395 How neural networks extrapolate: From feedforward to graph neural networks, 2021.
- 396 David A Nix and Andreas S Weigend. Estimating the mean and variance of the target probability
397 distribution. In *Proceedings of 1994 IEEE international conference on neural networks (ICNN'94)*,
398 volume 1, pages 55–60. IEEE, 1994.
- 399 Søren Hauberg. Only bayes should learn a manifold (on the estimation of differential geometric
400 structure from data), 2019.
- 401 Georgios Arvanitidis, Lars Kai Hansen, and Søren Hauberg. Latent space oddity: on the curvature of
402 deep generative models. *arXiv preprint arXiv:1710.11379*, 2017.
- 403 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*
404 *arXiv:1312.6114*, 2013.

- 405 Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and
406 approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014.
- 407 Armen Der Kiureghian and Ove Ditlevsen. Aleatory or epistemic? does it matter? *Structural safety*,
408 31(2):105–112, 2009.
- 409 Burr Settles. Active learning. *Synthesis lectures on artificial intelligence and machine learning*, 6(1):
410 1–114, 2012.
- 411 Jonas Moćkus. On bayesian methods for seeking the extremum. In *Optimization techniques IFIP*
412 *technical conference*, pages 400–404. Springer, 1975.
- 413 Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*
414 (*Adaptive Computation and Machine Learning*). The MIT Press, 2005.
- 415 Alex Graves. Practical variational inference for neural networks. *Advances in neural information*
416 *processing systems*, 24:2348–2356, 2011.
- 417 Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in
418 neural networks. *arXiv preprint arXiv:1505.05424*, 2015.
- 419 José Miguel Hernández-Lobato and Ryan Adams. Probabilistic backpropagation for scalable learning
420 of bayesian neural networks. In *International Conference on Machine Learning*, pages 1861–1869,
421 2015.
- 422 Leonard Hasenclever, Stefan Webb, Thibaut Lienart, Sebastian Vollmer, Balaji Lakshminarayanan,
423 Charles Blundell, and Yee Whye Teh. Distributed bayesian learning with stochastic natural gradient
424 expectation propagation and the posterior server. *The Journal of Machine Learning Research*, 18
425 (1):3744–3780, 2017.
- 426 Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In
427 *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 681–688,
428 2011.
- 429 Jost Tobias Springenberg, Aaron Klein, Stefan Falkner, and Frank Hutter. Bayesian optimization
430 with robust bayesian neural networks. *Advances in neural information processing systems*, 29:
431 4134–4142, 2016.
- 432 Leo Breiman. Bagging predictors. *Machine learning*, 24(2):123–140, 1996.
- 433 Jeppe Thagaard, Søren Hauberg, Bert van der Vegt, Thomas Ebstrup, Johan D. Hansen, and Anders B.
434 Dahl. Can you trust predictive uncertainty under real dataset shifts in digital pathology? In *Medical*
435 *Image Computing and Computer-Assisted Intervention (MICCAI)*, Lima, Peru, October 2020.
- 436 Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua
437 Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty?
438 evaluating predictive uncertainty under dataset shift. In *Advances in Neural Information Processing*
439 *Systems*, pages 13991–14002, 2019.
- 440 Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- 441 Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov.
442 Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine*
443 *learning research*, 15(1):1929–1958, 2014.
- 444 Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-of-distribution
445 image detection in neural networks. *arXiv preprint arXiv:1706.02690*, 2017.
- 446 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
447 examples. *arXiv preprint arXiv:1412.6572*, 2014a.
- 448 Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier
449 exposure. *arXiv preprint arXiv:1812.04606*, 2018.

- 450 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
451 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information*
452 *processing systems*, 27:2672–2680, 2014b.
- 453 Zihang Dai, Zhilin Yang, Fan Yang, William W Cohen, and Russ R Salakhutdinov. Good semi-
454 supervised learning that requires a bad gan. In *Advances in neural information processing systems*,
455 pages 6510–6520, 2017.
- 456 Andrew Stirn and David A Knowles. Variational variance: Simple and reliable predictive variance
457 parameterization. *arXiv preprint arXiv:2006.04910*, 2020.
- 458 Pierre-Alexandre Mattei and Jes Frellsen. Leveraging the exact likelihood of deep latent variable
459 models. In *Advances in Neural Information Processing Systems*, pages 3855–3866, 2018a.
- 460 Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin.
461 *Bayesian data analysis*. CRC press, 2013.
- 462 Martin Jørgensen. *Stochastic Representations with Gaussian Processes and Geometry*. PhD thesis,
463 Technical University of Denmark, 2020.
- 464 Elisa T Lee and John Wang. *Statistical methods for survival data analysis*, volume 476. John Wiley
465 & Sons, 2003.
- 466 Ashish Shrivastava, Tomas Pfister, Oncel Tuzel, Joshua Susskind, Wenda Wang, and Russell Webb.
467 Learning from simulated and unsupervised images through adversarial training. In *Proceedings of*
468 *the IEEE conference on computer vision and pattern recognition*, pages 2107–2116, 2017.
- 469 Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International*
470 *Conference on Machine Learning*, pages 1530–1538. PMLR, 2015.
- 471 George Papamakarios, Theo Pavlakou, and Iain Murray. Masked autoregressive flow for density
472 estimation. *arXiv preprint arXiv:1705.07057*, 2017.
- 473 A Philip Dawid. The well-calibrated bayesian. *Journal of the American Statistical Association*, 77
474 (379):605–610, 1982.
- 475 David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians.
476 *Journal of the American statistical Association*, 112(518):859–877, 2017.
- 477 Andrew YK Foong, Yingzhen Li, José Miguel Hernández-Lobato, and Richard E Turner. 'in-
478 between' uncertainty in bayesian neural networks. *arXiv preprint arXiv:1906.11537*, 2019.
- 479 Jakob D. Havtorn, Jes Frellsen, Søren Hauberg, and Lars Maaløe. Hierarchical vaes know what they
480 don't know, 2021.
- 481 P-A Mattei and J Frellsen. Refit your encoder when new data comes by. In *3rd NeurIPS workshop*
482 *on Bayesian Deep Learning*, 2018b.
- 483 Norman L Johnson, Samuel Kotz, and Narayanaswamy Balakrishnan. *Continuous univariate distri-*
484 *butions, volume 1, 2nd Edition*. John wiley & sons, 1994.
- 485 Christian Bauckhage. Computing the kullback-leibler divergence between two generalized gamma
486 distributions. *arXiv preprint arXiv:1401.6853*, 2014.

487 **Checklist**

- 488 1. For all authors...
- 489 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
490 contributions and scope? [Yes]
- 491 (b) Did you describe the limitations of your work? [Yes] Yes, we have included a separate
492 subsection on the limitations. See Section 4
- 493 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Covered
494 specifically in Section 4
- 495 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
496 them? [Yes] Yes.
- 497 2. If you are including theoretical results...
- 498 (a) Did you state the full set of assumptions of all theoretical results? [Yes] Our results are
499 obtained through empirical evidence. We do, however, include relevant derivations in
500 the supplementary material.
- 501 (b) Did you include complete proofs of all theoretical results? [N/A]
- 502 3. If you ran experiments...
- 503 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
504 perimental results (either in the supplemental material or as a URL)? [Yes] Code is
505 included in the supplementary material.
- 506 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
507 were chosen)? [Yes] Yes, every experiment is accompanied by a separate settings file
508 in yaml format.
- 509 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
510 ments multiple times)? [Yes] Each experiment was run at least in triplicate.
- 511 (d) Did you include the total amount of compute and the type of resources used (e.g., type
512 of GPUs, internal cluster, or cloud provider)? [Yes] Yes, in supplementary material.
- 513 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 514 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 515 (b) Did you mention the license of the assets? [Yes] We ran our experiments on standard,
516 well known datasets, reference the source, which itself includes licenses.
- 517 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
518 Our code is available in the supplementary material
- 519 (d) Did you discuss whether and how consent was obtained from people whose data you're
520 using/curating? [N/A] We did not run any experiments on personal data.
- 521 (e) Did you discuss whether the data you are using/curating contains personally identifiable
522 information or offensive content? [N/A]
- 523 5. If you used crowdsourcing or conducted research with human subjects...
- 524 (a) Did you include the full text of instructions given to participants and screenshots, if
525 applicable? [N/A]
- 526 (b) Did you describe any potential participant risks, with links to Institutional Review
527 Board (IRB) approvals, if applicable? [N/A]
- 528 (c) Did you include the estimated hourly wage paid to participants and the total amount
529 spent on participant compensation? [N/A]