
Laplacian Autoencoders for Learning Stochastic Representations

Marco Miani¹, Frederik Warburg¹,
 Pablo Moreno-Muñoz, Nicke Skaftø Detlefsen, Søren Hauberg
 {mmia, frwa, pabmo, nsde, sohau}@dtu.dk
 Technical University of Denmark

<https://github.com/FrederikWarburg/LaplaceAE>

Abstract

Established methods for unsupervised representation learning such as variational autoencoders produce none or poorly calibrated uncertainty estimates making it difficult to evaluate if learned representations are stable and reliable. In this work, we present a Bayesian autoencoder for unsupervised representation learning, which is trained using a novel variational lower-bound of the autoencoder evidence. This is maximized using Monte Carlo EM with a variational distribution that takes the shape of a Laplace approximation. We develop a new Hessian approximation that scales linearly with data size allowing us to model high-dimensional data. Empirically, we show that our Laplacian autoencoder estimates well-calibrated uncertainties in both latent and output space. We demonstrate that this results in improved performance across a multitude of downstream tasks.

1 Introduction

Unsupervised representation learning is a brittle matter. Consider the *classic autoencoder* (AE) (Rumelhart et al., 1986), which compresses data x to a low-dimensional representation z from which data is approximately reconstructed. The nonlinearity of the model implies that sometimes small changes to data x gives a large change in the latent representation z (and sometimes not). Likewise, for some data, reconstructions are of low quality, while for others it is near perfect. Unfortunately, the model does not have a built-in quantification of its uncertainty, and we cannot easily answer when the representation is reliable and accurately reflects data.

The celebrated *variational autoencoder* (VAE) (Kingma and Welling, 2014; Rezende et al., 2014) address this concern directly through an explicit likelihood model $p(x|z)$ and a variational approximation of the representation posterior $p(z|x)$. Both these distributions have parameters predicted by neural networks that act similarly to the encoder–decoder pair of the classic autoencoder.

But is the VAEs quantification of reliability reliable? To investigate, we fit a VAE with a two-dimensional latent representation to the MNIST dataset (Lecun et al., 1998), and illustrate the predicted uncertainty of $p(x|z)$ in Fig. 1a. The model learns to assign high uncertainty to low-level image features such as edges, but predicts its smallest values far away from the data distribution. Not only is such behavior counter-intuitive, it is also suboptimal in terms of data likelihood (Sec. 1.1). Retrospectively, this should not be surprising as the uncertainty levels away from the data are governed by the extrapolatory behavior of the neural network determining $p(x|z)$. This suggests that perhaps uncertainty should be a derived quantity rather than a predicted one.

¹ Denotes equal contribution.

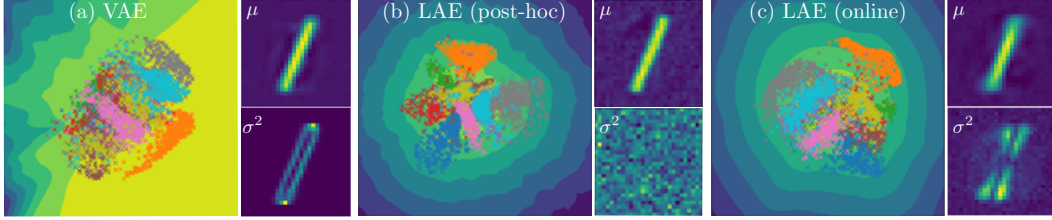


Figure 1: (a) The VAE learns to estimate high variance for low level image features such as edges, but fails at extrapolating uncertainties away from training data. (b) Applying post-hoc laplace to the AE setup shows much better extrapolating capabilities, but fails in estimating calibrated uncertainties in output space. (c) Our online, sampling-based optimization of a Laplacian autoencoder (LAE) gives well-behaved uncertainties in both latent and output space.

From a Bayesian perspective the natural solution is to form an (approximate) posterior over the weights of the neural networks. To investigate, we consider a state-of-the-art implementation of a *post-hoc* Laplace approximation (Daxberger et al., 2021) of the weight posterior. This amounts to training a regular autoencoder, and thereafter approximating the weight uncertainty with the Hessian of the loss (Sec. 1.1). Fig. 1b shows that uncertainty now grows, as intuitively expected, with the distance to the data distribution, but there seems to be little semantic structure in the uncertainty in output space. This suggests that while the post-hoc procedure is computationally attractive it is too simplistic.

In this paper we introduce a new framework for Bayesian autoencoders in unsupervised representation learning. Our method takes inspiration from the Laplace approximation to build a variational distribution on the neural network weights. We first apply existing post-hoc methods to the autoencoder domain, showcasing their lack of proper calibrated uncertainty. Secondly, we develop a new fast and memory-efficient Hessian approximation, which allows us to maximize a variational lower bound using Monte Carlo EM, such that model uncertainty is a key part of model training rather than estimated post-hoc. Fig. 1c gives an example of the corresponding uncertainty, which exhibits a natural and semantically meaningful behavior.

1.1 Background

The VAE is a latent variable model that parametrize the data density $p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$ using a latent variable (representation) \mathbf{z} . Here $p(\mathbf{z})$ is a, usually standard normal, prior over the representation, and $p(\mathbf{x}|\mathbf{z})$ is a likelihood with parameters predicted by a neural network. The nonlinearity of the likelihood parameters render the marginalization of \mathbf{z} intractable, and a variational lower bound of $p(\mathbf{x})$ is considered instead. To arrive at this, one first introduces a variational approximation $q(\mathbf{z}|\mathbf{x}) \approx p(\mathbf{z}|\mathbf{x})$ and write $p(\mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [p(\mathbf{x}|\mathbf{z})^{p(\mathbf{z})/q(\mathbf{z}|\mathbf{x})}]$. A lower bound on $p(\mathbf{x})$ then follows by direct application of Jensen’s inequality,

$$p(\mathbf{x}) \geq \mathcal{L}_{\text{VAE}}(\mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] - \text{KL}(q(\mathbf{z}|\mathbf{x})\|p(\mathbf{z})). \quad (1)$$

If we momentarily assume that $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mu(\mathbf{z}), \sigma^2(\mathbf{z}))$, we see that optimally $\sigma^2(\mathbf{z})$ should be as large as possible away from training data in order to increase $p(\mathbf{x})$ on the training data. Yet this is not the observed empirical behavior in Fig. 1a. Since the σ^2 network is left untrained away from training data, its predictions depends on extrapolation. In practice, σ^2 take fairly small values near training data (assuming the mean μ provides a reasonable data fit), and σ^2 likewise extrapolates to small values even if this is suboptimal in terms of data likelihood. Similar remarks hold for other likelihood models $p(\mathbf{x}|\mathbf{z})$ and encoder distributions $q(\mathbf{z}|\mathbf{x})$: *relying on neural network extrapolation to predict uncertainty does not work*.

The Laplace approximation (Laplace, 1774; MacKay, 1992) is an integral part of our proposed solution. In the context of Bayesian neural networks, we seek the weight-posterior $p(\theta|\mathcal{D}) \propto \exp(-\mathcal{L}(\mathcal{D};\theta))$, where θ are network weights, \mathcal{D} is the training data, and \mathcal{L} is the applied loss function interpreted as an unnormalized log-posterior. This is generally intractable and Laplace’s approximation (LA) amounts to a second-order Taylor expansion around a chosen weight vector θ^*

$$\log p(\theta|\mathcal{D}) = \mathcal{L}^* + (\theta - \theta^*)^\top \nabla \mathcal{L}^* + \frac{1}{2}(\theta - \theta^*)^\top \nabla^2 \mathcal{L}^* (\theta - \theta^*) + \mathcal{O}(\|\theta - \theta^*\|^3) \quad (2)$$

where we use the short-hand $\mathcal{L}^* = \mathcal{L}(\mathcal{D}; \theta^*)$. The approximation, thus, assumes that $p(\theta|\mathcal{D})$ is Gaussian. Note that when θ^* is a MAP estimate, the first order term vanishes and the second order term is negative semi-definite. We provide more details on the Laplace approximation in Appendix B. In practice, computing the full Hessian is too taxing both in terms of computation and memory, and various approximations are applied (Sec. 3).

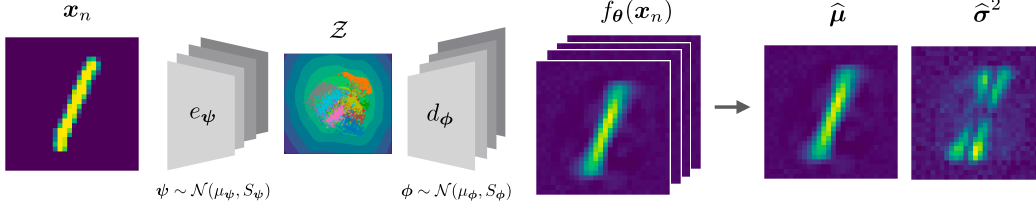


Figure 2: **Model overview.** We learn a distribution over parameters such that we can sample encoders e_ψ and decoders d_ϕ . This allows us to compute the empirical mean and variance in both the latent space z and the output space $f_\theta(\mathbf{x}_n) = d_\phi(e_\psi(\mathbf{x}_n))$.

2 Laplacian Autoencoders

We consider unsupervised representation learning from i.i.d. data $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$ consisting of observations $\mathbf{x}_n \in \mathbb{R}^D$. We also define a continuous latent space such that representations $\mathbf{z}_n \in \mathbb{R}^K$. Similar to AEs (Hinton and Salakhutdinov, 2006), we consider two neural networks $e_\psi: \mathbb{R}^D \rightarrow \mathbb{R}^K$ and $d_\phi: \mathbb{R}^K \rightarrow \mathbb{R}^D$, widely known as the *encoder* and *decoder*. These have parameters $\theta = \{\psi, \phi\}$. We refer to the composition of encoder and decoder as $f_\theta = d_\phi \circ e_\psi$.

Model overview. The autoencoder network structure implies that we model the data as being distributed on a K -dimensional manifold parametrized by θ . We then seek the distribution of the *reconstruction* $\mathbf{x}_{\text{rec}} = f_\theta(\mathbf{x})$ given observation \mathbf{x} , where the uncertainty comes from θ being unknown,

$$p(\mathbf{x}_{\text{rec}}|\mathbf{x}, f) = \mathbb{E}_{\theta \sim p(\theta|\mathbf{x}, f)} [p(\mathbf{x}_{\text{rec}}|\theta, \mathbf{x}, f)]. \quad (3)$$

Notice that we explicitly condition on f , which is the operator $\theta \mapsto f_\theta$, even if this is not stochastic; this choice will simplify notation at a later stage. Mimicking the standard autoencoder reconstruction loss, we set $p(\mathbf{x}_{\text{rec}}|\theta, \mathbf{x}, f) = \mathcal{N}(\mathbf{x}_{\text{rec}}|f_\theta(\mathbf{x}), \mathbb{I})$. Since $p(\theta|\mathbf{x}, f)$ is unknown, the reconstruction likelihood (3) is intractable, and approximations are in order. Similar to Blundell et al. (2015), we resort to a Gaussian approximation, but rather than learning the variance variationally, we opt for LA. This will allow us to sample NNs and deduce uncertainties in both latent and output space as illustrated in Fig. 2.

Intractable joint distribution. Any meaningful approximate posterior over θ should be similar to the marginal of the joint distribution $p(\theta, \mathbf{x}_{\text{rec}}|\mathbf{x}, f)$. This marginal is

$$p(\theta|\mathbf{x}, f) = \mathbb{E}_{\mathbf{x}_{\text{rec}} \sim p(\mathbf{x}_{\text{rec}}|\mathbf{x}, f)} [p(\theta|\mathbf{x}_{\text{rec}}, \mathbf{x}, f)] \quad (4)$$

which can be bounded on a log-scale using Jensen’s inequality,

$$\log p(\theta|\mathbf{x}, f) \geq \mathcal{L}_\theta = \mathbb{E}_{\mathbf{x}_{\text{rec}} \sim p(\mathbf{x}_{\text{rec}}|\mathbf{x}, f)} [\log p(\theta|\mathbf{x}_{\text{rec}}, \mathbf{x}, f)]. \quad (5)$$

Our first approximation is a LA of $p(\theta|\mathbf{x}, f) \approx q^t(\theta|\mathbf{x}, f) = \mathcal{N}(\theta|\theta^{(t)}, \mathbf{H}_t^{-1})$, where we postpone the details on how to acquire $\theta^{(t)}$ and \mathbf{H}_t . These will eventually be iteratively computed from the lower bound (5); hence the t index. Fig. 3 (a) illustrates the situation thus far: $p(\theta|\mathbf{x}, f)$ is approximately Gaussian, but the non-linearity of f gives $p(\mathbf{x}_{\text{rec}}|\mathbf{x}, f)$ a non-trivial density.

Linearization for gradient updates. Standard gradient-based learning can be viewed as a linearization of f in θ , i.e. for a loss $\mathcal{L} = l(f_\theta)$, the gradient is $\nabla_\theta \mathcal{L} = J_\theta f_\theta \nabla_f l(f)$, where $J_\theta f_\theta$ is the Jacobian of f . In a similar spirit, we linearize f in θ in order to arrive at a tractable approximation of $p(\mathbf{x}_{\text{rec}}|\mathbf{x}, f)$. Specifically, we perform a Taylor expansion around $\theta^{(t)}$

$$f_\theta(\mathbf{x}) = f_{\theta^{(t)}}(\mathbf{x}) + J_\theta f_{\theta^{(t)}}(\mathbf{x})(\theta - \theta^{(t)}) + \mathcal{O}(\|\theta - \theta^{(t)}\|^2) \quad (6)$$

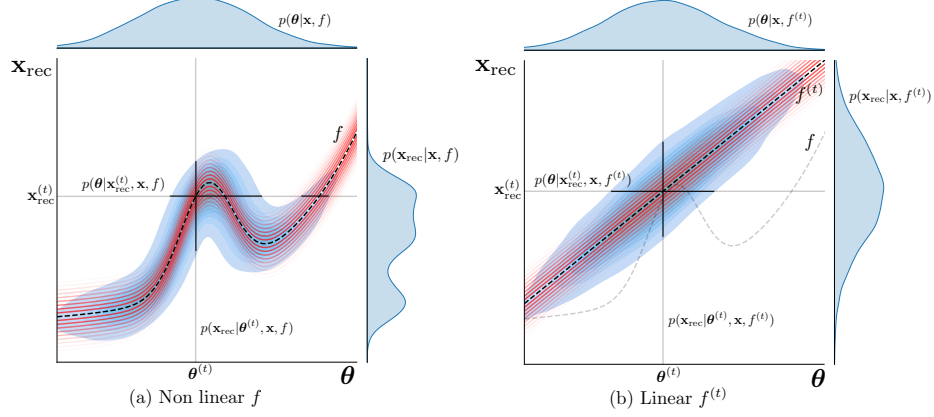


Figure 3: Illustrative example for fixed \mathbf{x} . The likelihood for a fixed $\theta^{(t)}$, shown by the columns, are assumed Gaussian $p(\mathbf{x}_{\text{rec}}|\theta^{(t)}, \mathbf{x}, f) = N(f_{\theta^{(t)}}(\mathbf{x}), \mathbb{I})$. We model the marginalised density $p(\theta|\mathbf{x}, f)$ (first axis) over parameters θ with Gaussians. With the additional assumption of linear f then $p(\mathbf{x}_{\text{rec}}|\mathbf{x}, f)$ (second axis) is Gaussian. This makes the joint distribution tractable.

where $f^{(t)}$ denote the associated first-order approximation. Under this approximation, the joint distribution $p(\theta, \mathbf{x}_{\text{rec}}|\mathbf{x}, f)$ becomes Gaussian, $p(\theta, \mathbf{x}_{\text{rec}}|\mathbf{x}, f^{(t)}) \approx \mathcal{N}(\theta, \mathbf{x}_{\text{rec}}|\mu^{(t)}, \Sigma^{(t)})$, with

$$\mu^{(t)} = \begin{pmatrix} \theta^{(t)} \\ f_{\theta^{(t)}}^{(t)}(\mathbf{x}) \end{pmatrix}, \text{ and } \Sigma^{(t)} = \begin{pmatrix} \mathbf{H}_t^{-1} & J_{\theta} f_{\theta^{(t)}}(\mathbf{x})^{\top} \\ J_{\theta} f_{\theta^{(t)}}(\mathbf{x}) & (J_{\theta} f_{\theta^{(t)}}(\mathbf{x})^{\top} \mathbf{H}_t J_{\theta} f_{\theta^{(t)}}(\mathbf{x}))^{-1} + \mathbb{I} \end{pmatrix}. \quad (7)$$

This approximation is illustrated in Fig. 3(b).

Iterative learning. With these approximations we can readily develop an iterative learning scheme which updates $q(\theta|\mathbf{x}, f)$. The mean of this approximate posterior can be updated according to a standard variational gradient step, $\theta^{(t+1)} = \theta^{(t)} + \lambda \nabla_{\theta} \mathcal{L}_{\mathbf{x}_{\text{rec}}}$, where

$$\mathcal{L}_{\mathbf{x}_{\text{rec}}} = \mathbb{E}_{\theta \sim q^t(\theta|\mathbf{x}, f)} [\log p(\mathbf{x}_{\text{rec}}|\theta, \mathbf{x}, f^{(t)})], \quad (8)$$

is a lower bound on Eq. 3, which we evaluate with a single Monte Carlo sample. The last missing piece is the covariance of $q^{t+1}(\theta|\mathbf{x}, f)$, which ideally should be the inverse of the Hessian of $\log p(\theta|\mathbf{x}, f)$ at $\theta^{(t+1)}$. Since this is intractable, we instead compute the Hessian of the lower bound (5)

$$\mathbf{H}_{t+1} = -\nabla_{\theta}^2 \mathcal{L}_{\theta} |_{\theta^{(t+1)}} = \mathbb{E}_{p(\mathbf{x}_{\text{rec}}|\mathbf{x}, f^{(t)})} \left[-\nabla_{\theta}^2 \log p(\mathbf{x}_{\text{rec}}|\theta, \mathbf{x}, f^{(t)}) - \nabla_{\theta}^2 \log q^t(\theta|\mathbf{x}, f) \right] |_{\theta^{(t+1)}} \quad (9)$$

The last term can be approximated by noting that $\nabla_{\theta}^2 \log q^t(\theta|\mathbf{x}, f) |_{\theta=\theta^{(t+1)}} = -\mathbf{H}_t + \mathcal{O}(\|\theta^{(t+1)} - \theta^{(t)}\|)$. To efficiently cope with the \mathcal{O} -term, we introduce a parameter α , such that the final approximation is

$$\mathbf{H}_{t+1} \approx (1 - \alpha) \mathbf{H}_t + \mathbb{E}_{p(\mathbf{x}_{\text{rec}}|\mathbf{x}, f)} \left[-\nabla_{\theta}^2 \log p(\mathbf{x}_{\text{rec}}|\theta, \mathbf{x}, f^{(t)}) \right] |_{\theta^{(t+1)}} \quad (10)$$

$$= (1 - \alpha) \mathbf{H}_t - J_{\theta} f_{\theta^{(t)}}^{(t)\top} \nabla_{\mathbf{x}_{\text{rec}}}^2 \log p(\mathbf{x}_{\text{rec}}|\theta^{(t+1)}, \mathbf{x}, f^{(t)}) J_{\theta} f_{\theta^{(t)}}^{(t)}. \quad (11)$$

The practical intuition is that when α is small, the network uncertainty decreases faster. The overall training procedure is summarized in Fig. 4. We initialize q^0 as a Gaussian with $\theta^{(0)} = 0$ and $\mathbf{H}_0 = \mathbb{I}$. We provide more details on the model, the linearization, and iterative learning in Appendix D.

Why not just...? The proposed training procedure may at first appear non-trivial, and it is reasonable to wonder if existing methods could be applied to similar results. Variational inference often achieves similar results to Laplace approximations, so could we use ‘Bayes by Backprop’ (Blundell et al., 2015) to get an alternative Gaussian approximate posterior over θ ? Similar to the supervised experiences of Jospin et al. (2020), we unfortunately found this approach too brittle to allow for practical model fitting. But then perhaps a post-hoc LA as proposed by Daxberger et al.

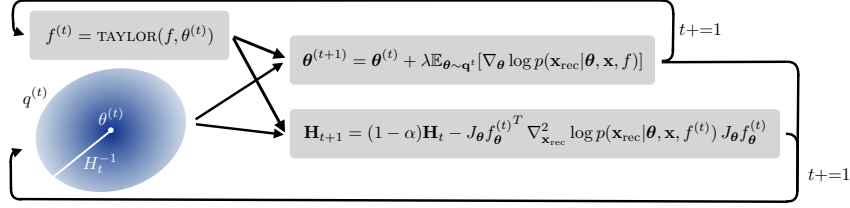


Figure 4: **Recursive training procedure.** Given a distribution q^t over parameters, and a linearized function $f^{(t)}$, compute first and second order derivatives to update the distribution on parameters.

(2021) for supervised learning? Empirically, we found it to be important to center the approximate posterior around a point where the Hessian provides useful uncertainty estimates. Our online training moves in this direction as the Hessian is part of the procedure, but this is not true for the post-hoc LA.

Conceptually, we argue that our approach, while novel, is not entirely separate from existing methods. Our reliance on lower bounds make the method an instance of variational inference (Jordan et al., 1999; Opper and Archambeau, 2009b), and we maximize the bounds using Monte Carlo EM (Cappé, 2009). We rely on a LA as our choice of variational distribution, which has also been explored by Park et al. (2019). Finally, we note that our linearization trick (6) bares great similarities to classic extended Kalman filtering (Gelb et al., 1974).

3 Scaling the Hessian to Large Images

The largest obstacle to apply LA in practice stem from the Hessian matrix. This matrix has a quadratic memory complexity in the number of network parameters, which very quickly exceeds the capabilities of available hardware. To counter this issue, several approximations have been proposed (Ritter et al., 2018; Botev et al., 2017; Martens and Grosse, 2015b) that improve the scaling w.r.t. to the number parameters. The currently most efficient Hessian implementations (Dangel et al., 2020; Daxberger et al., 2021) builds on the generalized Gauss-Newton (GGN) approximation of the Hessian

$$\nabla_{\theta_l}^2 \mathcal{L}(f_{\theta}(\mathbf{x})) \approx J_{\theta_l} f_{\theta}(\mathbf{x})^{\top} \cdot \nabla_{\mathbf{x}_{\text{rec}}}^2 \mathcal{L}(\mathbf{x}_{\text{rec}}) \cdot J_{\theta_l} f_{\theta}(\mathbf{x}), \quad (12)$$

which neglects second order derivatives of f w.r.t. the parameters. Besides, the computational benefits of this approximations, previous works on LA (Daxberger et al., 2021) relies on GGN to ensure that the Hessian is always semi-negative definite. In contrast, the model presented in Sec. 2 implies that GGN is no longer a practical and unprincipled trick, but rather the exact Hessian.

Albeit relying on first order derivatives, the layer-block-diagonal GGN, which assume that layers are independent to each other, scales quadratically with the *output* dimension of the considered neural network f . This lack of scaling is particularly detrimental for convolutional layers as these have low parameter counts, but potentially very high output dimensions.

Expanding $J_{\theta_l} f_{\theta}(\mathbf{x})$ with the chain rule, one realises that the Jacobian can be computed as a function of the Jacobian of the next layer. Fig. 5 illustrate that an intermediate quantity M , which is initialised as $\nabla_{\mathbf{x}_{\text{rec}}}^2 \mathcal{L}(\mathbf{x}_{\text{rec}})$, can be efficiently backpropagated through multiplication with the Jacobian w.r.t. input of each layer. This process leads to a **block diagonal** approximation of the Hessian as illustrated in Fig. 6(a). However, diagonal blocks are generally too large to store and invert. To combat this, each block can be further approximated by its **exact diagonal** (Daxberger et al., 2021) as depicted in Fig. 6(b). This scales linearly w.r.t. parameters, but still scales quadratically w.r.t. the output resolution (Tab. 1).

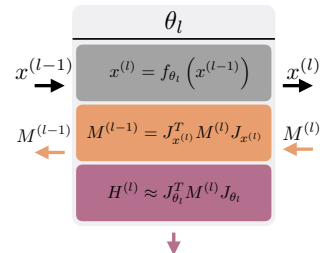


Figure 5: Forward pass of feature map x for layer l with parameters θ_l and **extended backward pass in which M is backpropagated to previous layers. Via the chain rule and M the Hessian of each layer can be computed efficiently.**

In order to scale our Laplacian autoencoder to high-dimensional data, we propose to approximate the diagonal of the Hessian rather than relying on exact computations. This is achieved by only backpropagating a diagonal form of M as illustrated in Fig. 6(c). This assumes that features from the same layer are uncorrelated, and consequently

APPROXIMATIONS	MEMORY	TIME
Block diag.	$\mathcal{O}(R_m^2 + W_s^2)$	$\mathcal{O}(R_s^2 + W_s^2)$
KFAC	$\mathcal{O}(R_m^2 + W_s)$	$\mathcal{O}(R_s^2 + W_s)$
Exact diag.	$\mathcal{O}(R_m^2 + W_s)$	$\mathcal{O}(R_s^2 + W_s)$
Approx. diag. (ours)	$\mathcal{O}(R_m + W_s)$	$\mathcal{O}(R_s + W_s)$
Mixed diag. (ours)	$\mathcal{O}(R_m + W_s)$	$\mathcal{O}(R_s + W_s)$

Table 1: **Memory & time complexity of Hessian approximations.** For an L -layer network, let $R_m = \max_{l=0 \dots L} |x^{(l)}|$, $R_s = \sum_l |x^{(l)}|$, $R_s^2 = \sum_l |x^{(l)}|^2$, and $W_s = \sum_l |\theta_l|$. Only our approximation scales linearly with both the output resolution and parameters.

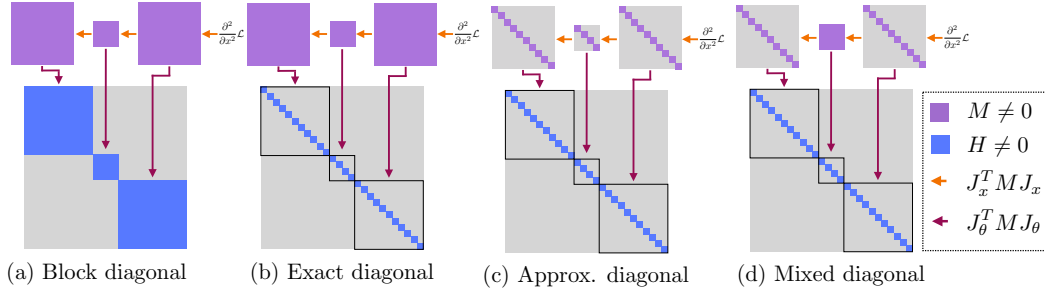


Figure 6: **Comparison of Hessian approximation methods.** Common approximations (a–b) scale quadratically with the output resolution. Our proposed approximate and mixed diagonal Hessians (c–d) scale linearly with the resolution. This is essential for scaling the LAE to large images.

has linear complexity in both time and memory with respect to the output dimension (Tab. 1). This makes it viable for our model.

We can further tailor this approximation to the autoencoder setting by leveraging the bottleneck architecture. We note that the quadratic scaling of the exact diagonal Hessian is less of an issue in the layers near the bottleneck than in the layers closer to the output space. We can therefore dynamically switch between our approximate diagonal and the exact one, depending on the feature dimension. This lessens the approximation error while remaining tractable in practice. We provide more details on the fast hessian computations in Appendix E.

4 Related Work

Deep generative models, and particularly the family of variational auto-encoders (VAEs) (Kingma and Welling, 2014; Rezende et al., 2014), address unsupervised representation learning from a probabilistic viewpoint by approximating the posterior over the representation space. Despite their widespread adoption, model parameters are still *deterministic* and sensible to ill-suited local minima, e.g. over-fitted to training data (Zhang et al., 2021), which may cause poor generalization. The Bayesian NNS favour inference over the NN weights for addressing such issues (MacKay, 1995; Neal, 1996). This approach deduce distributions on data space by learning posterior distributions on the parameter space. However, several shortcomings (Wenzel et al., 2020), e.g. expensive training, tuning and implementations, often limit their applicability to autoencoder-style models. Alternatively, other methods such as deep ensembles (Lakshminarayanan et al., 2017), stochastic weight averaging (SWAG) (Maddox et al., 2019) or Monte-Carlo dropout (Gal and Ghahramani, 2016) also promise Bayesian approximations to NN weight’s posterior, but at the cost of increased training time, poor empirical performance or limited Bayesian interpretation.

As demonstrated by Daxberger et al. (2021) LA is a scalable and well-behaved alternative to Bayesian NNS if used *post-hoc* to approximate the intractable posterior over the weights after *maximum-a-posteriori* (MAP) training in classification and regression. The general utility of LA has also motivated its use as an approximation to the marginal likelihood over NN weights. Recent methods, including Daxberger et al. (2021), have explored this path to find model hyperparameters (Immer et al., 2021) or learning invariances (Immer et al., 2022). However, its computational burden, for instance in Hessian matrices, has prescribed diagonal or Kronecker factored approximations (Ritter et al., 2018; Martens and Grosse, 2015a; Botev et al., 2017), which are now widely used for second-order optimization. We provide more details on the connections to existing hessian based methods Daxberger et al.

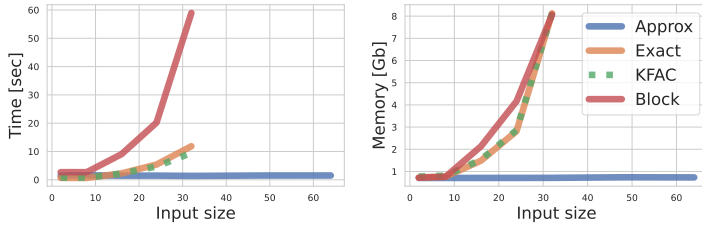


Figure 7: **Memory & time usage of Hessian approximations.** The exact and KFAC scales poorly with the image resolution. In contrast, our proposed approximate diagonal Hessian scales linearly.

(2021); Zhang et al. (2017); Khan et al. (2017), Bayes by Backpropagation Blundell et al. (2015) and Adam Kingma and Ba (2015a) in Appendix C.

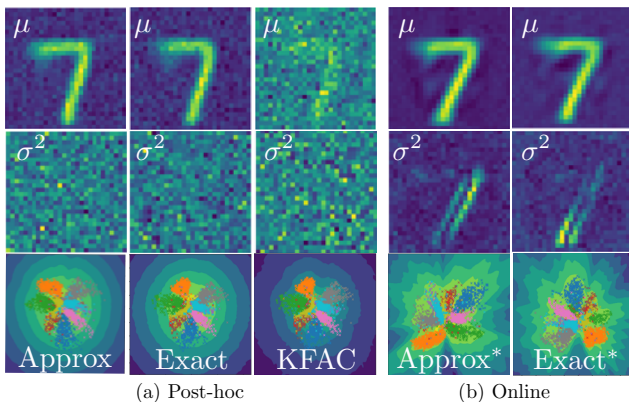
The *full* Bayesian perspective on VAE weights was first explored by Daxberger and Hernández-Lobato (2019) which we find similar in spirit to our work. In contrast to them (1) we follow a principled Bayesian derivation. (2) Neither do we depend on Hamiltonian Monte-Carlo sampling, which is generally hard to scale to efficient training.

5 Experiments

First we demonstrate the computational advantages of the proposed Hessian approximation, and then that our sampling-based training leads to well calibrated uncertainties that can be used for OOD detection, data imputation and semi-supervised learning. For all the downstream tasks we consider the following baselines: AE (Hinton and Salakhutdinov, 2006) with constant and learned variance, VAE (Rezende et al., 2014; Kingma and Welling, 2014), Monte-Carlo dropout AE (Gal and Ghahramani, 2016) and Ensembles of AE (Lakshminarayanan et al., 2017). We extend StochManDetlefsen et al. (2021) with the Hessian backpropagation for the approximate and mixed diagonals. The training code is implemented in PyTorch and available². Appendix A provides more details on the experimental setup.

Efficient Hessian Approximation. For practical applications, training time and memory usage of the Hessian approximation must be kept low. We here show that the proposed approximate diagonal Hessian is sufficient and even outperforms other approximations when combined with our online training.

Fig. 7 show the time and memory requirement for different approximation methods as a function of input size for a 5 layer convolutional network that preserve channel and input dimension. As baselines we use efficient implementations of the exact and KFAC approximation (Daxberger et al., 2021; Dangel et al., 2020). The exact diagonal approximation run out of memory for an $\sim 36 \times 36 \times 3$ image on a 11 Gb NVIDIA GeForce GTX 1080 Ti. In contrast, our approximate diagonal Hessian scales linearly with the resolution, which is especially beneficial for convolutional layers.



Hessian	$\log p(x) \uparrow$	MSE \downarrow
KFAC	-10065.02	125.58
Exact	-299.95	27.21
Approx	-242.24	26.63
Exact*	-26.10	26.06
Approx*	-26.11	26.08

Table 2: Online training (indicated by *) outperforms post-hoc LA. The approximate diagonal always has similar performance with the exact diagonal.

Figure 8: Mean and variance estimates of 100 sampled NN.

Tab. 2 shows that the exact or approximate Hessian diagonal gives similar performance. Using post-hoc LA results in good mean reconstructions (low MSE), but each sampled NN does not give

² <https://github.com/FrederikWarburg/LaplaceAE>

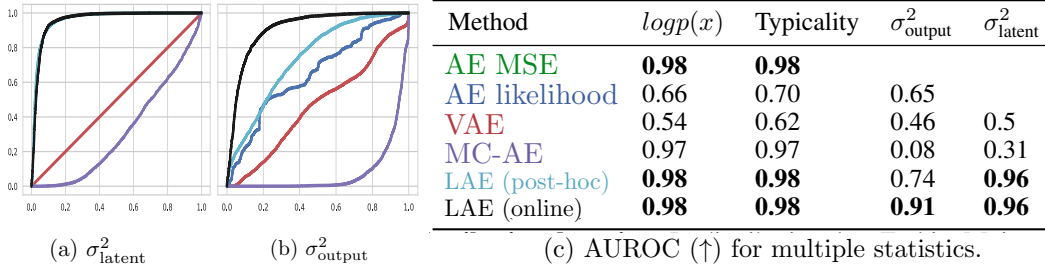


Figure 9: **Out of Distribution detection.** In-distribution data FashionMnist and OOD data Mnist. (a) and (b) shows the ROC curves for latent and output space uncertainties. (c) shows (AUROC \uparrow) for log-likelihood, typicality score, latent space σ_{latent} and output space σ_{output} uncertainties. Online LAE is able to discriminate between in and OOD using the deduced variances in latent and output space.

good reconstructions (low $\log p(x)$). Using our online training procedure results in a much higher log-likelihood. This indicate that every sampled NN predict good reconstructions.

Fig. 8 shows the latent representation, mean and variance of the reconstructions with the KFAC, exact and approximate diagonal for both post-hoc and online setup. Note that the online training makes the uncertainties more well fitted, both in latent and data space. These well-fitted uncertainties have several practical downstream applications, which we demonstrate next.

Out-of-Distribution (OOD) Detection capabilities are critical for identifying distributional shifts, outliers and irregular user inputs, which can hinder the propagation of erroneous decision in an automated system. We evaluate OOD performance on the commonly used benchmarks (Nalisnick et al., 2019b), where we use FASHIONMNIST (Xiao et al., 2017) as in-distribution and MNIST (Lecun et al., 1998) as OOD. Fig. 9 (c) shows that our online LAE outperforms existing models in both log-likelihood and Typicality score (Nalisnick et al., 2019a). This stems from the calibrated model uncertainties, which are exemplified in the models ability to detect OOD examples from the uncertainty deduced in latent and output space; see Fig. 9 (a,b) for ROC curves.

Fig. 10 shows distribution of the output variances for in- and OOD data. This illustrates that using LA improves OoD detection. Furthermore, the online training improves the model calibration.

Missing Data imputation. Another application of stochastic representation learning is to provide distributions over unobserved values (Rezende et al., 2014). In many application domains, sensor readings go missing, which we may mimic by letting parts of an image be unobserved. Rezende et al. (2014) show that we can then draw samples from the distribution of the entire image conditioned on the observed part, by imputing the missing pixels with noise and repeatedly encode and decode while keeping observed pixels fixed. Fig. 11 show samples using this procedure from a VAE, a post-hoc LAE and our online LAE, where we only observe the lower half of an MNIST image. This implies ambiguity about the original digit, e.g. the lower half of a “5” could be a “3” and similarly a “7” could be a “9”. Our LAE captures this ambiguity, which is exemplified by the multi-modal reconstructions from the sampled networks in Fig. 11. The baselines only capture unimodal reconstructions.

Capturing the ambiguity of partly missing data can improve downstream tasks such as the calibration of an image classifier. In Fig. 11 (c) we demonstrate how averaging the predictions of a simple classifier across reconstructions improves standard calibration metrics. This is because the classifier inherits the uncertainty and ambiguity in the learned representations. A deep ensemble of AEs perform similarly to ours, but comes at the cost of training and storing multiple models.

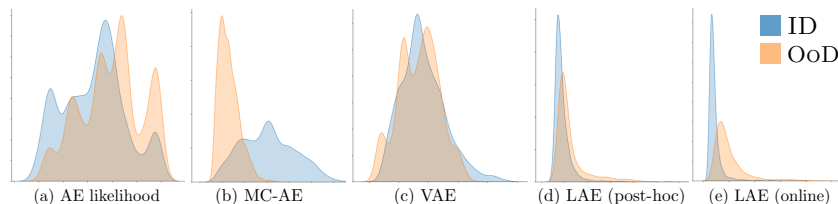


Figure 10: **Histograms of variance** for in- and OOD reconstructions in output space. Note that MC-AE separates the distributions well, but the model assign higher variance to the in- than OOD data.

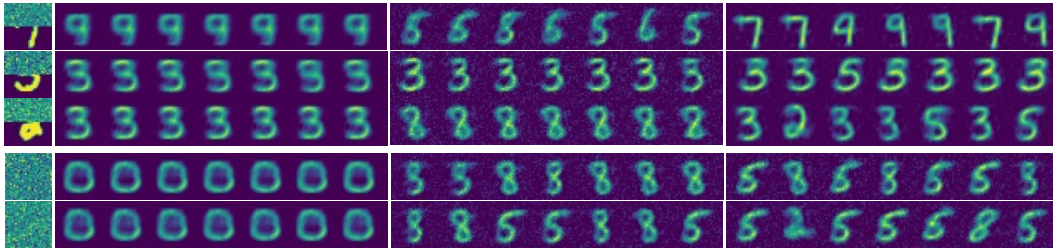


Figure 11: **Missing data imputation & generative capabilities.** Online training of LAE improves representational robustness. This is exemplified by the multimodal behavior in the data imputation (top rows) that accurately model the ambiguity in the data. The bottom two rows, shows that the LAE is able to generate crisp digits from random noise.

Method	MSE ↓	$\log p(x)$ ↑	Acc. ↑	ECE ↓	MCE ↓	RMSCE ↓
Classifier			0.53	0.16	0.25	0.18
VAE	104.73	-104.75	0.22	0.18	0.34	0.19
MC-AE	106.05	-106.05	0.45	0.28	0.38	0.29
AE Ensemble	96.23	-100.94	0.53	0.12	0.2	0.13
LAE (post-hoc)	101.62	-107.25	0.51	0.16	0.27	0.18
LAE (online)	99.59	-106.29	0.53	0.12	0.16	0.13

Table 3: Reconstruction quality measured by the MSE and log-likelihood for the data imputation. Our well-calibrated uncertainties propagates to the MNIST classifier and improves the calibration metrics ECE, MCE and RMSCE.

When the entire input image is missing, the imputation procedure can be seen as a sampling mechanism, such that our LAE can be viewed as a generative model. The bottom rows in Fig. 11 show that the LAE indeed does generate sharp images from a multi-modal distribution.

Semi-supervised learning combines a small amount of label data with a large amount of unlabeled data. The hope is that the structure in the unlabeled data can be used to infer properties of the data that cannot be extracted from a few labeled points. Embedding the same labeled data point multiple times using a stochastic representation scales up the amount of labeled data that is available during training.

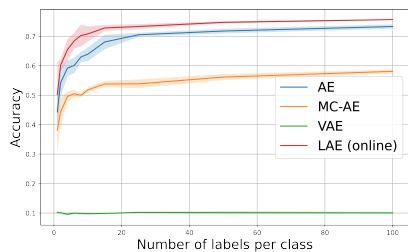


Figure 12: Accuracy as an function of the number of labels per class on MNIST.

Attribute	AE	VAE	MC-AE	LAE (online)
Arched Eyebrows	0.50	0.52	0.55	0.60
Attractive	0.52	0.50	0.49	0.53
Bald	0.98	0.98	0.98	0.98
Wearing Lipstick	0.52	0.49	0.50	0.54
Heavy Makeup	0.45	0.52	0.49	0.56
Overall	0.73	0.72	0.73	0.74

Table 4: Semi-supervised classification accuracy on CELEBA using only 10 labeled datapoints.

Fig. 12 shows the accuracy of a K -nearest neighbour classifier trained on different amounts of labeled data from the MNIST dataset. For all models with a stochastic encoder, we encode each labeled datapoint 100 times and repeat the experiment 5 times. When only few labels per class available (1-20) we clearly observe that our LAE model outperforms all other models, stochastic and deterministic. Increasing the number of labels beyond 100 per class makes the AE and LAE equal in their classification performance with the AE model eventually outperforming the LAE model.

In Tab. 4 we conduct a similar experiment on the CELEBA (Liu et al., 2015) facial dataset, where the the task is to predict 40 different binary labels per datapoint. When evaluating the overall accuracy of predicting all 40 facial attributes, we see no significant difference in performance. However, when we zoom in on specific facial attributes we gain a clear performance advantage over other models.

Limitations. Empirically, the LAE improvements are more significant for overparameterized networks. The additional capacity seems to help the optimizer find a local mode where a Gaussian fit is appropriate. It seems the regularization induced by marginalizing θ compensates for the added flexibility.

6 Conclusion

In this paper we have introduced a Bayesian autoencoder that is realized using Laplace approximations. Unlike current models, this Laplacian autoencoder produces well-behaved uncertainties in both latent and data space. We have proposed a novel variational lower-bound of the autoencoder evidence and an efficient way to compute its Hessian on high dimensional data that scales linearly with data size. Empirically, we demonstrate that our proposed model predicts reliable stochastic representations that are useful for a multitude of downstream tasks: out-of-distribution detection, missing data imputation and semi-supervised classification. Our work opens the way for fully Bayesian representation learning where we can marginalize the representation in downstream tasks. We find this to consistently improve performance.

Acknowledgments and Disclosure of Funding

This work was supported by research grants (15334, 42062) from VILLUM FONDEN. This project has also received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement 757360). This work was funded in part by the Novo Nordisk Foundation through the Center for Basic Machine Learning Research in Life Science (NNF20OC0062606).

— Supplementary Material —

The supplementary material is organised as follows. First, we give more technical details on the experiments. Second, we discuss the Laplace approximation and the “mean shift” issue encountered in non-local maxima. Third, we elaborate more on the difference between related works and our proposed method. Fourth, we present a much more thorough explanation of the proposed method accompanied with relevant proofs. Fifth, we present more details on the hessian derivations.

A Experimental details

All experiments was conducted on one of the following datasets: MNIST (Lecun et al., 1998), FASHIONMNIST (Xiao et al., 2017) and CELEBA (Liu et al., 2015). For training and testing the default splits were used. For validation, we sampled 5000 data points randomly from the training sets. All images were normalized to be in the $[0, 1]$ range and for the CELEBA dataset the images were resized to 64×64 additionally.

If nothing else is stated, the models were trained with the following configuration: we used Adam optimizer (Kingma and Ba, 2015b) with learning rate 0.001 and default Pytorch settings (Paszke et al., 2019). The learning rate was adjusted using *ReduceLROnPlateau* learning rate scheduler with parameters *factor=0.5* and *patience=5*, meaning that the learning rate was halved whenever 5 epochs had passed where the validation loss had not decreased. The mean squared error loss was used as the reconstruction loss in all models. Models where trained until convergence, defined as whenever the validation loss had not decreased for 8 epochs. Models trained on MNIST and FASHIONMNIST used a batch size of 64 and on CELEBA a batch size of 128 was used.

Model specific details:

- VAE: Models were trained with KL-scaling of 0.001. We use two encoders and two decoders, such that the model has twice the number of parameters compared to the other models.
- MC-AE: Models were trained with dropout between all trainable layers with probability $p = 0.2$. We keep the same dropout rate during testing.
- EMSEMBLE-AE: Each ensemble consists of 5 models, each initialised with a different random seed.
- LAE (POSTHOC): For experiments with linear layers we used the Laplace Redux (Daxberger et al., 2021) implementation. For convolutions, we found it necessary to use our proposed hessian approximation. We use a diagonal approximation of the hessian in all experiments. After fitting the hessian, we optimise for the prior precision using the marginal likelihood (Daxberger et al., 2021). We use 100 MC sampling in all experiments.
- LAE (ONLINE): We use the exact diagonal in experiments with linear layers and the mixed diagonal approximation in all experiments with convolutional layers. We use a hessian memory factor of 0.0001 and sample only 1 network per iteration. We found that it was not necessary to optimise for the prior precision when trained online.

A.1 Hessian approximation

We use a linear encoder-decoder with three layers in the encoder and decoder, TANH activation functions, and latent size 2. We choose this architecture as Laplace Redux (Daxberger et al., 2021) supports various hessian approximations for this simple network. We use Laplace Redux for all post-hoc experiments except for the approximate diagonal hessian.

A.2 Out-of-distribution

We use a convolutional encoder-decoder architecture. The encoder consisted of a CONV2D, TANH, MAXPOOL2D, CONV2D, TANH, MAXPOOL2D, LINEAR, TANH, LINEAR, TANH, LINEAR, where the decoder was mirrored but with nearest neighbour Upsampling rather than MAXPOOL2D. We used a latent size 2 in these experiments for all models.

A.3 Missing data imputation

To elaborate on the procedure, we reconstruct 5 samples from the half/fully masked image. For each of these reconstructions, we make 5 more reconstructions and take the average of these reconstructions. The intuition is that the first stage explores the multi-modal behaviour of the reconstructions. In the second stage, the uncertainty of the reconstructed digit is reduced, and each sample will reconstruct the same modality. By averaging over these modalities, we achieve a more crisp reconstruction. We use the same architecture as in the hessian approximation experiment.

A.4 Semi-supervised learning

For the experiments on MNIST, we use a single convergence checkpoint for each model. We use the same model architecture as in the hessian approximation. We did 5 repetitions for each model where we first sampled n labels from each of the 10 classes from the validation set, then embedded the $10 \times n$ data points into latent space, trained a KNN classifier on the embedded points and finally evaluated the accuracy of the classifier on the remaining validation set. This procedure was repeated for different values of n in the $[1, 100]$ range. For the stochastic encoders (VAE, MC-AE, LAE), we repeated the embedding step 100 times with the goal that the uncertainty could help the downstream classification. For the KNN classifier we use cross-validation ($K = 2$) to find the optimal number of nearest neighbours.

For the experiments on CELEBA we repeated the exact same experiments but with a fixed value if $n = 10$. Additionally, the classifier was changed to a multi-label version KNN-classifier to accommodate the multiple binary features in the dataset. For CELEBA we use a convolutional architecture. The encoder consists of 5 convolutional layers with TANH and MAXPOOL2D in between each parametric layer. We use a latent size of 64. The decoder mirrors the encoder, but we replace MAXPOOL2D with nearest neighbor upsampling.

B Laplace Approximation

Laplace approximation is an operator that maps local properties (derivatives) to global properties. The idea is to infer a density on every point based on the curvature in a single point. This is done through a Taylor expansion.

Given a vectorial space Θ of size D , let $P(\theta)$ be an arbitrary distribution on Θ and let $\theta^* \in \Theta$ be an arbitrary point. Consider the second order Taylor expansion of the log density around θ^*

$$\ln P(\theta) = \ln P(\theta^*) + \nabla_{\theta} \ln P(\theta^*)(\theta - \theta^*) + \frac{1}{2}(\theta - \theta^*)^{\top} \nabla_{\theta}^2 \ln P(\theta^*)(\theta - \theta^*) + \mathcal{O}(\|\theta - \theta^*\|^3),$$

where

$$[\nabla_{\theta} \ln P(\theta^*)]_i = \left. \frac{\partial}{\partial \theta_i} \ln P(\theta) \right|_{\theta=\theta^*} \quad [\nabla_{\theta}^2 \ln P(\theta^*)]_{ij} = \left. \frac{\partial^2}{\partial \theta_i \partial \theta_j} \ln P(\theta) \right|_{\theta=\theta^*}.$$

are the first and second order derivatives.

Note that if $P(\theta)$ is a Gaussian, then its log density is a second order polynomial and the second order Taylor expansion is exact. This implies that, starting from a Gaussian, the Laplace approximation can infer the full exact density just from the values of $\nabla_{\theta} \ln P$ and $\nabla_{\theta}^2 \ln P$ in a single point θ^* .

The main takeaway from the Laplace approximation is the strong, intrinsic, tie between the covariance matrix and the negative inverse of the hessian of the log probability: $-(\nabla_{\theta}^2 \ln P(\theta^*))^{-1}$.

B.1 If θ^* is a local maxima

If $\nabla_{\theta} \ln P(\theta^*) = 0$ the Taylor expansion consists of only two terms and the Laplace derivation is easier. We also present this in order to develop intuition, although this case is a subcase of the non-local maxima case. In the next section we will consider the more general setting.

Define the Gaussian

$$\text{LAPLACE}_{\max}(\theta^*; P) := \mathcal{N}(\theta | \mu = \theta^*, \sigma^2 = -(\nabla_{\theta}^2 \ln P(\theta^*))^{-1}), \quad (13)$$

which has density

$$Q(\boldsymbol{\theta}) := \frac{P(\boldsymbol{\theta}^*)}{Z_Q} e^{\frac{1}{2}(\boldsymbol{\theta}-\boldsymbol{\theta}^*)^\top (\nabla_{\boldsymbol{\theta}}^2 \ln P(\boldsymbol{\theta}^*))^{-1}(\boldsymbol{\theta}-\boldsymbol{\theta}^*)},$$

where $Z_Q = P(\boldsymbol{\theta}^*) \sqrt{(-2\pi)^D \det(\nabla_{\boldsymbol{\theta}}^2 \ln P(\boldsymbol{\theta}^*))}$ is the normalizing constant. Then, $Q(\boldsymbol{\theta})$ is a good approximation of $P(\boldsymbol{\theta})$ in the sense that

$$\begin{aligned} \ln Q(\boldsymbol{\theta}) + \ln Z_Q &= \ln P(\boldsymbol{\theta}^*) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top (\nabla_{\boldsymbol{\theta}}^2 \ln P(\boldsymbol{\theta}^*))^{-1}(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \cong \\ &\cong \ln P(\boldsymbol{\theta}^*) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top (\nabla_{\boldsymbol{\theta}}^2 \ln P(\boldsymbol{\theta}^*))^{-1}(\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \mathcal{O}(\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|^3) = \\ &= \ln P(\boldsymbol{\theta}). \end{aligned}$$

Notice that $\boldsymbol{\theta}^*$ is in a local maxima, which ensures that the hessian $\nabla_{\boldsymbol{\theta}}^2 \ln P(\boldsymbol{\theta}^*)$ is negative semi-definite. This in turn ensure that the normalizing constant Z_Q exists and that $\text{LAPLACE}(\boldsymbol{\theta}^*; P)$ is well defined.

B.2 If $\boldsymbol{\theta}^*$ is not a local maxima

In order to proceed to a similar derivation when $\nabla_{\boldsymbol{\theta}} \ln P(\boldsymbol{\theta}^*) \neq 0$, we first rearrange the terms in the Taylor expansion. For a more compact notation, we write $\nabla \ln P$ instead of $\nabla_{\boldsymbol{\theta}} \ln P(\boldsymbol{\theta}^*)$ and $\nabla^2 \ln P$ instead of $\nabla_{\boldsymbol{\theta}}^2 \ln P(\boldsymbol{\theta}^*)$.

$$\begin{aligned} \ln P(\boldsymbol{\theta}) &\cong \ln P(\boldsymbol{\theta}^*) + \nabla \ln P(\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla^2 \ln P(\boldsymbol{\theta} - \boldsymbol{\theta}^*) = \\ &= \ln P(\boldsymbol{\theta}^*) - \frac{1}{2} \nabla \ln P^\top \nabla^2 \ln P^{-1} \nabla \ln P \\ &\quad + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}^* + \nabla^2 \ln P^{-1} \nabla \ln P)^\top \nabla^2 \ln P(\boldsymbol{\theta} - \boldsymbol{\theta}^* + \nabla^2 \ln P^{-1} \nabla \ln P) = \\ &= \ln P(\boldsymbol{\theta}^*) - \frac{1}{2} \nabla \ln P^\top \nabla^2 \ln P^{-1} \nabla \ln P + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_1^*)^\top \nabla^2 \ln P(\boldsymbol{\theta} - \boldsymbol{\theta}_1^*), \end{aligned}$$

where we define the new point $\boldsymbol{\theta}_1^*$ as

$$\boldsymbol{\theta}_1^* = \boldsymbol{\theta}^* - (\nabla_{\boldsymbol{\theta}}^2 \ln P(\boldsymbol{\theta}^*))^{-1} \nabla_{\boldsymbol{\theta}} \ln P(\boldsymbol{\theta}^*). \quad (14)$$

Define the Gaussian

$$\text{LAPLACE}(\boldsymbol{\theta}^*; P) := \mathcal{N}(\boldsymbol{\theta} | \mu = \boldsymbol{\theta}_1^*, \sigma^2 = -(\nabla_{\boldsymbol{\theta}}^2 \ln P(\boldsymbol{\theta}^*))^{-1}), \quad (15)$$

which has density

$$Q(\boldsymbol{\theta}) := \frac{P(\boldsymbol{\theta}^*)}{Z_Q} e^{-\frac{1}{2} \nabla \ln P^\top \nabla^2 \ln P^{-1} \nabla \ln P + \frac{1}{2}(\boldsymbol{\theta}-\boldsymbol{\theta}_1^*)^\top (\nabla_{\boldsymbol{\theta}}^2 \ln P(\boldsymbol{\theta}^*))^{-1}(\boldsymbol{\theta}-\boldsymbol{\theta}_1^*)},$$

where $Z_Q = P(\boldsymbol{\theta}^*) e^{-\frac{1}{2} \nabla \ln P^\top \nabla^2 \ln P^{-1} \nabla \ln P} \sqrt{(-2\pi)^D \det(\nabla_{\boldsymbol{\theta}}^2 \ln P(\boldsymbol{\theta}^*))}$ is the normalizing constant. Then, $Q(\boldsymbol{\theta})$ is a good approximation of $P(\boldsymbol{\theta})$ in the sense that

$$\begin{aligned} \ln Q(\boldsymbol{\theta}) + \ln Z_Q &= \ln P(\boldsymbol{\theta}^*) - \frac{1}{2} \nabla \ln P^\top \nabla^2 \ln P^{-1} \nabla \ln P + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_1^*)^\top \nabla^2 \ln P(\boldsymbol{\theta} - \boldsymbol{\theta}_1^*) = \\ &= \ln P(\boldsymbol{\theta}^*) + \nabla \ln P(\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla^2 \ln P(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \cong \\ &\cong \ln P(\boldsymbol{\theta}^*) + \nabla \ln P(\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla^2 \ln P(\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \mathcal{O}(\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|^3) = \\ &= \ln P(\boldsymbol{\theta}). \end{aligned}$$

We highlight four points:

(1) In order to ensure that the normalizing constant Z_Q exists and consequently that $\text{LAPLACE}(\boldsymbol{\theta}^*; P)$ is well defined, the hessian $\nabla_{\boldsymbol{\theta}}^2 \ln P(\boldsymbol{\theta}^*)$ must be negative semi-definite. Differently from the locally maximal case, this is no longer guaranteed.

(2) The method is numerically unstable if $\nabla_{\boldsymbol{\theta}}^2 \ln P(\boldsymbol{\theta}^*)$ has eigenvalues close to 0. This has been empirically observed to commonly be the case (Sagun et al., 2016).

(3) We emphasize that the Taylor expansion is accurate around θ^* , but the Laplace Gaussian is centered in θ_1^* . This is often referred to as “mean shift” and implies that the sampled parameters from the normal distribution over weights are sampled far away from the actual mean.

(4) Now, assume $\ln P(\theta)$ to be the composition of a function $f(\theta)$ and a loss $l(f)$. If, as in our setting, f is linear and l is concave, then $\ln P(\theta)$ is concave, its hessian is guaranteed to be negative semi-definite, and the Laplace approximation is well-defined.

C Extended related work

Our proposed method’s update rule resemble different methods in the literature. This is rather interesting since, despite arriving at a similar algorithm, they follow different derivation. In this section, we seek to explain the nuances between existing methods and ours.

It is useful to recall our precision update rule

$$\mathbf{H}_{t+1} = (1 - \alpha)\mathbf{H}_t + J_{\theta} f^t(\mathbf{x})^{\top} \Sigma_{\text{data}}^{-1} J_{\theta} f^t(\mathbf{x}) \quad \mathbf{H}_0 = \Sigma_{\text{prior}}^{-1}. \quad (16)$$

To the best of our knowledge, a key conceptual difference with all related works in the literature is **Exact vs. Approximate Hessian**: The second term on the RHS, $J^{\top} \Sigma^{-1} J$ is identical in all the update formulations, but comes from different derivation. In the formulation of Daxberger et al. (2021) (and others) this is an approximation of the hessian, that *happens to guarantee negative definiteness*. In our formulation, this term comes from the linearization. This term is the exact hessian, i.e. no approximation, such that negative definiteness is implied by the linearization. This results in a more intuitive understanding of the linearization error.

C.1 Differences with Laplace Redux

Daxberger et al. (2021) specifically develop a post-hoc method, which assume access to a MAP parameter, nevertheless they also present an online training scheme. This appears to be grounded in ideas from second order optimization, while ours is closer linked to the probabilistic model.

Neglecting their prior optimization procedure, Daxberger et al. (2021) considers

$$\mathbf{H}_{t+1} = \Sigma_{\text{prior}}^{-1} + J_{\theta} f^t(\mathbf{x})^{\top} \Sigma_{\text{data}}^{-1} J_{\theta} f^t(\mathbf{x}). \quad (17)$$

We highlight two main differences.

Recursively Updating the Hessian vs. Using an Uninformed Prior: In our formulation the first term of the RHS is the previous precision, which we “discount” with a forgetting term α that comes from the error inherited by moving away from the previous linearization. In the formulation of Daxberger et al. (2021) this is a prior that is assumed to be in the form $\gamma^2 \mathbb{I}$, where the scalar γ is optimized at every step through maximization of the evidence (Eq. 6 in their paper). This evidence maximization comes with some drawbacks: (1) It tends to make the Laplace approximation overconfident to outliers. They partially address this issue by adding to the evidence an auxiliary term that depends on an OOD dataset, penalizing it (Eq. 12 in their Appendix). (2) Besides inducing the avoidable need for an OOD dataset, this technique have the same pitfall as VAEs, namely uncertainty should be a derived quantity, not a learned one.

Computing the hessian at every iteration vs. every epoch: We update the hessian estimate every iteration. In practice, Daxberger et al. (2021) updates the hessian every epoch. In principle, they could update the hessian more regularly, but to the best of our knowledge, they do not explore this path.

C.2 Differences with Variational Adaptive Newton Method

Khan et al. (2017) propose the following precision update rule

$$\mathbf{H}_{t+1} = \mathbf{H}_t + \alpha J_{\theta} f^t(\mathbf{x})^{\top} \Sigma_{\text{data}}^{-1} J_{\theta} f^t(\mathbf{x}). \quad (18)$$

The updates are very similar to ours, where the only difference is that the scaling is made on the Jacobian product instead of on the previous precision \mathbf{H}_t . Both their and our updates are “additive”

in the sense that the magnitude is increasing, for them strictly monotonically, for us on average depending on the magnitude of α . Thus, in both cases, the variance will approach 0 in the limit, modelling epistemic uncertainty disappearing for infinitely long training.

Their update rule comes from the Variational Optimization setting, that can be viewed as an instance of Variational Inference neglecting the KL term. They highlight strong similarities with Newton’s method.

C.3 Differences with Noisy Natural Gradient

Zhang et al. (2017) propose the precision update rule

$$\mathbf{H}_{t+1} = (1 - \alpha)\mathbf{H}_t + \alpha J_{\theta} f^t(\mathbf{x})^{\top} \Sigma_{\text{data}}^{-1} J_{\theta} f^t(\mathbf{x}). \quad (19)$$

Again, the update rule is very similar to ours, besides the scaling being applied to the Jacobian product. They use a convex sum, which makes the update rule “norm preserving”. Thus, these updates have very different asymptotic behaviour than ours.

Their update rule comes from applying, in the context of Variational Inference, natural gradient to the variational posterior distribution, instead of directly on the parameters space. Natural gradient is deeply connected with LAE and natural parameters highlight the importance of updating the precision matrix instead of the covariance.

They also extend their derivation to the KFAC approximation of the hessian (in place of the exact diagonal), this is made through the use of matrix variate Gaussian. Despite not considered in this work, a similar derivation is in principle feasible in our setting.

C.4 Connection with Adam

Both Zhang et al. (2017) and Khan et al. (2017) highlight strong similarities with the Adam method (Kingma and Ba, 2015a). We share these similarities and we highlight them too. The “connection point” is a noisy version of Adam. Zhang et al. (2017) describe this method and call it “Noisy Adam” (Algorithm 1 in their paper). The difference to the vanilla version is that at each step, instead of using the current parameter, they use a noisy version of it. The noise magnitude is the pseudo-second order term (v_t in the original Adam paper (Kingma and Ba, 2015b)).

We can then interpret Noisy Adam as an instance of our variational setting, specifically where the expectation estimate $\mathbb{E}_{q(\theta)}[\cdot]$ is made through Monte Carlo estimation with $N = 1$ samples. Having the methods in the same setting, we can compare them and highlight the two main differences.

First, we emphasize the difference is in the second order term. Similarly to Zhang et al. (2017) and Khan et al. (2017), we use the diagonal of the hessian (technically the diagonal of the GGN), while Adam uses the pointwise square of the gradient, which they call second raw moment and is intended as a cheaply computable approximation of the hessian.

Another difference is that Adam applies the square root. This is a minor point since, as pointed out by Zhang et al. (2017), this change may affect optimization performance, but does not change the fixed points.

C.5 Connection with Bayes by Backpropagation

Bayes by Backprop (Blundell et al., 2015) can be viewed as a sample-based approximation of Laplace. In order to show this relation, we recall two very powerful equations (Opper and Archambeau, 2009a). Specifically, let μ, Σ be the parameter of a Gaussian distribution $q(\theta) \sim \mathcal{N}(\mu, \Sigma)$, and let $V(\theta)$ be an arbitrary L^2 integrable function (the log-likelihood in our case). Then, we are interested in the derivatives of $\mathbb{E}_{\theta \sim q}[V(\theta)]$. By standard Fourier analysis and integration by part, we have

$$\nabla_{\mu} \mathbb{E}_{\theta \sim q}[V(\theta)] = \mathbb{E}_{\theta \sim q}[\nabla_{\theta} V(\theta)], \quad (20)$$

$$\nabla_{\Sigma} \mathbb{E}_{\theta \sim q}[V(\theta)] = \frac{1}{2} \mathbb{E}_{\theta \sim q}[\nabla_{\theta}^2 V(\theta)]. \quad (21)$$

While the first equation is somehow trivial, the second highlight a very deep relationship. The LHS can be rewritten as

$$\nabla_{\Sigma} \mathbb{E}_{\theta \sim q}[V(\theta)] = \nabla_{\Sigma} \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1)}[V(\mu + \epsilon \Sigma)] = \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1)}[\nabla_{\Sigma} V(\mu + \epsilon \Sigma)]. \quad (22)$$

We can recognize that inside the expectation on the RHS is exactly the update rule for the variance in the bayes by backprop method (Blundell et al., 2015). Thus, we can interpret bayes by backprop as a one-sample Monte Carlo estimation of that expectation value. The equations shows that in the limit of infinite samples, the bayes by backprop update step $\nabla_{\Sigma} V(\mu + \epsilon \Sigma)$ converges to the (averaged) Laplace approximation step $\mathbb{E}_{\theta \sim q}[\nabla_{\theta}^2 V(\theta)]$ up to a factor 2. This motivates both the power and the instability of bayes by backprop.

D Model

D.1 Overview

Let $\mathbb{X} = \mathbb{R}^D$ be the *data* space, let $\mathbb{Y} = \mathbb{R}^D$ be the *reconstruction* space and let $\Theta = \mathbb{R}^P$ be the *parameter* space. Let $\mathbb{F} = (\Theta \rightarrow (\mathbb{X} \rightarrow \mathbb{Y}))$ be the space of operators from Θ to the space of operators from \mathbb{X} to \mathbb{Y} . We will denote a function $f \in \mathbb{F}$ applied to a parameter $\theta \in \Theta$ as $f_{\theta} : \mathbb{X} \mapsto \mathbb{Y}$. This will represent, for example, a NN f with a specific set of parameter (weights) $\theta \in \Theta$, that maps some data $\mathbf{x} \in \mathbb{X}$ to some reconstruction $f_{\theta}(\mathbf{x}) = \mathbf{y} \in \mathbb{Y}$.

Let $\mathbb{X} \times \mathbb{Y} \times \Theta \times \mathbb{F}$ be a probability space. The only assumption we make on this space is

$$p(\mathbf{y}|\mathbf{x}, \theta, f) \sim \mathcal{N}(\mathbf{y}|\mu = f_{\theta}(\mathbf{x}), \sigma^2 = \Sigma) \quad \forall (\mathbf{x}, \theta, f) \in \mathbb{X} \times \Theta \times \mathbb{F} \quad (23)$$

where $\Sigma \in \mathfrak{M}(\mathbb{R}^D \times \mathbb{R}^D)$ is a *fixed* variance matrix. This is a common assumption for regression tasks, and is sometimes referred to as the “data noise” or “reconstruction error”. In this paper, we fix $\Sigma = \mathbb{I}$, but the derivations hold for the general case. With only this assumption, the distribution is undefined and multiple solutions can exist. Thus, we require more assumptions.

A dataset $\mathcal{D} = \{\mathbf{x}_n\}$ is a finite of infinite collection of data $\mathbf{x}_n \in \mathbb{X}$ that is assumed to follow a certain, fixed but unknown, distribution

$$\mathbf{x}_n \sim p(\mathbf{x}). \quad (24)$$

Sacrificing slim notation for the sake of clarity, we introduce an operator $\mathcal{I} : \mathbb{X} \mapsto \mathbb{Y}$. This represent the ideal reconstruction for a given \mathbf{x} . In the standard supervised setting this would be the operator (defined on the dataset only) that maps each data input to its label. In our unsupervised setting, where \mathbb{X} and \mathbb{Y} are *the same space*, the operator \mathcal{I} is simply the identity (Indeed they are *not* the same space, they are isomorphic spaces that we identify through the operator \mathcal{I}). Since \mathcal{I} is the identity, it is often neglected in literature, which can lead to unclear and potentially ambiguous Bayesian derivations. Thus, we choose to adopt this heavier, but more precise notation.

We assume access to a specific $f^{NN} \in \mathbb{F}$. Practically this will be our NN architecture, i.e. an operator that, given a set of parameter $\theta \in \Theta$ give rise to a function from \mathbb{X} to \mathbb{Y} . Having f^{NN} fixed, one may consider f not to be stochastic anymore, we choose to still explicitly condition on f in order to have a clearer notation in later stages. Note that, despite not being covered in this work, a proper stochastic derivation also on the NN architecture should be feasible.

D.2 Objective

The NN’s parameter optimization process in this full Bayesian probabilistic framework can be viewed as: given a fixed $f^{NN} \in \mathbb{F}$, namely the NN architecture, maximise the reconstruction probability of $\mathcal{I}(\mathbf{x}_n)$ over the dataset \mathcal{D}

$$\mathbb{E}_{\mathbf{x}_n \sim p(\mathbf{x})} \left[p(\mathbf{y}|\mathbf{x}_n, f^{NN}) \Big|_{\mathbf{y}=\mathcal{I}(\mathbf{x}_n)} \right] = \sum_{\mathbf{x}_n \in \mathcal{D}} p(\mathbf{y}|\mathbf{x}_n, f^{NN}) \Big|_{\mathbf{y}=\mathcal{I}(\mathbf{x}_n)}, \quad (25)$$

where the untractable $p(\mathbf{y}|\mathbf{x}_n, f^{NN})$ can be expanded in θ and thus related to our hypothesis (23) as

$$p(\mathbf{y}|\mathbf{x}_n, f^{NN}) = \mathbb{E}_{\theta \sim p(\theta|\mathbf{x}_n, f^{NN})} [p(\mathbf{y}|\mathbf{x}_n, \theta, f^{NN})]. \quad (26)$$

Notice that the only unfixed quantity is the distribution on the parameters, which we will optimize for. We are not interested in finding a datapoint-dependant distribution, but rather one that maximise all reconstructions at the same time, i.e. $p(\theta|f^{NN}) = p(\theta|\mathbf{x}_n, f^{NN})$. We can then frame Bayesian

optimization as: find a distribution on parameters such that

$$q(\boldsymbol{\theta}) \in \arg \max_{p(\boldsymbol{\theta}|f^{NN})} \sum_{\mathbf{x}_n \in \mathcal{D}} \mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta}|f^{NN})} \left[p(\mathbf{y}|\mathbf{x}_n, \boldsymbol{\theta}, f^{NN}) \Big|_{\mathbf{y}=\mathcal{I}(\mathbf{x}_n)} \right] \quad (27)$$

$$= \arg \max_{p(\boldsymbol{\theta}|f^{NN})} \sum_{\mathbf{x}_n \in \mathcal{D}} \mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta}|f^{NN})} \left[p(\mathcal{I}(\mathbf{x}_n) | \mathcal{N}(f_{\boldsymbol{\theta}}^{NN}(\mathbf{x}_n), \Sigma)) \right]. \quad (28)$$

Moreover, finding this optima in the space $\Delta(\Theta)$ of all distributions on Θ is not tractable. So, as commonly done, we restrict ourselves to the subset $\mathcal{G}(\Theta) \subset \Delta(\Theta)$ of Gaussians over Θ . Then, a *solution* in our context is

$$q(\boldsymbol{\theta}) \in \arg \max_{q \in \mathcal{G}(\Theta)} \sum_{\mathbf{x}_n \in \mathcal{D}} \mathbb{E}_{\boldsymbol{\theta} \sim q(\boldsymbol{\theta})} \left[p(\mathcal{I}(\mathbf{x}_n) | \mathcal{N}(f_{\boldsymbol{\theta}}^{NN}(\mathbf{x}_n), \Sigma)) \right]. \quad (29)$$

We emphasize that this solution has no guarantees of being unique, but we are interested in finding one of them.

D.3 Joint distribution for a fixed datapoint

Let us first get a better understanding of the joint distribution on $\mathbb{Y} \times \Theta$ conditional to a fixed datapoint $\mathbf{x} \in \mathbb{X}$ and a network architecture $f \in \mathbb{F}$

$$p(\mathbf{y}, \boldsymbol{\theta} | \mathbf{x}, f). \quad (30)$$

This distribution has two *marginals*

$$p(\mathbf{y} | \mathbf{x}, f) \quad (31)$$

$$p(\boldsymbol{\theta} | \mathbf{x}, f) \quad (32)$$

and two

$$p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{x}, f) \quad (33)$$

$$p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{x}, f). \quad (34)$$

These four quantities must satisfy the system of two ‘‘recursive’’ equations

$$\begin{aligned} p(\mathbf{y} | \mathbf{x}, f) &= \int_{\Theta} p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{x}, f) p(\boldsymbol{\theta} | \mathbf{x}, f) d\boldsymbol{\theta} = \\ &= \mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta} | \mathbf{x}, f)} [p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{x}, f)] \end{aligned} \quad (35)$$

$$\begin{aligned} p(\boldsymbol{\theta} | \mathbf{x}, f) &= \int_{\mathbb{Y}} p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{x}, f) p(\mathbf{y} | \mathbf{x}, f) d\mathbf{y} = \\ &= \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y} | \mathbf{x}, f)} [p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{x}, f)]. \end{aligned} \quad (36)$$

If these are satisfied then the joint is a well-defined distribution and we can apply Bayes rule

$$p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{x}, f) p(\boldsymbol{\theta} | \mathbf{x}, f) = p(\mathbf{y}, \boldsymbol{\theta} | \mathbf{x}, f) = p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{x}, f) p(\mathbf{y} | \mathbf{x}, f) \quad (37)$$

which in logarithmic form is

$$\log p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{x}, f) + \log p(\boldsymbol{\theta} | \mathbf{x}, f) = \log p(\boldsymbol{\theta} | \mathbf{y}, \mathbf{x}, f) + \log p(\mathbf{y} | \mathbf{x}, f). \quad (38)$$

We can factor in the assumptions. The ‘‘data noise’’ *assumption* gives us one of the two conditionals:

$$p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{x}, f) \sim \mathcal{N}(\mathbf{y} | \mu = f_{\boldsymbol{\theta}}(\mathbf{x}), \sigma^2 = \Sigma). \quad (39)$$

The ‘‘Gaussian parameter’’ *assumption* gives us one of the two marginals:

$$p(\boldsymbol{\theta} | \mathbf{x}, f) = q^t(\boldsymbol{\theta}) \sim \mathcal{N}(\boldsymbol{\theta} | \mu = \boldsymbol{\theta}_t, \sigma^2 = \mathbf{H}_t^{-1}). \quad (40)$$

With these in place, the joint distribution is uniquely defined. The other marginal, by Eq. 35, is

$$p(\mathbf{y} | \mathbf{x}, f) = \mathbb{E}_{\boldsymbol{\theta} \sim q^t(\boldsymbol{\theta})} [p(\mathbf{y} | \mathcal{N}(f_{\boldsymbol{\theta}}(\mathbf{x}), \Sigma))] \quad (41)$$

and the other conditional, by Bayes rule, is

$$p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{x}, f) = \frac{p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{x}, f)p(\boldsymbol{\theta}|\mathbf{x}, f)}{p(\mathbf{y}|\mathbf{x}, f)}. \quad (42)$$

Despite being uniquely defined, the integral in Eq. 41 is, with a general f , intractable, and so is the joint distribution.

But why do we even care? The intractability of $p(\mathbf{y}|\mathbf{x}, f)$ in Eq. 41 may at first glance appear irrelevant. This is the case, for instance, with bayes by backprop (Blundell et al., 2015) methods. They simply need access to the gradient of this quantity. For this purpose a simple Monte Carlo estimate of the expectation is enough.

On the other hand, we are interested in recovering a meaningful distribution on parameters. This imply that we aim at using Eq. 41 to enforce that Eq. 36 holds. For this purpose we need access to the the density $p(\mathbf{y}|\mathbf{x}, f)$, so a Monte Carlo estimate of Eq. 41 is not enough.

D.4 Linear f

Theorem 1. *Given the data noise assumption from Eq. 39*

$$p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{x}, f) \sim \mathcal{N}(\mathbf{y}|\mu = f_{\boldsymbol{\theta}}(\mathbf{x}), \sigma^2 = \Sigma), \quad (43)$$

given the Gaussian parameter assumption from Eq. 40 for some $\boldsymbol{\theta}_t, \mathbf{H}_t$

$$p(\boldsymbol{\theta}|\mathbf{x}, f) = q^t(\boldsymbol{\theta}) \sim \mathcal{N}(\boldsymbol{\theta}|\mu = \boldsymbol{\theta}_t, \sigma^2 = \mathbf{H}_t^{-1}), \quad (44)$$

assume that f is linear in $\boldsymbol{\theta}$, i.e.

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = f_0(\mathbf{x}) + Jf(\mathbf{x})\boldsymbol{\theta} \quad \forall \boldsymbol{\theta} \in \Theta, \forall \mathbf{x} \in \mathbb{X}. \quad (45)$$

Then the joint distribution is Gaussian itself

$$p(\boldsymbol{\theta}, \mathbf{y}|\mathbf{x}, f) \sim \mathcal{N}((\boldsymbol{\theta}, \mathbf{y})|\mu_t, \Sigma_t) \quad (46)$$

where

$$\mu_t = \begin{pmatrix} \boldsymbol{\theta}_t \\ f_{\boldsymbol{\theta}_t}(\mathbf{x}) \end{pmatrix}, \text{ and } \Sigma_t = \begin{pmatrix} \mathbf{H}_t^{-1} & \Sigma Jf(\mathbf{x})^\top \\ Jf(\mathbf{x})\Sigma & (Jf(\mathbf{x})^\top \mathbf{H}_t Jf(\mathbf{x}))^{-1} + \Sigma \end{pmatrix}.$$

Proof. With the further assumption of linearity of f , we can explicitly carry out the integral in the expectation in Eq. 41

$$\begin{aligned} p(\mathbf{y}|\mathbf{x}, f) &= \mathbb{E}_{\boldsymbol{\theta} \sim q^t(\boldsymbol{\theta})} [p(\mathbf{y}|\mathcal{N}(f_{\boldsymbol{\theta}}(\mathbf{x}), \Sigma))] = \\ &= \int_{\Theta} p(\mathbf{y}|\mathcal{N}(f_{\boldsymbol{\theta}}(\mathbf{x}), \Sigma)) q^t(\boldsymbol{\theta}) d\boldsymbol{\theta} = \\ &= \int_{\Theta} p(\mathbf{y}|\mathcal{N}(f_{\boldsymbol{\theta}}(\mathbf{x}), \Sigma)) p(\boldsymbol{\theta}|\mathcal{N}(\boldsymbol{\theta}_t, \mathbf{H}_t^{-1})) d\boldsymbol{\theta} = \\ &= \int_{\Theta} p(\mathbf{y}|\mathcal{N}(f_0(\mathbf{x}) + Jf(\mathbf{x})\boldsymbol{\theta}, \Sigma)) p(\boldsymbol{\theta}|\mathcal{N}(\boldsymbol{\theta}_t, \mathbf{H}_t^{-1})) d\boldsymbol{\theta} = \\ &= p(\mathbf{y}|\mathcal{N}(f_{\boldsymbol{\theta}_t}(\mathbf{x}), (Jf(\mathbf{x})^\top \mathbf{H}_t Jf(\mathbf{x}))^{-1} + \Sigma)). \end{aligned}$$

We emphasize that as a consequence $\nabla_{\boldsymbol{\theta}}^2 \log p(\mathbf{y}|\mathbf{x}, f)$ is not dependent on $\boldsymbol{\theta}_t$. \square

Having Theorem 1 in place, we can go back to our original problem. We need to deal with a fixed *non-linear* architecture f^{NN} . We can exploit the Theorem 1 by defining f^t : a linearization of f^{NN} with a first order Taylor expansion around $\boldsymbol{\theta}_t$

$$\begin{aligned} f_{\boldsymbol{\theta}}^t(\mathbf{x}) &:= \text{TAYLOR}(f^{NN}, \boldsymbol{\theta}_t)(\mathbf{x}) = \\ &= f_{\boldsymbol{\theta}_t}^{NN}(\mathbf{x}) + J_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}^{NN}(\mathbf{x})(\boldsymbol{\theta} - \boldsymbol{\theta}_t) \end{aligned} \quad (47)$$

and it holds that

$$f_{\boldsymbol{\theta}}^{NN}(\mathbf{x}) = f_{\boldsymbol{\theta}}^t(\mathbf{x}) + \mathcal{O}(\|\boldsymbol{\theta} - \boldsymbol{\theta}_t\|^2). \quad (48)$$

Recalling Eq. 23 both for f^{NN} and for f^t

$$p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}, f^{NN}) \sim \mathcal{N}(\mathbf{y}|\mu = f_{\boldsymbol{\theta}}^{NN}(\mathbf{x}), \sigma^2 = \Sigma) \quad (49)$$

$$p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}, f^t) \sim \mathcal{N}(\mathbf{y}|\mu = f_{\boldsymbol{\theta}}^t(\mathbf{x}), \sigma^2 = \Sigma) \quad (50)$$

that, together with Eq. 48, imply

$$\mathcal{N}(\mathbf{y}|\mu = f_{\boldsymbol{\theta}}^{NN}(\mathbf{x}), \sigma^2 = \Sigma) \sim \mathcal{N}(\mathbf{y}|\mu = f_{\boldsymbol{\theta}}^t(\mathbf{x}) + \mathcal{O}(\|\boldsymbol{\theta} - \boldsymbol{\theta}_t\|^2), \sigma^2 = \Sigma), \quad (51)$$

where we can interpret the unknown $\mathcal{O}(\|\boldsymbol{\theta} - \boldsymbol{\theta}_t\|^2)$ as $\boldsymbol{\theta}$ -dependent noise. More specifically, calling $\gamma > 0$ the scalar constant of the \mathcal{O} -term, we assume that

$$\mathcal{O}(\|\boldsymbol{\theta} - \boldsymbol{\theta}_t\|^2) \approx \epsilon(\boldsymbol{\theta} - \boldsymbol{\theta}_t)^2 \quad \text{where } \epsilon \sim \mathcal{N}(0, \gamma\mathbb{I}) \quad (52)$$

and thus, from Eq. 51, we have

$$\mathcal{N}(\mathbf{y}|\mu = f_{\boldsymbol{\theta}}^{NN}(\mathbf{x}), \sigma^2 = \Sigma) \sim \mathcal{N}(\mathbf{y}|\mu = f_{\boldsymbol{\theta}}^t(\mathbf{x}), \sigma^2 = \Sigma + \gamma\|\boldsymbol{\theta} - \boldsymbol{\theta}_t\|^2\mathbb{I}). \quad (53)$$

Despite at this point integrals are not analytically tractable, and thus a proper proof is not feasible, the intuition is that this increased variance reflects in an increased variance in $p(\boldsymbol{\theta}|\mathbf{x}, f^{NN})$

$$\nabla_{\boldsymbol{\theta}}^2 \log p(\boldsymbol{\theta}|\mathbf{x}, f^{NN})|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t+1}} \approx \nabla_{\boldsymbol{\theta}}^2 \log p(\boldsymbol{\theta}|\mathbf{x}, f^t)|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t+1}} + \gamma\|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|^2, \quad (54)$$

where we introduce an hyperparameter $\alpha > 0$ to cope with this added variance. If we then assume the Jacobian $J_{\boldsymbol{\theta}} f^{NN}$ to be a Lipschitz function, then the Lipschitz constant is an upper bound on γ , as follows from the Taylor expansion of Eq. (47). If this Lipschitz constant is smaller than the inverse of the gradient step $1/\|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|$ (which is not an unreasonable assumption for the gradient ascent to be stable) we have

$$\gamma\|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|^2 \approx \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\| \quad (55)$$

that gives us a plausible order of magnitude for choosing the hyperparameter α .

Motivation: During training, we produce a sequence of Gaussians $q^t(\boldsymbol{\theta}) \sim \mathcal{N}(\boldsymbol{\theta}_t, \mathbf{H}_t^{-1})$ that we assume to be the distribution $p(\boldsymbol{\theta}|\mathbf{x}, f^t)$, at every step $t \geq 0$. This distribution q^t is then used for (1) Gaussian derivation in the linear case, for (2) Monte Carlo sampling in the update rule $\boldsymbol{\theta}_t \rightarrow \boldsymbol{\theta}_{t+1}$ and (3), as second order derivative, for update rule $\mathbf{H}_t \rightarrow \mathbf{H}_{t+1}$.

Moreover, given that this distribution q^t is our "best guess so far", we assume it to be also the distribution $p(\boldsymbol{\theta}|\mathbf{x}, f^{NN})$. This, being f^{NN} not linear, (1) cannot be used for Gaussian derivation, (2) can reasonably be used for sampling (and thus we derive the improved update rule Eq. (64)), and (3) can somehow be used as second order derivative (and thus we derive the improved update rule Eq. (65)), but the latter requires some more care. That is why we introduce the parameter α .

D.5 Iterative learning

Our learning method produces a sequence of Gaussians

$$q^t(\boldsymbol{\theta}) \sim \mathcal{N}(\boldsymbol{\theta}|\mu = \boldsymbol{\theta}_t, \sigma^2 = \mathbf{H}_t^{-1}) \quad (56)$$

and a sequence of linearized functions

$$f^t = \text{TAYLOR}(f^{NN}, \boldsymbol{\theta}_t) \quad (57)$$

for every $t \geq 0$.

Initialization is trivially done using a Gaussian prior on the parameters.

$$\boldsymbol{\theta}_0 = \boldsymbol{\theta}^{\text{prior}} \quad \mathbf{H}_0 = (\Sigma^{\text{prior}})^{-1} \quad (58)$$

Iterative step is made in two steps. First, having access to $\boldsymbol{\theta}_t$, we "generate" the linearization f^t . Practically this is equivalent to computing the two quantities $f_{\boldsymbol{\theta}_t}^{NN}(\mathbf{x})$ and $J_{\boldsymbol{\theta}} f_{\boldsymbol{\theta}_t}^{NN}(\mathbf{x})$, that, together with the value $\boldsymbol{\theta}_t$ are actually equivalent to "generating" f^t , as Eq. 47 shows.

Second, we compute the Gaussian parameters $\boldsymbol{\theta}_{t+1}$ and \mathbf{H}_{t+1} .

Recalling our aim of maximizing the quantity Eq. (29), update on $q(\cdot)$ means, $\boldsymbol{\theta}_t \rightarrow \boldsymbol{\theta}_{t+1}$, is *ideally* made through gradient ascent steps on $p(\mathbf{y}|\mathbf{x}, f)|_{\mathbf{y}=\mathcal{I}(\mathbf{x})}$. As this is intractable, we instead do gradient steps on the lower bound \mathcal{L}_y of (the log of) Eq. 35

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \lambda \nabla_{\boldsymbol{\theta}} \mathcal{L}_y \Big|_{\mathbf{y}=\mathcal{I}(\mathbf{x})} \quad (59)$$

where

$$\mathcal{L}_y = \mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta}|\mathbf{x}, f^t)} [\log p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{x}, f^t)] \leq \log \mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta}|\mathbf{x}, f^t)} [p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{x}, f^t)] \quad (60)$$

and so

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} \mathcal{L}_y \Big|_{\mathbf{y}=\mathcal{I}(\mathbf{x})} &= \mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta}|\mathbf{x}, f^t)} \left[\nabla_{\boldsymbol{\theta}} \log p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{x}, f^t) \Big|_{\mathbf{y}=\mathcal{I}(\mathbf{x})} \right] = \\ &= \mathbb{E}_{\boldsymbol{\theta} \sim q^t(\boldsymbol{\theta})} \left[\nabla_{\boldsymbol{\theta}} \log p(\mathcal{I}(\mathbf{x})|\mathcal{N}(f_{\boldsymbol{\theta}}^t(\mathbf{x}), \Sigma)) \right]. \end{aligned} \quad (61)$$

Recalling the Laplace approximation Eq. (15), the negative precision, $-\mathbf{H}_{t+1}$, is *ideally* set to be the the hessian of the log probability $p(\boldsymbol{\theta}|\mathbf{x}, f)$, evaluated in $\boldsymbol{\theta}_{t+1}$. As this is intractable we instead set it to the hessian of the lower bound \mathcal{L}_{θ} of (the log of) Eq. (36)

$$\mathbf{H}_{t+1} = -\nabla_{\boldsymbol{\theta}}^2 \mathcal{L}_{\theta} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t+1}} \quad (62)$$

where

$$\mathcal{L}_{\theta} = \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, f^t)} [\log p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{x}, f^t)] \leq \log \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, f^t)} [p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{x}, f^t)] \quad (63)$$

and so

$$\begin{aligned} \nabla_{\boldsymbol{\theta}}^2 \mathcal{L}_{\theta} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t+1}} &= \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, f^t)} \left[\nabla_{\boldsymbol{\theta}}^2 \log p(\boldsymbol{\theta}|\mathbf{y}, \mathbf{x}, f^t) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t+1}} \right] \\ &\text{via Eq. (38)} \\ &= \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, f^t)} \left[\nabla_{\boldsymbol{\theta}}^2 \log p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{x}, f^t) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t+1}} + \nabla_{\boldsymbol{\theta}}^2 \log p(\boldsymbol{\theta}|\mathbf{x}, f^t) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t+1}} + \right. \\ &\quad \left. - \nabla_{\boldsymbol{\theta}}^2 \log p(\mathbf{y}|\mathbf{x}, f^t) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t+1}} \right] \\ &\text{via Theorem 1} \\ &= \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, f^t)} \left[\nabla_{\boldsymbol{\theta}}^2 \log p(\mathbf{y}|\boldsymbol{\theta}, \mathbf{x}, f^t) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t+1}} + \nabla_{\boldsymbol{\theta}}^2 \log p(\boldsymbol{\theta}|\mathbf{x}, f^t) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t+1}} \right] \\ &\text{via the chain rule Eq. (67)} \\ &= \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, f^t)} \left[J_{\boldsymbol{\theta}} f^t(\mathbf{x})^\top \nabla_{\mathbf{y}}^2 \log p(\mathbf{y}|\boldsymbol{\theta}_{t+1}, \mathbf{x}, f^t) J_{\boldsymbol{\theta}} f^t(\mathbf{x}) + \nabla_{\boldsymbol{\theta}}^2 \log p(\boldsymbol{\theta}|\mathbf{x}, f^t) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t+1}} \right] \\ &\text{via HP Eq. (23)} \\ &= \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, f^t)} \left[-J_{\boldsymbol{\theta}} f^t(\mathbf{x})^\top \Sigma^{-1} J_{\boldsymbol{\theta}} f^t(\mathbf{x}) + \nabla_{\boldsymbol{\theta}}^2 \log p(\boldsymbol{\theta}|\mathbf{x}, f^t) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t+1}} \right] \\ &= -J_{\boldsymbol{\theta}} f^t(\mathbf{x})^\top \Sigma^{-1} J_{\boldsymbol{\theta}} f^t(\mathbf{x}) + \nabla_{\boldsymbol{\theta}}^2 \log p(\boldsymbol{\theta}|\mathbf{x}, f^t) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t+1}} \\ &= -J_{\boldsymbol{\theta}} f^t(\mathbf{x})^\top \Sigma^{-1} J_{\boldsymbol{\theta}} f^t(\mathbf{x}) + \nabla_{\boldsymbol{\theta}}^2 \log q^t(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{t+1}} \\ &= -J_{\boldsymbol{\theta}} f^t(\mathbf{x})^\top \Sigma^{-1} J_{\boldsymbol{\theta}} f^t(\mathbf{x}) - \mathbf{H}_t. \end{aligned}$$

D.5.1 Improved update rule

As said, the update $\boldsymbol{\theta}_t \rightarrow \boldsymbol{\theta}_{t+1}$ is *ideally* made through gradient ascent steps on $p(\mathbf{y}|\mathbf{x}, f)|_{\mathbf{y}=\mathcal{I}(\mathbf{x})}$, but we instead use the tractable lower bound with f^t . We can perform the same derivation using f^{NN} in place of f^t . Assuming $p(\boldsymbol{\theta}|\mathbf{x}, f^{NN}) \sim q^t(\boldsymbol{\theta})$ for sampling, leads to the improved update rule

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \lambda \mathbb{E}_{\boldsymbol{\theta} \sim q^t(\boldsymbol{\theta})} \left[\nabla_{\boldsymbol{\theta}} \log p(\mathcal{I}(\mathbf{x})|\mathcal{N}(f_{\boldsymbol{\theta}}^{NN}(\mathbf{x}), \Sigma)) \right]. \quad (64)$$

Similarly, the negative precision $-\mathbf{H}_{t+1}$ is *ideally* set to be the hessian of the log probability $p(\boldsymbol{\theta}|\mathbf{x}, f)$, but instead we use the tractable lower bound with f^t . Here we cannot perform the same derivation,

since Theorem 1 does not hold anymore. Instead, we can rely on the estimate Eq. (54) to improve the term

$$\nabla_{\theta}^2 \log p(\theta | \mathbf{x}, f^t) \Big|_{\theta=\theta_{t+1}} = -\mathbf{H}_t \quad \longrightarrow \quad \nabla_{\theta}^2 \log p(\theta | \mathbf{x}, f^{NN}) \Big|_{\theta=\theta_{t+1}} \approx -(1-\alpha)\mathbf{H}_t,$$

and this leads to the improved update rule

$$\mathbf{H}_{t+1} = (1-\alpha)\mathbf{H}_t + J_{\theta} f^t(\mathbf{x})^{\top} \Sigma^{-1} J_{\theta} f^t(\mathbf{x}). \quad (65)$$

E Fast Hessian

We are interested in computing the hessian of a loss function. For this purpose, the Jacobian of the NN w.r.t. parameters plays a crucial role. In this section we develop a better understanding of this object, we derive the backpropagation (also used by the BackPack library (Dangel et al., 2020)) and finally we explain our approximated backpropagation that allows linear scaling.

E.1 Jacobian of a Neural Network

Let us first define some terminology that we will need for chain rule derivations. A NN is a composition of l functions $f := f_{L_l} \circ f_{L_{l-1}} \circ \dots \circ f_{L_2} \circ f_{L_1}$:

$$x_0 \xrightarrow{f_{L_1}} x_1 \longrightarrow \dots \longrightarrow x_{i-1} \xrightarrow{f_{L_i}} x_i \longrightarrow \dots \longrightarrow x_{l-1} \xrightarrow{f_{L_l}} x_l$$

where there are parametric and non-parametric function f_{L_i} . We here highlight two common parametric functions and a common non-parametric function.

Parametric function such as a linear layer

$$x_i = f_{L_i}(x_{i-1}) = \phi_{L_i} x_{i-1} \quad \text{where } \phi_{L_i} \in \mathfrak{M}(|x_i|, |x_{i-1}|)$$

or convolution

$$x_i = f_{L_i}(x_{i-1}) = \text{conv}_{\text{feat}=\phi_{L_i}}(x_{i-1}) \quad \text{where } \phi_{L_i} \in \mathfrak{M}(\text{in channel, out channel, feat height, feat width})$$

Non-parametric such as activation functions $L_i = \tanh, \text{ReLU} \dots$

$$x_i = f_{L_i}(x_{i-1}) \quad \text{where } |x_i| = |x_{i-1}|$$

What is conventionally called *layer* is actually a composition of two function: a linear function and an activation function. For sake of clarity in our derivation we do not adopt this convention and we use the word “layer” to indicate singular “function” component of the NN.

Let us now consider a NN with w parametric layers and $l-w$ activation layers

$$x_l = f_{\phi}(x_0) = f_{L_l} \circ \dots \circ f_{L_1}(x_0) \quad \text{where } \phi = (\phi_1, \dots, \phi_w),$$

and define the bijection

$$\mathcal{W}: \begin{array}{ccc} \{1, \dots, w\} & \longrightarrow & \{i \mid \text{s.t. } L_i \text{ is parametric layer}\} \subseteq \{1, \dots, l\} \\ p & \longmapsto & i \text{ s.t. } L_i \text{ has parameters } \phi_p \end{array}$$

from the subset of parametric layers to the corresponding index in $\phi = (\phi_1, \dots, \phi_w)$.

The Jacobian w.r.t. the parameters $J_{\phi} f_{\phi}(x_0) \in \mathfrak{M}(|x_l|, |\phi_1| + \dots + |\phi_w|)$ is a matrix with number-of-output $|x_l|$ rows and number-of-parameters $|\phi|$ columns. For reference, just storing this matrix can exceed memory limits even with the smallest autoencoder working on MNIST.

There are two ways of looking at this matrix: (1) **row by row**, that is output by output or (2) **block-of-columns by block-of-columns**, that is layer by layer.

E.2 Jacobian per output

Each row of the Jacobian correspond to the gradient w.r.t. the parameters of an element of the output

$$J_{\phi} f_{\phi}(x_0) = \begin{pmatrix} \nabla_{\phi} [f_{\phi}(x_0)]_1 \\ \vdots \\ \nabla_{\phi} [f_{\phi}(x_0)]_{|x_l|} \end{pmatrix}$$

This can be computed by defining $|x_l|$ loss functions

$$\text{loss}_k(x_l) := [x_l]_k \quad \text{for } k = 1, \dots, |x_l|$$

and backpropagating each of those to obtain one line at a time. The disadvantage with this formulation is that we cannot reuse computation for one loss function to improve computation of other loss functions (they are *independent*). Moreover, we need to store all these rows at the same time in order to compute $J^\top J$, which is computationally intractable.

E.3 Jacobian per layer

Each column of the Jacobian is the derivative of the output vector w.r.t. a single parameter. We can then group the parameters (i.e. columns) layer by layer

$$J_{\phi} f_{\phi}(x_0) = \left(\begin{array}{c|c|c} J_{\phi_1} f_{\phi}(x_0) & \dots & J_{\phi_w} f_{\phi}(x_0) \end{array} \right)$$

where $J_{\phi_p} f_{\phi}(x_0) \in \mathfrak{M}(|x_l|, |\phi_p|)$. Let us focus on the computation of $J_{\phi_p} f_{\phi}(x_0)$ for a fixed layer $p = 1, \dots, w$. First notice that the parameters ϕ_p in $f_{\phi} = f_{L_l} \circ \dots \circ f_{L_1}$ only appear in $f_{L_{\mathcal{W}(p)}}$ and so

$$\frac{\partial f_{L_i}}{\partial \phi_p}(x_{i-1}) = 0 \text{ if } i \neq \mathcal{W}(p).$$

Chain rule (informal)

$$\begin{aligned} \frac{\partial f_{\phi}(x_0)}{\partial \phi_p} &= \frac{\partial x_l}{\partial \phi_p} = \\ &= \frac{\partial x_l}{\partial x_{l-1}} \frac{\partial x_{l-1}}{\partial \phi_p} = \\ &= \underbrace{\frac{\partial x_l}{\partial x_{l-1}}}_{\substack{\text{layer } l \\ \text{w.r.t. input}}} \underbrace{\frac{\partial x_{l-1}}{\partial x_{l-2}}}_{\substack{\text{layer } l-1 \\ \text{w.r.t. input}}} \dots \underbrace{\frac{\partial x_{\mathcal{W}(p)+1}}{\partial x_{\mathcal{W}(p)}}}_{\substack{\text{layer } \mathcal{W}(p)+1 \\ \text{w.r.t. input}}} \underbrace{\frac{\partial x_{\mathcal{W}(p)}}{\partial \phi_p}}_{\substack{\text{layer } \mathcal{W}(p) \\ \text{w.r.t. parameters}}} \end{aligned}$$

Chain rule (formal)

$$J_{\phi_p} f_{\phi}(x_0) = \left(\prod_{k=l}^{\mathcal{W}(p)+1} J_{x_{k-1}} f_{L_k}(x_{k-1}) \right) J_{\phi_p} f_{L_{\mathcal{W}(p)}}(x_{\mathcal{W}(p)-1}) \quad (66)$$

The intuition for the chain rule is that the Jacobian $J_{\phi_p} f_{\phi}(x_0)$ is the composition of the Jacobians w.r.t. *input* of subsequent layers times the Jacobian w.r.t. *parameters* of the specific layer. Thus, we can reuse computation for one layer to improve computation of other layer, specifically the product of Jacobians w.r.t. input. Moreover, we can compute $J_p^\top J_p$ layer by layer without ever storing the full Jacobian.

E.3.1 Jacobian of a layer w.r.t. to input

The Jacobian of a standard **linear layer** w.r.t. to the input is

$$J_{x_{p-1}} f_{\phi_p}(x_{p-1}) = \phi_p$$

and this remains the same also in the case with a bias. The Jacobian of a **convolutional layer** w.r.t. to the input is

$$J_{x_{p-1}} \text{conv}_{\text{feat}=\phi_p}(x_{p-1}) = \mathcal{M}(\text{conv}_{\text{feat}=\phi_p})$$

The Jacobian of the activation function depends on the specific choice. Recall that, for each layer i , $x_i \in \mathbb{R}^{|x_i|}$ is a vector

$$x_i = ([x_i]_k)_{k=1, \dots, |x_i|}$$

where $[x_i]_k \in \mathbb{R}$ is the value in position k of the vector x_i .

If L_i is **tanh**

$$[x_i]_k = [f_{L_i}(x_{i-1})]_k = \mathbf{tanh}([x_{i-1}]_k) = \frac{e^{[x_{i-1}]_k} - e^{-[x_{i-1}]_k}}{e^{[x_{i-1}]_k} + e^{-[x_{i-1}]_k}} \quad \text{for } k = 1, \dots, |x_{i-1}|$$

then

$$[J_{x_{i-1}} f_{L_i}(x_{i-1})]_{kj} = \delta_{kj} (1 - (\mathbf{tanh}([x_{i-1}]_k))^2) = \delta_{kj} (1 - [x_i]_k^2) \quad \text{for } k, j = 1, \dots, |x_{i-1}|.$$

If L_i is **ReLU**

$$[x_i]_k = [f_{L_i}(x_{i-1})]_k = \mathbf{ReLU}([x_{i-1}]_k) = \max(0, [x_{i-1}]_k) \quad \text{for } k = 1, \dots, |x_{i-1}|$$

then

$$[J_{x_{i-1}} f_{L_i}(x_{i-1})]_{kj} = \delta_{kj} (1 \text{ if } [x_{i-1}]_k > 0 \text{ else } 0) \quad \text{for } k, j = 1, \dots, |x_{i-1}|.$$

E.3.2 Jacobian of a layer w.r.t. to parameters

The Jacobian of a standard linear layer w.r.t. to the parameters is

$$J_{\phi_i} f_{\phi_i}(x_{i-1}) = \mathbb{I}_{|x_i|} \otimes x_{i-1} \in \mathfrak{M}(|x_i|, |\phi_i|)$$

and in the case with bias $b_i \in \mathbb{R}^{x_i}$ the Jacobian is

$$J_{\phi_i, b_i} f_{\phi_i, b_i}(x_{i-1}) = \mathbb{I}_{|x_i|} \otimes [x_{i-1}, 1] \in \mathfrak{M}(|x_i|, |\phi_i| + |b_i|)$$

The Jacobian of a convolutional layer w.r.t. to the parameters is

$$J_{\phi_i} \mathbf{conv}_{\text{feat}=\phi_i}(x_{i-1}) = J_{\phi_i} \mathbf{conv}_{\text{feat}=\text{rev}(x_{i-1})}^T(\phi_i) = \mathcal{M}(\mathbf{conv}_{\text{feat}=\text{rev}(x_{i-1})})^T \in \mathfrak{M}(|x_i|, |\phi_i|)$$

and in the case with bias $b_i \in \mathbb{R}^{o_i}$ the Jacobian is

$$J_{\phi_i} \mathbf{conv}_{\text{feat}=\phi_i, \text{bias}=b_i}(x_{i-1}) = (\mathcal{M}(\mathbf{conv}_{\text{feat}=\text{rev}(x_{i-1})})^T | I) \in \mathfrak{M}(|x_i|, |\phi_i| + |b_i|)$$

E.4 Hessian of a Neural Network

Consider a function $\mathcal{L} : \mathbb{R}^{|x_i|} \rightarrow \mathbb{R}$ from the output of the NN to scalar value. This later will be interpreted as loss or likelihood, but for now let us stick to the general case.

We are interested in the hessian of this scalar value w.r.t. the parameters of the NN

$$\nabla_{\phi}^2 (\mathcal{L}(f_{\phi}(x_0))) \in \mathfrak{M} \left(\sum_{p=1}^w |\phi_p|, \sum_{p=1}^w |\phi_p| \right).$$

Similarly to the previous section, it is convenient to see this matrix as block matrices, separated layer wise

$$\nabla_{\phi}^2 \mathcal{L}(f_{\phi}(x_0)) = \begin{pmatrix} \nabla_{\phi_1}^2 \mathcal{L}(f_{\phi}(x_0)) & \frac{\partial^2}{\partial \phi_1 \partial \phi_2} \mathcal{L}(f_{\phi}(x_0)) & \dots & \frac{\partial^2}{\partial \phi_1 \partial \phi_w} \mathcal{L}(f_{\phi}(x_0)) \\ \frac{\partial^2}{\partial \phi_2 \partial \phi_1} \mathcal{L}(f_{\phi}(x_0)) & \nabla_{\phi_2}^2 \mathcal{L}(f_{\phi}(x_0)) & & \\ \vdots & & \ddots & \\ \frac{\partial^2}{\partial \phi_w \partial \phi_1} \mathcal{L}(f_{\phi}(x_0)) & & & \nabla_{\phi_w}^2 \mathcal{L}(f_{\phi}(x_0)) \end{pmatrix}$$

First common assumption is to consider layers to be independent to each other, i.e.

$$\frac{\partial^2}{\partial \phi_i \partial \phi_j} \mathcal{L}(f_{\phi}(x_0)) = 0 \quad \forall i \neq j$$

Let us now fix a layer $p = 1, \dots, w$ and focus on a single diagonal block.

$$\nabla_{\phi_p}^2 \mathcal{L}(f_{\phi}(x_0)) \in \mathfrak{M}(|\phi_p|, |\phi_p|).$$

According to the chain rule

$$\nabla_{\phi_p}^2 \mathcal{L}(f_\phi(x_0)) = \underbrace{J_{\phi_p} f_\phi(x_0)^T \cdot \nabla_{x_l}^2 \mathcal{L}(x_l) \cdot J_{\phi_p} f_\phi(x_0)}_{=:G(\phi)} + \sum_{i=1}^{|x_l|} [\nabla_{x_l} \mathcal{L}(x_l)]_i \cdot \nabla_{\phi_p}^2 [f_\phi(x_0)]_i \quad (67)$$

The second term of the RHS is equal to 0 if the model perfectly fits the dataset, $\nabla_{x_l} \mathcal{L}(x_l) = 0$, OR if f is linear in the parameters, $H_{\phi_p} [f_\phi(x_0)]_i = 0$.

The first term of the RHS of Eq. 67, $G(\phi)$, is in literature referred to as Generalized Gauss Newton (GGN) matrix. It can be computed efficiently thanks to the view of the Jacobian as layer by layer. Using equation (66), the expression for the approximated hessian w.r.t. to ϕ_p is then

$$\begin{aligned} G(\phi) &= J_{\phi_p} f_\phi(x_0)^T \cdot H_{x_l} \mathcal{L}(x_l) \cdot J_{\phi_p} f_\phi(x_0) = \\ &= J_{\phi_p} f_{L_{\mathcal{W}(p)}}^T \left(\prod_{k=\mathcal{W}(p)+1}^l J_{x_k} f_{L_k}^T \right) H_{x_l} \mathcal{L}(x_l) \left(\prod_{k=\mathcal{W}(p)+1}^{k=l} J_{x_k} f_{L_k} \right) J_{\phi_p} f_{L_{\mathcal{W}(p)}} \end{aligned}$$

and from this we can build an efficient backpropagation algorithm.

Algorithm 1 Algorithm for $J_\phi f^T \cdot \nabla^2 \mathcal{L} \cdot J_\phi f$

```

M = ∇xl2 ℒ(xl)
for k = l, l - 1, ..., 1 do
  if Lk is parametric with φp (i.e. k = ℳ(p)): then
    Hp = Jφp fLk⊤ · M · Jφp fLk
  end if
  M = Jxk fLk⊤ · M · Jxk fLk
end for
return (H1, ..., Hw)

```

As it is written, each H_p is a matrix $|\phi_p| \times |\phi_p|$ so we technically obtain the hessian in a block-diagonal form

$$\begin{pmatrix} H_1 & & 0 \\ & \ddots & \\ 0 & & H_w \end{pmatrix}$$

if we are interested in the diagonal only, we can construct that by concatenation of the diagonals for each H_p .

E.4.1 Fast approximated version of the Algorithm

The idea is to backpropagate only the diagonal of the matrix M , neglecting all the non-diagonal elements. In a single backpropagation step we have

$$M' = J_{x_p} f_{L_p}^{\top} \cdot M \cdot J_{x_p} f_{L_p}$$

In order to backpropagate the diagonal only, we need to use the operator

$$\text{diag}(M) \mapsto \text{diag}(M')$$

For linear layers and activation function this operator is trivial. For convolutional layer it turns out that this operator is itself a convolution

$$\text{diag}(M') = \text{conv}_{\text{feat}=\phi_p^{(2)}}(\text{diag}(M))$$

where the kernel tensor $\phi_p^{(2)}$ is the pointwise square of the kernel tensor ϕ_p .

E.5 Hessian of a Reconstruction Loss

In the previous Algorithm, the backpropagated quantity is initialized as the hessian of the loss w.r.t. the output of the NN $\nabla_{x_l}^2 \mathcal{L}(x_l)$, or, in an equivalent but more compact notation $\nabla_f^2 \mathcal{L}(f)$. The value of this hessian clearly depends on the specific choice of loss function \mathcal{L} .

The most common choice of likelihood for regression is the Gaussian distribution, while for classification it is the Bernoulli distribution. The Gaussian log likelihood is

$$\mathcal{L}(f) := \log p(x|\mu = f_\theta(x), \sigma^2 = \sigma_d^2) = -\frac{1}{2\sigma_d^2} \|x - f_\theta(x)\|^2 - \log(\sqrt{2\pi}\sigma_d) \quad (68)$$

and its hessian is identity scaled with σ_d

$$\nabla_f^2 \log p(x|\mu = f_\theta(x), \sigma^2 = \sigma_d^2) = -(\sigma_d)^{-2} \mathbb{I} \quad (69)$$

The bernoulli log likelihood is

$$\mathcal{L}(f) := \log p(c|f_\theta(x)) = \log[\text{softmax}(f_\theta(x))]_c = [f_\theta(x)]_c - \log\left(\sum_i e^{[f_\theta(x)]_i}\right) \quad (70)$$

and its hessian can be written in terms of the vector $\pi = \text{softmax}(f_\theta(x))$ of predicted probabilities

$$\nabla_f^2 \log p(c|f_\theta(x)) = -\nabla_f^2 \log\left(\sum_i e^{[f_\theta(x)]_i}\right) = \text{diag}(\pi) - \pi\pi^T \quad (71)$$

We highlight that both Hessians are independent on the label, and thus the GGN is equal to the Fisher matrix, this is true every time $p(y|f_\theta(x))$ is an exponential family distribution with natural parameters $f_\theta(x)$. In this paper, we focus on the Gaussian log likelihood, but we emphasize that the method is not limited to this distribution.

References

- C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra. Weight uncertainty in neural networks. In *Proceedings of the 32nd International Conference on Machine Learning - Volume 37, ICML'15*, page 1613–1622. JMLR.org, 2015.
- A. Botev, H. Ritter, and D. Barber. Practical Gauss-Newton optimisation for deep learning. In *International Conference on Machine Learning (ICML)*, pages 557–565. PMLR, 2017.
- O. Cappé. Online sequential monte carlo em algorithm. In *2009 IEEE/SP 15th Workshop on Statistical Signal Processing*, pages 37–40. IEEE, 2009.
- F. Dangel, F. Kunstner, and P. Hennig. BackPACK: Packing more into backprop. In *International Conference on Learning Representations (ICLR)*, 2020.
- E. Daxberger and J. M. Hernández-Lobato. Bayesian variational autoencoders for unsupervised out-of-distribution detection. *arXiv preprint arXiv:1912.05651*, 2019.
- E. Daxberger, A. Kristiadi, A. Immer, R. Eschenhagen, M. Bauer, and P. Hennig. Laplace Redux - Effortless Bayesian deep learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 34, 2021.
- N. S. Detlefsen, A. Pouplin, C. W. Feldager, C. Geng, D. Kalatzis, H. Hauschultz, M. G. Duque, F. Warburg, M. Miani, and S. Hauberg. Stochman. *GitHub. Note: <https://github.com/MachineLearningLifeScience/stochman/>*, 2021.
- Y. Gal and Z. Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *International Conference on Machine Learning (ICML)*, volume 48, pages 1050–1059, New York, New York, USA, 20–22 Jun 2016. PMLR.
- A. Gelb et al. *Applied optimal estimation*. MIT press, 1974.
- G. E. Hinton and R. R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006.
- A. Immer, M. Bauer, V. Fortuin, G. Rätsch, and K. M. Emtiyaz. Scalable marginal likelihood estimation for model selection in deep learning. In *International Conference on Machine Learning (ICML)*, pages 4563–4573. PMLR, 2021.
- A. Immer, T. F. van der Ouderaa, V. Fortuin, G. Rätsch, and M. van der Wilk. Invariance learning in deep neural networks with differentiable Laplace approximations. *arXiv preprint arXiv:2202.10638*, 2022.
- M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. An introduction to variational methods for graphical models. *Machine Learning*, 37(2):183–233, 1999.
- L. V. Jospin, W. Buntine, F. Boussaid, H. Laga, and M. Bennamoun. Hands-on bayesian neural networks – a tutorial for deep learning users, 2020. URL <https://arxiv.org/abs/2007.06823>.
- M. E. Khan, W. Lin, V. Tangkaratt, Z. Liu, and D. Nielsen. Variational adaptive-newton method for explorative learning, 2017. URL <https://arxiv.org/abs/1711.05560>.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In Y. Bengio and Y. LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015a. URL <http://arxiv.org/abs/1412.6980>.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015b.
- D. P. Kingma and M. Welling. Auto-encoding variational Bayes. In *International Conference on Learning Representations (ICLR)*, 2014.
- B. Lakshminarayanan, A. Pritzel, and C. Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems (NIPS)*, volume 30, 2017.
- P. S. Laplace. Mémoire sur la probabilité des causes par les événements. *Mémoire de l'Académie Royale des Sciences*, 1774.
- Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. doi: 10.1109/5.726791.

- Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- D. J. C. MacKay. Bayesian interpolation. *Neural Computation*, 4(3):415–447, 05 1992.
- D. J. C. MacKay. Probable networks and plausible predictions: A review of practical Bayesian methods for supervised neural networks. *Network: Computation in Neural Systems*, 6(3):469, 1995.
- W. J. Maddox, T. Garipov, P. Izmailov, D. P. Vetrov, and A. G. Wilson. A simple baseline for bayesian uncertainty in deep learning. *Neurips*, abs/1902.02476, 2019.
- J. Martens and R. Grosse. Optimizing neural networks with Kronecker-factored approximate curvature. In *International Conference on Machine Learning (ICML)*, pages 2408–2417. PMLR, 2015a.
- J. Martens and R. B. Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International Conference on Machine Learning (ICML)*, 2015b.
- E. Nalisnick, A. Matsukawa, Y. W. Teh, and B. Lakshminarayanan. Detecting out-of-distribution inputs to deep generative models using typicality, 2019a. URL <https://arxiv.org/abs/1906.02994>.
- E. T. Nalisnick, A. Matsukawa, Y. W. Teh, D. Görür, and B. Lakshminarayanan. Do deep generative models know what they don’t know? In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019b. URL <https://openreview.net/forum?id=H1xwNhCcYm>.
- R. M. Neal. *Bayesian leaning for neural networks*. University of Toronto, 1996.
- M. Opper and C. Archambeau. The Variational Gaussian Approximation Revisited. *Neural Computation*, 21(3): 786–792, 03 2009a. ISSN 0899-7667. doi: 10.1162/neco.2008.08-07-592. URL <https://doi.org/10.1162/neco.2008.08-07-592>.
- M. Opper and C. Archambeau. The variational Gaussian approximation revisited. *Neural Computation*, 21(3): 786–792, 2009b.
- Y. Park, C. Kim, and G. Kim. Variational Laplace autoencoders. In *International Conference on Machine Learning (ICML)*, pages 5032–5041. PMLR, 2019.
- A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 8024–8035. 2019.
- D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning (ICML)*, volume 32 of *Proceedings of Machine Learning Research*, pages 1278–1286. PMLR, 22–24 Jun 2014.
- H. Ritter, A. Botev, and D. Barber. A scalable Laplace approximation for neural networks. In *International Conference on Learning Representations (ICLR)*, volume 6, 2018.
- D. Rumelhart, G. Hinton, and R. Williams. Learning internal representations by error propagation. In *Parallel Distributed Processing: Explorations in the Microstructure of Cognition, vol. 1: Foundations*, pages 318–362. 1986.
- L. Sagun, L. Bottou, and Y. LeCun. Eigenvalues of the hessian in deep learning: Singularity and beyond. *arXiv preprint arXiv:1611.07476*, 2016.
- F. Wenzel, K. Roth, B. Veeling, J. Swiatkowski, L. Tran, S. Mandt, J. Snoek, T. Salimans, R. Jenatton, and S. Nowozin. How good is the Bayes posterior in deep neural networks really? In *International Conference on Machine Learning (ICML)*, pages 10248–10259, 2020.
- H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017.
- G. Zhang, S. Sun, D. Duvenaud, and R. B. Grosse. Noisy natural gradient as variational inference. *CoRR*, 2017. URL <http://arxiv.org/abs/1712.02390>.
- M. Zhang, P. Hayes, and D. Barber. Generalization gap in amortized inference. *Workshop on Bayesian Deep Learning @ NeurIPS*, 2021.