

# VARIATIONAL DEEP LEARNING VIA IMPLICIT REGULARIZATION

**Jonathan Wenger**<sup>†</sup>  
Columbia University

**Beau Coker**<sup>†</sup>  
Columbia University

**Juraj Marusic**  
Columbia University

**John P. Cunningham**  
Columbia University

## ABSTRACT

Modern deep learning models generalize remarkably well in-distribution, despite being overparametrized and trained with little to no *explicit* regularization. Instead, current theory credits *implicit* regularization imposed by the choice of architecture, hyperparameters, and optimization procedure. However, deep neural networks can be surprisingly non-robust, resulting in overconfident predictions and poor out-of-distribution generalization. Bayesian deep learning addresses this via model averaging, but typically requires significant computational resources as well as carefully elicited priors to avoid overriding the benefits of implicit regularization. Instead, in this work, we propose to regularize variational neural networks solely by relying on the *implicit bias of (stochastic) gradient descent*. We theoretically characterize this inductive bias in overparametrized linear models as generalized variational inference and demonstrate the importance of the choice of parametrization. Empirically, our approach demonstrates strong in- and out-of-distribution performance without additional hyperparameter tuning and with minimal computational overhead.

## 1 INTRODUCTION

The success of deep learning across many application domains is, on the surface, remarkable, given that deep neural networks are usually overparameterized and trained with little to no *explicit* regularization. The generalization properties observed in practice have been explained by *implicit* regularization instead, resulting from the choice of architecture [1], hyperparameters [2, 3], and optimizer [4–10]. Notably, the corresponding inductive biases often require no additional computation, in contrast to enforcing a desired inductive bias through explicit regularization.

In the last two decades, there has been an increasing focus on improving the reliability and robustness of deep learning models via (approximately) Bayesian approaches [11] to improve performance on out-of-distribution data [12], in continual learning [13], and in sequential decision-making [14]. However, despite its promise, in practice, Bayesian deep learning can suffer from issues with prior elicitation [15], can be challenging to scale [16], and explicit regularization via a prior combined with approximate inference may result in pathological inductive biases and uncertainty [17–20].

In this work, we demonstrate both theoretically and empirically how to exploit the implicit bias of optimization for approximate inference in probabilistic neural networks, thus regularizing training implicitly rather than explicitly via the prior. This not only narrows the gap to how standard neural networks are trained, but also reduces the computational overhead of training compared to variational inference. More specifically, we propose to learn a variational distribution over the weights of a deep neural network by maximizing the *expected* log-likelihood in analogy to training via maximum likelihood in the standard case. However, in contrast to variational Bayes, there is *no explicit regularization* via a Kullback-Leibler divergence to the prior. Surprisingly, we show theoretically and empirically that training this way does not cause uncertainty to collapse away from the training data, if initialized and parametrized correctly. More so, for overparametrized linear models we rigorously characterize the implicit bias of SGD as generalized variational inference with a 2-Wasserstein regularizer penalizing deviations from the prior. Figure 1 illustrates our approach on a toy example.

<sup>†</sup>Equal contribution.

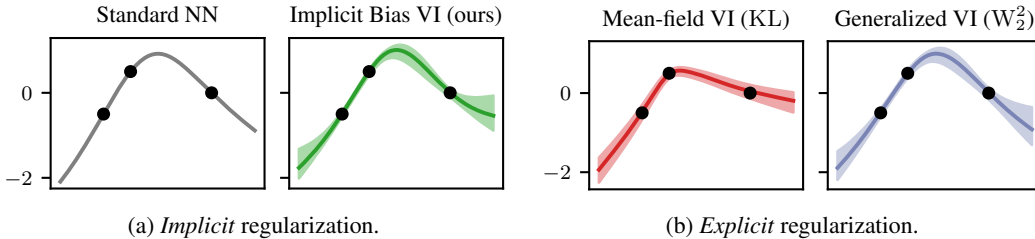


Figure 1: *Variational deep learning via implicit regularization.* Neural networks generalize well without explicit regularization due to implicit regularization from the architecture and optimization. We can exploit this implicit bias for variational deep learning, removing the computational overhead of explicit regularization and narrowing the gap to deep learning practice. As illustrated for a two-hidden layer MLP and proven rigorously for overparameterized linear models in Theorems 1 and 2, the implicit bias of (S)GD in variational networks (see (a)) can be understood as generalized variational inference with a 2-Wasserstein regularizer (see (b)). This differs from the standard ELBO objective with a KL divergence to the prior as used for example in mean-field VI (see (b)).

**Contributions** In this work, we propose a new approach to Bayesian deep learning that generalizes robustly by exploiting the implicit regularization of (stochastic) gradient descent. We fully characterize this implicit bias for regression (Theorem 1) and binary classification (Theorem 2) in overparameterized linear models, generalizing results for non-probabilistic models and drawing a rigorous connection to generalized Bayesian inference. We also demonstrate the importance of the parametrization for the inductive bias and its impact on hyperparameter choice. In several benchmarks, we demonstrate competitive performance to state-of-the-art baselines for Bayesian deep learning, at minimal computational overhead compared to standard neural networks. Finally, we provide an open-source implementation of our approach as a standalone library: [inferno](#).

## 2 BACKGROUND

Given a training dataset  $(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$  of input-output pairs, supervised learning seeks a function  $f_{\mathbf{w}}(\mathbf{x})$  to predict the corresponding output  $y(\mathbf{x})$  for a test input  $\mathbf{x}$ . The parameters  $\mathbf{w} \in \mathbb{R}^P$  of the function are typically trained via empirical risk minimization, i.e.

$$\mathbf{w}_* \in \arg \min_{\mathbf{w}} \ell_r(\mathbf{w}) \quad \text{with} \quad \ell_r(\mathbf{w}) = \ell(\mathbf{y}, f_{\mathbf{w}}(\mathbf{X})) + \lambda r(\mathbf{w}), \quad (1)$$

where the loss  $\ell(\mathbf{y}, f_{\mathbf{w}}(\mathbf{X}))$  encourages fitting the training data and the regularizer  $r(\mathbf{w})$ , given some  $\lambda > 0$ , discourages overfitting, which can lead to poor generalization on test data.

**Implicit Bias of Optimization** One remarkable observation in deep learning is that training overparameterized neural networks ( $P > N$ ) with gradient descent without *explicit* regularization can nonetheless lead to effective (in-distribution) generalization, despite there being many global minima of the loss corresponding to functions  $f_{\mathbf{w}}$  which achieve zero training error [21]. This can be explained by the optimizer, initialization, and parametrization *implicitly* regularizing the optimization problem  $\arg \min_{\mathbf{w}} \ell(\mathbf{y}, f_{\mathbf{w}}(\mathbf{X}))$ , thereby preferring certain global minima [e.g. 4, 5, 7, 22, 23]. Nonetheless, deep neural networks can be surprisingly brittle when predicting *out-of-distribution*, often displaying overconfidence and a significant drop in generalization performance.

**Bayesian Deep Learning** Approximate Bayesian techniques like the Laplace approximation [24–26], stochastic weight averaging [27, 28], deep ensembles [29], and variational approaches [30–33] attempt to address the aforementioned shortcomings of deep learning by learning a distribution over functions as opposed to merely a point estimate. The idea being that a weighted combination of models, all of which achieve low training error, generalizes more robustly while at the same time providing uncertainty quantification.

**Variational Inference** In Bayesian inference this weighted combination is defined by the posterior distribution  $p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}) \propto p(\mathbf{y} \mid \mathbf{X}, \mathbf{w})p(\mathbf{w})$  over weights, induced by a likelihood  $p(\mathbf{y} \mid \mathbf{w})$  and

a choice of prior  $p(\mathbf{w})$  that expresses an explicit preference for some models over others. Approximating the posterior with  $q_\theta(\mathbf{w}) \approx p(\mathbf{w} \mid \mathbf{X}, \mathbf{y})$  by maximizing a lower bound to the log-evidence leads to the following variational optimization problem [34]:

$$\boldsymbol{\theta}_* \in \arg \min_{\boldsymbol{\theta}} \ell_r(\boldsymbol{\theta}) \quad \text{s.t.} \quad \ell_r(\boldsymbol{\theta}) = \mathbb{E}_{q_\theta(\mathbf{w})}(-\log p(\mathbf{y} \mid \mathbf{w})) + \text{KL}(q_\theta(\mathbf{w}) \parallel p(\mathbf{w})). \quad (2)$$

Equation (2) is an instance of the empirical risk minimization objective in Equation (1), with the key difference that one optimizes over variational parameters  $\boldsymbol{\theta}$  of a family of distributions  $q_\theta(\mathbf{w}) \in \mathcal{Q}$ . If that family includes the posterior,  $q_\theta(\mathbf{w}) = p(\mathbf{w} \mid \mathbf{X}, \mathbf{y})$  is the unique global minimum. In the case of a potentially misspecified prior or likelihood, the variational formulation (2) can be generalized to arbitrary loss functions  $\ell$  and statistical distances  $D$  to the prior [35–37], such that

$$\ell_r(\boldsymbol{\theta}) = \mathbb{E}_{q_\theta(\mathbf{w})}(\ell(\mathbf{y}, f_{\mathbf{w}}(\mathbf{X}))) + \lambda D(q_\theta, p). \quad (3)$$

### 3 VARIATIONAL DEEP LEARNING VIA IMPLICIT REGULARIZATION

Our overarching goal is to enable deep neural networks to generalize robustly out-of-distribution without sacrificing their in-distribution performance, at minimal computational overhead. We approach this goal within the framework of Bayesian deep learning, by learning a distribution over neural networks  $f_{\mathbf{w}}$ , induced by a parametrized variational distribution  $q_\theta(\mathbf{w})$  over its weights. However, rather than approximating the Bayesian posterior, which trades off training error against an explicit, a priori preference for certain models, we enforce that *all models have zero training error* while using *implicit* regularization to weight them. Doing so preserves the implicit regularization of the optimizer, which determines the generalization performance of neural networks to a substantial degree, rather than purely relying on explicit regularization induced by the prior. Importantly, this approach leads to robust out-of-distribution generalization, while providing uncertainty quantification at small computational overhead over standard deep learning.

#### 3.1 TRAINING VIA THE EXPECTED LOSS

We propose to train a variational neural network defined by an architecture  $f_{\mathbf{w}}$  and a variational distribution over weights  $q_\theta(\mathbf{w})$  by *minimizing the expected loss*  $\bar{\ell}(\boldsymbol{\theta})$  in analogy to how deep neural networks are usually trained. In other words, the optimal variational parameters are given by

$$\boldsymbol{\theta}_* \in \arg \min_{\boldsymbol{\theta}} \underbrace{\mathbb{E}_{q_\theta(\mathbf{w})}(\ell(\mathbf{y}, f_{\mathbf{w}}(\mathbf{X})))}_{:=\bar{\ell}(\boldsymbol{\theta})} + \lambda D(q_\theta, p). \quad (4)$$

At first glance removing the divergence term from the variational objective in Eq. (3) seems problematic because the new objective is clearly minimized when the variational distribution is a point mass at the minimum loss solution, i.e.  $q_{\boldsymbol{\theta}_*}(\mathbf{w}) = \delta_{\mathbf{w}_*}(\mathbf{w})$  where  $\mathbf{w}_* \in \arg \min_{\mathbf{w}} \ell(\mathbf{y}, f_{\mathbf{w}}(\mathbf{X}))$ . This seemingly defeats the point of a Bayesian deep learning framework, given that there is no variability in predictions on test data. Moreover, the new objective no longer involves a prior distribution, ostensibly removing the ability to manually favor some models over others entirely. The key to understanding our approach is that, in the overparameterized setting, a point mass is only one of many optima corresponding to distributions  $q_{\boldsymbol{\theta}_*}(\mathbf{w})$ , and it is the implicit bias of the optimization procedure that chooses among them. As we will see, if one trains an overparametrized linear model via the expected loss using (stochastic) gradient descent, this implicit bias can be explicitly characterized to depend on the initialization.

#### 3.2 IMPLICIT BIAS OF SGD AS GENERALIZED VARIATIONAL INFERENCE

Assume we train an overparametrized linear model with a Gaussian variational family via the expected loss. For an appropriate learning rate sequence, (stochastic) gradient descent converges to a global minimum  $\boldsymbol{\theta}_*^{\text{GD}} \in \arg \min_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})$  of the training objective. As we show in Section 4, if SGD is *initialized to the prior*, i.e.  $q_{\boldsymbol{\theta}_0}(\mathbf{w}) = p(\mathbf{w})$ , its implicit bias can be understood as selecting the distribution over models with zero training error that is closest to the prior in 2-Wasserstein distance:

$$q_{\boldsymbol{\theta}_*^{\text{GD}}} = \arg \min_{q_\theta} W_2^2(q_\theta, p) \quad \text{s.t. } \boldsymbol{\theta} \in \arg \min_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta})$$

Therefore, we can interpret the implicit bias of (S)GD when training a variational linear model as performing *generalized variational inference*. More precisely, the above is equivalent to  $q_{\theta^{\text{GD}}}$  minimizing the objective in Equation (3) for a certain regularization strength, but with a regularizer that is *not* a KL divergence as it would be for standard variational inference, but rather a 2-Wasserstein distance to the prior. This characterization directly generalizes results for (non-probabilistic) models, where the implicit bias of SGD selects minima that are close to the initialization in Euclidean distance [5, 21]. We therefore call our method **Implicit Bias Variational Inference (IBVI)**. From a practical perspective, by exploiting the implicit regularization of SGD, rather than performing generalized variational inference directly, we no longer need to compute the regularizer explicitly or allocate memory for the prior hyperparameters.<sup>1</sup>

Section 4 provides a detailed version of the regression result introduced here and proves a similar result for binary classification. Our experiments in Section 5 focus on the application to deep neural networks, where we generally expect the implicit regularization to be more complex.

### 3.3 COMPUTATIONAL EFFICIENCY

In practice, we minibatch the expected loss both over training data and parameter samples  $w_m$  drawn from the variational distribution  $q_{\theta}(w)$  such that

$$\bar{\ell}(\theta) = \mathbb{E}_{q_{\theta}(w)}(\ell(\mathbf{y}, f_w(\mathbf{X}))) \approx \frac{1}{N_b M} \sum_{n=1}^{N_b} \sum_{m=1}^M \ell(\mathbf{y}_n, f_{w_m}(\mathbf{x}_n)). \quad (5)$$

The training cost is primarily determined by two factors. The number of parameter samples  $M$  we draw for each evaluation of the objective, and the variational family, which determines the number of additional parameters of the model and the cost for sampling a set of parameters in each forward pass. We wish to keep the overhead compared to a vanilla deep neural network as small as possible.

**Training With A Single Parameter Sample ( $M = 1$ )** When drawing fewer parameter samples  $w_m$  the training objective in Eq. (5) becomes noisier similar to using a smaller batch size. This is concerning since the optimization procedure may not converge given this additional noise. However, one can *train with a single parameter sample only*, simply by reducing the learning rate appropriately, as we show experimentally in Figure 2 and Section S3.2. Therefore, given a set of sampled parameters, the cost of a forward and backward pass is identical to a standard neural network (up to the overhead of the covariance parameters). When using fewer parameter samples in the expected loss, training is unstable unless the learning rate is chosen sufficiently small. For a fixed number of optimizer steps this decreases performance, but either training for more steps, or using momentum closes this gap.

**Variational Family and Covariance Structure** We choose a Gaussian variational distribution  $q_{\theta}(w)$  over (a subset of the) weights of the neural network. While at first glance this may seem restrictive, there is ample evidence that variational families in deep neural networks do not need to be complex to be expressive [38, 39]. In fact, in analogy to deep feedforward NNs with ReLU activations being universal approximators [40], one can show that Bayesian neural networks with ReLU activations and at least one Gaussian hidden layer are universal conditional distribution approximators, meaning they can approximate any continuous conditional distribution arbitrarily well [39]. As we show in Section 4, training an overparametrized linear model with SGD via the expected loss amounts to generalized variational inference *if the covariance is factorized*, i.e.  $\Sigma = \mathbf{S}\mathbf{S}^T$  where  $\mathbf{S} \in \mathbb{R}^{P \times R}$  is a dense matrix with rank  $R \leq P$ . The implicit bias of SGD for arbitrary parametrizations of the covariance matrix remains an open problem. Throughout our experiments we use Gaussian layers with factorized covariances for all architectures.

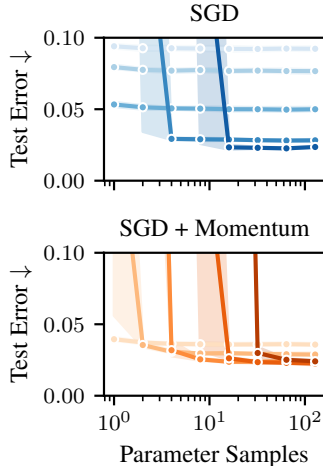


Figure 2: *Training with a single parameter sample given a small enough learning rate.* Lighter color shades correspond to smaller learning rates. See also Section S3.2.

<sup>1</sup>We only need them to initialize the optimizer after which we can free up the allocated memory.

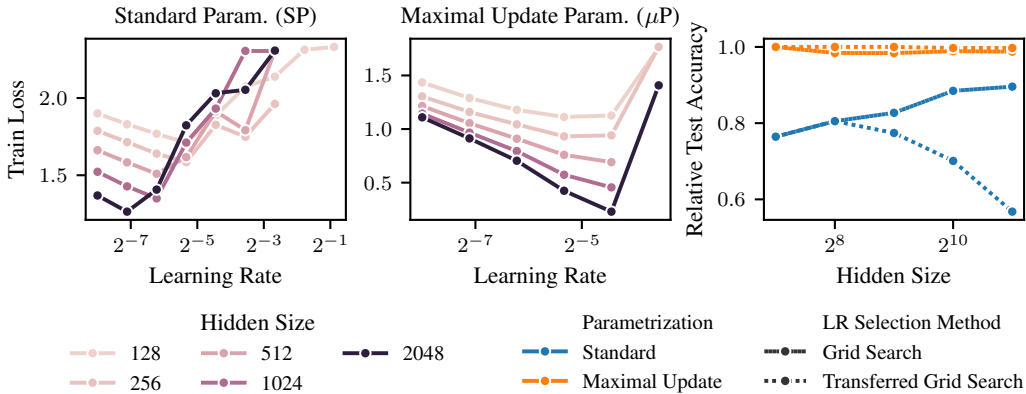


Figure 3: *Hyperparameter Transfer*. When scaling the size of a neural network, one has to re-tune the hyperparameters, such as the learning rate, when using the standard parametrization (SP). The same is true for probabilistic networks as we show here on CIFAR-10 (left). However, when using our proposed extension of the maximal update parametrization ( $\mu P$ ) [41] to probabilistic networks, one can tune the learning rate on a small model and achieve optimal generalization for larger models by “transferring” the optimal learning rate from a smaller model (center and right).

### 3.4 PARAMETRIZATION, FEATURE LEARNING AND HYPERPARAMETER TRANSFER

The inductive bias of SGD depends on the initialization and choice of *parametrization*, a bijective map  $\rho : \Theta' \rightarrow \Theta$  reparametrizing a (variational) model such that  $f_{\theta} \equiv f_{\rho(\theta')}$ . When training deep neural networks, it is not unusual to use layer-specific learning rates. These can be absorbed into the weights of the model and the initialization, meaning they effectively just define a different parametrization [Lemma J.1, 41]. While parameterization is well-studied for non-probabilistic deep learning, it has been identified as one of the “grand challenges of Bayesian computation” [42].

In deep learning, the “standard parameterization” (SP) initializes the weights of a neural network randomly from a distribution with variance  $\propto 1/\text{fan\_in}$  (e.g., as in Kaiming initialization, the PyTorch default) and makes no further adjustments to the forward pass or learning rate. In contrast, the maximal update parametrization ( $\mu P$ ) [43] ensures feature learning even as the width of the network tends to infinity. In addition, under  $\mu P$ , hyperparameters like the learning rate, can be tuned on a small model and transferred to a large-scale model [41].

Given our interpretation of training via the expected loss as generalized variational inference with a prior implied by the parametrization and initialization, a natural question is whether we can extend  $\mu P$  to the variational setting and thus inherit its inductive bias. In the probabilistic setting, feature learning occurs when the *distribution* over hidden units changes from initialization. At any point during training, the  $i$ th hidden unit in layer  $l$  is a function of four random variables: the variational mean and covariance parameters  $(\mu, \mathbf{S})$ , Gaussian noise  $\mathbf{z}$ , and the previous layer hidden units:

$$\mathbf{h}_i^{(l)}(\mathbf{x}) = \mathbf{W}_i \mathbf{h}^{(l-1)}(\mathbf{x}) = (\mu_i + \mathbf{S}_i \mathbf{z}) \mathbf{h}^{(l-1)}(\mathbf{x}). \quad (6)$$

The parameters are random because of the stochasticity in the initialization and/or optimization procedure, while the noise is randomly drawn during each forward pass. Since the  $\mathbf{S}_i \mathbf{z}$  term is a sum over  $R$  terms, where  $R$  is the rank of  $\mathbf{S} \in \mathbb{R}^{P \times R}$ , applying the central limit theorem we propose scaling this term by  $R^{-1/2}$  and then applying  $\mu P$  to the mean and covariance parameters. In practice, we implement the scaling via an adjustment to the covariance initialization and learning rate. Section S2 in the supplement provides empirical investigating of this scaling, demonstrating feature learning in the last hidden layer as the width is increased.

Figure 3 demonstrates that our proposed maximal update parametrization enables hyperparameter transfer in a probabilistic model. We train two-hidden-layer MLPs on CIFAR10, using a low rank covariance in the final two layers. Under standard parametrization (left panel), the learning rate that results in the smallest training loss decreases with hidden size. In contrast, under  $\mu P$  (middle panel), it remains the same across hidden sizes. The right panel of Fig. 3 demonstrates the practical

implications for model selection. For each parametrization and each hidden size  $D$ , we select the learning rate based on a grid search. In “transferred grid search” we do a grid search using the smallest model (hidden size 128) and transfer the best validating learning rate to the hidden size  $D$  model, whereas in “grid search” we perform the grid search on the hidden size  $D$  model. Relative to the test accuracy of the best performing model across learning rate and parametrization, we see that (a)  $\mu$ P outperforms SP, though the gap decreases with hidden size, and (b) the transfer strategy works well for  $\mu$ P but poorly for SP once the hidden size exceeds 256.

### 3.5 RELATED WORK

Variational inference in the context of Bayesian deep learning has seen rapid development in recent years [30–33, 44, 45]. Using a Wasserstein regularizer [37] in the context of generalized VI [36] is arguably most related to our work, given our theoretical results. Structure in the variational parameters has always played an important role for computational reasons [38, 46, 47] and often only a few layers are treated probabilistically [39], with some methods only considering the last layer, effectively treating the neural network as a feature extractor [48, 49]. The Laplace approximation if applied in the last-layer also falls under this category, which has the advantage that it can be applied post-hoc [13, 24–26, 50–54]. Deep ensembles repeat the standard training process using multiple random initializations [29, 55] and have been linked to Bayesian methods [56, 57] with certain caveats [58, 59]. While we use SGD only to optimize the variational parameters and arguably average over samples by using momentum, SGD has also been used widely to directly approximate samples from a target distribution [27, 56, 60, 61], a popular example being stochastic weight averaging (SWA) [27, 28]. Our theoretical analysis extends recent developments on the implicit bias of overparameterized linear models [4, 5, 7] to the probabilistic setting. For classification, works have focused on convergence rates [6], SGD [7], SGD with momentum [8], and the multiclass setting [10]. Results on the implicit bias of neural network training [22] often assume large widths [9, 62–65] allowing similar arguments as for linear models. The former is exemplified by the neural tangent parametrization, under which neural networks behave like kernel methods in the infinite width limit [66]. Yang et al. [41, 43, 64, 65] developed an alternative parameterization that still admits feature learning in the infinite width limit, which we extended to the case of variational networks.

## 4 THEORETICAL ANALYSIS

Consider an overparameterized linear model with a Gaussian prior, trained via the expected loss using (stochastic) gradient descent. We show that, in both regression (Theorem 1) and binary classification (Theorem 2), our approach can be understood as generalized variational inference with a 2-Wasserstein regularizer, which penalizes deviation from the prior among models with zero training error. Theorems 1 and 2 recover analogous results for non-probabilistic models [4, 5, 21].

### 4.1 LINEAR REGRESSION

#### Theorem 1 (Implicit Bias in Regression)

Let  $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}$  be an overparameterized linear model with  $P > N$ . Define a Gaussian prior  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_0, \mathbf{S}_0 \mathbf{S}_0^\top)$  and likelihood  $p(\mathbf{y} | \mathbf{w}) = \mathcal{N}(\mathbf{y}; f_{\mathbf{w}}(\mathbf{X}), \sigma^2 \mathbf{I})$  and assume a variational family  $q_{\boldsymbol{\theta}}(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}, \mathbf{S} \mathbf{S}^\top)$  with  $\boldsymbol{\theta} = (\boldsymbol{\mu}, \mathbf{S})$  such that  $\boldsymbol{\mu} \in \mathbb{R}^P$  and  $\mathbf{S} \in \mathbb{R}^{P \times R}$  where  $R \leq P$ . If the learning rate sequence  $(\eta_t)_t$  is chosen such that the limit point  $\boldsymbol{\theta}_*^{\text{GD}} = \lim_{t \rightarrow \infty} \boldsymbol{\theta}_t^{\text{GD}}$  identified by gradient descent, initialized at  $\boldsymbol{\theta}_0 = (\boldsymbol{\mu}_0, \mathbf{S}_0)$ , is a (global) minimizer of the expected log-likelihood  $\bar{\ell}(\boldsymbol{\theta})$ , then

$$\boldsymbol{\theta}_*^{\text{GD}} \in \underset{\substack{\boldsymbol{\theta}=(\boldsymbol{\mu}, \mathbf{S}) \\ \text{s.t. } \boldsymbol{\theta} \in \arg \min \bar{\ell}(\boldsymbol{\theta})}}{\arg \min}} W_2^2(q_{\boldsymbol{\theta}}, p). \quad (7)$$

Further, this also holds in the case of stochastic gradient descent and when using momentum.

*Proof.* See Section S1.1.1. □

Theorem 1 states that, among those variational parameters which minimize the expected loss, SGD (with momentum) converges to the unique variational distribution which is closest in 2-Wasserstein

distance to the prior. This characterization of the implicit regularization of SGD as generalized variational inference differs from a standard ELBO objective (2) in VI via the choice of regularizer. Since the variational parameters minimize the expected loss in Equation (7), all samples from the predictive distribution interpolate the training data (see Figure 1(b), right panel), the same way a standard neural network would. In contrast, when training with a KL regularizer, the uncertainty does not collapse at the training data (see Figure 1(b), left panel). In fact, a KL regularizer would diverge to infinity for a Gaussian with vanishing variance. Now, for test points that are increasingly out-of-distribution, i.e. less aligned with the span of the training data, the variational predictive matches the prior predictive more closely. Interestingly,  $q_{\theta^{\text{GD}}}$  is equal to the distribution over weights of an ensemble of linear models initialized from the prior and trained independently (see Section S1.1.3). Next, we prove a similar result for binary classification.

## 4.2 BINARY CLASSIFICATION OF LINEARLY SEPARABLE DATA

Consider a binary classification problem with labels  $y_n \in \{-1, 1\}$ , a linear model  $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}$ , and a variational distribution  $q_{\theta}(\mathbf{w})$  with variational parameters  $\theta$ . The expected empirical loss is  $\bar{\ell}(\theta) = \sum_{n \in [N]} \mathbb{E}_{q_{\theta}(\mathbf{w})}(\ell(y_n \mathbf{x}_n^\top \mathbf{w}))$ . We assume without loss of generality<sup>2</sup> that all labels are positive, i.e.  $y_n = 1$  for all  $n$ , and that the dataset is linearly separable.

**Assumption 1** The dataset is *linearly separable*:  $\exists \mathbf{w} \in \mathbb{R}^P$  such that  $\forall n : \mathbf{w}^\top \mathbf{x}_n > 0$ .

For an overparametrized linear model, if  $\mathbf{X} \in \mathbb{R}^{N \times P}$  has full row rank the dataset is guaranteed to be linearly separable.<sup>3</sup> Define the solution to the hard margin SVM, the  $L_2$  max margin vector as

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^P} \|\mathbf{w}\|_2^2 \quad \text{s.t.} \quad \mathbf{w}^\top \mathbf{x}_n \geq 1, \quad (8)$$

and the set of *support vectors*  $\mathcal{S} = \arg \min_{n \in [N]} \mathbf{x}_n^\top \hat{\mathbf{w}}$  indexing the data points on the margin. We make the following additional assumption which is satisfied with high probability under mild assumptions on the training data distribution and degree of overparametrization [67, 68].

**Assumption 2** The SVM support vectors span the dataset:  $\text{span}(\{\mathbf{x}_n\}_{n \in [N]}) = \text{span}(\{\mathbf{x}_n\}_{n \in \mathcal{S}})$ .

We can now characterize the implicit bias in the case of binary classification.

### Theorem 2 (Implicit Bias in Binary Classification)

Let  $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}$  be an (overparametrized) linear model and define a Gaussian prior  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_0, \mathbf{S}_0 \mathbf{S}_0^\top)$ . Assume a variational distribution  $q_{\theta}(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}, \mathbf{S} \mathbf{S}^\top)$  over the weights  $\mathbf{w} \in \mathbb{R}^P$  with variational parameters  $\theta = (\boldsymbol{\mu}, \mathbf{S})$  such that  $\mathbf{S} \in \mathbb{R}^{P \times R}$  and  $R \leq P$ . Assume we are using the exponential loss  $\ell(u) = \exp(-u)$  and optimize the expected empirical loss  $\bar{\ell}(\theta)$  via gradient descent initialized at the prior, i.e.  $\theta_0 = (\boldsymbol{\mu}_0, \mathbf{S}_0)$ , with a sufficiently small learning rate  $\eta$ . Then for almost any dataset which is linearly separable (Assumption 1) and for which the support vectors span the data (Assumption 2), the rescaled gradient descent iterates (rGD)

$$\theta_t^{\text{rGD}} = (\boldsymbol{\mu}_t^{\text{rGD}}, \mathbf{S}_t^{\text{rGD}}) = \left( \frac{1}{\log(t)} \boldsymbol{\mu}_t^{\text{GD}} + \mathbf{P}_{\text{null}(\mathbf{X})} \boldsymbol{\mu}_0, \mathbf{S}_t^{\text{GD}} \right) \quad (9)$$

converge to a limit point  $\theta_*^{\text{rGD}} = \lim_{t \rightarrow \infty} \theta_t^{\text{rGD}}$  for which it holds that

$$\theta_*^{\text{rGD}} \in \arg \min_{\substack{\theta = (\boldsymbol{\mu}, \mathbf{S}) \\ \text{s.t. } \theta \in \Theta_*}} W_2^2(q_{\theta}, p), \quad (10)$$

where the feasible set  $\Theta_* = \{(\boldsymbol{\mu}, \mathbf{S}) \mid \mathbf{P}_{\text{range}(\mathbf{X}^\top)} \boldsymbol{\mu} = \hat{\mathbf{w}} \text{ and } \forall n : \text{Var}_{q_{\theta}}(f_{\mathbf{w}}(\mathbf{x}_n)) = 0\}$  consists of mean parameters which, if projected onto the training data, are equivalent to the  $L_2$  max margin vector and covariance parameters such that there is no uncertainty at training data.

*Proof.* See Section S1.2. □

Theorem 2 states that the mean parameters  $\boldsymbol{\mu}_t$  converge to the  $L_2$  max-margin vector  $\hat{\mathbf{w}}$  in the span of the training data, i.e. the data manifold, and there uncertainty collapses to zero. This is analogous

<sup>2</sup>This is not a restriction since we can always absorb the sign into the inputs, such that  $\mathbf{x}'_n := y_n \mathbf{x}_n$ .

<sup>3</sup>We can always choose  $\mathbf{w} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{1}$ , i.e. the weights linearly interpolating  $\mathbf{y} = \mathbf{1} = (1, \dots, 1)^\top$ .

to the regression case, where zero training loss enforces interpolation of the training data. In the null space of the training data, i.e. off of the data manifold, the model falls back on the prior as enforced by the 2-Wasserstein distance. The assumption of an exponential loss is standard in the literature and we expect this to extend to (binary) cross-entropy in the same way it does in results for standard neural networks [4, 6–8, 10]. Similarly, we conjecture that Theorem 2 can be extended to SGD with momentum [cf. 7, 8]. While Theorem 2 is similar to Theorem 1, there are some subtle differences. First, the feasible set for the minimization problem in Equation (10) is not the set of minima of the expected loss. This is because the exponential function does not have an optimum in contrast to a quadratic function. However, the sequence of variational parameters identified by gradient descent still satisfies  $\lim_{t \rightarrow \infty} \bar{\ell}(\boldsymbol{\theta}_t) = 0$ . Second, without transformation of the mean parameters, the exponential loss results in the mean parameters being unbounded. This necessitates the transformation in Equation (9) as we explain in detail in Section S1.3.

## 5 EXPERIMENTS

We benchmark the *generalization* and *robustness* of our approach, **Implicit Bias VI (IBVI)**, against standard neural networks and several baselines for uncertainty quantification, namely **Temperature Scaling (TS)** [69], **Laplace approximation (LA-GS) & (LA-ML)** [24–26], **Weight-Space VI (WSVI)** [30, 31], **SWA-Gaussian (SWAG)** [28], and **Deep Ensembles (DE)** [29], on a set of standard benchmark datasets for image classification and robustness to input corruptions. We use convolutional architectures (LeNet5 [70] or ResNet34 [71]), which, for all datasets but MNIST, are initialized with pretrained weights except for the input and output layer. All models are trained with SGD with momentum  $\gamma = 0.9$  and a batch size of  $N_b = 128$  for 200 epochs in single precision on an NVIDIA GH200 GPU. Results shown are averaged across five random seeds. A detailed description of the datasets, metrics, models and training can be found in Section S3. An implementation of our method can be found at:

<https://github.com/inferno-ml/inferno>

**In-Distribution Generalization and Uncertainty Quantification** In order to assess the in-distribution generalization, we measure the test error, negative log-likelihood (NLL), and calibration error (ECE) on MNIST, CIFAR10, CIFAR100 and TinyImageNet. As Figure 4 shows for CIFAR100, and Figure S10 for all datasets, the test error for post-hoc methods (**TS**, **LA-GS**, **LA-ML**) is unchanged. As expected, **SWAG** and **IBVI** perform similarly with only **Ensembles** providing an increase in accuracy, but at substantial memory overhead compared to most other approaches. Similarity of **IBVI** to **Ensembles** is perhaps expected in light of their equivalence for linear models (see Proposition S1). In-distribution uncertainty quantification measured in terms of NLL is improved substantially by **TS**, **DE**, and **IBVI**, with only **LA** and **WSVI** showing occasional worsening of NLL compared to the base model. The full results in Figure S10 show that **TS**, **DE**, and **IBVI** consistently are also the best calibrated. As described in Section 3.3, for **IBVI** we train with a single sample only and a probabilistic input and output layer with low-rank covariance, reducing the memory overhead compared to a standard neural network to as little as  $\approx 10\%$  with similar training time (see Figure 4). See Section S3.3.2 for the full experimental results including different parametrizations (SP vs  $\mu$ P).

**Robustness to Input Corruptions** We evaluate the robustness of the different models on MNISTC [72], CIFAR10C, CIFAR100C, and TinyImageNetC [73]. These are corrupted versions of the original datasets, where the images are modified via a set of 15 corruptions, such as impulse noise, blur, pixelation, etc. We selected the maximum severity for each corruption and averaged the performance across all. As expected, the performance of all models drops compared to the in-distribution performance measured on the standard test sets as Figure 5 shows. Besides **DE** which consistently show lower test error, also **IBVI** shows improved accuracy on corrupted data compared to all other approaches. When using the maximal update parametrization, **SWAG** shows good accuracy on the two larger datasets (see Figure S12). **TS**, **DE**, and **IBVI** perform consistently well in terms of uncertainty quantification (both for NLL and ECE) across all datasets, with **LA-ML** being somewhat competitive in terms of NLL. However, compared to the in-distribution setting **IBVI** has better uncertainty quantification than **Ensembles** across all datasets.

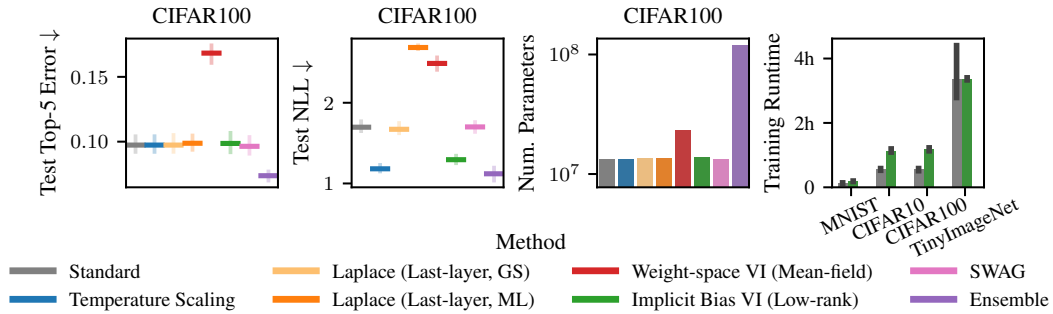


Figure 4: *In-distribution generalization and uncertainty quantification.* Implicit Bias VI (IBVI) has similar test error to other Bayesian deep learning approaches and achieves competitive uncertainty quantification on in-distribution data. While ensembles have improved accuracy, they come at an additional memory overhead. Training a probabilistic model via IBVI has only a minor computational overhead during training, both in time and memory, over standard deep learning.

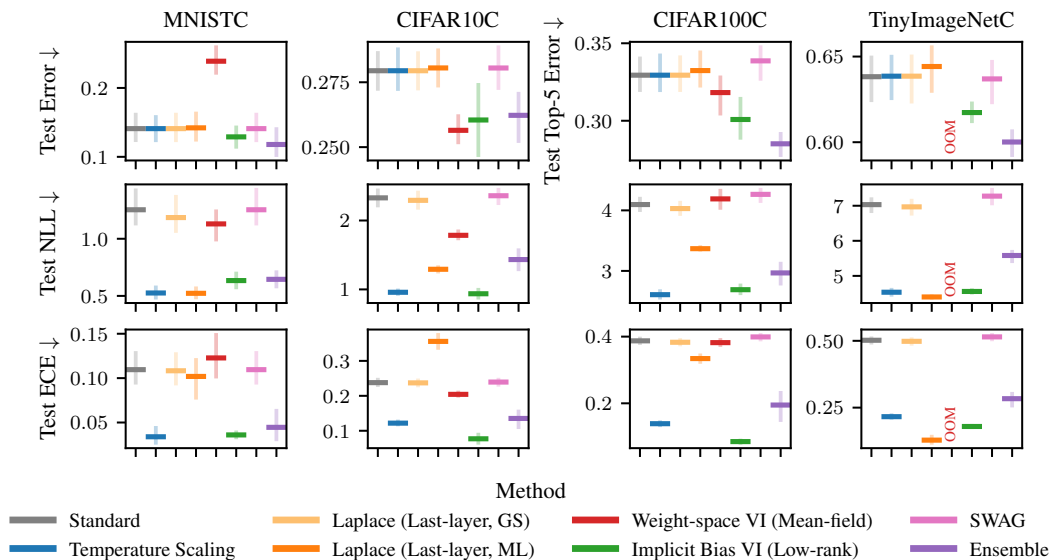


Figure 5: *Generalization on robustness benchmark problems.* When comparing different methods for Bayesian deep learning with regards to robustness to 15 different input corruptions, our approach, Implicit Bias VI, consistently has competitive uncertainty quantification across different datasets and metrics without sacrificing accuracy compared to a non-probabilistic network.

**Limitations** Compared to standard neural networks, when training via Implicit Bias VI, we observed that often lower learning rates were necessary due to the additional stochasticity in the objective (see also Section 3.3). While this does not have a significant impact on generalization, it sometimes requires slightly more epochs to achieve similar in-distribution performance to standard neural networks. Effectively, early in training it takes a bit more time for IBVI to become sufficiently certain about those features which are critical for in-distribution performance. This also means that folk knowledge on learning rate settings for specific architectures may not immediately transfer. In the experiments we train models with probabilistic in- and output layers with our approach, but we have so far not explored other covariance structures or where in the network probabilistic layers are most beneficial. While there is theoretical evidence that even just a single probabilistic hidden layer may be sufficient [39], we believe there is potential for improvement. Beyond the prior induced by the choice of parametrization, we did not experiment with more informative or learned priors, which could potentially give significant performance improvements on certain tasks [15].

## 6 CONCLUSION

In this paper, we demonstrated how to improve the robustness of deep neural networks while quantifying predictive uncertainty by exploiting the implicit regularization of (stochastic) gradient descent. We rigorously characterized this implicit bias for an overparametrized linear model and showed that our approach is equivalent to generalized variational inference with a 2-Wasserstein regularizer at reduced computational cost. We demonstrated the importance of parameterization and how it impacts the inductive bias via the initialization — thus conferring desirable properties such as learning rate transfer. Lastly, we empirically demonstrated competitive performance with state-of-the-art methods for Bayesian deep learning on a set of in- and out-of-distribution benchmarks with minimal computational overhead over standard deep learning. In principle, our approach is not restricted to Gaussian variational families and should seamlessly extend to location-scale families, which could further improve performance. Finally, it would be interesting to explore connections between Implicit Bias VI and Bayesian deep learning in function-space [e.g., 37, 54, 74–76].

## ACKNOWLEDGMENTS

JW, BC, JM and JPC are supported by the Gatsby Charitable Foundation (GAT3708), the Simons Foundation (542963), the NSF AI Institute for Artificial and Natural Intelligence (ARNI: NSF DBI 2229929) and the Kavli Foundation. This work used the DeltaAI system at the National Center for Supercomputing Applications through allocations CIS250340 and CIS250292 from the Advanced Cyberinfrastructure Coordination Ecosystem: Services & Support (ACCESS) program, which is supported by U.S. National Science Foundation grants #2138259, #2138286, #2138307, #2137603, and #2138296. The authors would like to thank Hanna Dettki for valuable input, which significantly improved this paper.

## REFERENCES

- [1] M. Goldblum, M. Finzi, K. Rowan, and A. G. Wilson. “The No Free Lunch Theorem, Kolmogorov Complexity, and the Role of Inductive Biases in Machine Learning”. In: *International Conference on Machine Learning (ICML)*. 2024. DOI: [10.48550/arXiv.2304.05366](https://arxiv.org/abs/2304.05366) (cit. on p. 1).
- [2] M. S. Nacson, R. Mulayoff, G. Ongie, T. Michaeli, and D. Soudry. “The Implicit Bias of Minima Stability in Multivariate Shallow ReLU Networks”. In: *International Conference on Learning Representations (ICLR)*. 2023. DOI: [10.48550/arXiv.2306.17499](https://arxiv.org/abs/2306.17499) (cit. on p. 1).
- [3] R. Mulayoff, T. Michaeli, and D. Soudry. “The Implicit Bias of Minima Stability: A View from Function Space”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2021. URL: <https://proceedings.neurips.cc/paper/2021/hash/944a5ae3483ed5c1e10bbccb7942a279-Abstract.html> (cit. on p. 1).
- [4] D. Soudry, E. Hoffer, M. S. Nacson, S. Gunasekar, and N. Srebro. “The Implicit Bias of Gradient Descent on Separable Data”. In: *Journal of Machine Learning Research (JMLR)* (2018). DOI: [10.48550/arXiv.1710.10345](https://arxiv.org/abs/1710.10345) (cit. on pp. 1, 2, 6, 8, 23, 24, 28, 32).

- [5] S. Gunasekar, J. Lee, D. Soudry, and N. Srebro. “Characterizing Implicit Bias in Terms of Optimization Geometry”. In: *International Conference on Machine Learning (ICML)*. 2018. DOI: [10.48550/arXiv.1802.08246](https://doi.org/10.48550/arXiv.1802.08246) (cit. on pp. 1, 2, 4, 6, 23).
- [6] M. S. Nacson, J. D. Lee, S. Gunasekar, P. H. P. Savarese, N. Srebro, and D. Soudry. “Convergence of Gradient Descent on Separable Data”. In: *International Conference on Artificial Intelligence and Statistics (AISTATS)*. 2019. DOI: [10.48550/arXiv.1803.01905](https://doi.org/10.48550/arXiv.1803.01905) (cit. on pp. 1, 6, 8).
- [7] M. S. Nacson, N. Srebro, and D. Soudry. “Stochastic Gradient Descent on Separable Data: Exact Convergence with a Fixed Learning Rate”. In: *International Conference on Artificial Intelligence and Statistics (AISTATS)*. 2019. DOI: [10.48550/arXiv.1806.01796](https://doi.org/10.48550/arXiv.1806.01796) (cit. on pp. 1, 2, 6, 8).
- [8] B. Wang, Q. Meng, H. Zhang, R. Sun, W. Chen, Z.-M. Ma, and T.-Y. Liu. “Does Momentum Change the Implicit Regularization on Separable Data?” In: *Advances in Neural Information Processing Systems (NeurIPS)* (2022) (cit. on pp. 1, 6, 8).
- [9] H. Jin and G. Montúfar. *Implicit Bias of Gradient Descent for Mean Squared Error Regression with Two-Layer Wide Neural Networks*. arXiv:2006.07356 [stat]. May 2023. DOI: [10.48550/arXiv.2006.07356](https://doi.org/10.48550/arXiv.2006.07356) (cit. on pp. 1, 6).
- [10] H. Ravi, C. Scott, D. Soudry, and Y. Wang. “The Implicit Bias of Gradient Descent on Separable Multiclass Data”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2024. DOI: [10.48550/arXiv.2411.01350](https://doi.org/10.48550/arXiv.2411.01350) (cit. on pp. 1, 6, 8).
- [11] T. Papamarkou, M. Skoularidou, K. Palla, L. Aitchison, J. Arbel, D. Dunson, M. Filippone, V. Fortuin, P. Hennig, J. M. Hernández-Lobato, A. Hubin, A. Immer, T. Karaletsos, M. E. Khan, A. Kristiadi, Y. Li, S. Mandt, C. Nemeth, M. A. Osborne, T. G. J. Rudner, D. Rügamer, Y. W. Teh, M. Welling, A. G. Wilson, and R. Zhang. “Position: Bayesian Deep Learning is Needed in the Age of Large-Scale AI”. In: *International Conference on Machine Learning (ICML)*. 2024. DOI: [10.48550/arXiv.2402.00809](https://doi.org/10.48550/arXiv.2402.00809) (cit. on p. 1).
- [12] D. Tran, J. Liu, M. W. Dusenberry, D. Phan, M. Collier, J. Ren, K. Han, Z. Wang, Z. Mariet, H. Hu, N. Band, T. G. J. Rudner, K. Singhal, Z. Nado, J. v. Amersfoort, A. Kirsch, R. Jenatton, N. Thain, H. Yuan, K. Buchanan, K. Murphy, D. Sculley, Y. Gal, Z. Ghahramani, J. Snoek, and B. Lakshminarayanan. *Plex: Towards Reliability using Pretrained Large Model Extensions*. July 15, 2022. DOI: [10.48550/arXiv.2207.07411](https://doi.org/10.48550/arXiv.2207.07411). arXiv: [2207.07411](https://arxiv.org/abs/2207.07411) [cs]. URL: <http://arxiv.org/abs/2207.07411> (visited on 05/16/2025) (cit. on p. 1).
- [13] H. Ritter, A. Botev, and D. Barber. “Online Structured Laplace Approximations For Overcoming Catastrophic Forgetting”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2018. DOI: [10.48550/arXiv.1805.07810](https://doi.org/10.48550/arXiv.1805.07810) (cit. on pp. 1, 6).
- [14] Y. L. Li, T. G. J. Rudner, and A. G. Wilson. “A Study of Bayesian Neural Network Surrogates for Bayesian Optimization”. In: *International Conference on Learning Representations (ICLR)*. 2024. DOI: [10.48550/arXiv.2305.20028](https://doi.org/10.48550/arXiv.2305.20028) (cit. on p. 1).
- [15] V. Fortuin. “Priors in Bayesian Deep Learning: A Review”. In: *International Statistical Review* 90.3 (2022), pp. 563–591. DOI: [10.1111/insr.12502](https://doi.org/10.1111/insr.12502) (cit. on pp. 1, 10).
- [16] P. Izmailov, S. Vikram, M. D. Hoffman, and A. G. Wilson. “What Are Bayesian Neural Network Posteriors Really Like?” In: *International Conference on Machine Learning (ICML)*. 2021. DOI: [10.48550/arXiv.2104.14421](https://doi.org/10.48550/arXiv.2104.14421) (cit. on p. 1).
- [17] B. Adlam, J. Snoek, and S. L. Smith. *Cold Posteriors and Aleatoric Uncertainty*. July 31, 2020. DOI: [10.48550/arXiv.2008.00029](https://doi.org/10.48550/arXiv.2008.00029). arXiv: [2008.00029](https://arxiv.org/abs/2008.00029) [stat]. URL: <http://arxiv.org/abs/2008.00029> (visited on 05/15/2025) (cit. on p. 1).
- [18] T. Cinquin, A. Immer, M. Horn, and V. Fortuin. “Pathologies in priors and inference for Bayesian transformers”. In: *NeurIPS Bayesian Deep Learning Workshop*. 2021. DOI: [10.48550/arXiv.2110.04020](https://doi.org/10.48550/arXiv.2110.04020) (cit. on p. 1).
- [19] B. Coker, W. P. Bruinsma, D. R. Burt, W. Pan, and F. Doshi-Velez. “Wide Mean-Field Bayesian Neural Networks Ignore the Data”. In: *International Conference on Artificial Intelligence and Statistics (AISTATS)*. 2022. DOI: [10.48550/arXiv.2202.11670](https://doi.org/10.48550/arXiv.2202.11670) (cit. on p. 1).
- [20] A. Y. K. Foong, D. R. Burt, Y. Li, and R. E. Turner. “On the Expressiveness of Approximate Inference in Bayesian Neural Networks”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2020. DOI: [10.48550/arXiv.1909.00719](https://doi.org/10.48550/arXiv.1909.00719) (cit. on p. 1).

- [21] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals. “Understanding deep learning requires rethinking generalization”. In: *International Conference on Learning Representations (ICLR)*. 2017. DOI: [10.48550/arXiv.1611.03530](https://doi.org/10.48550/arXiv.1611.03530) (cit. on pp. 2, 4, 6).
- [22] G. Vardi. “On the Implicit Bias in Deep-Learning Algorithms”. In: *Commun. ACM* 66.6 (May 2023), pp. 86–93. DOI: [10.1145/3571070](https://doi.org/10.1145/3571070) (cit. on pp. 2, 6).
- [23] B. Vasudeva, P. Deora, and C. Thrampoulidis. *Implicit Bias and Fast Convergence Rates for Self-attention*. 2024. DOI: [10.48550/arXiv.2402.05738](https://doi.org/10.48550/arXiv.2402.05738) (cit. on p. 2).
- [24] D. J. C. MacKay. “A Practical Bayesian Framework for Backpropagation Networks”. In: *Neural Computation* 4 (1992). ISSN: 0899-7667, 1530-888X. DOI: [10.1162/neco.1992.4.3.448](https://doi.org/10.1162/neco.1992.4.3.448) (cit. on pp. 2, 6, 8).
- [25] H. Ritter, A. Botev, and D. Barber. “A Scalable Laplace Approximation for Neural Networks”. In: *International Conference on Learning Representations (ICLR)*. 2018 (cit. on pp. 2, 6, 8).
- [26] E. Daxberger, A. Kristiadi, A. Immer, R. Eschenhagen, M. Bauer, and P. Hennig. “Laplace Redux – Effortless Bayesian Deep Learning”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2021. DOI: [10.48550/arXiv.2106.14806](https://doi.org/10.48550/arXiv.2106.14806) (cit. on pp. 2, 6, 8, 41).
- [27] P. Izmailov, D. Podoprikin, T. Garipov, D. Vetrov, and A. G. Wilson. “Averaging Weights Leads to Wider Optima and Better Generalization”. In: *Conference on Uncertainty in Artificial Intelligence (UAI)*. 2018. URL: <https://arxiv.org/abs/1803.05407v3> (cit. on pp. 2, 6).
- [28] W. Maddox, T. Garipov, P. Izmailov, D. Vetrov, and A. G. Wilson. “A Simple Baseline for Bayesian Uncertainty in Deep Learning”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2019. DOI: [10.48550/arXiv.1902.02476](https://doi.org/10.48550/arXiv.1902.02476) (cit. on pp. 2, 6, 8, 42).
- [29] B. Lakshminarayanan, A. Pritzel, and C. Blundell. “Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2017. DOI: [10.48550/arXiv.1612.01474](https://doi.org/10.48550/arXiv.1612.01474). URL: <http://arxiv.org/abs/1612.01474> (cit. on pp. 2, 6, 8, 42).
- [30] A. Graves. “Practical Variational Inference for Neural Networks”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2011. URL: [https://papers.nips.cc/paper\\_files/paper/2011/hash/7eb3c8be3d411e8ebfab08eba5f49632-Abstract.html](https://papers.nips.cc/paper_files/paper/2011/hash/7eb3c8be3d411e8ebfab08eba5f49632-Abstract.html) (cit. on pp. 2, 6, 8, 42).
- [31] C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra. “Weight Uncertainty in Neural Networks”. In: *International Conference on Machine Learning (ICML)*. 2015. DOI: [10.48550/arXiv.1505.05424](https://doi.org/10.48550/arXiv.1505.05424) (cit. on pp. 2, 6, 8, 42).
- [32] K. Osawa, S. Swaroop, A. Jain, R. Eschenhagen, R. E. Turner, R. Yokota, and M. E. Khan. “Practical Deep Learning with Bayesian Principles”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2019. DOI: [10.48550/arXiv.1906.02506](https://doi.org/10.48550/arXiv.1906.02506) (cit. on pp. 2, 6).
- [33] Y. Shen, N. Daheim, B. Cong, P. Nickl, G. M. Marconi, C. Bazan, R. Yokota, I. Gurevych, D. Cremers, M. E. Khan, and T. Möllenhoff. “Variational Learning is Effective for Large Deep Networks”. In: *International Conference on Machine Learning (ICML)*. 2024. DOI: [10.48550/arXiv.2402.17641](https://doi.org/10.48550/arXiv.2402.17641) (cit. on pp. 2, 6).
- [34] A. Zellner. “Optimal Information Processing and Bayes’s Theorem”. In: *The American Statistician* 42.4 (1988), pp. 278–280. DOI: [10.2307/2685143](https://doi.org/10.2307/2685143) (cit. on p. 3).
- [35] P. G. Bissiri, C. Holmes, and S. Walker. “A General Framework for Updating Belief Distributions”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 78.5 (Nov. 2016), pp. 1103–1130. ISSN: 1369-7412, 1467-9868. DOI: [10.1111/rssb.12158](https://doi.org/10.1111/rssb.12158) (cit. on p. 3).
- [36] J. Knoblauch, J. Jewson, and T. Damoulas. “An Optimization-centric View on Bayes’ Rule: Reviewing and Generalizing Variational Inference”. In: *Journal of Machine Learning Research (JMLR)* 23.132 (2022), pp. 1–109. ISSN: 1533-7928. URL: <http://jmlr.org/papers/v23/19-1047.html> (cit. on pp. 3, 6).
- [37] V. D. Wild, R. Hu, and D. Sejdinovic. “Generalized Variational Inference in Function Spaces: Gaussian Measures meet Bayesian Deep Learning”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. Oct. 2022. DOI: [10.48550/arXiv.2205.06342](https://doi.org/10.48550/arXiv.2205.06342) (cit. on pp. 3, 6, 10).

- [38] S. Farquhar, L. Smith, and Y. Gal. “Liberty or Depth: Deep Bayesian Neural Nets Do Not Need Complex Weight Posterior Approximations”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2020. DOI: [10.48550/arXiv.2002.03704](https://doi.org/10.48550/arXiv.2002.03704). URL: <http://arxiv.org/abs/2002.03704> (cit. on pp. 4, 6).
- [39] M. Sharma, S. Farquhar, E. Nalisnick, and T. Rainforth. “Do Bayesian Neural Networks Need To Be Fully Stochastic?”. In: *International Conference on Artificial Intelligence and Statistics (AISTATS)*. 2023. DOI: [10.48550/arXiv.2211.06291](https://doi.org/10.48550/arXiv.2211.06291) (cit. on pp. 4, 6, 10).
- [40] B. Hanin and M. Sellke. *Approximating Continuous Functions by ReLU Nets of Minimal Width*. arXiv:1710.11278 [stat]. Mar. 2018. DOI: [10.48550/arXiv.1710.11278](https://doi.org/10.48550/arXiv.1710.11278). URL: <http://arxiv.org/abs/1710.11278> (cit. on p. 4).
- [41] G. Yang, E. J. Hu, I. Babuschkin, S. Sidor, X. Liu, D. Farhi, N. Ryder, J. Pachocki, W. Chen, and J. Gao. “Tensor Programs V: Tuning Large Neural Networks via Zero-Shot Hyperparameter Transfer”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2021. DOI: [10.48550/arXiv.2203.03466](https://doi.org/10.48550/arXiv.2203.03466) (cit. on pp. 5, 6, 36, 37).
- [42] A. Bhattacharya, A. Linero, and C. J. Oates. “Grand Challenges in Bayesian Computation”. In: *Bulletin of the International Society for Bayesian Analysis (ISBA)* 31.3 (Sept. 2024). DOI: [10.48550/arXiv.2410.00496](https://doi.org/10.48550/arXiv.2410.00496) (cit. on p. 5).
- [43] G. Yang and E. J. Hu. “Tensor Programs IV: Feature Learning in Infinite-Width Neural Networks”. In: *International Conference on Machine Learning (ICML)*. 2021. DOI: [10.48550/arXiv.2011.14522](https://doi.org/10.48550/arXiv.2011.14522) (cit. on pp. 5, 6, 34).
- [44] G. Zhang, S. Sun, D. Duvenaud, and R. Grosse. *Noisy Natural Gradient as Variational Inference*. Feb. 26, 2018. DOI: [10.48550/arXiv.1712.02390](https://doi.org/10.48550/arXiv.1712.02390). arXiv: [1712.02390](https://arxiv.org/abs/1712.02390)[cs]. URL: <http://arxiv.org/abs/1712.02390> (visited on 05/15/2025) (cit. on p. 6).
- [45] M.-N. Tran, N. Nguyen, D. Nott, and R. Kohn. *Bayesian Deep Net GLM and GLMM*. May 25, 2018. DOI: [10.48550/arXiv.1805.10157](https://doi.org/10.48550/arXiv.1805.10157). arXiv: [1805.10157](https://arxiv.org/abs/1805.10157)[stat]. URL: <http://arxiv.org/abs/1805.10157> (visited on 05/15/2025) (cit. on p. 6).
- [46] C. Louizos and M. Welling. *Structured and Efficient Variational Deep Learning with Matrix Gaussian Posteriors*. June 23, 2016. DOI: [10.48550/arXiv.1603.04733](https://doi.org/10.48550/arXiv.1603.04733). arXiv: [1603.04733](https://arxiv.org/abs/1603.04733)[stat]. URL: <http://arxiv.org/abs/1603.04733> (visited on 05/15/2025) (cit. on p. 6).
- [47] A. Mishkin, F. Kunstner, D. Nielsen, M. Schmidt, and M. E. Khan. *SLANG: Fast Structured Covariance Approximations for Bayesian Deep Learning with Natural Gradient*. Jan. 12, 2019. DOI: [10.48550/arXiv.1811.04504](https://doi.org/10.48550/arXiv.1811.04504). arXiv: [1811.04504](https://arxiv.org/abs/1811.04504)[cs]. URL: <http://arxiv.org/abs/1811.04504> (visited on 05/15/2025) (cit. on p. 6).
- [48] J. Harrison, J. Willes, and J. Snoek. “Variational Bayesian Last Layers”. In: *International Conference on Learning Representations (ICLR)*. Apr. 2024. DOI: [10.48550/arXiv.2404.11599](https://doi.org/10.48550/arXiv.2404.11599) (cit. on p. 6).
- [49] J. Z. Liu, Z. Lin, S. Padhy, D. Tran, T. Bedrax-Weiss, and B. Lakshminarayanan. “Simple and Principled Uncertainty Estimation with Deterministic Deep Learning via Distance Awareness”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. Oct. 2020. DOI: [10.48550/arXiv.2006.10108](https://doi.org/10.48550/arXiv.2006.10108) (cit. on p. 6).
- [50] M. E. Khan, A. Immer, E. Abedi, and M. Korzepa. “Approximate Inference Turns Deep Networks into Gaussian Processes”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2019. DOI: [10.48550/arXiv.1906.01930](https://doi.org/10.48550/arXiv.1906.01930) (cit. on p. 6).
- [51] A. Immer, M. Korzepa, and M. Bauer. “Improving predictions of Bayesian neural nets via local linearization”. In: *International Conference on Artificial Intelligence and Statistics (AISTATS)*. 2021 (cit. on p. 6).
- [52] E. Daxberger, E. Nalisnick, J. U. Allingham, J. Antorán, and J. M. Hernández-Lobato. “Bayesian Deep Learning via Subnetwork Inference”. In: *International Conference on Machine Learning (ICML)*. 2021. DOI: [10.48550/arXiv.2010.14689](https://doi.org/10.48550/arXiv.2010.14689) (cit. on p. 6).
- [53] A. Kristiadi, A. Immer, R. Eschenhagen, and V. Fortuin. “Promises and Pitfalls of the Linearized Laplace in Bayesian Optimization”. In: *Advances in Approximate Bayesian Inference (AABI)*. 2023. DOI: [10.48550/arXiv.2304.08309](https://doi.org/10.48550/arXiv.2304.08309) (cit. on p. 6).
- [54] T. Cinquin, M. Pförtner, V. Fortuin, P. Hennig, and R. Bamler. “FSP-Laplace: Function-Space Priors for the Laplace Approximation in Bayesian Deep Learning”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. Oct. 2024. DOI: [10.48550/arXiv.2407.13711](https://doi.org/10.48550/arXiv.2407.13711). URL: <http://arxiv.org/abs/2407.13711> (cit. on pp. 6, 10).

- [55] S. Fort, H. Hu, and B. Lakshminarayanan. *Deep Ensembles: A Loss Landscape Perspective*. June 25, 2020. DOI: [10.48550/arXiv.1912.02757](https://doi.org/10.48550/arXiv.1912.02757). arXiv: [1912.02757](https://arxiv.org/abs/1912.02757) [stat]. URL: <http://arxiv.org/abs/1912.02757> (visited on 05/15/2025) (cit. on p. 6).
- [56] A. G. Wilson and P. Izmailov. “Bayesian Deep Learning and a Probabilistic Perspective of Generalization”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2020. DOI: [10.48550/arXiv.2002.08791](https://doi.org/10.48550/arXiv.2002.08791) (cit. on p. 6).
- [57] V. D. Wild, S. Ghalebikesabi, D. Sejdinovic, and J. Knoblauch. “A Rigorous Link between Deep Ensembles and (Variational) Bayesian Methods”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2023. DOI: [10.48550/arXiv.2305.15027](https://doi.org/10.48550/arXiv.2305.15027) (cit. on p. 6).
- [58] T. Abe, E. K. Buchanan, G. Pleiss, R. Zemel, and J. P. Cunningham. “Deep Ensembles Work, But Are They Necessary?” In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2022. DOI: [10.48550/arXiv.2202.06985](https://doi.org/10.48550/arXiv.2202.06985) (cit. on p. 6).
- [59] N. Dern, J. P. Cunningham, and G. Pleiss. *Theoretical Limitations of Ensembles in the Age of Overparameterization*. arXiv:2410.16201 [stat]. Oct. 2024. DOI: [10.48550/arXiv.2410.16201](https://doi.org/10.48550/arXiv.2410.16201) (cit. on p. 6).
- [60] C. Mingard, G. Valle-Pérez, J. Skalse, and A. A. Louis. “Is SGD a Bayesian sampler? Well, almost.” In: *Journal of Machine Learning Research (JMLR)* (2020) (cit. on p. 6).
- [61] J. A. Lin, J. Antorán, S. Padhy, D. Janz, J. M. Hernández-Lobato, and A. Terenin. “Sampling from Gaussian Process Posteriors using Stochastic Gradient Descent”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2023. DOI: [10.48550/arXiv.2306.11589](https://doi.org/10.48550/arXiv.2306.11589) (cit. on p. 6).
- [62] J. Lee, L. Xiao, S. S. Schoenholz, Y. Bahri, R. Novak, J. Sohl-Dickstein, and J. Pennington. “Wide Neural Networks of Any Depth Evolve as Linear Models Under Gradient Descent”. In: *Journal of Statistical Mechanics: Theory and Experiment* 2020.12 (2020). DOI: [10.1088/1742-5468/abc62b](https://doi.org/10.1088/1742-5468/abc62b) (cit. on p. 6).
- [63] J. Lai, M. Xu, R. Chen, and Q. Lin. *Generalization Ability of Wide Neural Networks on  $\mathbb{R}$* . Feb. 12, 2023. DOI: [10.48550/arXiv.2302.05933](https://doi.org/10.48550/arXiv.2302.05933). arXiv: [2302.05933](https://arxiv.org/abs/2302.05933) [stat]. URL: <http://arxiv.org/abs/2302.05933> (visited on 05/15/2025) (cit. on p. 6).
- [64] G. Yang. “Tensor Programs I: Wide Feedforward or Recurrent Neural Networks of Any Architecture are Gaussian Processes”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2019. DOI: [10.48550/arXiv.1910.12478](https://doi.org/10.48550/arXiv.1910.12478) (cit. on p. 6).
- [65] G. Yang. *Tensor Programs II: Neural Tangent Kernel for Any Architecture*. 2020. DOI: [10.48550/arXiv.2006.14548](https://doi.org/10.48550/arXiv.2006.14548) (cit. on p. 6).
- [66] A. Jacot, F. Gabriel, and C. Hongler. “Neural Tangent Kernel: Convergence and Generalization in Neural Networks”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2018. DOI: [10.48550/arXiv.1806.07572](https://doi.org/10.48550/arXiv.1806.07572) (cit. on p. 6).
- [67] V. Muthukumar, A. Narang, V. Subramanian, M. Belkin, D. Hsu, and A. Sahai. “Classification vs regression in overparameterized regimes: Does the loss function matter?” In: *Journal of Machine Learning Research (JMLR)* (Oct. 2021). DOI: [10.48550/arXiv.2005.08054](https://doi.org/10.48550/arXiv.2005.08054) (cit. on p. 7).
- [68] D. Hsu, V. Muthukumar, and J. Xu. *On the proliferation of support vectors in high dimensions*. 2022. DOI: [10.48550/arXiv.2009.10670](https://doi.org/10.48550/arXiv.2009.10670). URL: <http://arxiv.org/abs/2009.10670> (cit. on p. 7).
- [69] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. “On Calibration of Modern Neural Networks”. In: *International Conference on Machine Learning (ICML)*. 2017. DOI: [10.48550/arXiv.1706.04599](https://doi.org/10.48550/arXiv.1706.04599) (cit. on pp. 8, 33, 41).
- [70] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. “Gradient-based learning applied to document recognition”. In: *Proceedings of the IEEE* 86.11 (Nov. 1998), pp. 2278–2324. ISSN: 1558-2256. DOI: [10.1109/5.726791](https://doi.org/10.1109/5.726791). URL: <https://ieeexplore.ieee.org/document/726791> (cit. on pp. 8, 39, 40).
- [71] K. He, X. Zhang, S. Ren, and J. Sun. “Deep Residual Learning for Image Recognition”. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2016, pp. 770–778. ISBN: 978-1-4673-8851-1. DOI: [10.1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90). URL: <http://ieeexplore.ieee.org/document/7780459/> (cit. on pp. 8, 41).

- [72] N. Mu and J. Gilmer. “MNIST-C: A Robustness Benchmark for Computer Vision”. In: *ICML Workshop on Uncertainty and Robustness in Deep Learning*. June 2019. DOI: [10.48550/arXiv.1906.02337](https://doi.org/10.48550/arXiv.1906.02337). URL: <http://arxiv.org/abs/1906.02337> (cit. on pp. 8, 39).
- [73] D. Hendrycks and T. Dietterich. “Benchmarking Neural Network Robustness to Common Corruptions and Perturbations”. In: *International Conference on Learning Representations (ICLR)*. 2019. DOI: [10.48550/arXiv.1903.12261](https://doi.org/10.48550/arXiv.1903.12261). URL: <http://arxiv.org/abs/1903.12261> (cit. on pp. 8, 39).
- [74] D. R. Burt, S. W. Ober, A. Garriga-Alonso, and M. van der Wilk. *Understanding Variational Inference in Function-Space*. Nov. 2020. DOI: [10.48550/arXiv.2011.09421](https://doi.org/10.48550/arXiv.2011.09421) (cit. on p. 10).
- [75] S. Qiu, T. G. J. Rudner, S. Kapoor, and A. G. Wilson. “Should We Learn Most Likely Functions or Parameters?” In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2023. DOI: [10.48550/arXiv.2311.15990](https://doi.org/10.48550/arXiv.2311.15990) (cit. on p. 10).
- [76] T. G. J. Rudner, Z. Chen, Y. W. Teh, and Y. Gal. “Tractable Function-Space Variational Inference in Bayesian Neural Networks”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2023. DOI: [10.48550/arXiv.2312.17199](https://doi.org/10.48550/arXiv.2312.17199) (cit. on p. 10).
- [77] S. P. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge University Press, 2004. ISBN: 978-0-521-83378-3 (cit. on p. 17).
- [78] Y. Nesterov. “A method for solving the convex programming problem with convergence rate  $O(\frac{1}{k^2})$ ”. In: *Dokl Akad Nauk SSSR* 269 (1983), p. 543 (cit. on p. 19).
- [79] B. T. Polyak. “Some methods of speeding up the convergence of iteration methods”. In: *USSR Computational Mathematics and Mathematical Physics* 4.5 (1964), pp. 1–17. DOI: [10.1016/0041-5553\(64\)90137-5](https://doi.org/10.1016/0041-5553(64)90137-5) (cit. on p. 19).
- [80] A. Krizhevsky et al. *Learning multiple layers of features from tiny images*. Tech. rep. 2009 (cit. on p. 39).
- [81] Y. Le and X. Yang. “Tiny ImageNet Visual Recognition Challenge”. In: *Stanford CS 231N* (2015). URL: <http://cs231n.stanford.edu/tiny-imagenet-200.zip> (cit. on p. 39).
- [82] P. Goyal, P. Dollár, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He. *Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour*. Tech. rep. 2018. URL: <http://arxiv.org/abs/1706.02677> (cit. on p. 40).
- [83] S. L. Smith and Q. V. Le. “A Bayesian Perspective on Generalization and Stochastic Gradient Descent”. In: *International Conference on Learning Representations (ICLR)*. 2018. DOI: [10.48550/arXiv.1710.06451](https://doi.org/10.48550/arXiv.1710.06451) (cit. on p. 40).
- [84] S. L. Smith, P.-J. Kindermans, C. Ying, and Q. V. Le. “Don’t Decay the Learning Rate, Increase the Batch Size”. In: *International Conference on Learning Representations (ICLR)*. 2018. DOI: [10.48550/arXiv.1711.00489](https://doi.org/10.48550/arXiv.1711.00489) (cit. on p. 40).
- [85] T. maintainers and contributors. *TorchVision: PyTorch’s Computer Vision library*. <https://github.com/pytorch/vision>. 2016 (cit. on p. 41).

## SUPPLEMENTARY MATERIAL

This supplementary material contains additional results and proofs for all theoretical statements. References referring to sections, equations or theorem-type environments within this document are prefixed with ‘S’, while references to, or results from, the main paper are stated as is.

<b>S1 Theoretical Results</b>	<b>16</b>
S1.1 Overparametrized Linear Regression . . . . .	17
S1.1.1 Characterization of Implicit Bias (Proof of Theorem 1) . . . . .	17
S1.1.2 Non-Asymptotic Error Analysis . . . . .	20
S1.1.3 Connection to Ensembles . . . . .	22
S1.2 Binary Classification of Linearly Separable Data . . . . .	23
S1.2.1 Preliminaries . . . . .	24
S1.2.2 Gradient Flow for the Expected Loss . . . . .	24
S1.2.3 Complete Proof of Theorem 2 . . . . .	27
S1.3 NLL Overfitting and the Need for (Temperature) Scaling . . . . .	32
<b>S2 Parametrization, Feature Learning and Hyperparameter Transfer</b>	<b>33</b>
S2.1 Definitions of Stability and Feature Learning . . . . .	34
S2.2 Initialization Scaling for a Linear Network . . . . .	34
S2.3 Proposed Scaling . . . . .	36
S2.4 Details on Hyperparameter Transfer Experiment . . . . .	37
<b>S3 Experiments</b>	<b>38</b>
S3.1 Setup and Details . . . . .	38
S3.1.1 Datasets . . . . .	39
S3.1.2 Metrics . . . . .	39
S3.2 Time and Memory-Efficient Training . . . . .	40
S3.3 In- and Out-of-distribution Generalization . . . . .	40
S3.3.1 Architectures, Training, and Methods . . . . .	40
S3.3.2 In-Distribution Generalization and Uncertainty Quantification . . . . .	43
S3.3.3 Robustness to Input Corruptions . . . . .	44
S3.3.4 Comparison to Generalized VI with 2-Wasserstein Regularization . . . . .	44

## S1 THEORETICAL RESULTS

**Lemma S1**

Let  $q(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ ,  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$  such that  $\boldsymbol{\mu}, \boldsymbol{\mu}_0 \in \mathbb{R}^P$ ,  $\boldsymbol{\Sigma}, \boldsymbol{\Sigma}_0 \in \mathbb{R}^{P \times P}$  positive semi-definite and let  $\mathbf{V}_A \in \mathbb{R}^{P \times N}$ ,  $\mathbf{V}_B \in \mathbb{R}^{P \times (P-N)}$  be matrices with pairwise orthonormal columns that together define an orthonormal basis of  $\mathbb{R}^P$ , i.e. for  $\mathbf{V} = [\mathbf{V}_A \ \mathbf{V}_B]$  it holds that  $\mathbf{V}\mathbf{V}^\top = \mathbf{V}^\top\mathbf{V} = \mathbf{I}$  and  $\text{span}(\mathbf{V}) = \mathbb{R}^P$ . Assume further that

$$\mathbf{V}_A^\top \boldsymbol{\Sigma} \mathbf{V}_A = \mathbf{0}, \quad (\text{S11})$$

then the squared 2-Wasserstein distance is given by

$$W_2^2(q, p) = \|\mathbf{V}_A^\top \boldsymbol{\mu} - \mathbf{V}_A^\top \boldsymbol{\mu}_0\|_2^2 + W_2^2(\mathcal{N}(\mathbf{V}_B^\top \boldsymbol{\mu}, \mathbf{V}_B^\top \boldsymbol{\Sigma} \mathbf{V}_B), \mathcal{N}(\mathbf{V}_B^\top \boldsymbol{\mu}_0, \mathbf{V}_B^\top \boldsymbol{\Sigma}_0 \mathbf{V}_B)) + C, \quad (\text{S12})$$

where the constant  $C$  is independent of  $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .

*Proof.* Consider the matrix

$$\mathbf{V}^\top \boldsymbol{\Sigma} \mathbf{V} = \begin{bmatrix} \mathbf{0}_{N \times N} & \mathbf{V}_A^\top \boldsymbol{\Sigma} \mathbf{V}_B \\ \mathbf{V}_B^\top \boldsymbol{\Sigma} \mathbf{V}_A & \mathbf{V}_B^\top \boldsymbol{\Sigma} \mathbf{V}_B \end{bmatrix}.$$

Since  $\mathbf{V}^\top \boldsymbol{\Sigma} \mathbf{V}$  is symmetric positive semi-definite, its off-diagonal block  $\mathbf{V}_A^\top \boldsymbol{\Sigma} \mathbf{V}_B$  satisfies

$$(\mathbf{I} - \mathbf{0}\mathbf{0}^\dagger) \mathbf{V}_A^\top \boldsymbol{\Sigma} \mathbf{V}_B = \mathbf{0} \iff \mathbf{V}_A^\top \boldsymbol{\Sigma} \mathbf{V}_B = \mathbf{0}$$

by Boyd and Vandenberghe [A5.5, 77]. Therefore, we have

$$\mathbf{V}^\top \Sigma \mathbf{V} = \begin{bmatrix} \mathbf{0}_{N \times N} & \mathbf{V}_A^\top \Sigma \mathbf{V}_B \\ \mathbf{V}_B^\top \Sigma \mathbf{V}_A & \mathbf{V}_B^\top \Sigma \mathbf{V}_B \end{bmatrix} = \begin{bmatrix} \mathbf{0}_{N \times N} & \mathbf{0}_{N \times (P-N)} \\ \mathbf{0}_{(P-N) \times N} & \mathbf{V}_B^\top \Sigma \mathbf{V}_B \end{bmatrix}. \quad (\text{S13})$$

The squared 2-Wasserstein distance between  $q(\mathbf{w})$  and  $p(\mathbf{w})$  is given by

$$W_2^2(q, p) = \|\boldsymbol{\mu} - \boldsymbol{\mu}_0\|_2^2 + \text{tr}(\Sigma - 2(\Sigma^{\frac{1}{2}} \Sigma_0 \Sigma^{\frac{1}{2}})^{\frac{1}{2}} + \Sigma_0).$$

For the squared norm term it holds by unitary invariance of  $\|\cdot\|_2$  that

$$\|\boldsymbol{\mu} - \boldsymbol{\mu}_0\|_2^2 = \|\mathbf{V}^\top (\boldsymbol{\mu} - \boldsymbol{\mu}_0)\|_2^2 = \left\| \begin{bmatrix} \mathbf{V}_A^\top (\boldsymbol{\mu} - \boldsymbol{\mu}_0) \\ \mathbf{V}_B^\top (\boldsymbol{\mu} - \boldsymbol{\mu}_0) \end{bmatrix} \right\|_2^2 = \|\mathbf{V}_A^\top \boldsymbol{\mu} - \mathbf{V}_A^\top \boldsymbol{\mu}_0\|_2^2 + \|\mathbf{V}_B^\top \boldsymbol{\mu} - \mathbf{V}_B^\top \boldsymbol{\mu}_0\|_2^2.$$

Now for the trace term we have that

$$\begin{aligned} & \text{tr}(\mathbf{V} \mathbf{V}^\top (\Sigma - 2(\Sigma^{\frac{1}{2}} \Sigma_0 \Sigma^{\frac{1}{2}})^{\frac{1}{2}} + \Sigma_0)) \\ &= \text{tr}(\mathbf{V}^\top \Sigma \mathbf{V}) - 2 \text{tr}(\mathbf{V}^\top (\Sigma^{\frac{1}{2}} \Sigma_0 \Sigma^{\frac{1}{2}})^{\frac{1}{2}} \mathbf{V}) + \text{tr}(\mathbf{V}^\top \Sigma_0 \mathbf{V}) \\ &= \text{tr}(\mathbf{V}_A^\top \Sigma \mathbf{V}_A) + \text{tr}(\mathbf{V}_B^\top \Sigma \mathbf{V}_B) + \text{tr}(\mathbf{V}_A^\top \Sigma_0 \mathbf{V}_A) + \text{tr}(\mathbf{V}_B^\top \Sigma_0 \mathbf{V}_B) - 2 \text{tr}(\mathbf{V}^\top (\Sigma^{\frac{1}{2}} \Sigma_0 \Sigma^{\frac{1}{2}})^{\frac{1}{2}} \mathbf{V}) \\ &\stackrel{\pm c}{=} \text{tr}(\mathbf{V}_B^\top \Sigma \mathbf{V}_B) + \text{tr}(\mathbf{V}_B^\top \Sigma_0 \mathbf{V}_B) - 2 \text{tr}(\mathbf{V}^\top (\Sigma^{\frac{1}{2}} \Sigma_0 \Sigma^{\frac{1}{2}})^{\frac{1}{2}} \mathbf{V}) \end{aligned} \quad (\text{S14})$$

where we used Eq. (S11) and  $\stackrel{\pm c}{=}$  denotes equality up to constants independent of  $(\boldsymbol{\mu}, \Sigma)$ .

Now by Eq. (S13), we have that  $\Sigma = \mathbf{V}_B \mathbf{M} \mathbf{V}_B^\top$  for  $\mathbf{M} = \mathbf{V}_B^\top \Sigma \mathbf{V}_B$  and its unique principal square root is given by  $\Sigma^{\frac{1}{2}} = \mathbf{V}_B \mathbf{M}^{\frac{1}{2}} \mathbf{V}_B^\top$  since

$$(\mathbf{V}_B \mathbf{M}^{\frac{1}{2}} \mathbf{V}_B^\top)(\mathbf{V}_B \mathbf{M}^{\frac{1}{2}} \mathbf{V}_B^\top) = \mathbf{V}_B \mathbf{M}^{\frac{1}{2}} \mathbf{I}_{(P-N) \times (P-N)} \mathbf{M}^{\frac{1}{2}} \mathbf{V}_B^\top = \Sigma.$$

It also holds that the unique principal square root

$$(\Sigma^{\frac{1}{2}} \Sigma_0 \Sigma^{\frac{1}{2}})^{\frac{1}{2}} = \mathbf{V}_B (\mathbf{M}^{\frac{1}{2}} \mathbf{V}_B^\top \Sigma_0 \mathbf{V}_B \mathbf{M}^{\frac{1}{2}})^{\frac{1}{2}} \mathbf{V}_B^\top$$

since direct calculation gives

$$\begin{aligned} & (\mathbf{V}_B (\mathbf{M}^{\frac{1}{2}} \mathbf{V}_B^\top \Sigma_0 \mathbf{V}_B \mathbf{M}^{\frac{1}{2}})^{\frac{1}{2}} \mathbf{V}_B^\top)(\mathbf{V}_B (\mathbf{M}^{\frac{1}{2}} \mathbf{V}_B^\top \Sigma_0 \mathbf{V}_B \mathbf{M}^{\frac{1}{2}})^{\frac{1}{2}} \mathbf{V}_B^\top) \\ &= \mathbf{V}_B \mathbf{M}^{\frac{1}{2}} \mathbf{V}_B^\top \Sigma_0 \mathbf{V}_B \mathbf{M}^{\frac{1}{2}} \mathbf{V}_B^\top = \Sigma^{\frac{1}{2}} \Sigma_0 \Sigma^{\frac{1}{2}}. \end{aligned}$$

Therefore we have that

$$\text{tr}(\mathbf{V}^\top (\Sigma^{\frac{1}{2}} \Sigma_0 \Sigma^{\frac{1}{2}})^{\frac{1}{2}} \mathbf{V}) = \text{tr}(\mathbf{V}^\top \mathbf{V}_B (\mathbf{M}^{\frac{1}{2}} \mathbf{V}_B^\top \Sigma_0 \mathbf{V}_B \mathbf{M}^{\frac{1}{2}})^{\frac{1}{2}} \mathbf{V}_B^\top \mathbf{V}) = \text{tr}((\mathbf{M}^{\frac{1}{2}} \mathbf{V}_B^\top \Sigma_0 \mathbf{V}_B \mathbf{M}^{\frac{1}{2}})^{\frac{1}{2}}).$$

Putting it all together we obtain

$$\begin{aligned} W_2^2(q, p) &\stackrel{\pm c}{=} \|\mathbf{V}_A^\top \boldsymbol{\mu} - \mathbf{V}_A^\top \boldsymbol{\mu}_0\|_2^2 + \|\mathbf{V}_B^\top \boldsymbol{\mu} - \mathbf{V}_B^\top \boldsymbol{\mu}_0\|_2^2 + \text{tr}(\mathbf{V}_B^\top \Sigma \mathbf{V}_B) + \text{tr}(\mathbf{V}_B^\top \Sigma_0 \mathbf{V}_B) - 2 \text{tr}(\mathbf{V}^\top (\Sigma^{\frac{1}{2}} \Sigma_0 \Sigma^{\frac{1}{2}})^{\frac{1}{2}} \mathbf{V}) \\ &= \|\mathbf{V}_A^\top \boldsymbol{\mu} - \mathbf{V}_A^\top \boldsymbol{\mu}_0\|_2^2 + \|\mathbf{V}_B^\top \boldsymbol{\mu} - \mathbf{V}_B^\top \boldsymbol{\mu}_0\|_2^2 + \text{tr}(\mathbf{V}_B^\top \Sigma \mathbf{V}_B) + \text{tr}(\mathbf{V}_B^\top \Sigma_0 \mathbf{V}_B) - 2 \text{tr}((\mathbf{M}^{\frac{1}{2}} \mathbf{V}_B^\top \Sigma_0 \mathbf{V}_B \mathbf{M}^{\frac{1}{2}})^{\frac{1}{2}}) \\ &= \|\mathbf{V}_A^\top \boldsymbol{\mu} - \mathbf{V}_A^\top \boldsymbol{\mu}_0\|_2^2 + W_2^2(\mathcal{N}(\mathbf{V}_B^\top \boldsymbol{\mu}, \mathbf{V}_B^\top \Sigma \mathbf{V}_B), \mathcal{N}(\mathbf{V}_B^\top \boldsymbol{\mu}_0, \mathbf{V}_B^\top \Sigma_0 \mathbf{V}_B)) \end{aligned}$$

which completes the proof.  $\square$

## S1.1 OVERPARAMETRIZED LINEAR REGRESSION

### S1.1.1 CHARACTERIZATION OF IMPLICIT BIAS (PROOF OF THEOREM 1)

#### Theorem 1 (Implicit Bias in Regression)

Let  $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}$  be an overparametrized linear model with  $P > N$ . Define a Gaussian prior  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_0, \mathbf{S}_0 \mathbf{S}_0^\top)$  and likelihood  $p(\mathbf{y} | \mathbf{w}) = \mathcal{N}(\mathbf{y}; f_{\mathbf{w}}(\mathbf{X}), \sigma^2 \mathbf{I})$  and assume a variational family  $q_{\boldsymbol{\theta}}(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}, \mathbf{S} \mathbf{S}^\top)$  with  $\boldsymbol{\theta} = (\boldsymbol{\mu}, \mathbf{S})$  such that  $\boldsymbol{\mu} \in \mathbb{R}^P$  and  $\mathbf{S} \in \mathbb{R}^{P \times P}$  where  $R \leq P$ . If the learning rate sequence  $(\eta_t)_t$  is chosen such that the limit point  $\boldsymbol{\theta}_*^{\text{GD}} = \lim_{t \rightarrow \infty} \boldsymbol{\theta}_t^{\text{GD}}$  identified by gradient descent, initialized at  $\boldsymbol{\theta}_0 = (\boldsymbol{\mu}_0, \mathbf{S}_0)$ , is a (global) minimizer of the expected log-likelihood  $\bar{\ell}(\boldsymbol{\theta})$ , then

$$\boldsymbol{\theta}_*^{\text{GD}} \in \underset{\boldsymbol{\theta} = (\boldsymbol{\mu}, \mathbf{S})}{\arg \min} W_2^2(q_{\boldsymbol{\theta}}, p). \quad (7)$$

s.t.  $\boldsymbol{\theta} \in \arg \min \bar{\ell}(\boldsymbol{\theta})$

Further, this also holds in the case of stochastic gradient descent and when using momentum.

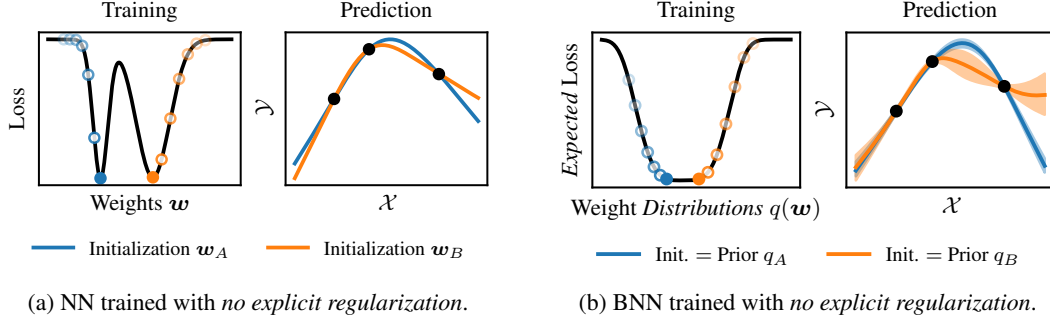


Figure S1: *Implicit regularization in standard neural networks versus in probabilistic networks*. Left panels: A neural network trained without explicit regularization can converge to different global minima of the loss. Optimization of the weights will implicitly regularize towards one or the other. Right panels: Analogously, there are multiple distributions over neural networks that are global minima of the *expected* loss. Optimization of the *distribution* over the weights will implicitly regularize towards one or the other. Our approach uses this implicit regularization instead of an explicit regularization to a prior.

*Proof.* Let  $\theta_\star = (\boldsymbol{\mu}_\star, \boldsymbol{S}_\star)$  be a minimizer of  $\bar{\ell}(\boldsymbol{\theta})$ . By assumption it holds that the expected negative log-likelihood is equal to the following non-negative loss function up to an additive constant:

$$\begin{aligned}\bar{\ell}(\boldsymbol{\theta}) &= \mathbb{E}_{q_{\boldsymbol{\theta}}(\boldsymbol{w})}(\ell(\boldsymbol{y}, f_{\boldsymbol{w}}(\boldsymbol{X}))) = \mathbb{E}_{q_{\boldsymbol{\theta}}(\boldsymbol{w})}(-\log p(\boldsymbol{y} | \boldsymbol{w})) \\ &\stackrel{\pm c}{=} \frac{1}{2\sigma^2} \mathbb{E}_{q_{\boldsymbol{\theta}}(\boldsymbol{w})}(\|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{w}\|_2^2) \\ &= \frac{1}{2\sigma^2} (\|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\mu}\|_2^2 + \text{tr}(\boldsymbol{X}\boldsymbol{\Sigma}\boldsymbol{X}^\text{T})) \geq 0,\end{aligned}$$

where  $\boldsymbol{\Sigma} = \boldsymbol{S}\boldsymbol{S}^\text{T}$  and non-negativity follows from  $\boldsymbol{\Sigma}$  being symmetric positive semi-definite. Therefore any (global) minimizer  $\boldsymbol{\theta}_\star = (\boldsymbol{\mu}_\star, \boldsymbol{\Sigma}_\star)$  necessarily satisfies

$$\|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\mu}_\star\|_2^2 = 0, \quad (\text{S15})$$

$$\text{tr}(\boldsymbol{X}\boldsymbol{\Sigma}_\star\boldsymbol{X}^\text{T}) = 0. \quad (\text{S16})$$

Let  $\boldsymbol{V} = [\boldsymbol{V}_{\text{range}} \quad \boldsymbol{V}_{\text{null}}] \in \mathbb{R}^{P \times P}$  be the orthonormal matrix of right singular vectors of  $\boldsymbol{X} = \boldsymbol{U}\boldsymbol{\Lambda}\boldsymbol{V}^\text{T}$ , where  $\boldsymbol{V}_{\text{range}} \in \mathbb{R}^{P \times N}$  and  $\boldsymbol{V}_{\text{null}} \in \mathbb{R}^{P \times (P-N)}$ . Since  $\boldsymbol{X} \in \mathbb{R}^{N \times P}$  and we are in the overparametrized regime, i.e.  $P > N$ , the optimal mean parameter decomposes into the least-squares solution and a null space contribution

$$\boldsymbol{\mu}_\star = \boldsymbol{V}_{\text{range}}\boldsymbol{u}_\star + \boldsymbol{V}_{\text{null}}\boldsymbol{z} = \boldsymbol{X}^\dagger\boldsymbol{y} + \boldsymbol{V}_{\text{null}}\boldsymbol{z}. \quad (\text{S17})$$

Furthermore, it holds for positive semi-definite  $\boldsymbol{\Sigma} \in \mathbb{R}^{P \times P}$  that

$$\begin{aligned}0 \leq \text{tr}(\boldsymbol{X}\boldsymbol{\Sigma}\boldsymbol{X}^\text{T}) &= \text{tr}(\boldsymbol{U}\boldsymbol{\Lambda}\boldsymbol{V}^\text{T}\boldsymbol{\Sigma}\boldsymbol{V}\boldsymbol{\Lambda}\boldsymbol{U}^\text{T}) = \text{tr}(\boldsymbol{\Lambda}\boldsymbol{V}^\text{T}\boldsymbol{\Sigma}\boldsymbol{V}\boldsymbol{\Lambda}) \\ &= \text{tr}\left([\boldsymbol{\Lambda}_{N \times N} \quad \mathbf{0}] \begin{bmatrix} \boldsymbol{V}_{\text{range}}^\text{T}\boldsymbol{\Sigma}\boldsymbol{V}_{\text{range}} & * \\ * & * \end{bmatrix} \begin{bmatrix} \boldsymbol{\Lambda}_{N \times N} \\ \mathbf{0} \end{bmatrix}\right) \\ &= \text{tr}(\boldsymbol{\Lambda}_{N \times N}\boldsymbol{V}_{\text{range}}^\text{T}\boldsymbol{\Sigma}\boldsymbol{V}_{\text{range}}\boldsymbol{\Lambda}_{N \times N}) \\ &= \sum_{i=1}^N \lambda_i^2 [\boldsymbol{V}_{\text{range}}^\text{T}\boldsymbol{\Sigma}\boldsymbol{V}_{\text{range}}]_{ii}\end{aligned}$$

where  $\lambda_i^2 > 0$  are the squared singular values of  $\boldsymbol{X}$ , which are strictly positive since  $\text{rank}(\boldsymbol{X}) = N$ . Therefore using Equation (S16) any global minimizer necessarily satisfies  $[\boldsymbol{V}_{\text{range}}^\text{T}\boldsymbol{\Sigma}_\star\boldsymbol{V}_{\text{range}}]_{ii} = 0$  for  $i \in \{1, \dots, N\}$ . Now since  $\boldsymbol{V}_{\text{range}}^\text{T}\boldsymbol{\Sigma}_\star\boldsymbol{V}_{\text{range}}$  is symmetric positive semi-definite and its diagonal is zero, so is its trace and therefore the sum of its non-negative eigenvalues is necessarily zero. Thus all eigenvalues are zero and therefore

$$\boldsymbol{V}_{\text{range}}^\text{T}\boldsymbol{\Sigma}\boldsymbol{V}_{\text{range}} = \mathbf{0}. \quad (\text{S18})$$

Now by Lemma S1 we have that the squared 2-Wasserstein distance between  $q_{\theta_*}(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$  and the initialization  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$  is given up to a constant independent of  $(\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$  by

$$\begin{aligned} W_2(q_{\theta_*}, p) &\stackrel{\pm c}{=} \left\| \mathbf{V}_{\text{range}}^\top \boldsymbol{\mu}_* - \mathbf{V}_{\text{range}}^\top \boldsymbol{\mu}_0 \right\|_2^2 + W_2^2(\mathcal{N}(\mathbf{V}_{\text{null}}^\top \boldsymbol{\mu}_*, \mathbf{V}_{\text{null}}^\top \boldsymbol{\Sigma}_* \mathbf{V}_{\text{null}}), \mathcal{N}(\mathbf{V}_{\text{null}}^\top \boldsymbol{\mu}_0, \mathbf{V}_{\text{null}}^\top \boldsymbol{\Sigma}_0 \mathbf{V}_{\text{null}})) \\ &= \left\| \mathbf{X}^\dagger \mathbf{y} - \mathbf{V}_{\text{range}}^\top \boldsymbol{\mu}_0 \right\|_2^2 + W_2^2(\mathcal{N}(\mathbf{V}_{\text{null}}^\top \boldsymbol{\mu}_*, \mathbf{V}_{\text{null}}^\top \boldsymbol{\Sigma}_* \mathbf{V}_{\text{null}}), \mathcal{N}(\mathbf{V}_{\text{null}}^\top \boldsymbol{\mu}_0, \mathbf{V}_{\text{null}}^\top \boldsymbol{\Sigma}_0 \mathbf{V}_{\text{null}})) \\ &\stackrel{\pm c}{=} W_2^2(\mathcal{N}(\mathbf{V}_{\text{null}}^\top \boldsymbol{\mu}_*, \mathbf{V}_{\text{null}}^\top \boldsymbol{\Sigma}_* \mathbf{V}_{\text{null}}), \mathcal{N}(\mathbf{V}_{\text{null}}^\top \boldsymbol{\mu}_0, \mathbf{V}_{\text{null}}^\top \boldsymbol{\Sigma}_0 \mathbf{V}_{\text{null}})). \end{aligned}$$

Therefore among variational distributions  $q_{\theta_*}$  with parameters  $\theta_*$  that minimize the expected loss  $\bar{\ell}(\theta)$ , any such  $\theta_*$  that minimizes the squared 2-Wasserstein distance to the prior satisfies

$$\underbrace{(\mathbf{V}_{\text{null}}^\top \boldsymbol{\mu}_*)}_{=: \mathbf{z}}, \underbrace{(\mathbf{V}_{\text{null}}^\top \boldsymbol{\Sigma}_* \mathbf{V}_{\text{null}})}_{=: \mathbf{M}} = (\mathbf{V}_{\text{null}}^\top \boldsymbol{\mu}_0, \mathbf{V}_{\text{null}}^\top \boldsymbol{\Sigma}_0 \mathbf{V}_{\text{null}}). \quad (\text{S19})$$

**(Stochastic) Gradient Descent** It remains to show that (stochastic) gradient descent identifies a minimum of the expected loss  $\bar{\ell}(\theta)$ , such that the above holds. By assumption we have for the loss on a batch  $\mathbf{X}_b$  of data that

$$\begin{aligned} \bar{\ell}(\theta) &= \mathbb{E}_{q_{\theta}(\mathbf{w})}(\ell(\mathbf{y}_b, f_{\mathbf{w}}(\mathbf{X}_b))) = \mathbb{E}_{q_{\theta}(\mathbf{w})}(-\log p(\mathbf{y}_b | \mathbf{w})) \\ &\stackrel{\pm c}{=} \frac{1}{2\sigma^2} (\|\mathbf{y}_b - \mathbf{X}_b \boldsymbol{\mu}\|_2^2 + \text{tr}(\mathbf{X}_b \boldsymbol{\Sigma} \mathbf{X}_b^\top)). \end{aligned}$$

Therefore, at convergence of (stochastic) gradient descent the variational parameters  $\theta_\infty = (\boldsymbol{\mu}_\infty, \mathbf{S}_\infty)$  are given by

$$\boldsymbol{\mu}_\infty = \boldsymbol{\mu}_0 - \sum_{t=1}^{\infty} \eta_t \nabla_{\boldsymbol{\mu}} \bar{\ell}_b(\boldsymbol{\theta}_{t-1}) = \boldsymbol{\mu}_0 + \sum_{t=1}^{\infty} \frac{\eta_t}{\sigma^2} \mathbf{X}_b^\top (\mathbf{y}_b - \mathbf{X}_b \boldsymbol{\mu}_{t-1})$$

as well as

$$\mathbf{S}_\infty = \mathbf{S}_0 - \sum_{t=1}^{\infty} \eta_t \nabla_{\mathbf{S}} \bar{\ell}_b(\boldsymbol{\theta}_{t-1}) = \mathbf{S}_0 - \sum_{t=1}^{\infty} \frac{\eta_t}{\sigma^2} \mathbf{X}_b^\top \mathbf{X}_b \mathbf{S}_{t-1}$$

and therefore

$$\begin{aligned} \mathbf{z}_\infty &= \mathbf{V}_{\text{null}}^\top \boldsymbol{\mu}_\infty = \mathbf{V}_{\text{null}}^\top \boldsymbol{\mu}_0 + \sum_{t=1}^{\infty} \frac{\eta_t}{\sigma^2} \mathbf{V}_{\text{null}}^\top \underbrace{\mathbf{X}_b^\top (\mathbf{y}_b - \mathbf{X}_b \boldsymbol{\mu}_{t-1})}_{\in \text{range}(\mathbf{X}_b^\top)} = \mathbf{V}_{\text{null}}^\top \boldsymbol{\mu}_0 \\ \mathbf{V}_{\text{null}}^\top \mathbf{S}_\infty &= \mathbf{V}_{\text{null}}^\top \mathbf{S}_0 - \sum_{t=1}^{\infty} \frac{\eta_t}{\sigma^2} \mathbf{V}_{\text{null}}^\top \underbrace{\mathbf{X}_b^\top \mathbf{X}_b \mathbf{S}_{t-1}}_{\text{columns} \in \text{range}(\mathbf{X}_b^\top)} = \mathbf{V}_{\text{null}}^\top \mathbf{S}_0 \end{aligned}$$

where we used continuity of linear maps between finite-dimensional spaces. It follows that

$$\mathbf{M}_\infty = \mathbf{V}_{\text{null}}^\top \boldsymbol{\Sigma}_\infty \mathbf{V}_{\text{null}} = \mathbf{V}_{\text{null}}^\top \mathbf{S}_\infty \mathbf{S}_\infty^\top \mathbf{V}_{\text{null}} = \mathbf{V}_{\text{null}}^\top \mathbf{S}_0 \mathbf{S}_0^\top \mathbf{V}_{\text{null}} = \mathbf{V}_{\text{null}}^\top \boldsymbol{\Sigma}_0 \mathbf{V}_{\text{null}}.$$

Therefore any limit point of (stochastic) gradient descent that minimizes the expected log-likelihood also minimizes the 2-Wasserstein distance to the prior, since  $\theta_\infty$  satisfies Equation (S19).

**Momentum** In case we are using (stochastic) gradient descent with momentum, the updates are given by

$$\begin{aligned} \boldsymbol{\mu}_{t+1} &= \boldsymbol{\mu}_t + \gamma_t \Delta \boldsymbol{\mu}_t - \eta_t \nabla_{\boldsymbol{\mu}} \bar{\ell}_b(\boldsymbol{\theta}_t + \alpha_t \Delta \boldsymbol{\theta}_t) \\ \mathbf{S}_{t+1} &= \mathbf{S}_t + \gamma_t \Delta \mathbf{S}_t - \eta_t \nabla_{\mathbf{S}} \bar{\ell}_b(\boldsymbol{\theta}_t + \alpha_t \Delta \boldsymbol{\theta}_t) \end{aligned} \quad (\text{S20})$$

where

$$\Delta \boldsymbol{\theta}_t = \begin{pmatrix} \Delta \boldsymbol{\mu}_t \\ \Delta \mathbf{S}_t \end{pmatrix} = \boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1}, \quad \Delta \boldsymbol{\theta}_0 = \mathbf{0}.$$

for parameters  $\gamma_t, \alpha_t \geq 0$ , which includes Nesterov's acceleration ( $\gamma_t = \alpha_t$ ) [78] and heavy ball momentum ( $\alpha_t = 0$ ) [79].

To prove that the updates of the variational parameters are always orthogonal to the null space of  $\mathbf{X}_b$ , we proceed by induction. The base case is trivial since  $\Delta\boldsymbol{\theta}_0 = \mathbf{0}$ . Assume now that  $\mathbf{V}_{\text{null}}^\top \Delta\boldsymbol{\mu}_t = \mathbf{0}$  and  $\mathbf{V}_{\text{null}}^\top \Delta\mathbf{S}_t = \mathbf{0}$ , then by Equation (S20), we have

$$\begin{aligned}\mathbf{V}_{\text{null}}^\top \Delta\boldsymbol{\mu}_{t+1} &= \mathbf{V}_{\text{null}}^\top (\boldsymbol{\mu}_{t+1} - \boldsymbol{\mu}_t) = \gamma_t \mathbf{V}_{\text{null}}^\top \Delta\boldsymbol{\mu}_t - \eta_t \mathbf{V}_{\text{null}}^\top \nabla_{\boldsymbol{\mu}} \bar{\ell}_b(\boldsymbol{\theta}_t + \alpha_t \Delta\boldsymbol{\theta}_t) = \mathbf{0} \\ \mathbf{V}_{\text{null}}^\top \Delta\mathbf{S}_{t+1} &= \mathbf{V}_{\text{null}}^\top (\mathbf{S}_{t+1} - \mathbf{S}_t) = \gamma_t \mathbf{V}_{\text{null}}^\top \Delta\mathbf{S}_t - \eta_t \mathbf{V}_{\text{null}}^\top \nabla_{\mathbf{S}} \bar{\ell}_b(\boldsymbol{\theta}_t + \alpha_t \Delta\boldsymbol{\theta}_t) = \mathbf{0}\end{aligned}$$

where we used the induction hypothesis and the fact that the gradients are orthogonal to the null space as shown earlier.

Therefore by the same argument as above we have that  $\boldsymbol{\theta}_\infty$  computed via (stochastic) gradient descent with momentum satisfies Equation (S19), which directly implies Theorem 1.  $\square$

### S1.1.2 NON-ASYMPTOTIC ERROR ANALYSIS

#### Theorem S3 (Non-Asymptotic Error of Gradient Flow)

Let  $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}$  be a linear model. Define a prior  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_0, \mathbf{S}_0 \mathbf{S}_0^\top)$  and assume noise-free observations  $y(\cdot) = f_{\mathbf{w}}(\cdot)$  for  $\mathbf{w} \sim p(\mathbf{w})$ . Further, define a variational distribution  $q_{\boldsymbol{\theta}}(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}, \mathbf{S} \mathbf{S}^\top)$  with  $\boldsymbol{\theta} = (\boldsymbol{\mu}, \mathbf{S})$  such that  $\boldsymbol{\mu} \in \mathbb{R}^P$  and  $\mathbf{S} \in \mathbb{R}^{P \times R}$  where  $R \leq P$ . Let  $\boldsymbol{\theta}(t) = (\boldsymbol{\mu}(t), \mathbf{S}(t))$  be the variational parameters at time  $t \geq 0$  given by the gradient flow of the expected loss

$$\dot{\boldsymbol{\theta}}(t) = -\nabla_{\boldsymbol{\theta}} \bar{\ell}(\boldsymbol{\theta}(t)) \quad (\text{S21})$$

initialized at  $\boldsymbol{\theta}(0) = (\boldsymbol{\mu}_0, \mathbf{S}_0)$ . Then the expected squared error of the mean prediction

$$\mathbb{E}_{\left(\begin{smallmatrix} \mathbf{y} \\ y_{\text{test}} \end{smallmatrix}\right)} \left( (y_{\text{test}} - f_{\boldsymbol{\mu}(t)}(\mathbf{x}_{\text{test}}))^2 \right) = \text{Var}_{\mathbf{w} \sim q_{\boldsymbol{\theta}(t)}}(f_{\mathbf{w}}(\mathbf{x}_{\text{test}})) \quad (\text{S22})$$

at any test point  $\mathbf{x}_{\text{test}} \in \mathbb{R}^P$ . In other words, assuming the training and test data are drawn from the prior predictive, the predictive error of  $f_{\boldsymbol{\mu}(t)}(\cdot)$  at any time  $t \geq 0$  is exactly quantified by the predictive uncertainty of the variational distribution, not only at initialization and in the limit  $t \rightarrow \infty$ .

*Proof.* The dynamics of the variational parameters as defined by the gradient flow in Equation (S21) are given by

$$\begin{aligned}\dot{\boldsymbol{\mu}}(t) &= -\nabla_{\boldsymbol{\mu}} \bar{\ell}(\boldsymbol{\mu}(t)) = \mathbf{X}^\top (\mathbf{y} - \mathbf{X} \boldsymbol{\mu}(t)) = -\mathbf{X}^\top \mathbf{X} (\boldsymbol{\mu}(t) - \mathbf{w}) = \frac{d}{dt} (\boldsymbol{\mu}(t) - \mathbf{w}), \\ \dot{\mathbf{S}}(t) &= -\nabla_{\mathbf{S}} \bar{\ell}(\mathbf{S}(t)) = -\mathbf{X}^\top \mathbf{X} \mathbf{S}(t).\end{aligned}$$

Since these dynamics are matrix differential equations, the mean and covariance parameters as a function of time are given by

$$\boldsymbol{\mu}(t) = \mathbf{w} + e^{-\mathbf{X}^\top \mathbf{X} t} (\boldsymbol{\mu}_0 - \mathbf{w}), \quad (\text{S23})$$

$$\mathbf{S}(t) = e^{-\mathbf{X}^\top \mathbf{X} t} \mathbf{S}_0. \quad (\text{S24})$$

Thus the expected predictive error at time step  $t \geq 0$  is given by

$$\begin{aligned}\mathbb{E}_{\left(\begin{smallmatrix} \mathbf{y} \\ y_{\text{test}} \end{smallmatrix}\right)} \left( \|y_{\text{test}} - f_{\boldsymbol{\mu}(t)}(\mathbf{x}_{\text{test}})\|_2^2 \right) &= \mathbb{E}_{\left(\begin{smallmatrix} \mathbf{X} \\ \mathbf{x}_{\text{test}} \end{smallmatrix}\right)} \left( \|y_{\text{test}} - \mathbf{x}_{\text{test}}^\top \boldsymbol{\mu}(t)\|_2^2 \right) \\ &= \mathbb{E}_{\mathbf{w}} \left( \left\| \mathbf{x}_{\text{test}}^\top \mathbf{w} - \mathbf{x}_{\text{test}}^\top \left( \mathbf{w} + e^{-\mathbf{X}^\top \mathbf{X} t} (\boldsymbol{\mu}_0 - \mathbf{w}) \right) \right\|_2^2 \right) \\ &= \mathbb{E}_{\mathbf{w}} \left( \left\| \mathbf{x}_{\text{test}}^\top e^{-\mathbf{X}^\top \mathbf{X} t} (\boldsymbol{\mu}_0 - \mathbf{w}) \right\|_2^2 \right)\end{aligned}$$

where we used Equation (S23). We have since  $\mathbb{E}(\mathbf{w}) = \boldsymbol{\mu}_0$ , that the above

$$\begin{aligned}&= \text{tr} \left( \text{Cov}(\mathbf{w} - \boldsymbol{\mu}_0) e^{-\mathbf{X}^\top \mathbf{X} t} \mathbf{x}_{\text{test}} \mathbf{x}_{\text{test}}^\top e^{-\mathbf{X}^\top \mathbf{X} t} \right) \\ &= \text{tr} \left( \mathbf{x}_{\text{test}}^\top e^{-\mathbf{X}^\top \mathbf{X} t} \mathbf{S}_0 \mathbf{S}_0^\top e^{-\mathbf{X}^\top \mathbf{X} t} \mathbf{x}_{\text{test}} \right)\end{aligned}$$

$$\begin{aligned}
&= \text{tr}(\mathbf{x}_{\text{test}}^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_{\text{test}}) \\
&= \text{Var}_{\mathbf{w} \sim q_{\theta(t)}}(f_{\mathbf{w}}(\mathbf{x}_{\text{test}}))
\end{aligned}$$

where we used Equation (S24) in the second-to-last equality. This completes the proof.  $\square$

**Theorem S4** (Non-Asymptotic Error of SGD)

Let  $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}$  be a linear model. Define a prior  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_0, \mathbf{S}_0 \mathbf{S}_0^\top)$  and assume noise-free observations  $y(\cdot) = f_{\mathbf{w}}(\cdot)$  for  $\mathbf{w} \sim p(\mathbf{w})$ . Further, define a variational distribution  $q_{\theta}(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}, \mathbf{S} \mathbf{S}^\top)$  with  $\theta = (\boldsymbol{\mu}, \mathbf{S})$  such that  $\boldsymbol{\mu} \in \mathbb{R}^P$  and  $\mathbf{S} \in \mathbb{R}^{P \times R}$  where  $R \leq P$ . Assume the expected loss is given by  $\bar{\ell}(\theta) = \mathbb{E}_{q_{\theta}(\mathbf{w})} \left( \frac{1}{2} \|\mathbf{y} - \mathbf{X} \mathbf{w}\|_2^2 \right)$  and let  $\theta(t) = (\boldsymbol{\mu}(t), \mathbf{S}(t))$  be the variational parameters at step  $t$  of (stochastic) gradient descent with learning rate sequence  $(\eta_t)_t$ , initialized at  $\theta(0) = (\boldsymbol{\mu}_0, \mathbf{S}_0)$ . Then the expected squared error of the mean prediction

$$\mathbb{E}_{\binom{\mathbf{y}}{y_{\text{test}}}} \left( (y_{\text{test}} - f_{\boldsymbol{\mu}(t)}(\mathbf{x}_{\text{test}}))^2 \right) = \text{Var}_{\mathbf{w} \sim q_{\theta(t)}}(f_{\mathbf{w}}(\mathbf{x}_{\text{test}})) \quad (\text{S25})$$

at any test point  $\mathbf{x}_{\text{test}} \in \mathbb{R}^P$ . In other words, assuming the training and test data are drawn from the prior predictive, the predictive error of  $f_{\boldsymbol{\mu}(t)}(\cdot)$  at any optimization step  $t$  is exactly quantified by the predictive uncertainty of the variational distribution.

Further, if the learning rate  $\eta_t \leq \frac{1}{\lambda_{\max}(\mathbf{X}_t^\top \mathbf{X}_t)}$  for all steps  $t$ , then

$$\text{tr}(\text{Cov}_{\mathbf{w} \sim q_{\theta(t+1)}}(\mathbf{w})) \leq \text{tr}(\text{Cov}_{\mathbf{w} \sim q_{\theta(t)}}(\mathbf{w})), \quad (\text{S26})$$

i.e. uncertainty about the parameters decreases monotonically during optimization.

*Proof.* The expected loss is given up to an additive constant by

$$\bar{\ell}(\theta) = \mathbb{E}_{q_{\theta}(\mathbf{w})}(\ell(\mathbf{y}, f_{\mathbf{w}}(\mathbf{X}))) \stackrel{+c}{=} \frac{1}{2} (\|\mathbf{y} - \mathbf{X} \boldsymbol{\mu}\|_2^2 + \text{tr}(\mathbf{X} \mathbf{S} \mathbf{S}^\top \mathbf{X}^\top)).$$

Now let  $(\mathbf{X}_t, \mathbf{y}_t)$  be the minibatch at step  $t \geq 1$ . Then it holds that

$$f_{\boldsymbol{\mu}(t)}(\mathbf{x}_{\text{test}}) - y_{\text{test}} = \mathbf{x}_{\text{test}}^\top (\boldsymbol{\mu}(t) - \mathbf{w}). \quad (\text{S27})$$

Further, the mean parameters identified by SGD are given by

$$\begin{aligned}
\boldsymbol{\mu}(t) - \mathbf{w} &= \boldsymbol{\mu}(t-1) - \mathbf{w} - \eta_t \nabla_{\boldsymbol{\mu}} \bar{\ell}(\theta(t-1)) \\
&= \boldsymbol{\mu}(t-1) - \mathbf{w} - \eta_t \mathbf{X}_t^\top (\mathbf{X}_t \boldsymbol{\mu}(t-1) - \mathbf{y}_t) \\
&= \boldsymbol{\mu}(t-1) - \mathbf{w} - \eta_t \mathbf{X}_t^\top \mathbf{X}_t (\boldsymbol{\mu}(t-1) - \mathbf{w}) \\
&= (\mathbf{I} - \eta_t \mathbf{X}_t^\top \mathbf{X}_t) (\boldsymbol{\mu}(t-1) - \mathbf{w}) \\
&= \prod_{j=1}^t (\mathbf{I} - \eta_j \mathbf{X}_j^\top \mathbf{X}_j) (\boldsymbol{\mu}(0) - \mathbf{w}) \\
&= \mathbf{B}_t (\boldsymbol{\mu}_0 - \mathbf{w})
\end{aligned}$$

where we defined  $\mathbf{B}_t = \prod_{j=1}^t (\mathbf{I} - \eta_j \mathbf{X}_j^\top \mathbf{X}_j)$ . The covariance parameters are given by

$$\begin{aligned}
\mathbf{S}(t) &= \mathbf{S}(t-1) - \eta_t \nabla_{\mathbf{S}} \bar{\ell}(\theta(t-1)) \\
&= \mathbf{S}(t-1) - \eta_t \mathbf{X}_t^\top \mathbf{X}_t \mathbf{S}(t-1) \\
&= (\mathbf{I} - \eta_t \mathbf{X}_t^\top \mathbf{X}_t) \mathbf{S}(t-1) \\
&= \prod_{j=1}^t (\mathbf{I} - \eta_j \mathbf{X}_j^\top \mathbf{X}_j) \mathbf{S}(0) \\
&= \mathbf{B}_t \mathbf{S}_0
\end{aligned}$$

Therefore the predictive error at step  $t \in \{0, 1, \dots\}$  is given by

$$\mathbb{E}_{\binom{\mathbf{y}}{y_{\text{test}}}} \left( \|y_{\text{test}} - f_{\boldsymbol{\mu}(t)}(\mathbf{x}_{\text{test}})\|_2^2 \right) = \mathbb{E}_{\binom{\mathbf{X}}{\mathbf{x}_{\text{test}}}} \left( \|y_{\text{test}} - \mathbf{x}_{\text{test}}^\top \boldsymbol{\mu}(t)\|_2^2 \right)$$

$$\begin{aligned}
&= \mathbb{E}_{\mathbf{w}} \left( \|\mathbf{x}_{\text{test}}^\top (\boldsymbol{\mu}(t) - \mathbf{w})\|_2^2 \right) \\
&= \mathbb{E}_{\mathbf{w}} \left( \|\mathbf{x}_{\text{test}}^\top \mathbf{B}_t (\boldsymbol{\mu}_0 - \mathbf{w})\|_2^2 \right).
\end{aligned}$$

We have since  $\mathbb{E}(\boldsymbol{\mu}_0 - \mathbf{w}) = \mathbf{0}$ , that the above

$$\begin{aligned}
&= \text{tr}(\mathbf{B}_t^\top \mathbf{x}_{\text{test}} \mathbf{x}_{\text{test}}^\top \mathbf{B}_t \text{Cov}(\mathbf{w} - \boldsymbol{\mu}_0)) \\
&= \text{tr}(\mathbf{x}_{\text{test}}^\top \mathbf{B}_t \mathbf{S}_0 \mathbf{S}_0^\top \mathbf{B}_t^\top \mathbf{x}_{\text{test}}) \\
&= \text{tr}(\mathbf{x}_{\text{test}}^\top \mathbf{S}(t) \mathbf{S}(t)^\top \mathbf{x}_{\text{test}}) \\
&= \text{Var}_{\mathbf{w} \sim q_{\theta(t)}}(f_{\mathbf{w}}(\mathbf{x}_{\text{test}})).
\end{aligned}$$

This proves Equation (S25).

To prove the second statement, we begin by showing that  $\mathbf{I} - \eta_t \mathbf{X}_t^\top \mathbf{X}_t$  has a spectrum in the interval  $[0, 1]$ . We have by Weyl's theorem, since  $\mathbf{I}$  and  $\mathbf{C}_{t+1} := -\eta_{t+1} \mathbf{X}_{t+1}^\top \mathbf{X}_{t+1}$  are hermitian, that

$$\begin{aligned}
&\lambda_p(\mathbf{I}) + \lambda_{\min}(\mathbf{C}_{t+1}) \leq \lambda_p(\mathbf{I} + \mathbf{C}_{t+1}) \leq \lambda_p(\mathbf{I}) + \lambda_{\max}(\mathbf{C}_{t+1}) \\
\iff &1 - \eta_{t+1} \lambda_{\max}(\mathbf{X}_{t+1}^\top \mathbf{X}_{t+1}) \leq \lambda_p(\mathbf{I} + \mathbf{C}_{t+1}) \leq 1 - \eta_{t+1} \lambda_{\min}(\mathbf{X}_{t+1}^\top \mathbf{X}_{t+1}) \\
\implies &1 - \frac{\lambda_{\max}(\mathbf{X}_{t+1}^\top \mathbf{X}_{t+1})}{\lambda_{\max}(\mathbf{X}_{t+1}^\top \mathbf{X}_{t+1})} \leq \lambda_p(\mathbf{I} + \mathbf{C}_{t+1}) \leq 1 \\
&\iff 0 \leq \lambda_p(\mathbf{I} + \mathbf{C}_{t+1}) \leq 1
\end{aligned}$$

where we used the assumption on the learning rate that  $\forall t : \eta_t \leq \frac{1}{\lambda_{\max}(\mathbf{X}_t^\top \mathbf{X}_t)}$ . Now by von Neumann's trace inequality, it holds that

$$\begin{aligned}
\text{tr}(\text{Cov}_{\mathbf{w} \sim q_{\theta(t+1)}}(\mathbf{w})) &= \text{tr}((\mathbf{I} - \eta_{t+1} \mathbf{X}_{t+1}^\top \mathbf{X}_{t+1}) \mathbf{S}_t \mathbf{S}_t^\top (\mathbf{I} - \eta_{t+1} \mathbf{X}_{t+1}^\top \mathbf{X}_{t+1})^\top) \\
&= \text{tr}(\mathbf{S}_t \mathbf{S}_t^\top (\mathbf{I} - \eta_{t+1} \mathbf{X}_{t+1}^\top \mathbf{X}_{t+1}) (\mathbf{I} - \eta_{t+1} \mathbf{X}_{t+1}^\top \mathbf{X}_{t+1})) \\
&\leq \sum_{p=1}^P \lambda_p(\mathbf{S}_t \mathbf{S}_t^\top) \lambda_p((\mathbf{I} - \eta_{t+1} \mathbf{X}_{t+1}^\top \mathbf{X}_{t+1})^2) \\
&= \sum_{p=1}^P \lambda_p(\mathbf{S}_t \mathbf{S}_t^\top) \lambda_p((\mathbf{I} - \eta_{t+1} \mathbf{X}_{t+1}^\top \mathbf{X}_{t+1}))^2 \\
&\leq \sum_{p=1}^P \lambda_p(\mathbf{S}_t \mathbf{S}_t^\top) \\
&= \text{tr}(\text{Cov}_{\mathbf{w} \sim q_{\theta(t)}}(\mathbf{w})).
\end{aligned}$$

□

### S1.1.3 CONNECTION TO ENSEMBLES

#### Proposition S1 (Connection to Ensembles)

Consider an ensemble of overparametrized linear models  $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}$  initialized with weights drawn from the prior  $\mathbf{w}_0^{(i)} \sim \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_0, \mathbf{S}_0 \mathbf{S}_0^\top)$ . Assume each model is trained independently to convergence via (S)GD such that  $\mathbf{w}_\star^{(i)} = \arg \min_{\mathbf{w}} \ell(\mathbf{y}, f_{\mathbf{w}}(\mathbf{X}))$ . Then the distribution over the weights of the trained ensemble  $q_{\text{Ens}}(\mathbf{w})$  is equal to the variational approximation  $q_{\theta_\star}(\mathbf{w})$  learned via (S)GD initialized at the prior hyperparameters  $\boldsymbol{\theta}_0 = (\boldsymbol{\mu}_0, \mathbf{S}_0)$ , i.e.

$$q_{\text{Ens}}(\mathbf{w}) = q_{\theta_\star^{\text{GD}}}(\mathbf{w}). \quad (\text{S28})$$

*Proof.* The parameters  $\mathbf{w}_\infty^{(i)}$  of the (independently) trained ensemble members identified via (stochastic) gradient descent are given by

$$\mathbf{w}_\infty^{(i)} = \arg \min_{\mathbf{w} \in F} \|\mathbf{w} - \mathbf{w}_0^{(i)}\|_2$$

where  $F = \{\mathbf{w} \in \mathbb{R}^P \mid f_{\mathbf{w}}(\mathbf{X}) = \mathbf{X}\mathbf{w} = \mathbf{y}\}$  is the set of interpolating solutions [5, Sec. 2.1]. Since we can write  $F$  equivalently via the minimum norm solution and an arbitrary null space contribution, s.t.  $F = \{\mathbf{w} = \mathbf{X}^\dagger \mathbf{y} + \mathbf{w}_{\text{null}} \mid \mathbf{w}_{\text{null}} \in \text{null}(\mathbf{X})\}$  we have

$$\begin{aligned} &= \mathbf{X}^\dagger \mathbf{y} + \arg \min_{\mathbf{w}_{\text{null}} \in \text{null}(\mathbf{X})} \|\mathbf{w}_{\text{null}} - (\mathbf{w}_0^{(i)} - \mathbf{X}^\dagger \mathbf{y})\|_2 \\ &= \mathbf{X}^\dagger \mathbf{y} + \text{proj}_{\text{null}(\mathbf{X})} \left( \mathbf{w}_0^{(i)} - \underbrace{\mathbf{X}^\dagger \mathbf{y}}_{\in \text{range}(\mathbf{X}^\top)} \right) \end{aligned}$$

where we used the characterization of an orthogonal projection onto a linear subspace as the (unique) closest point in the subspace. Finally, we use that the minimum norm solution is in the range space of the data and rewrite the projection in matrix form, s.t.

$$= \mathbf{X}^\dagger \mathbf{y} + \mathbf{P}_{\text{null}} \mathbf{w}_0^{(i)}.$$

Therefore the distribution over the parameters  $\mathbf{w}_\infty^{(i)}$  of the ensemble members computed via (S)GD with initial parameters  $\mathbf{w}_0 \sim \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_0, \mathbf{S}_0 \mathbf{S}_0^\top)$  is given by

$$q_{\text{Ens}}(\mathbf{w}) = \mathcal{N} \left( \mathbf{w}; \underbrace{\mathbf{X}^\dagger \mathbf{y} + \mathbf{P}_{\text{null}} \boldsymbol{\mu}_0}_{=\boldsymbol{\mu}_{\text{Ens}}}, \underbrace{\mathbf{P}_{\text{null}} \mathbf{S}_0 \mathbf{S}_0^\top \mathbf{P}_{\text{null}}}_{=\mathbf{S}_{\text{Ens}}} \right).$$

Now the expected negative log-likelihood of the distribution over the parameters of the trained ensemble members  $q_{\text{Ens}}(\mathbf{w})$  with hyperparameters  $\boldsymbol{\theta}_{\text{Ens}} = (\boldsymbol{\mu}_{\text{Ens}}, \mathbf{S}_{\text{Ens}})$  is

$$\bar{\ell}(\boldsymbol{\theta}_{\text{Ens}}) \stackrel{\pm c}{=} \frac{1}{2\sigma^2} (\|\mathbf{y} - \mathbf{X} \boldsymbol{\mu}_{\text{Ens}}\|_2^2 + \text{tr}(\mathbf{X} \mathbf{S}_{\text{Ens}} \mathbf{S}_{\text{Ens}}^\top \mathbf{X}^\top)) = 0$$

and therefore  $\boldsymbol{\theta}_{\text{Ens}}$  is a minimizer of the expected log-likelihood. Further it holds that

$$\mathbf{z} = \mathbf{V}_{\text{null}}^\top (\mathbf{P}_{\text{null}} \boldsymbol{\mu}_0) = \mathbf{V}_{\text{null}}^\top \boldsymbol{\mu}_0$$

$$\mathbf{M} = \mathbf{V}_{\text{null}}^\top (\mathbf{P}_{\text{null}} \mathbf{S}_0) (\mathbf{P}_{\text{null}} \mathbf{S}_0)^\top \mathbf{V}_{\text{null}} = \mathbf{V}_{\text{null}}^\top \mathbf{S}_0 \mathbf{S}_0^\top \mathbf{V}_{\text{null}} = \mathbf{V}_{\text{null}}^\top \boldsymbol{\Sigma}_0 \mathbf{V}_{\text{null}}$$

and thus by Equation (S19), the distribution of the trained ensemble parameters minimizes the 2-Wasserstein distance to the prior distribution, i.e.

$$q_{\text{Ens}} = \arg \min_{q(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}, \boldsymbol{\Sigma})} W_2^2(q(\mathbf{w}), \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)).$$

Combining this with the characterization of the variational posterior in Theorem 1 proves the claim.  $\square$

## S1.2 BINARY CLASSIFICATION OF LINEARLY SEPARABLE DATA

In this subsection we provide proofs of claims from Section 4.2. We begin with presenting some preliminary results from Soudry et al. [4] which will be used throughout the proof. Next, we will analyze the gradient flow of the expected loss. We extend the results for the gradient flow to gradient descent and derive the characterization of the implicit bias, completing the proof of Theorem 2.

### Theorem 2 (Implicit Bias in Binary Classification)

Let  $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}$  be an (overparametrized) linear model and define a Gaussian prior  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_0, \mathbf{S}_0 \mathbf{S}_0^\top)$ . Assume a variational distribution  $q_{\boldsymbol{\theta}}(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}, \mathbf{S} \mathbf{S}^\top)$  over the weights  $\mathbf{w} \in \mathbb{R}^P$  with variational parameters  $\boldsymbol{\theta} = (\boldsymbol{\mu}, \mathbf{S})$  such that  $\mathbf{S} \in \mathbb{R}^{P \times R}$  and  $R \leq P$ . Assume we are using the exponential loss  $\ell(u) = \exp(-u)$  and optimize the expected empirical loss  $\bar{\ell}(\boldsymbol{\theta})$  via gradient descent initialized at the prior, i.e.  $\boldsymbol{\theta}_0 = (\boldsymbol{\mu}_0, \mathbf{S}_0)$ , with a sufficiently small learning rate  $\eta$ . Then for almost any dataset which is linearly separable (Assumption 1) and for which the support vectors span the data (Assumption 2), the rescaled gradient descent iterates (rGD)

$$\boldsymbol{\theta}_t^{\text{rGD}} = (\boldsymbol{\mu}_t^{\text{rGD}}, \mathbf{S}_t^{\text{rGD}}) = \left( \frac{1}{\log(t)} \boldsymbol{\mu}_t^{\text{GD}} + \mathbf{P}_{\text{null}}(\mathbf{X}) \boldsymbol{\mu}_0, \mathbf{S}_t^{\text{GD}} \right) \quad (9)$$

converge to a limit point  $\boldsymbol{\theta}_*^{\text{rGD}} = \lim_{t \rightarrow \infty} \boldsymbol{\theta}_t^{\text{rGD}}$  for which it holds that

$$\boldsymbol{\theta}_*^{\text{rGD}} \in \arg \min_{\substack{\boldsymbol{\theta} = (\boldsymbol{\mu}, \mathbf{S}) \\ \text{s.t. } \boldsymbol{\theta} \in \Theta_*}} W_2^2(q_{\boldsymbol{\theta}}, p), \quad (10)$$

where the feasible set  $\Theta_* = \{(\boldsymbol{\mu}, \mathbf{S}) \mid \mathbf{P}_{\text{range}(\mathbf{X}^\top)} \boldsymbol{\mu} = \hat{\mathbf{w}} \text{ and } \forall n : \text{Var}_{q_{\boldsymbol{\theta}}}(f_{\mathbf{w}}(\mathbf{x}_n)) = 0\}$  consists of mean parameters which, if projected onto the training data, are equivalent to the  $L_2$  max margin vector and covariance parameters such that there is no uncertainty at training data.

### S1.2.1 PRELIMINARIES

Recall that the expected loss is given by

$$\bar{\ell}(\boldsymbol{\theta}) = \sum_{n=1}^N \mathbb{E}_{q_{\boldsymbol{\theta}}(\boldsymbol{w})}(\ell(y_n \boldsymbol{x}_n^T \boldsymbol{w})), \quad (\text{S29})$$

and specifically, for the exponential loss, we have

$$\bar{\ell}(\boldsymbol{\theta}) = \bar{\ell}(\boldsymbol{\mu}, \boldsymbol{S}) = \sum_{n=1}^N \exp(-\boldsymbol{x}_n^T \boldsymbol{\mu} + \frac{1}{2} \boldsymbol{x}_n^T \boldsymbol{S} \boldsymbol{S}^T \boldsymbol{x}_n). \quad (\text{S30})$$

Throughout these proofs, for any mean parameter iterate  $\boldsymbol{\mu}_t$ , we define the residual as

$$\boldsymbol{r}_t = \boldsymbol{\mu}_t - \hat{\boldsymbol{w}} \log t - \tilde{\boldsymbol{w}} \quad (\text{S31})$$

where  $\hat{\boldsymbol{w}}$  is the solution to the hard margin SVM, and  $\tilde{\boldsymbol{w}}$  is the vector which satisfies

$$\forall n \in \mathcal{S} : \eta \exp(-\boldsymbol{x}_n^T \tilde{\boldsymbol{w}}) = \alpha_n, \quad (\text{S32})$$

where weights  $\alpha_n$  are defined through the KKT conditions on the hard margin SVM problem, i.e.

$$\hat{\boldsymbol{w}} = \sum_{n \in \mathcal{S}} \alpha_n \boldsymbol{x}_n. \quad (\text{S33})$$

In Lemma 12 (Appendix B) of Soudry et al. [4], it is shown that, for almost any dataset, there are no more than  $P$  support vectors and  $\alpha_n \neq 0, \forall n \in \mathcal{S}$ . Furthermore, we denote the minimum margin to a non-support vector as:

$$\kappa = \min_{n \notin \mathcal{S}} \boldsymbol{x}_n^T \hat{\boldsymbol{w}} > 1. \quad (\text{S34})$$

Finally, we define  $\boldsymbol{P}_{\mathcal{S}} \in \mathbb{R}^{P \times P}$  as the orthogonal projection matrix to the subspace spanned by the support vectors, and  $\bar{\boldsymbol{P}}_{\mathcal{S}} = \boldsymbol{I} - \boldsymbol{P}_{\mathcal{S}}$  as the complementary projection.

### S1.2.2 GRADIENT FLOW FOR THE EXPECTED LOSS

Similar as in Soudry et al. [4], we begin by studying the gradient flow dynamics, i.e. taking the continuous time limit of gradient descent:

$$\dot{\boldsymbol{\theta}}_t = -\nabla \bar{\ell}(\boldsymbol{\theta}_t), \quad (\text{S35})$$

which can be written componentwise as:

$$\dot{\boldsymbol{\mu}}_t = -\nabla_{\boldsymbol{\mu}} \bar{\ell}(\boldsymbol{\mu}_t, \boldsymbol{S}_t) = \sum_{n=1}^N \exp\left(-\boldsymbol{\mu}_t^T \boldsymbol{x}_n + \frac{1}{2} \boldsymbol{x}_n^T \boldsymbol{S}_t \boldsymbol{S}_t^T \boldsymbol{x}_n\right) \boldsymbol{x}_n \quad (\text{S36})$$

$$\dot{\boldsymbol{S}}_t = -\nabla_{\boldsymbol{S}} \bar{\ell}(\boldsymbol{\mu}_t, \boldsymbol{S}_t) = -\sum_{n=1}^N \exp\left(-\boldsymbol{\mu}_t^T \boldsymbol{x}_n + \frac{1}{2} \boldsymbol{x}_n^T \boldsymbol{S}_t \boldsymbol{S}_t^T \boldsymbol{x}_n\right) \boldsymbol{x}_n \boldsymbol{x}_n^T \boldsymbol{S}_t. \quad (\text{S37})$$

We begin by showing that the total uncertainty, as measured by the Frobenius norm of the covariance factor, is bounded during the gradient flow dynamics. To that end, we derive the following dynamics:

$$\frac{d}{dt} \frac{1}{2} \|\boldsymbol{S}_t\|_F^2 = \text{tr}(\boldsymbol{S}_t^T \dot{\boldsymbol{S}}_t) = -\sum_{n=1}^N \exp\left(-\boldsymbol{\mu}_t^T \boldsymbol{x}_n + \frac{1}{2} \boldsymbol{x}_n^T \boldsymbol{S}_t \boldsymbol{S}_t^T \boldsymbol{x}_n\right) \|\boldsymbol{x}_n^T \boldsymbol{S}_t\|^2 \leq 0, \quad (\text{S38})$$

and therefore

$$\|\boldsymbol{S}_t\|_F^2 \leq \|\boldsymbol{S}_0\|_F^2. \quad (\text{S39})$$

Finally, by Cauchy-Schwarz inequality, we have that

$$\|\boldsymbol{S}_t \boldsymbol{S}_t^T\|_F \leq \|\boldsymbol{S}_t\|_F^2 \leq \|\boldsymbol{S}_0\|_F^2. \quad (\text{S40})$$

We continue by studying the convergence behavior of the mean parameter  $\boldsymbol{\mu}_t$ .

**Mean parameter** Our goal is to show that  $\|\mathbf{r}_t\|$  is bounded. Equation (S31) implies that

$$\dot{\mathbf{r}}_t = \dot{\boldsymbol{\mu}}_t - \frac{1}{t}\dot{\boldsymbol{w}} = -\nabla_{\boldsymbol{\mu}}\bar{\ell}(\boldsymbol{\mu}_t, \mathbf{S}_t) - \frac{1}{t}\dot{\boldsymbol{w}}. \quad (\text{S41})$$

This in turn implies that

$$\begin{aligned} \frac{1}{2}\frac{d}{dt}\|\mathbf{r}_t\|^2 &= \dot{\mathbf{r}}_t^\top \mathbf{r}_t \\ &= \sum_{n=1}^N \exp\left(-\boldsymbol{\mu}_t^\top \mathbf{x}_n + \frac{1}{2}\mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n\right) \mathbf{x}_n^\top \mathbf{r}_t - \frac{1}{t}\dot{\boldsymbol{w}}^\top \mathbf{r}_t \\ &= \sum_{n \in \mathcal{S}} \exp\left(-\log(t)\hat{\boldsymbol{w}}^\top \mathbf{x}_n - \tilde{\boldsymbol{w}}^\top \mathbf{x}_n + \frac{1}{2}\mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n - \mathbf{x}_n^\top \mathbf{r}_t\right) \mathbf{x}_n^\top \mathbf{r}_t - \frac{1}{t}\dot{\boldsymbol{w}}^\top \mathbf{r}_t \\ &\quad + \sum_{n \notin \mathcal{S}} \exp\left(-\log(t)\hat{\boldsymbol{w}}^\top \mathbf{x}_n - \tilde{\boldsymbol{w}}^\top \mathbf{x}_n + \frac{1}{2}\mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n - \mathbf{x}_n^\top \mathbf{r}_t\right) \mathbf{x}_n^\top \mathbf{r}_t \\ &= \left[\frac{1}{t} \sum_{n \in \mathcal{S}} \exp(-\tilde{\boldsymbol{w}}^\top \mathbf{x}_n) \left(\exp\left(-\mathbf{x}_n^\top \mathbf{r}_t + \frac{1}{2}\mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n\right) - 1\right) \mathbf{x}_n^\top \mathbf{r}_t\right] \\ &\quad + \left[\sum_{n \notin \mathcal{S}} \left(\frac{1}{t}\right)^{\hat{\boldsymbol{w}}^\top \mathbf{x}_n} \exp\left(-\tilde{\boldsymbol{w}}^\top \mathbf{x}_n + \frac{1}{2}\mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n\right) \exp(-\mathbf{x}_n^\top \mathbf{r}_t) \mathbf{x}_n^\top \mathbf{r}_t\right]. \end{aligned} \quad (\text{S42})$$

where in last line we used the fact that  $\hat{\boldsymbol{w}}^\top \mathbf{x}_n = 1$  for  $n \in \mathcal{S}$  as in (S32), and that  $\sum_{n \in \mathcal{S}} \exp(-\mathbf{x}_n^\top \tilde{\boldsymbol{w}}) \mathbf{x}_n = \hat{\boldsymbol{w}}$  as in (S33). We begin by examining the first bracket, studying three possible cases for each of the summands. First, note that if  $\mathbf{x}_n^\top \mathbf{r}_t \leq 0$ , then since  $\frac{1}{2}\mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n \geq 0$ , we have that

$$\left(\exp\left(-\mathbf{x}_n^\top \mathbf{r}_t + \frac{1}{2}\mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n\right) - 1\right) \mathbf{x}_n^\top \mathbf{r}_t \leq 0. \quad (\text{S43})$$

Next, by defining  $B := \|\mathbf{S}_0\|_F^2 \max_n \|\mathbf{x}_n\|_2$ , if  $0 < \mathbf{x}_n^\top \mathbf{r}_t < \frac{B}{2}$ , we have that

$$\left|\left(\exp\left(-\mathbf{x}_n^\top \mathbf{r}_t + \frac{1}{2}\mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n\right) - 1\right) \mathbf{x}_n^\top \mathbf{r}_t\right| < \left(\exp\left(\frac{B}{2}\right) - 1\right) \frac{B}{2}, \quad (\text{S44})$$

and if  $\mathbf{x}_n^\top \mathbf{r}_t \geq \frac{B}{2}$ , we have that

$$\left(\exp\left(-\mathbf{x}_n^\top \mathbf{r}_t + \frac{1}{2}\mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n\right) - 1\right) \mathbf{x}_n^\top \mathbf{r}_t \leq 0. \quad (\text{S45})$$

Finally, for arbitrary  $\epsilon \in \max\{B, 1\}$ , if  $|\mathbf{x}_n^\top \mathbf{r}_t| \geq \epsilon$ , we have that

$$\left(\exp\left(-\mathbf{x}_n^\top \mathbf{r}_t + \frac{1}{2}\mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n\right) - 1\right) \mathbf{x}_n^\top \mathbf{r}_t \leq \left(\exp\left(-\frac{B}{2}\right) - 1\right) \epsilon < 0. \quad (\text{S46})$$

Furthermore, let  $\gamma_* = \min_{n \in \mathcal{S}} \tilde{\boldsymbol{w}}^\top \mathbf{x}_n$  and  $\gamma^* = \max_{n \in \mathcal{S}} \tilde{\boldsymbol{w}}^\top \mathbf{x}_n$ . Now, by taking  $\epsilon \geq \max\{B, 1\}$  large enough such that

$$\left|\exp(-\gamma^*) \left(\exp\left(-\frac{B}{2}\right) - 1\right) \epsilon\right| \geq |\mathcal{S}| \exp(-\gamma_*) \left(\exp\left(\frac{B}{2}\right) - 1\right) \frac{B}{2}, \quad (\text{S47})$$

if there exists a support vector  $n \in \mathcal{S}$  such that  $|\mathbf{x}_n^\top \mathbf{r}_t| \geq \epsilon$ , then

$$\frac{1}{t} \sum_{n \in \mathcal{S}} \exp(-\tilde{\boldsymbol{w}}^\top \mathbf{x}_n) \left(\exp\left(-\mathbf{x}_n^\top \mathbf{r}_t + \frac{1}{2}\mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n\right) - 1\right) \mathbf{x}_n^\top \mathbf{r}_t \leq 0. \quad (\text{S48})$$

The idea of this is that if there exists a support vector such that  $|\mathbf{x}_n^\top \mathbf{r}_t|$  is sufficiently big, then the first bracket in Eq. (S42) is negative.

On the other hand, for the second bracket in Eq. (S42), note that for  $n \notin \mathcal{S}$ , we have that  $\mathbf{x}_n^\top \hat{\mathbf{w}} \geq \kappa$ , and hence

$$\begin{aligned} & \sum_{n \notin \mathcal{S}} \left(\frac{1}{t}\right)^{\hat{\mathbf{w}}^\top \mathbf{x}_n} \exp\left(-\hat{\mathbf{w}}^\top \mathbf{x}_n + \frac{1}{2} \mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n\right) \exp(-\mathbf{x}_n^\top \mathbf{r}_t) \mathbf{x}_n^\top \mathbf{r}_t \\ & \leq \frac{1}{t^\kappa} \exp\left(\frac{1}{2} \|\mathbf{S}_0\|_F^2 \max_n \mathbf{x}_n^\top \mathbf{x}_n\right) \sum_{n \notin \mathcal{S}} \exp(-\hat{\mathbf{w}}^\top \mathbf{x}_n) = \mathcal{O}\left(\frac{1}{t^\kappa}\right), \end{aligned} \quad (\text{S49})$$

where in the last line we used that  $ze^{-z} \leq 1, \forall z \in \mathbb{R}$  and fact that  $\|\mathbf{S}_t \mathbf{S}_t^\top\|_F \leq \|\mathbf{S}_0\|_F^2 < \infty$ .

We will now combine the results from above to show that the residual  $\mathbf{r}_t$  is bounded in the following way: if there exists a support vector  $n \in \mathcal{S}$  such that  $|\mathbf{x}_n^\top \mathbf{r}_t| \geq \epsilon$  for big enough  $\epsilon > 0$ , then  $\frac{1}{2} \frac{d}{dt} \|\mathbf{r}_t\|^2 = \mathcal{O}(t^{-\kappa})$ . If such a support vector does not exist at time  $t$ , we will show that  $\mathbf{r}_t$  is contained inside a compact set. To that end, if  $\|\mathbf{P}_\mathcal{S} \mathbf{r}_t\| \geq \epsilon_1$ , we have that

$$\max_{n \in \mathcal{S}} |\mathbf{x}_n^\top \mathbf{r}_t|^2 \geq \frac{1}{|\mathcal{S}|} \sum_{n \in \mathcal{S}} |\mathbf{x}_n^\top \mathbf{P}_\mathcal{S} \mathbf{r}_t|^2 = \frac{1}{|\mathcal{S}|} \|\mathbf{X}_\mathcal{S}^\top \mathbf{P}_\mathcal{S} \mathbf{r}_t\|^2 \geq \frac{1}{|\mathcal{S}|} \sigma_{\min}^2(\mathbf{X}_\mathcal{S}) \epsilon_1^2, \quad (\text{S50})$$

where in the first inequality we used the fact that  $\mathbf{P}_\mathcal{S}^\top \mathbf{x}_n = \mathbf{x}_n$  for  $n \in \mathcal{S}$ . Hence by choosing  $\epsilon_1$  such that  $\sigma_{\min}^2(\mathbf{X}_\mathcal{S}) \epsilon_1^2 / |\mathcal{S}| = \epsilon^2$ , where the  $\epsilon$  is chosen in Eq. (S47), we have that

$$\|\mathbf{P}_\mathcal{S} \mathbf{r}_t\| \geq \epsilon_1 \Rightarrow \frac{1}{2} \frac{d}{dt} \|\mathbf{r}_t\|^2 = \mathcal{O}(t^{-\kappa}). \quad (\text{S51})$$

On the other hand, if  $\|\mathbf{P}_\mathcal{S} \mathbf{r}_t\| \leq \epsilon_1$ , recall that

$$\mathbf{r}_t = (\boldsymbol{\mu}_t - \boldsymbol{\mu}_0) + \boldsymbol{\mu}_0 - \hat{\mathbf{w}} \log t - \tilde{\mathbf{w}}, \quad (\text{S52})$$

and since all updates to the mean parameter are in the space spanned by the support vectors (Assumption 2), we have that

$$\bar{\mathbf{P}}_\mathcal{S} \mathbf{r}_t = \bar{\mathbf{P}}_\mathcal{S} \boldsymbol{\mu}_0 - \bar{\mathbf{P}}_\mathcal{S} \tilde{\mathbf{w}}. \quad (\text{S53})$$

We can now conclude that

$$\|\mathbf{P}_\mathcal{S} \mathbf{r}_t\| \leq \epsilon_1 \Rightarrow \|\mathbf{r}_t\| \leq \|\mathbf{P}_\mathcal{S} \mathbf{r}_t\| + \|\bar{\mathbf{P}}_\mathcal{S} \mathbf{r}_t\| \leq \epsilon_1 + \|\bar{\mathbf{P}}_\mathcal{S} \boldsymbol{\mu}_0\| + \|\bar{\mathbf{P}}_\mathcal{S} \tilde{\mathbf{w}}\| < \infty. \quad (\text{S54})$$

Finally, combining the results from Eq. (S49) and Eq. (S54), recalling that  $\kappa > 1$ , we have that  $\|\mathbf{r}_t\|$  is bounded for all  $t > 0$ . This completes the first part of the proof and shows that

$$\boldsymbol{\mu}_t = \hat{\mathbf{w}} \log t + \tilde{\mathbf{w}} + \mathbf{r}_t = \hat{\mathbf{w}} \log t + \mathcal{O}(1), \quad (\text{S55})$$

and in particular

$$\lim_{t \rightarrow \infty} \frac{\boldsymbol{\mu}_t}{\|\boldsymbol{\mu}_t\|} = \frac{\hat{\mathbf{w}}}{\|\hat{\mathbf{w}}\|}. \quad (\text{S56})$$

We proceed by showing that the limit covariance parameter vanishes in the span of the support vectors.

**Covariance parameter** We begin by substituting the definition of the residual  $\mathbf{r}_t$  (S31) into the gradient flow dynamics for the covariance factor  $\mathbf{S}_t$ :

$$\begin{aligned} \dot{\mathbf{S}}_t &= -\nabla_{\mathbf{S}} \bar{\ell}(\boldsymbol{\mu}_t, \mathbf{S}_t) \\ &= -\sum_{n=1}^N \exp(-\boldsymbol{\mu}_t^\top \mathbf{x}_n + \frac{1}{2} \mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n) \mathbf{x}_n \mathbf{x}_n^\top \mathbf{S}_t. \end{aligned} \quad (\text{S57})$$

Next, we split the sum into contributions from support vectors and non-support vectors. For  $n \in \mathcal{S}$ , we use the property  $\mathbf{x}_n^\top \hat{\mathbf{w}} = 1$ ; for  $n \notin \mathcal{S}$ , the margin is strictly larger than one, which introduces higher-order decay in  $t$ :

$$\begin{aligned} \dot{\mathbf{S}}_t &= -\sum_{n \in \mathcal{S}} \frac{1}{t} \exp(-\tilde{\mathbf{w}}^\top \mathbf{x}_n - \mathbf{r}_t^\top \mathbf{x}_n) \exp(\frac{1}{2} \mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n) \mathbf{x}_n \mathbf{x}_n^\top \mathbf{S}_t \\ &\quad - \sum_{n \notin \mathcal{S}} \left(\frac{1}{t}\right)^{\mathbf{x}_n^\top \hat{\mathbf{w}}} \exp(-\tilde{\mathbf{w}}^\top \mathbf{x}_n - \mathbf{r}_t^\top \mathbf{x}_n) \exp(\frac{1}{2} \mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n) \mathbf{x}_n \mathbf{x}_n^\top \mathbf{S}_t. \end{aligned} \quad (\text{S58})$$

Since  $\mathbf{r}_t$  is bounded (from the previous part of the proof), the exponential prefactor is uniformly bounded away from zero. We formalize this by defining

$$C := \min_{n \in [N]} \min_{t \geq 0} \exp(-\tilde{\mathbf{w}}^\top \mathbf{x}_n - \mathbf{r}_t^\top \mathbf{x}_n) > 0. \quad (\text{S59})$$

We also let  $\sigma_{\min}$  denote the smallest non-zero eigenvalue of the matrix  $\sum_{n \in \mathcal{S}} \mathbf{x}_n \mathbf{x}_n^\top$ . Finally, to measure the size of  $\mathbf{S}_t$  restricted to the support-vector subspace, we define

$$\Delta_t := \text{tr}(\mathbf{P}_\mathcal{S} \mathbf{S}_t \mathbf{S}_t^\top \mathbf{P}_\mathcal{S}).$$

We now compute the derivative of  $\Delta_t$  over time. Differentiating and substituting the dynamics of  $\mathbf{S}_t$  yields

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} \Delta_t &= \text{tr}(\mathbf{P}_\mathcal{S} \dot{\mathbf{S}}_t \mathbf{S}_t^\top \mathbf{P}_\mathcal{S}) \\ &= -\frac{1}{t} \sum_{n \in \mathcal{S}} \exp(-\tilde{\mathbf{w}}^\top \mathbf{x}_n - \mathbf{r}_t^\top \mathbf{x}_n) \exp\left(\frac{1}{2} \mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n\right) \text{tr}(\mathbf{P}_\mathcal{S} \mathbf{x}_n \mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{P}_\mathcal{S}) \\ &\quad + \mathcal{O}\left(\frac{1}{t^\kappa}\right). \end{aligned} \quad (\text{S60})$$

At this point we use two facts: 1. from (S59), the exponential prefactor is bounded below by  $C > 0$ , 2. from the definition of  $\sigma_{\min}$ , we can control the quadratic form  $\sum_{n \in \mathcal{S}} \mathbf{x}_n \mathbf{x}_n^\top$ . Applying both gives

$$\frac{1}{2} \frac{d}{dt} \Delta_t \leq -\frac{C \sigma_{\min}}{t} \Delta_t + \mathcal{O}\left(\frac{1}{t^\kappa}\right). \quad (\text{S61})$$

Finally, by Grönwall's lemma, there exists a constant  $K > 0$  and a starting time  $t_0 > 0$  such that

$$\Delta_t \leq \Delta_{t_0} \left(\frac{t}{t_0}\right)^{-2C \sigma_{\min}} + \frac{K}{2C \sigma_{\min} + \kappa - 1} t^{-(\kappa-1)}, \quad \forall t \geq t_0. \quad (\text{S62})$$

Since both  $|\mathcal{S}| C \sigma_{\min} > 0$  and  $\kappa > 1$ , we conclude that  $\Delta_t \rightarrow 0$  as  $t \rightarrow \infty$ . In words: the covariance factor vanishes when projected onto the span of the support vectors, i.e.

$$\forall n \in \mathcal{S} : \lim_{t \rightarrow \infty} \mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n = 0, \quad (\text{S63})$$

as claimed.  $\square$

### S1.2.3 COMPLETE PROOF OF THEOREM 2

We will now extend the results for the gradient flow to gradient descent and then use these results to characterize the implicit bias of gradient descent as generalized variational inference.

Throughout this proof, let

$$\mathbf{A}_t = \sum_{n=1}^N \exp\left(-\boldsymbol{\mu}_t^\top \mathbf{x}_n + \frac{1}{2} \mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n\right) \mathbf{x}_n \mathbf{x}_n^\top \quad (\text{S64})$$

be a positive definite matrix at iteration  $t$ . We begin the section with a few lemmata which will be used throughout the proof.

#### Lemma S2

Suppose that we start gradient descent from  $(\boldsymbol{\mu}_0, \mathbf{S}_0)$ . If  $\eta < \lambda_{\max}(\mathbf{A}_0)^{-1}$ , then for the gradient descent iterates

$$\mathbf{S}_{t+1} = \mathbf{S}_t - \eta \nabla_{\mathbf{S}} \bar{\ell}(\boldsymbol{\mu}_t, \mathbf{S}_t), \quad (\text{S65})$$

we have that  $\|\mathbf{S}_t\|_F \leq \|\mathbf{S}_0\|_F$  for all  $t \geq 0$ .

*Proof.* First, note that the gradient descent update for the covariance factor is given by

$$\mathbf{S}_{t+1} = \mathbf{S}_t (\mathbf{I} - \eta \mathbf{A}_t), \quad (\text{S66})$$

and hence we have that

$$\|\mathbf{S}_{t+1}\|_F = \|\mathbf{S}_t (\mathbf{I} - \eta \mathbf{A}_t)\|_F \leq \|\mathbf{S}_t\|_F \|(\mathbf{I} - \eta \mathbf{A}_t)\|_2. \quad (\text{S67})$$

Now, since  $\eta \leq \lambda_{\max}(\mathbf{A}_0)^{-1} \leq \lambda_{\max}(\mathbf{A}_t)^{-1}$  for all  $t \geq 0$  and noting that  $\mathbf{A}_t \succeq 0$ , we have that

$$\|(\mathbf{I} - \eta\mathbf{A}_t)\|_2 \leq 1, \quad (\text{S68})$$

and therefore

$$\|\mathbf{S}_{t+1}\|_F \leq \|\mathbf{S}_t\|_F. \quad (\text{S69})$$

Finally, we can conclude that  $\|\mathbf{S}_t\|_F \leq \|\mathbf{S}_0\|_F$  for all  $t \geq 0$ , as required.  $\square$

### Lemma S3

Suppose that we start gradient descent from  $(\boldsymbol{\mu}_0, \mathbf{S}_0)$ . If  $\eta < \lambda_{\max}(\mathbf{A}_0)^{-1}$ , then for the gradient descent iterates

$$\boldsymbol{\mu}_{t+1} = \boldsymbol{\mu}_t - \eta \nabla_{\boldsymbol{\mu}} \bar{\ell}(\boldsymbol{\mu}_t, \mathbf{S}_t), \quad (\text{S70})$$

we have that  $\sum_{u=0}^{\infty} \|\nabla_{\boldsymbol{\mu}} \bar{\ell}(\boldsymbol{\mu}_u, \mathbf{S}_u)\|^2 < \infty$ . Consequently, we also have that  $\lim_{t \rightarrow \infty} \|\nabla_{\boldsymbol{\mu}} \bar{\ell}(\boldsymbol{\mu}_t, \mathbf{S}_t)\|^2 = 0$ .

*Proof.* Note that our loss function is not globally smooth in  $\boldsymbol{\mu}$ . However, if we initialize at  $(\boldsymbol{\mu}_0, \mathbf{S}_0)$ , the gradient descent iterates with  $\eta < \lambda_{\max}(\mathbf{A}_0)^{-1}$  maintain bounded local smoothness. The statement now follows directly from Lemma 10 in Soudry et al. [4].  $\square$

### Lemma S4

By choosing  $\epsilon_1$  as in Eq. (S51), if  $\|\mathbf{P}_S \mathbf{r}_t\| \geq \epsilon_1$ , we have that

$$(\mathbf{r}_{t+1} - \mathbf{r}_t)^\top \mathbf{r}_t \leq \mathcal{O}\left(\frac{1}{t^\kappa}\right) + \mathcal{O}\left(\frac{1}{t^2}\right) \|\mathbf{r}_t\|. \quad (\text{S71})$$

If  $\|\mathbf{P}_S \mathbf{r}_t\| < \epsilon_1$ , there exists a constant  $C$  such that

$$(\mathbf{r}_{t+1} - \mathbf{r}_t)^\top \mathbf{r}_t \leq C. \quad (\text{S72})$$

*Proof.* We follow similar steps as in the gradient flow case. It holds that

$$\begin{aligned} & (\mathbf{r}_{t+1} - \mathbf{r}_t)^\top \mathbf{r}_t \\ &= (-\eta \nabla_{\boldsymbol{\mu}}(\boldsymbol{\mu}_t, \mathbf{S}_t) - \hat{\mathbf{w}}(\log(t+1) - \log(t)))^\top \mathbf{r}_t \\ &= \eta \sum_{n=1}^N \exp\left(-\boldsymbol{\mu}_t^\top \mathbf{x}_n + \frac{1}{2} \mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n\right) \mathbf{x}_n^\top \mathbf{r}_t - \hat{\mathbf{w}}^\top \mathbf{r}_t \log(1+t^{-1}) \\ &= \hat{\mathbf{w}}^\top \mathbf{r}_t (t^{-1} - \log(1+t^{-1})) + \eta \sum_{n \notin \mathcal{S}} \exp\left(-\boldsymbol{\mu}_t^\top \mathbf{x}_n + \frac{1}{2} \mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n\right) \mathbf{x}_n^\top \mathbf{r}_t \\ &+ \eta \sum_{n \in \mathcal{S}} \left[ -\frac{1}{t} \exp(-\tilde{\mathbf{w}}^\top \mathbf{x}_n) + \exp\left(-\boldsymbol{\mu}_t^\top \mathbf{x}_n + \frac{1}{2} \mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n\right) \right] \mathbf{x}_n^\top \mathbf{r}_t, \end{aligned} \quad (\text{S73})$$

where in the last equality we used Equation (S33) to expand  $\hat{\mathbf{w}}^\top \mathbf{r}_t$ . Furthermore, we can bound all four terms as follows, beginning with the first term:

$$\hat{\mathbf{w}}^\top \mathbf{r}_t (t^{-1} - \log(1+t^{-1})) \leq \|\mathbf{r}_t\| \mathcal{O}\left(\frac{1}{t^2}\right), \quad (\text{S74})$$

where we used that  $\log(1+t^{-1}) = t^{-1} + \mathcal{O}(t^{-2})$ . For the second term, using the same argument as in Equation (S49), we derive that

$$\eta \sum_{n \notin \mathcal{S}} \exp\left(-\boldsymbol{\mu}_t^\top \mathbf{x}_n + \frac{1}{2} \mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n\right) \mathbf{x}_n^\top \mathbf{r}_t \leq \mathcal{O}\left(\frac{1}{t^\kappa}\right). \quad (\text{S75})$$

For the third item, from Eq. (S48) and Eq. (S50), we have that  $\|\mathbf{P}_S \mathbf{r}_t\| \geq \epsilon_1$  implies that

$$\eta \sum_{n \in \mathcal{S}} \left[ -\frac{1}{t} \exp(-\tilde{\mathbf{w}}^\top \mathbf{x}_n) + \exp\left(-\boldsymbol{\mu}_t^\top \mathbf{x}_n + \frac{1}{2} \mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n\right) \right] \mathbf{x}_n^\top \mathbf{r}_t \leq 0. \quad (\text{S76})$$

The first result follows from combining the above three inequalities.

Next, if  $\|\mathbf{P}_S \mathbf{r}_t\| < \epsilon_1$ , by defining  $B := \|\mathbf{S}_0\|_F^2$ , following the steps in Eq. (S44), we have that

$$\eta \sum_{n \notin S} \exp\left(-\boldsymbol{\mu}_t^\top \mathbf{x}_n + \frac{1}{2} \mathbf{x}_n^\top \mathbf{S}_t \mathbf{x}_n\right) \mathbf{x}_n^\top \mathbf{r}_t \leq \eta |\mathcal{S}| \left(\exp\left(\frac{B}{2}\right) - 1\right) \frac{B}{2}, \quad (\text{S77})$$

and hence, combining this with Assumption 2 which implies that  $\mathbf{r}_t$  is bounded as in Eq. (S54), one can find a constant  $C$  such that

$$(\mathbf{r}_{t+1} - \mathbf{r}_t)^\top \mathbf{r}_t \leq C. \quad (\text{S78})$$

□

## Proof of Theorem 2

*Proof.* As in the simple version of the proof, we begin by considering the convergence behavior of the mean parameter  $\boldsymbol{\mu}_t$ .

**Mean parameter** Our goal is again to show that  $\|\mathbf{r}_t\|$  is bounded. To that end, we will provide an upper bound to the following equation

$$\|\mathbf{r}_{t+1}\|^2 = \|\mathbf{r}_{t+1} - \mathbf{r}_t\|^2 + 2(\mathbf{r}_{t+1} - \mathbf{r}_t)^\top \mathbf{r}_t + \|\mathbf{r}_t\|^2. \quad (\text{S79})$$

First, consider the first term in the above equation:

$$\begin{aligned} & \|\mathbf{r}_{t+1} - \mathbf{r}_t\|^2 \\ &= \|\boldsymbol{\mu}_{t+1} - \hat{\mathbf{w}} \log(t+1) - \tilde{\mathbf{w}} - \boldsymbol{\mu}_t + \hat{\mathbf{w}} \log(t) + \tilde{\mathbf{w}}\|^2 \\ &= \|\eta \nabla_{\boldsymbol{\mu}} \bar{\ell}(\boldsymbol{\mu}_t, \mathbf{S}_t) - \hat{\mathbf{w}} \log(1+t^{-1})\|^2 \\ &\leq 2 \left[ \eta^2 \|\nabla_{\boldsymbol{\mu}} \bar{\ell}(\boldsymbol{\mu}_t, \mathbf{S}_t)\|^2 + \|\hat{\mathbf{w}}\|^2 \log^2(1+t^{-1}) \right] \\ &\leq 2 \left[ \eta^2 \|\nabla_{\boldsymbol{\mu}} \bar{\ell}(\boldsymbol{\mu}_t, \mathbf{S}_t)\|^2 + \|\hat{\mathbf{w}}\|^2 t^{-2} \right] \end{aligned} \quad (\text{S80})$$

where in the first inequality we used the standard inequality that  $(x+y)^2 \leq 2(x^2+y^2)$ , and in the second inequality we used the fact that  $\log(1+x) \leq x$  for  $x \geq 0$ . Now, from Lemma S3 and the fact that  $t^{-2}$  is summable, we conclude that there exists  $C_1 < \infty$  such that

$$\sum_{t=1}^{\infty} \|\mathbf{r}_{t+1} - \mathbf{r}_t\|^2 \leq C_1 < \infty. \quad (\text{S81})$$

Next, for the second term, recall that in Lemma S4 we showed that if  $\|\mathbf{P}_S \mathbf{r}_t\| \geq \epsilon_1$ , then, for some constants  $C_2, C_3 < \infty$ , we have that, eventually

$$(\mathbf{r}_{t+1} - \mathbf{r}_t)^\top \mathbf{r}_t \leq C_2 \frac{1}{t^\kappa} + C_3 \frac{1}{t^2} \|\mathbf{r}_t\|, \quad (\text{S82})$$

and that if  $\|\mathbf{P}_S \mathbf{r}_t\| < \epsilon_1$ , then there exists a constant  $C_4 < \infty$  such that

$$(\mathbf{r}_{t+1} - \mathbf{r}_t)^\top \mathbf{r}_t \leq C_4. \quad (\text{S83})$$

We will show that when  $\|\mathbf{P}_S \mathbf{r}_t\| < \epsilon_1$ , the residual  $\mathbf{r}_t$  is contained in a compact set, and when  $\|\mathbf{P}_S \mathbf{r}_t\| \geq \epsilon_1$ , the residual  $\mathbf{r}_t$  can't escape to infinity. We now formally show this claim.

Let  $S_1$  be the first time such that  $\|\mathbf{P}_S \mathbf{r}_t\| \geq \epsilon_1$ , if such a time does not exist, we are done since the support vectors span the data and hence  $\|\mathbf{r}_t\|$  is bounded. Now, let  $T_1$  be the first time after  $S_1$  such that  $\|\mathbf{P}_S \mathbf{r}_t\| < \epsilon_1$ , where we allow  $T_1 = \infty$  if such a time does not exist. Continuing in this manner, we define the sequences  $S_1 < T_1 < S_2 < T_2 < \dots$ , where we allow  $T_i = \infty$  for some  $i$ .

We proceed by showing that  $\|\mathbf{r}_t\|$  is uniformly bounded on each of the intervals  $[S_i, T_i)$ . To that end, note that for  $t \in [S_i, T_i)$ , we have that

$$\|\mathbf{r}_{t+1}\|^2 - \|\mathbf{r}_t\|^2 \leq 2C_2 \frac{1}{t^\kappa} + 2C_3 \frac{1}{t^2} \|\mathbf{r}_t\| + \|\mathbf{r}_{t+1} - \mathbf{r}_t\|^2, \quad (\text{S84})$$

and hence, using the fact that  $\kappa > 1$ , by the discrete version of Grönwall's lemma, that

$$\max_{t \in [S_i, T_i)} (\|\mathbf{r}_t\|^2 - \|\mathbf{r}_{S_i}\|^2) \leq K, \quad (\text{S85})$$

for some constant  $K < \infty$  independent of  $i$ . Furthermore, we also know from Eq. (S83) that

$$\|\mathbf{r}_{S_i}\| \leq \epsilon_1 + 2C_4 + \|\mathbf{r}_{S_i} - \mathbf{r}_{S_{i-1}}\|^2 \leq \epsilon_1 + 2C_4 + \max_{t \geq 0} \|\mathbf{r}_{t+1} - \mathbf{r}_t\|^2 < \infty, \quad (\text{S86})$$

showing that the first jump outside the  $\epsilon_1$ -ball is bounded. Combining the two results, we conclude that  $\|\mathbf{r}_t\|$  is uniformly bounded on each of the intervals  $[S_i, T_i)$ .

Finally, by noting that the support vectors span the data, we have that  $\|\mathbf{r}_t\|$  is uniformly bounded on each of the intervals  $[T_i, S_{i+1})$ . Combining the two results, we conclude that  $\|\mathbf{r}_t\|$  is uniformly bounded for all  $t \geq 0$  and hence we have that

$$\lim_{t \rightarrow \infty} \frac{\boldsymbol{\mu}_t}{\|\boldsymbol{\mu}_t\|} = \frac{\hat{\boldsymbol{w}}}{\|\hat{\boldsymbol{w}}\|} \quad (\text{S87})$$

and the following lemma.

**Lemma S5**

For the mean parameter  $\boldsymbol{\mu}_t$ , we have that

$$\boldsymbol{\mu}_t = \log(t)\hat{\boldsymbol{w}} + \mathcal{O}(1). \quad (\text{S88})$$

*Proof.* This follows immediately from the definition of the residual in Equation (S31):

$$\boldsymbol{\mu}_t = \hat{\boldsymbol{w}} \log t + \mathbf{r}_t + \tilde{\boldsymbol{w}}_t,$$

and the fact that  $\mathbf{r}_t$  and  $\tilde{\boldsymbol{w}}_t$  are bounded as we showed above.  $\square$

We continue with the analysis of the covariance parameter over optimization iterations.

**Covariance parameter** As before, let  $\Delta_t = \text{tr}(\mathbf{P}_S \mathbf{S}_t \mathbf{S}_t^\top \mathbf{P}_S)$  be the trace of the projection of the covariance parameter on the space of support vectors in  $\mathcal{S}$ . By following the ideas from the gradient flow case, we have the following dynamics:

$$\begin{aligned} \Delta_{t+1} &= \text{tr}(\mathbf{P}_S (\mathbf{I} - \eta \mathbf{A}_t) \mathbf{S}_t \mathbf{S}_t^\top (\mathbf{I} - \eta \mathbf{A}_t)^\top \mathbf{P}_S) \\ &= \text{tr}(\mathbf{P}_S \mathbf{S}_t \mathbf{S}_t^\top \mathbf{P}_S) - 2\eta \text{tr}(\mathbf{P}_S \mathbf{S}_t \mathbf{S}_t^\top \mathbf{A}_t \mathbf{P}_S) + \eta^2 \text{tr}(\mathbf{P}_S \mathbf{A}_t \mathbf{S}_t \mathbf{S}_t^\top \mathbf{A}_t \mathbf{P}_S) \\ &\leq \Delta_t - \frac{2\eta}{t} C \sigma_{\min} \text{tr}(\mathbf{P}_S \mathbf{S}_t \mathbf{S}_t^\top \mathbf{P}_S) + \mathcal{O}\left(\frac{1}{t^\kappa}\right) + \mathcal{O}\left(\frac{1}{t^2}\right) \\ &= \Delta_t - \frac{2\eta}{t} C \sigma_{\min} \Delta_t + \mathcal{O}\left(\frac{1}{t^\kappa}\right) + \mathcal{O}\left(\frac{1}{t^2}\right), \end{aligned} \quad (\text{S89})$$

where we used the same arguments as in Equation (S60) to derive the last inequality, in addition to noting that  $\lambda_{\max}(\mathbf{A}_t^2) \leq \mathcal{O}\left(\frac{1}{t^2}\right)$  in order to bound the last term. Hence, we can write

$$\Delta_{t+1} - \Delta_t \leq -\frac{2\eta}{t} C \sigma_{\min} \Delta_t + \mathcal{O}\left(\frac{1}{t^\kappa}\right) + \mathcal{O}\left(\frac{1}{t^2}\right). \quad (\text{S90})$$

Again, by the discrete version of Grönwall's lemma, we derive the equivalent result to Eq. (S62). Now, noting that  $\sum_t \frac{1}{t}$  diverges, the fact that  $\kappa > 1$  and  $\eta C \sigma_{\min} > 0$ , we conclude that  $\Delta_t$  converges to zero. This implies that the covariance parameter converges to zero in the span of the support vectors, i.e.

$$\forall n \in \mathcal{S} : \lim_{t \rightarrow \infty} \mathbf{x}_n^\top \mathbf{S}_t \mathbf{S}_t^\top \mathbf{x}_n = 0, \quad (\text{S91})$$

as desired.

**Characterization as Generalized Variational Inference** As a final step we need to show that the solution identified by gradient descent if appropriately transformed identifies the minimum 2-Wasserstein solution in the feasible set. Define the feasible set

$$\Theta_* = \{(\boldsymbol{\mu}, \mathbf{S}) \mid \mathbf{P}_S \boldsymbol{\mu} = \hat{\boldsymbol{w}} \quad \text{and} \quad \forall n \in \mathcal{S} : \text{Var}_{q_\theta}(f_{\boldsymbol{w}}(\mathbf{x}_n)) = 0\} \quad (\text{S92})$$

$$= \{(\boldsymbol{\mu}, \mathbf{S}) \mid \mathbf{P}_S \boldsymbol{\mu} = \hat{\boldsymbol{w}} \quad \text{and} \quad \forall n \in \mathcal{S} : \mathbf{x}_n^\top \mathbf{S} \mathbf{S}^\top \mathbf{x}_n = 0\} \quad (\text{S93})$$

and the variational parameters identified by rescaled gradient descent as

$$\boldsymbol{\theta}_*^{\text{rGD}} = \lim_{t \rightarrow \infty} \boldsymbol{\theta}_t^{\text{rGD}} = \lim_{t \rightarrow \infty} \left( \frac{1}{\log(t)} \boldsymbol{\mu}_t + \mathbf{P}_{\text{null}(\mathbf{X})} \boldsymbol{\mu}_0, \mathbf{S}_t \right). \quad (\text{S94})$$

It holds by Lemma S5 that

$$\mathbf{P}_S \boldsymbol{\mu}_*^{\text{rGD}} = \mathbf{P}_S \left( \lim_{t \rightarrow \infty} \frac{1}{\log(t)} \boldsymbol{\mu}_t \right) + \mathbf{0} = \mathbf{P}_S \hat{\boldsymbol{w}} = \hat{\boldsymbol{w}} \quad (\text{S95})$$

and additionally by Equation (S91) we have for all  $n \in \mathcal{S}$  that

$$\mathbf{x}_n^\top \mathbf{S}_*^{\text{rGD}} (\mathbf{S}_*^{\text{rGD}})^\top \mathbf{x}_n = \lim_{t \rightarrow \infty} \mathbf{x}_n^\top \mathbf{S}_t (\mathbf{S}_t)^\top \mathbf{x}_n = 0. \quad (\text{S96})$$

Therefore, the limit point  $\boldsymbol{\theta}_*^{\text{rGD}}$  of rescaled gradient descent is in the feasible set. It remains to show that it is also a minimizer of the 2-Wasserstein distance to the prior / initialization. We will first show a more general result that does not require Assumption 2.

To that end define  $(\mathbf{V}_S \quad \mathbf{V}_{\mathbf{X} \perp \mathcal{S}} \quad \mathbf{V}_{\text{null}(\mathbf{X})}) \in \mathbb{R}^{P \times P}$  where  $\mathbf{V}_S \in \mathbb{R}^{P \times P_S}$  is an orthonormal basis of the span of the support vectors  $\text{range}(\mathbf{X}_S^\top)$ ,  $\mathbf{V}_{\mathbf{X} \perp \mathcal{S}} \in \mathbb{R}^{P \times (N - P_S)}$  an orthonormal basis of its orthogonal complement in  $\text{range}(\mathbf{X}^\top)$  and  $\mathbf{V}_{\text{null}(\mathbf{X})} \in \mathbb{R}^{P \times (P - N)}$  the corresponding orthonormal basis of the null space  $\text{null}(\mathbf{X})$  of the data. Let  $\mathbf{V} = (\mathbf{V}_S \quad \mathbf{V}_{\text{null}(\mathbf{X})}) \in \mathbb{R}^{P \times (P - N + P_S)}$  and define the projected variational distribution and prior onto the span of the support vectors and the null space of the data as

$$q_\theta^{\text{proj}}(\tilde{\boldsymbol{w}}) = \mathcal{N}(\tilde{\boldsymbol{w}}; \mathbf{P}_V \boldsymbol{\mu}, \mathbf{P}_V \boldsymbol{\Sigma} \mathbf{P}_V^\top) = \mathcal{N}(\tilde{\boldsymbol{w}}; \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}) \quad (\text{S97})$$

$$p^{\text{proj}}(\tilde{\boldsymbol{w}}) = \mathcal{N}(\tilde{\boldsymbol{w}}; \mathbf{P}_V \boldsymbol{\mu}_0, \mathbf{P}_V \boldsymbol{\Sigma}_0 \mathbf{P}_V^\top) = \mathcal{N}(\tilde{\boldsymbol{w}}; \tilde{\boldsymbol{\mu}}_0, \tilde{\boldsymbol{\Sigma}}_0) \quad (\text{S98})$$

where  $\tilde{\boldsymbol{w}} \in \mathbb{R}^{P - N + P_S}$ . Now earlier we showed that the limit point of rescaled gradient descent is in the feasible set, defined in Equation (S94), and thus the same holds for the projected limit point of rescaled gradient descent, i.e.

$$(\tilde{\boldsymbol{\mu}}_*^{\text{rGD}}, \tilde{\boldsymbol{\Sigma}}_*^{\text{rGD}}) \in \Theta_* \quad (\text{S99})$$

in particular

$$\mathbf{P}_S \tilde{\boldsymbol{\mu}}_*^{\text{rGD}} = \mathbf{P}_S \boldsymbol{\mu}_*^{\text{rGD}} = \hat{\boldsymbol{w}}, \quad (\text{S100})$$

$$\forall n \in \mathcal{S} : \quad \mathbf{x}_n^\top \tilde{\boldsymbol{\Sigma}}_*^{\text{rGD}} (\tilde{\boldsymbol{\Sigma}}_*^{\text{rGD}})^\top \mathbf{x}_n = \mathbf{x}_n^\top \mathbf{S}_*^{\text{rGD}} (\mathbf{S}_*^{\text{rGD}})^\top \mathbf{x}_n = 0. \quad (\text{S101})$$

Therefore, we have for all  $n \in \mathcal{S}$  that

$$0 = \mathbf{x}_n^\top \tilde{\boldsymbol{\Sigma}}_*^{\text{rGD}} (\tilde{\boldsymbol{\Sigma}}_*^{\text{rGD}})^\top \mathbf{x}_n = \|(\tilde{\boldsymbol{\Sigma}}_*^{\text{rGD}})^\top \mathbf{x}_n\|_2^2 \iff (\tilde{\boldsymbol{\Sigma}}_*^{\text{rGD}})^\top \mathbf{x}_n = \mathbf{0} \quad (\text{S102})$$

$$\iff (\tilde{\boldsymbol{\Sigma}}_*^{\text{rGD}})^\top \mathbf{V}_S = \mathbf{0} \quad (\text{S103})$$

and thus  $\mathbf{V}_S^\top \tilde{\boldsymbol{\Sigma}}_*^{\text{rGD}} (\tilde{\boldsymbol{\Sigma}}_*^{\text{rGD}})^\top \mathbf{V}_S = \mathbf{0}$ . Therefore by Lemma S1 it holds for the squared 2-Wasserstein distance between the projected limit point of rescaled gradient descent and the projected prior that

$$\begin{aligned} W_2^2(q_{\tilde{\boldsymbol{\theta}}_*}^{\text{proj}}, p^{\text{proj}}) &\stackrel{\pm c}{=} \|\mathbf{V}_S^\top \tilde{\boldsymbol{\mu}} - \mathbf{V}_S^\top \tilde{\boldsymbol{\mu}}_0\|_2^2 + W_2^2\left(\mathcal{N}\left(\mathbf{V}_{\text{null}}^\top \tilde{\boldsymbol{\mu}}, \mathbf{V}_{\text{null}}^\top \tilde{\boldsymbol{\Sigma}} \mathbf{V}_{\text{null}}\right), \mathcal{N}\left(\mathbf{V}_{\text{null}}^\top \tilde{\boldsymbol{\mu}}_0, \mathbf{V}_{\text{null}}^\top \tilde{\boldsymbol{\Sigma}}_0 \mathbf{V}_{\text{null}}\right)\right) \\ &= \left\| \begin{pmatrix} \mathbf{V}_S^\top \tilde{\boldsymbol{\mu}} - \mathbf{V}_S^\top \tilde{\boldsymbol{\mu}}_0 \\ \mathbf{0} \end{pmatrix} \right\|_2^2 + W_2^2\left(\mathcal{N}\left(\mathbf{V}_{\text{null}}^\top \tilde{\boldsymbol{\mu}}, \mathbf{V}_{\text{null}}^\top \tilde{\boldsymbol{\Sigma}} \mathbf{V}_{\text{null}}\right), \mathcal{N}\left(\mathbf{V}_{\text{null}}^\top \tilde{\boldsymbol{\mu}}_0, \mathbf{V}_{\text{null}}^\top \tilde{\boldsymbol{\Sigma}}_0 \mathbf{V}_{\text{null}}\right)\right) \\ &= \left\| \mathbf{V} \begin{pmatrix} \mathbf{V}_S^\top \tilde{\boldsymbol{\mu}} - \mathbf{V}_S^\top \tilde{\boldsymbol{\mu}}_0 \\ \mathbf{0} \end{pmatrix} \right\|_2^2 + W_2^2\left(\mathcal{N}\left(\mathbf{V}_{\text{null}}^\top \tilde{\boldsymbol{\mu}}, \mathbf{V}_{\text{null}}^\top \tilde{\boldsymbol{\Sigma}} \mathbf{V}_{\text{null}}\right), \mathcal{N}\left(\mathbf{V}_{\text{null}}^\top \tilde{\boldsymbol{\mu}}_0, \mathbf{V}_{\text{null}}^\top \tilde{\boldsymbol{\Sigma}}_0 \mathbf{V}_{\text{null}}\right)\right) \end{aligned}$$

$$\begin{aligned}
&= \|P_S \tilde{\boldsymbol{\mu}} - P_S \tilde{\boldsymbol{\mu}}_0\|_2^2 + W_2^2\left(\mathcal{N}\left(\mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\mu}}, \mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\Sigma}} \mathbf{V}_{\text{null}}\right), \mathcal{N}\left(\mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\mu}}_0, \mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\Sigma}}_0 \mathbf{V}_{\text{null}}\right)\right) \\
&= \|\hat{\boldsymbol{w}} - P_S \tilde{\boldsymbol{\mu}}_0\|_2^2 + W_2^2\left(\mathcal{N}\left(\mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\mu}}, \mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\Sigma}} \mathbf{V}_{\text{null}}\right), \mathcal{N}\left(\mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\mu}}_0, \mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\Sigma}}_0 \mathbf{V}_{\text{null}}\right)\right) \\
&\stackrel{\pm c}{=} W_2^2\left(\mathcal{N}\left(\mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\mu}}, \mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\Sigma}} \mathbf{V}_{\text{null}}\right), \mathcal{N}\left(\mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\mu}}_0, \mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\Sigma}}_0 \mathbf{V}_{\text{null}}\right)\right)
\end{aligned}$$

where we used that  $P_S \tilde{\boldsymbol{\mu}} = \hat{\boldsymbol{w}}$  for any  $(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{S}})$  in the feasible set  $\Theta_*$ . Therefore it suffices to show that the projected solution  $\tilde{\boldsymbol{\theta}}_*^{\text{rGD}}$  minimizes

$$W_2^2\left(\mathcal{N}\left(\mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\mu}}, \mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\Sigma}} \mathbf{V}_{\text{null}}\right), \mathcal{N}\left(\mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\mu}}_0, \mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\Sigma}}_0 \mathbf{V}_{\text{null}}\right)\right) \geq 0. \quad (\text{S104})$$

We have using the definition of the iterates in Equation (9) that

$$\mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\mu}}_*^{\text{rGD}} = \mathbf{V}_{\text{null}}^T P_V \left( \lim_{t \rightarrow \infty} \frac{1}{\log(t)} \boldsymbol{\mu}_t + P_{\text{null}(\mathbf{X})} \boldsymbol{\mu}_0 \right) \quad (\text{S105})$$

$$= \mathbf{V}_{\text{null}}^T (\hat{\boldsymbol{w}} + P_{\text{null}(\mathbf{X})} \boldsymbol{\mu}_0) = \mathbf{V}_{\text{null}}^T \boldsymbol{\mu}_0 \quad (\text{S106})$$

where we used  $\hat{\boldsymbol{w}} \in \text{range}(\mathbf{X}_S^T)$ . Further, it holds for the gradient of the expected loss (S30) with respect to the covariance factor parameters that

$$\mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{S}}_*^{\text{rGD}} = \mathbf{V}_{\text{null}}^T P_V \mathbf{S}_*^{\text{rGD}} = \mathbf{V}_{\text{null}}^T \mathbf{S}_*^{\text{rGD}} = \mathbf{V}_{\text{null}}^T \left( \mathbf{S}_0 - \underbrace{\sum_{t=1}^{\infty} \eta_t \nabla_{\mathbf{S}} \bar{\ell}(\boldsymbol{\mu}_t, \mathbf{S}_t)}_{\in \text{range}(\mathbf{X}^T)} \right) \quad (\text{S107})$$

$$= \mathbf{V}_{\text{null}}^T \mathbf{S}_0 = \mathbf{V}_{\text{null}}^T P_V \mathbf{S}_0 = \mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{S}}_0. \quad (\text{S108})$$

Therefore we have that

$$W_2^2\left(\mathcal{N}\left(\mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\mu}}_*^{\text{rGD}}, \mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\Sigma}}_*^{\text{rGD}} \mathbf{V}_{\text{null}}\right), \mathcal{N}\left(\mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\mu}}_0, \mathbf{V}_{\text{null}}^T \tilde{\boldsymbol{\Sigma}}_0 \mathbf{V}_{\text{null}}\right)\right) = 0 \quad (\text{S109})$$

and thus the projected variational parameters  $\tilde{\boldsymbol{\theta}}_*^{\text{rGD}}$  are both feasible (S99) and minimize the squared 2-Wasserstein distance to the projected initialization / prior (S104). This completes the proof for the generalized version of Theorem 2 without Assumption 2, which we state here for convenience.

### Lemma S6

Given the assumptions of Theorem 2, except for Assumption 2 meaning the support vectors  $\mathbf{X}_S$  do not necessarily span the data, it holds for the limit point of rescaled gradient descent that

$$\boldsymbol{\theta}_*^{\text{rGD}} \in \arg \min_{\substack{\boldsymbol{\theta}=(\boldsymbol{\mu}, \mathbf{S}) \\ \text{s.t. } \boldsymbol{\theta} \in \Theta_*}} W_2^2\left(q_{\boldsymbol{\theta}}^{\text{proj}}, p^{\text{proj}}\right). \quad (\text{S110})$$

If in addition Assumption 2 holds, i.e. the support vectors span the training data  $\mathbf{X}$ , such that

$$\text{span}(\{\mathbf{x}_n\}_{n \in [N]}) = \text{span}(\{\mathbf{x}_n\}_{n \in \mathcal{S}}), \quad (\text{S111})$$

then the orthogonal complement of the support vectors in  $\text{range}(\mathbf{X}^T)$  has dimension  $N - P_S = 0$  and thus the projection  $P_V = I_{P \times P}$  is the identity and therefore

$$q_{\boldsymbol{\theta}}^{\text{proj}} = q_{\boldsymbol{\theta}} \quad \text{and} \quad p^{\text{proj}} = p. \quad (\text{S112})$$

This completes the proof of Theorem 2. □

### S1.3 NLL OVERFITTING AND THE NEED FOR (TEMPERATURE) SCALING

In Theorem 2, we assume we rescale the mean parameters. This is because the exponential loss can be made arbitrarily small for a mean vector that is aligned with the  $L_2$  max-margin vector simply by increasing its magnitude. In fact, the sequence of mean parameters identified by gradient descent diverges to infinity at a logarithmic rate  $\boldsymbol{\mu}_t^{\text{GD}} \approx \log(t) \hat{\boldsymbol{w}}$  as we show<sup>4</sup> in Lemma S5 and illustrate in Figure S2 (right panel).

<sup>4</sup>This has been observed previously in the deterministic case (see Theorem 3 of Soudry et al. [4]) and thus naturally also appears in our probabilistic extension.

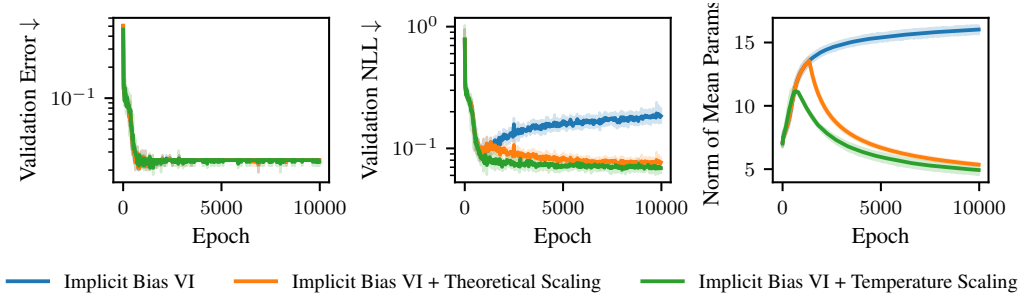


Figure S2: *NLL overfitting in classification due to implicit bias of the mean parameters.* As shown here for a two-hidden layer neural network on synthetic data, when training with vanilla SGD the mean parameters diverge to infinity  $\|\boldsymbol{\mu}_t\|_2 \approx \mathcal{O}(\log(t))$  (right) and thus the classifier will eventually overfit in terms of negative log-likelihood (left and middle). Rescaling the GD iterates as in Theorem 2 or using temperature scaling [69] avoids overfitting.

This bias of the mean parameters towards the max-margin solution does not impact the train loss or validation error, but leads to overfitting in terms of validation NLL (see Figure S2) as long as there is at least one misclassified datapoint  $\boldsymbol{x}$ , since then the (average) validation NLL is given by

$$\begin{aligned} \bar{\ell}(\boldsymbol{\theta}_t^{\text{GD}}) &= \mathbb{E}_{q_{\theta^{\text{GD}}}(w)}(\exp(-y\boldsymbol{x}^\top \boldsymbol{w})) = \exp(\boldsymbol{x}^\top \boldsymbol{\mu}_t^{\text{GD}} + \frac{1}{2}\boldsymbol{x}^\top \boldsymbol{S}_t^{\text{GD}}(\boldsymbol{S}_t^{\text{GD}})^\top \boldsymbol{x}) \\ &\approx \exp(\log(t)\boldsymbol{x}^\top \hat{\boldsymbol{w}} + \frac{1}{2}\boldsymbol{x}^\top \boldsymbol{S}_t^{\text{GD}}(\boldsymbol{S}_t^{\text{GD}})^\top \boldsymbol{x}) \rightarrow \infty \quad \text{as } t \rightarrow \infty. \end{aligned} \quad (\text{S113})$$

However, by rescaling the mean parameters as we do in Theorem 2, this can be prevented as Figure S2 (middle panel) illustrates for a two-hidden layer neural network on synthetic data. Such overfitting in terms of NLL has been studied extensively empirically with the perhaps most common remedy being Temperature Scaling (TS) [69]. As we show empirically in Figure S2, instead of using the theoretical rescaling, using temperature scaling performs very well, especially in the non-asymptotic regime, which is why we also adopt it for our experiments in Section 5.

The aforementioned divergence of the mean parameters to infinity also explains the need for the projection of the prior mean parameters in Equation (9), since any bias from the initialization vanishes in the limit of infinite training. At first glance the additional projection seems computationally prohibitive for anything but a zero mean prior, but close inspection of the implicit bias of the covariance parameters  $\boldsymbol{S}$  in Theorem 2 shows that at convergence

$$\forall n : \text{Var}_{q_{\theta}}(f_w(\boldsymbol{x}_n)) = \boldsymbol{x}_n^\top \boldsymbol{S} \boldsymbol{S}^\top \boldsymbol{x}_n = 0 \implies \text{range}(\boldsymbol{S}) \subset \text{null}(\boldsymbol{X}) \quad (\text{S114})$$

Meaning we can approximate a basis of the null space of the training data by computing a QR decomposition of the covariance factor in  $\mathcal{O}(PR^2)$  once at the end of training. For  $R = P$  the inclusion becomes an equality and the projection can be computed exactly.

## S2 PARAMETRIZATION, FEATURE LEARNING AND HYPERPARAMETER TRANSFER

**Notation** For this section we need a more detailed neural network notation. Denote an  $L$ -hidden layer, width- $D$  feedforward neural network by  $f(\boldsymbol{x}) \in \mathbb{R}_{\text{out}}^D$ , with inputs  $\boldsymbol{x} \in \mathbb{R}^{D_{\text{in}}}$ , weights  $\boldsymbol{W}^{(l)}$ , pre-activations  $\boldsymbol{h}^{(l)}(\boldsymbol{x}) \in \mathbb{R}^{D^{(l)}}$ , and post-activations (or “features”)  $\boldsymbol{g}^{(l)}(\boldsymbol{x}) \in \mathbb{R}^{D^{(l)}}$ . That is,  $\boldsymbol{h}^{(1)}(\boldsymbol{x}) = \boldsymbol{W}^{(1)}\boldsymbol{x}$  and, for  $l \in 1, \dots, L-1$ ,

$$\boldsymbol{g}^{(l)}(\boldsymbol{x}) = \phi(\boldsymbol{h}^{(l)}(\boldsymbol{x})), \quad \boldsymbol{h}^{(l+1)}(\boldsymbol{x}) = \boldsymbol{W}^{(l+1)}\boldsymbol{g}^{(l)}(\boldsymbol{x}),$$

and the network output is given by  $f(\boldsymbol{x}) = \boldsymbol{W}^{(L+1)}\boldsymbol{g}^{(L)}(\boldsymbol{x})$ , where  $\phi(\bullet)$  is an activation function.

For convenience, we may abuse notation and write  $\boldsymbol{h}^{(0)}(\boldsymbol{x}) = \boldsymbol{x}$  and  $\boldsymbol{h}^{(L+1)}(\boldsymbol{x}) = f(\boldsymbol{x})$ . Throughout we use  $\bullet^{(l)}$  to indicate the layer, subscript  $\bullet_t$  to indicate the training time (i.e., epoch),  $\Delta\bullet_t = \bullet_t - \bullet_0$  to indicate the change since initialization, and  $[\bullet]_i$ ,  $[\bullet]_{ij}$  to indicate the component within a vector or matrix.

### S2.1 DEFINITIONS OF STABILITY AND FEATURE LEARNING

The following definitions extend those of Yang and Hu [43] to the variational setting.

**Definition S1** (*bc scaling*)

In layer  $l$ , the variational parameters are initialized as

$$[\boldsymbol{\mu}_0^{(l)}]_i \sim \mathcal{N}\left(0, D^{-2b^{(l)}}\right), \quad [\mathbf{S}_0^{(l)}]_{ij} \sim \mathcal{N}\left(0, D^{-2\tilde{b}^{(l)}}\right)$$

and the learning rates for the mean and covariance parameters, respectively, are set to

$$\eta^{(l)} = \eta D^{-c^{(l)}}, \quad \tilde{\eta}^{(l)} = \eta D^{-\tilde{c}^{(l)}}.$$

The hyperparameter  $\eta$  represents a global learning rate that can be tuned, as for example in the hyperparameter transfer experiment from Section 3.4.

For the next two definitions, let  $m_r(X) = \mathbb{E}_{\mathbf{z}}((X - \mathbb{E}_{\mathbf{z}}(X))^r)$  denote the  $r$ th central moment of a random variable  $X$  with respect to  $\mathbf{z}$ , which represents all reparameterization noise in the random variable  $X$ . All Landau notation in Section S2 refers to asymptotic behavior in width  $D$  in probability over reparameterization noise  $\mathbf{z}$ . We say that a vector sequence  $\{\mathbf{v}_D\}_{D=1}^{\infty}$ , where each  $\mathbf{v}_D \in \mathbb{R}^D$ , is  $\mathcal{O}(D^{-a})$  if the scalar sequence  $\{\sqrt{\frac{1}{D}}\|\mathbf{v}_D\|^2\}_{D=1}^{\infty} = \{\text{RMSE}(\mathbf{v}_D)\}_{D=1}^{\infty}$  is  $\mathcal{O}(D^{-a})$ .

**Definition S2** (Stability of Moment  $r$ )

A neural network is *stable in moment  $r$* , if all of the following hold for all  $\mathbf{x}$  and  $l \in \{1, \dots, L\}$ .

1. At initialization ( $t = 0$ ):

(a) The pre- and post-activations are  $\Theta(1)$ :

$$m_r(\mathbf{h}_0^{(l)}(\mathbf{x})), m_r(\mathbf{g}_0^{(l)}(\mathbf{x})) = \Theta(1).$$

(b) The function is  $\mathcal{O}(1)$ :

$$m_r(f_0(\mathbf{x})) = \mathcal{O}(1).$$

2. At any point during training  $t > 0$ :

(a) The change from initialization in the pre- and post-activations are  $\mathcal{O}(1)$ :

$$\Delta m_r(\mathbf{h}_t^{(l)}(\mathbf{x})), \Delta m_r(\mathbf{g}_t^{(l)}(\mathbf{x})) = \mathcal{O}(1).$$

(b) The function is  $\mathcal{O}(1)$ :

$$m_r(f_t(\mathbf{x})) = \mathcal{O}(1).$$

**Definition S3** (Feature Learning of Moment  $r$ )

*Feature learning* occurs in moment  $r$  in layer  $l$  if, for any  $t > 0$ , the change from initialization is  $\Omega(1)$ :

$$\Delta m_r(\mathbf{g}_t^{(l)}(\mathbf{x})) = \Omega(1).$$

As we will see later, Figure S5 and Figure S6 investigate feature learning for the first two moments.

### S2.2 INITIALIZATION SCALING FOR A LINEAR NETWORK

In this section we illustrate how the initialization scaling  $\{(b^{(l)}, \tilde{b}^{(l)})\}$  can be chosen for stability. For simplicity, we consider a linear feedforward network of width  $D$  evaluated on a single input  $\mathbf{x} \in \mathbb{R}_{\text{in}}^D$ . We assume a Gaussian variational family that factorizes across layers. This implies the hidden units evolve as  $\mathbf{h}_t^{(l+1)} = \mathbf{W}_t^{(l+1)} \mathbf{h}_t^{(l)}$  and the weights are linked to the variational parameters by  $\text{vec}(\mathbf{W}_t^{(l)}) = \boldsymbol{\mu}_t^{(l)} + \mathbf{S}_t^{(l)} \mathbf{z}$ .

Therefore, the mean and variance of the  $i$ th component hidden units in layer  $l \in \{1, \dots, L+1\}$ , where  $i \in 1 \dots, D^{(l)}$ , are given by

$$\mathbb{E}_{\mathbf{z}}([\mathbf{h}_t^{(l)}]_i) = [\boldsymbol{\mu}_t^{(l)}]_i^{\top} \mathbb{E}_{\mathbf{z}}(\mathbf{h}_t^{(l-1)})$$

$$\text{Var}_z\left([\mathbf{h}_t^{(l)}]_i\right) = [\boldsymbol{\mu}_t^{(l)}]_I^\top \mathbf{C}_t^{(l-1)} [\boldsymbol{\mu}_t^{(l)}]_I + \text{tr}([\mathbf{S}_t^{(l)}]_{I,:}^\top \mathbf{A}_t^{(l-1)} [\mathbf{S}_t^{(l)}]_{I,:}),$$

where  $I = \{iD^{(l-1)}, \dots, (i+1)D^{(l-1)}\}$  and the second moment and covariance of layer- $l$  hidden units are denoted by

$$\begin{aligned} \mathbf{A}_t^{(l)} &= \mathbb{E}_z\left(\mathbf{h}_t^{(l)} \mathbf{h}_t^{(l)\top}\right) \\ \mathbf{C}_t^{(l)} &= \mathbf{A}_t^{(l)} - \mathbb{E}_z\left(\mathbf{h}_t^{(l)}\right) \mathbb{E}_z\left(\mathbf{h}_t^{(l)}\right)^\top. \end{aligned}$$

**Mean** We start with the mean of the hidden units, which conveniently depends only on the mean variational parameters and the previous layer hidden units.

$$\begin{aligned} \mathbb{E}_z\left([\mathbf{h}_0^{(l)}]_i\right) &= \sum_{j=1}^{D^{(l-1)}} [\boldsymbol{\mu}_0^{(l)}]_{I_j} \mathbb{E}_z\left([\mathbf{h}_0^{(l-1)}]_j\right) \\ &= \mathcal{O}\left(\sqrt{D^{(l-1)}} \cdot D^{-b^{(l)}} \cdot 1\right) \\ &= \begin{cases} \mathcal{O}\left(D^{-b^{(1)}}\right) & l = 1 \\ \mathcal{O}\left(D^{-(b^{(l)} - \frac{1}{2})}\right) & l \in \{2, \dots, L+1\}. \end{cases} \end{aligned}$$

Therefore, we require  $b^{(1)} \geq 0$  and  $b^{(l)} \geq \frac{1}{2}$  for  $l \in \{2, \dots, L+1\}$ .

**Variance** Next we examine the variance of hidden units. Consider the first term, which represents the contribution of the mean parameters.

$$\begin{aligned} [\boldsymbol{\mu}_0^{(l)}]_I^\top \mathbf{C}_0^{(l-1)} [\boldsymbol{\mu}_0^{(l)}]_I &= \sum_{j=1}^{D^{(l-1)}} [\boldsymbol{\mu}_0^{(l)}]_{I_j}^2 [\mathbf{C}_0^{(l-1)}]_{j,j} + \sum_{j \neq j'}^{D^{(l-1)}} [\boldsymbol{\mu}_0^{(l)}]_{I_j} [\mathbf{C}_0^{(l-1)}]_{j,j'} [\boldsymbol{\mu}_0^{(l)}]_{I_{j'}}, \\ &= \mathcal{O}\left(D^{(l-1)} \cdot D^{-2b^{(l)}} \cdot 1\right) + \mathcal{O}\left(\sqrt{D^{(l-1)}(D^{(l-1)} - 1)} \cdot D^{-b^{(l)}} \cdot 1 \cdot D^{-b^{(l)}}\right) \\ &= \mathcal{O}\left(D^{(l-1)} \cdot D^{-2b^{(l)}}\right) \\ &= \begin{cases} \mathcal{O}\left(D^{-2b^{(1)}}\right) & l = 1 \\ \mathcal{O}\left(D^{-(2b^{(l)} - 1)}\right) & l \in \{2, \dots, L+1\}. \end{cases} \end{aligned}$$

Therefore, we require  $b^{(1)} \geq 0$  and  $b^{(l)} \geq \frac{1}{2}$  for  $l \in \{2, \dots, L+1\}$ . Notice these are the same requirements as above for the mean of the hidden units. We summarize the scaling for the mean parameters as

$$b^{(l)} \geq \begin{cases} 0 & l = 1 \\ \frac{1}{2} & l \in \{2, \dots, L+1\}. \end{cases} \quad (\text{S115})$$

Now consider the second term in the variance of the hidden units. Assume the rank scales with the input and output dimension of a layer as  $R^{(l)} = (D^{(l-1)}D^{(l)})^{p^{(l)}}$ , where  $p^{(l)} \in [0, 1]$ .

$$\begin{aligned} \text{tr}([\mathbf{S}_0^{(l)}]_{I,:}^\top \mathbf{A}_0^{(l-1)} [\mathbf{S}_0^{(l)}]_{I,:}) &= \sum_{r=1}^{R^{(l)}} [\mathbf{S}_0^{(l)}]_{I,r}^\top \mathbf{A}_0^{(l-1)} [\mathbf{S}_0^{(l)}]_{I,r} \\ &= \sum_{r=1}^{R^{(l)}} \left( \sum_{j=1}^{D^{(l-1)}} [\mathbf{S}_0^{(l)}]_{I_j,r}^2 [\mathbf{A}_0^{(l-1)}]_{j,j} + \sum_{j \neq j'}^{D^{(l-1)}} [\mathbf{S}_0^{(l)}]_{I_j,r} [\mathbf{A}_0^{(l-1)}]_{j,j'} [\mathbf{S}_0^{(l)}]_{I_{j'},r} \right) \\ &= \mathcal{O}\left(R^{(l)} D^{(l-1)} \cdot D^{-2\bar{b}^{(l)}} \cdot 1\right) + \mathcal{O}\left(\sqrt{R^{(l)} D^{(l-1)} (D^{(l-1)} - 1)} \cdot D^{-\bar{b}^{(l)}} \cdot 1 \cdot D^{-\bar{b}^{(l)}}\right) \\ &= \mathcal{O}\left(R^{(l)} D^{(l-1)} D^{-2\bar{b}^{(l)}}\right) \end{aligned}$$

$$= \begin{cases} \mathcal{O}\left(D^{-(2\tilde{b}^{(1)}-p^{(1)})}\right) & l = 1 \\ \mathcal{O}\left(D^{-(2\tilde{b}^{(l)}-1-2p^{(l)})}\right) & l \in \{2, \dots, L\} \\ \mathcal{O}\left(D^{-(2\tilde{b}^{(L+1)}-1-p^{(L+1)})}\right) & l = L + 1. \end{cases}$$

Therefore we require  $\tilde{b}^{(0)} \geq \frac{p^{(1)}}{2}$ ,  $\tilde{b}^{(l)} \geq \frac{1}{2} + p^{(l)}$  for  $l \in \{2, \dots, L\}$ , and  $\tilde{b}^{(L+1)} \geq \frac{1}{2} + \frac{p^{(L+1)}}{2}$ . Notice we can write these conditions in terms of the mean scaling as

$$\tilde{b}^{(l)} \geq b^{(l)} + \begin{cases} \frac{p^{(l)}}{2} & l = 1 \\ p^{(l)} & l \in \{2, \dots, L\} \\ \frac{p^{(l)}}{2} & l = L + 1. \end{cases} \quad (\text{S116})$$

### S2.3 PROPOSED SCALING

The previous section derives the necessary conditions for stability at initialization. Recall from Section 3.4 that we propose scaling the contribution of the covariance parameters to the forward pass, i.e. the  $\mathbf{S}\mathbf{z}$  term, by  $R^{-1/2}$  since each element in the term is a sum over  $R$  random variables, where  $R$  is the rank of  $\mathbf{S}$ . In the more detailed notation of this section, the proposed scaling implies the forward pass in a linear layer is given by

$$[\mathbf{h}_t^{(l)}]_i = [\mathbf{W}_t]_{:,i} \mathbf{h}_t^{(l-1)} = \left( [\boldsymbol{\mu}_t^{(l)}]_I + R^{-1/2} [\mathbf{S}_t^{(l)}]_I \mathbf{z}^{(l)} \right) \mathbf{h}_t^{(l-1)}. \quad (\text{S117})$$

In practice, rather than scaling  $[\mathbf{S}_t^{(l)}]_I \mathbf{z}^{(l)}$  by  $R^{-1/2}$  in the forward pass, we apply Lemma J.1 from Yang et al. [41] to instead scale the initialization by  $R^{-1/2}$  and, in SGD, the learning rate by  $R^{-1}$ . Scaling by the rank allows treating the mean and covariance parameters as if they were weights parameterized by  $\mu\text{P}$  in a non-probabilistic network, inheriting any scaling that has already been derived for that architecture.

From Table 3 of Yang et al. [41], we therefore scale the mean parameters as

$$b^{(l)} = \begin{cases} 0 & l = 1 \\ 1/2 & l \in \{2, \dots, L\} \\ 1 & l = L + 1 \end{cases} \quad \text{and} \quad c^{(l)} = \begin{cases} -1 & l = 1 \\ 0 & l \in \{2, \dots, L\} \\ 1 & l = L + 1. \end{cases} \quad (\text{S118})$$

Assuming  $R^{(l)} = (D^{(l-1)}D^{(l)})^{p^{(l)}}$  as before, where  $p^{(l)} \in [0, 1]$ , we scale the covariance parameters as

$$\tilde{b}^{(l)} = b^{(l)} + \begin{cases} \frac{p^{(l)}}{2} & l = 1 \\ p^{(l)} & l \in \{2, \dots, L\} \\ \frac{p^{(l)}}{2} & l = L + 1 \end{cases} \quad \text{and} \quad \tilde{c}^{(l)} = c^{(l)} + \begin{cases} p^{(l)} & l = 1 \\ 2p^{(l)} & l \in \{2, \dots, L\} \\ p^{(l)} & l = L + 1. \end{cases} \quad (\text{S119})$$

By comparing to Equations S115 and S116, we see the mean and covariance parameters in all but the output layer are initialized as large as possible while still maintaining stability. The output layer parameters scale to zero faster, since, as in  $\mu\text{P}$  for the weights of non-probabilistic networks, we set  $b^{(L+1)}$  to 1 instead of  $1/2$ .

Note that in Section S2.2 we did not consider input and output dimensions that scaled with the width  $D$  for simplicity. For our experiments, we take the exact  $\mu\text{P}$  initialization and learning rate scaling from Yang et al. [41] — which includes, for example, a `1/fan_in` scaling in the input layer — for the means and then make the rank adjustment for the covariance parameters as described above.

We investigate the proposed scaling in Figures S4 and S5. We train two-hidden-layer ( $L = 2$ ) MLPs of hidden sizes 8, 16, 32, and 64 on a single observation  $(x, y) = (1, 1)$  using a squared error loss. We use SGD with a learning rate of 0.05. For the variational networks, we assume a multivariate Gaussian variational family with a full rank covariance.

Figures S3 and S4 show the RMSE of the change in the hidden units from initialization,  $\Delta \mathbf{g}_t^{(l)}(x) = \mathbf{g}_t^{(l)}(x) - \mathbf{g}_0^{(l)}(x)$ , as a function of the hidden size. The RMSE of the hidden units *at* initialization,

$\mathbf{g}_0^{(l)}$  is also shown in blue. Each panel corresponds to a layer of the network, so the first two panels correspond to features  $\mathbf{g}_t^{(1)}(x)$  and  $\mathbf{g}_t^{(2)}(x)$ , respectively, while the third panel corresponds to the output of the network,  $\mathbf{g}_t^{(3)}(x) = f_t(x)$ . The difference between the figures is the parameterization. Figure S3 uses standard parameterization (SP) while Figure S4 uses maximal update parameterization ( $\mu$ P). We observe that (a) the features change more under  $\mu$ P than SP and (b) training is more stable across hidden sizes under  $\mu$ P than SP, especially for smaller networks.

Figures S5 and S6 show the analogous results for a variational network. The top row shows the change in the mean of the hidden units, while the bottom row shows the change in the standard deviation. As in the non-probabilistic case, we observe that (a) both the mean and standard deviation of the features change more under  $\mu$ P than SP and (b) training is more stable across hidden sizes under  $\mu$ P than SP, especially for smaller networks.

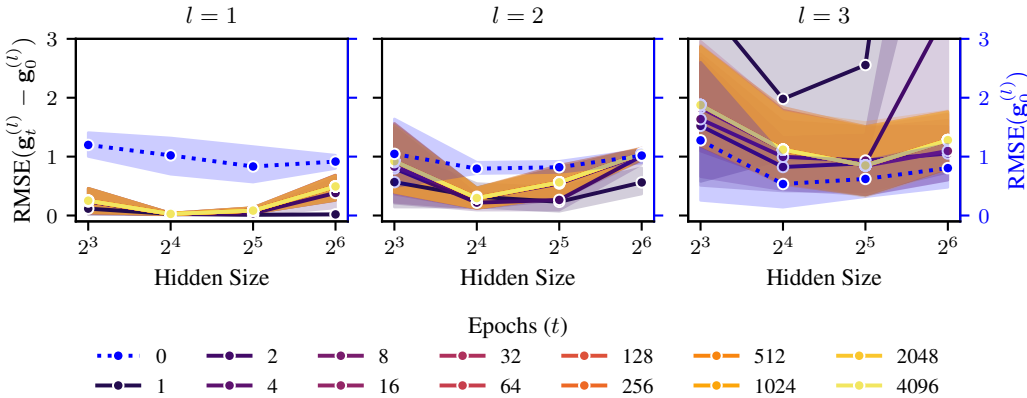


Figure S3: *MLP, Standard Parameterization*. RMSE of the change in the hidden units and, in blue, their initial values. Shaded region represents 95% confidence interval over 5 random initializations. The MLP is trained under SP.

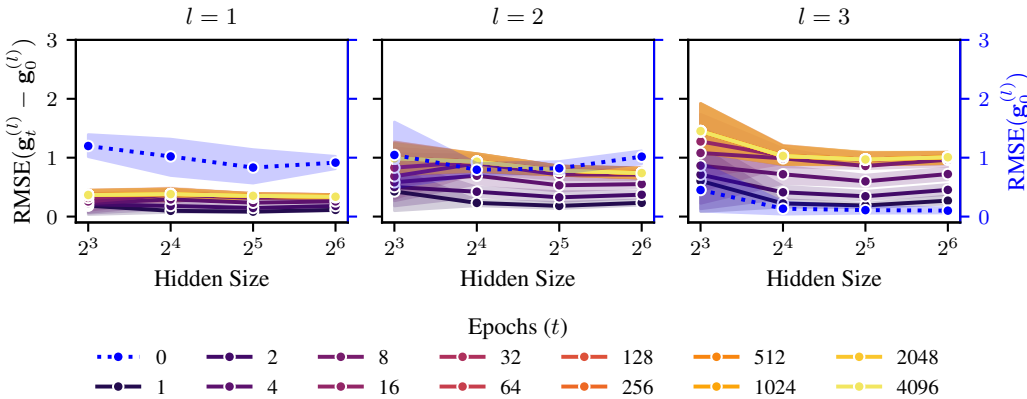


Figure S4: *MLP, Maximal Update Parameterization*. RMSE of the change in the hidden units and, in blue, their initial values. Shaded region represents 95% confidence interval over 5 random initializations. The MLP is trained under  $\mu$ P.

#### S2.4 DETAILS ON HYPERPARAMETER TRANSFER EXPERIMENT

As discussed in Section 3.4 we train two-hidden-layer MLPs of width 128, 256, 512, 1024, and 2048 on CIFAR-10. For comparability to Figure 3 in Tensor Programs V [41] we use the same hyperpa-

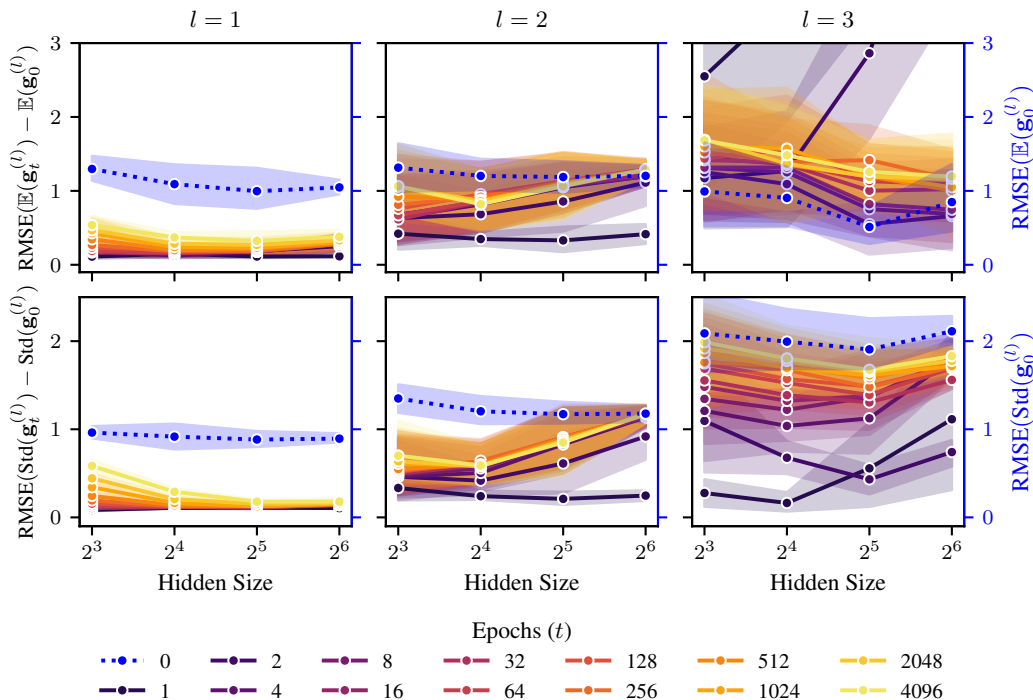


Figure S5: *Variational MLP, Standard Parameterization*. RMSE of the change in the hidden units and, in blue, their initial values. Shaded region represents 95% confidence interval over 5 random initializations. The variational MLP is trained under SP with a full rank covariance in each layer.

rameters but applied to the mean parameters.<sup>5</sup> For the input layer, we scale the mean parameters at initialization by a factor of 16 and in the forward pass by a factor of 1/16. For the output layer, we scale the mean parameters by 0.0 at initialization and by 32.0 in the forward pass. We use 20 epochs, batch size 64, and a grid of global learning rates ranging from  $2^{-8}$  to  $2^0$  with cosine annealing during training. For the grid search results shown in the right panel of Figure 3, we use validation NLL for model selection and then evaluate the relative test error compared to the best performing model for that width across parameterizations and learning rates.

### S3 EXPERIMENTS

This section outlines in more detail the experimental setup, including datasets (Section S3.1.1), metrics (Section S3.1.2), architectures, the training setup and method details (Section S3.3.1). It also contains additional experiments to the ones in the main paper (Sections S3.2, S3.3.2 and S3.3.3).

#### S3.1 SETUP AND DETAILS

In all of our experiments we used the following datasets and metrics.

<sup>5</sup>Specifically, we used the hyperparameters as indicated here: <https://github.com/microsoft/mup/blob/main/examples/MLP/demo.ipynb>

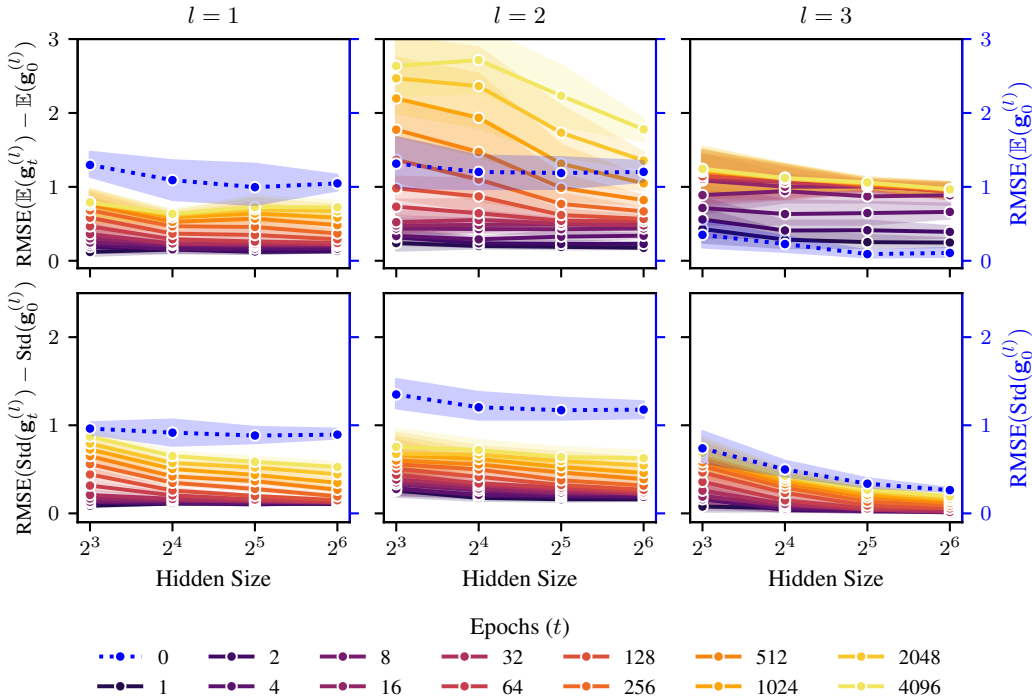


Figure S6: *Variational MLP, Maximal Update Parametrization*. RMSE of the change in the hidden units and, in blue, their initial values. Shaded region represents 95% confidence interval over 5 random initializations. The variational MLP is trained under  $\mu\text{P}$  with a full rank covariance in each layer.

### S3.1.1 DATASETS

Table S1: *Benchmark datasets used in our experiments*. All corrupted datasets are only intended for evaluation and thus only have test sets consisting of 15 different corruptions of the original test set.

Dataset	$N$	$N_{\text{test}}$	$D_{\text{in}}$	$C$	Train / Validation Split
MNIST [70]	60 000	10 000	$28 \times 28$	10	(0.9, 0.1)
CIFAR-10 [80]	50 000	10 000	$3 \times 32 \times 32$	10	(0.9, 0.1)
CIFAR-100 [80]	50 000	10 000	$3 \times 32 \times 32$	100	(0.9, 0.1)
TinyImageNet [81]	100 000	10 000	$3 \times 64 \times 64$	200	(0.9, 0.1)
MNIST-C [72]	-	150 000	$28 \times 28$	10	-
CIFAR-10-C [73]	-	150 000	$3 \times 32 \times 32$	10	-
CIFAR-100-C [73]	-	150 000	$3 \times 32 \times 32$	100	-
TinyImageNet-C [73]	-	150 000	$3 \times 64 \times 64$	200	-

### S3.1.2 METRICS

**Accuracy** The (top-k) accuracy is defined as

$$\text{Accuracy}_k(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{N_{\text{test}}} \sum_{n=1}^{N_{\text{test}}} \mathbb{1}_{(y_n \in \hat{y}_n^{1:k})}. \tag{S120}$$

**Negative Log-Likelihood (NLL)** The (normalized) negative log likelihood for classification is given by

$$\text{NLL}(\mathbf{y}, \hat{\mathbf{y}}) = -\frac{1}{N_{\text{test}}} \sum_{n=1}^{N_{\text{test}}} \log \hat{\mathbf{p}}_{\hat{y}_n}, \tag{S121}$$

where  $\hat{p}_{\hat{y}_n}$  is the probability a model assigns to the predicted class  $\hat{y}_n$ .

**Expected Calibration Error (ECE)** The expected calibration error measures how well a model is calibrated, i.e. how closely the predicted class probability matches the accuracy of the model. Assume the predicted probabilities of the model on the test set are binned into a given binning of the unit interval. Compute the accuracy  $a_j$  and average predicted probability  $\hat{p}_j$  of each bin, then the expected calibration error is given by

$$\text{ECE} = \sum_{j=1}^J b_j |a_j - \hat{p}_j|, \quad (\text{S122})$$

where  $b_j$  is the fraction of datapoints in bin  $j \in \{1, \dots, J\}$ .

### S3.2 TIME AND MEMORY-EFFICIENT TRAINING

To keep the time and memory overhead low during training, we would like to draw as few samples of the parameters as possible to evaluate the training objective  $\bar{\ell}(\theta)$ . Drawing  $M$  parameter samples for the loss increases the time and memory overhead of a forward and backward pass  $M$  times (disregarding parallelism). Therefore it is paramount for efficiency to use as few parameter samples as possible, ideally  $M = 1$ .

When drawing fewer samples from the variational distribution, the variance in the training loss and gradients increases. In practice this means one has to potentially choose a smaller learning rate to still achieve good performance. This is analogous to the previously observed linear relationship  $N_b \propto \eta$  between the optimal batch size  $N_b$  and learning rate  $\eta$  [e.g., 82–84]. Figure S7 shows this relationship between the number of parameter samples used for training and the learning rate on MNIST for a two-hidden layer MLP of width 128.

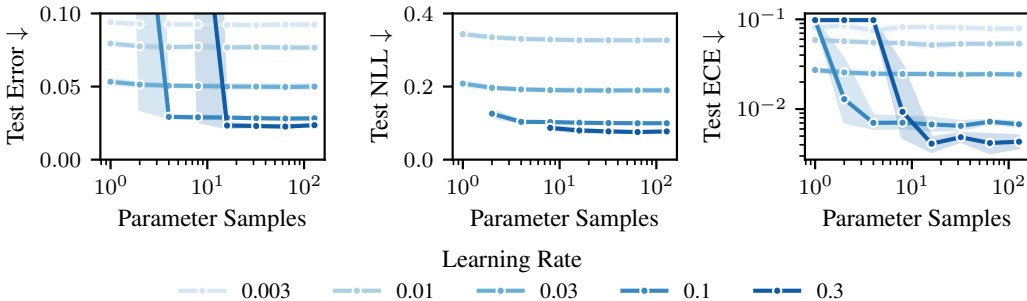


Figure S7: *Generalization versus number of parameter samples.* For a fixed number of epochs and batch size, fewer samples require a smaller learning rate. For a fixed learning rate, generalization performance quickly plateaus with more parameter samples.

As Figure S8 shows, when using momentum, generalization performance tends to increase, but only if either the number of samples is increased, or the learning rate is decreased accordingly. A similar relationship between noise in the objective and the use of momentum has previously been observed by Smith and Le [83], which propose and empirically verify a scaling law for the optimal batch size  $N_b \propto \frac{\eta}{1-\gamma}$  as a function of the momentum parameter  $\gamma > 0$ .

### S3.3 IN- AND OUT-OF-DISTRIBUTION GENERALIZATION

This section recounts details of the methods we benchmark in Section 5, how they are trained and additional experimental results.

#### S3.3.1 ARCHITECTURES, TRAINING, AND METHODS

**Architectures** We use convolutional architectures for all experiments in Section 5. For MNIST, we use a standard LeNet-5 [70] with ReLU activations. For CIFAR-10, CIFAR-100 and TinyIma-

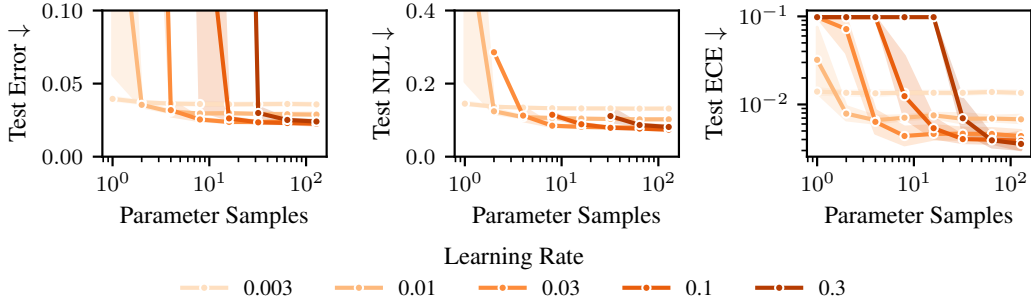


Figure S8: *Generalization versus number of parameter samples when using momentum.* Using momentum improves generalization performance, but when using fewer parameter samples, a smaller learning rate is necessary than for vanilla SGD as predicted by ??.

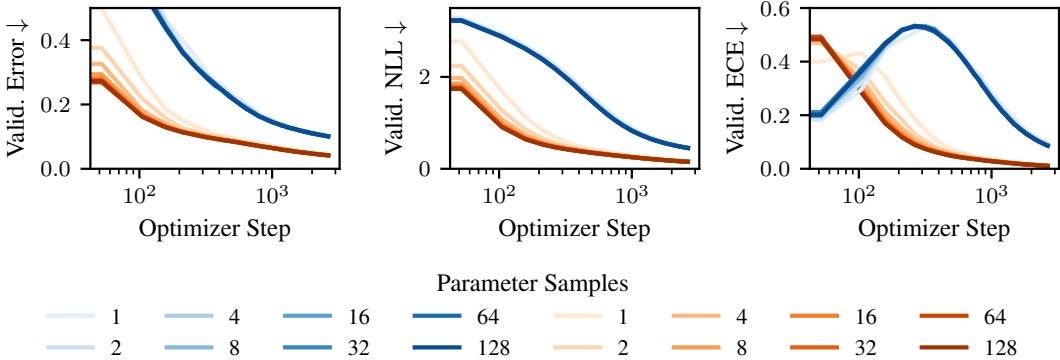


Figure S9: *Validation error during training for different numbers of parameter samples.* The difference in generalization error between different number of parameter samples vanishes with more optimization steps both for SGD (—) and when using momentum (—), if the learning rate is sufficiently small (in this example  $\eta = 0.003$ ).

geNet we use a ResNet-34 [71] where the first layer is a 2D convolution with `kernel_size=3`, `stride=1` and `padding=1` to account for the image resolution of CIFAR and TinyImageNet and the normalization layers are `GroupNorm` layers. We use pretrained weights from ImageNet for all but the first and last layer of the ResNets from `torchvision` [85] and fully finetune all parameters during training.

**Training** We train all models using SGD with momentum ( $\gamma = 0.9$ ) with batch size  $N_b = 128$  and learning rate  $\eta = 0.005$  for 200 epochs. We do not use a learning rate scheduler since we found that neither cosine annealing nor learning rate warm-up improved the results.

**Temperature Scaling** [69] For temperature scaling we optimize the scalar temperature parameter in the last layer on the validation set via the L-BFGS implementation in `torch` with an initial learning rate  $\eta_{rs} = 0.1$ , a maximum number of 100 iterations per optimization step and `history_size=100`.

**Laplace Approximation (Last-Layer, GS + ML)** [26] As recommended by Daxberger et al. [26] we use a post-hoc KFAC last-layer Laplace approximation with a GGN approximation to the Hessian. We tune the hyperparameters post-hoc using type-II maximum likelihood (ML). As an alternative we also do a grid search (GS) for the prior scale, which we found to be somewhat more robust in our experiments. Finally, we compute the predictive using an (extended) probit approximation. Our implementation of the Laplace approximation is a thin wrapper of `laplace` [26] and we use its default hyperparameters throughout.

**Weight-space VI (Mean-field) [30, 31]** For variational inference, we used a mean-field variational family and trained via an ELBO objective with a weighting of the Kullback-Leibler regularization term to the prior. We chose a unit-variance Gaussian prior with mean that was set to the pretrained weights, except for the in- and output layer which had zero mean. We found that using a KL weight and more than a single sample (here  $M = 8$ ) was necessary to achieve competitive performance. The KL weight was chosen to be inversely proportional to the number of parameters of the model, for which we observed better performance than a KL weight that was independent of the architecture. At test time we compute the predictive by averaging logits using 32 samples.

**Implicit Bias VI [ours]** For all architectures in Section 5 we use a Gaussian in- and output layer with a low-rank covariance ( $R = 10, 20$ ). We train with a single parameter sample  $M = 1$  throughout and do temperature scaling at the end of training on the validation set with the same settings as when just performing temperature scaling. We do temperature scaling in classification due to the specific form of the implicit bias in classification as described in Section S1.3. Since IBVI trains by optimizing a minibatch approximation of the expected negative log-likelihood (an average over log-probabilities with respect to parameter samples), we also average log-probabilities at test-time to compute the predictive distribution over class probabilities. Although we did not see a significant difference between averaging log-probabilities, probabilities or logits. Like for WSVI we use 32 samples at test time.

**SWAG [28]** We used a slightly modified implementation of SWAG based on `torch-uncertainty` and the original implementation by Maddox et al. [28]. The beginning of the averaging cycle set to half the number of total epochs and a cycle length of one, i.e. SWAG updates happen every epoch. For all other hyperparameters we use the default settings.

**Deep Ensembles [29]** We use five ensemble members initialized and trained independently. We compute the predictive by averaging the predicted probabilities of the ensemble members in line with standard practice [29]. We did not see a significant difference in performance between averaging logits or averaging class probabilities.

## S3.3.2 IN-DISTRIBUTION GENERALIZATION AND UNCERTAINTY QUANTIFICATION

The full results from the in-distribution generalization experiment in Section 5 can be found in Figure S10. The same experiment but done in the Maximal Update parametrization is depicted in Figure S11. When finetuning a pretrained model, we found that on some datasets (CIFAR-100, TinyImageNet)  $\mu\text{P}$  resulted in somewhat lower performance, contrary to the results in Section 3.4, where we trained from scratch. This suggests that, when pretraining, there may be a modification to the parametrization that could improve generalization.

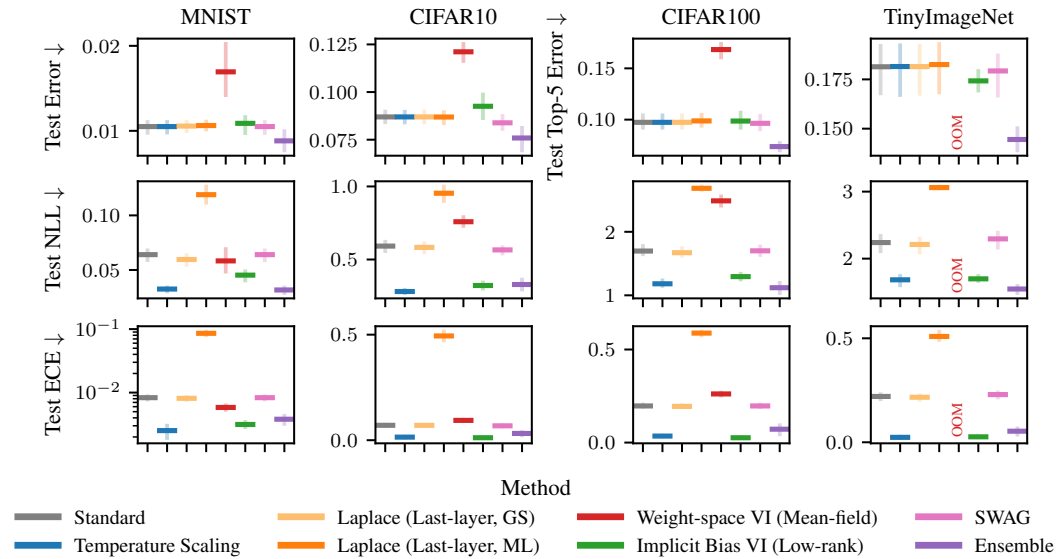


Figure S10: *In-distribution generalization and uncertainty quantification (Standard parametrization).*

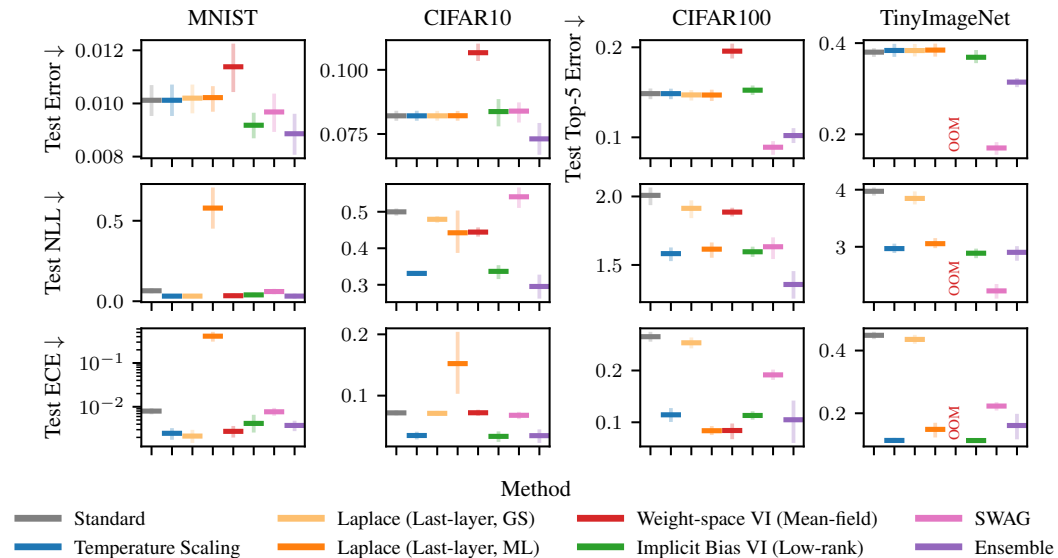


Figure S11: *In-distribution generalization and uncertainty quantification (Maximal Update parametrization).*

### S3.3.3 ROBUSTNESS TO INPUT CORRUPTIONS

Besides the benchmark in Figure S11, we also evaluated the models trained using the Maximal Update parametrization on the corrupted datasets. The results can be found in Figure S12.

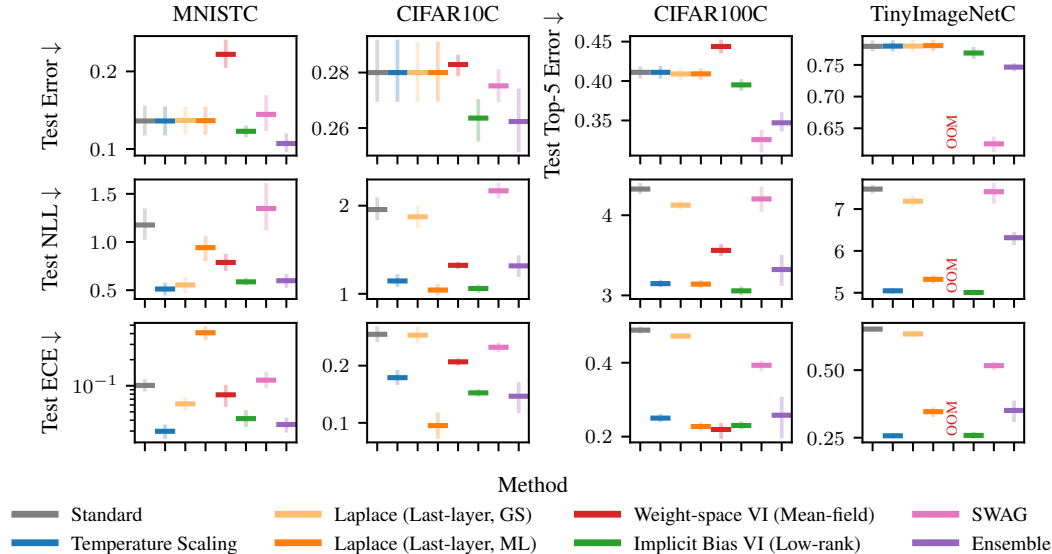


Figure S12: Generalization on robustness benchmark problems (Maximal Update parametrization).

### S3.3.4 COMPARISON TO GENERALIZED VI WITH 2-WASSERSTEIN REGULARIZATION

Theorems 1 and 2 characterize the implicit bias of gradient descent for an overparametrized linear model as a preference for distributions minimizing the expected loss, which are closest in 2-Wasserstein distance to the initialization. Given this characterization, by the KKT conditions there exists a Lagrange multiplier  $\lambda \geq 0$  such that the optimal variational parameters  $\theta_*^{\text{GD}}$  define a stationary point of the following unconstrained optimization objective:

$$\bar{\ell}_r(\theta) = \bar{\ell}(\theta) + \lambda W_2^2(q_\theta, p). \tag{S123}$$

In other words, Implicit Bias VI is equivalent to Generalized VI (GVI) with a 2-Wasserstein regularizer and some regularization strength  $\lambda \geq 0$  for overparametrized linear models.

**Experiment Results** To understand the difference in performance between IBVI and Generalized VI with a 2-Wasserstein regularizer for deep neural networks, we trained models via the GVI objective in Equation (S123) for different regularization strengths  $\lambda \geq 0$  with the same setup as in Section 5. The results on in-distribution test data can be found in Figure S13 and the results for corrupted test data are in Figure S14. Both on in- and out-of-distribution data GVI performs similar or worse than IBVI for all regularization strengths we tested in terms of test error. IBVI and GVI perform roughly similar in terms of uncertainty quantification with GVI only performing better for regularization strengths that harm accuracy.

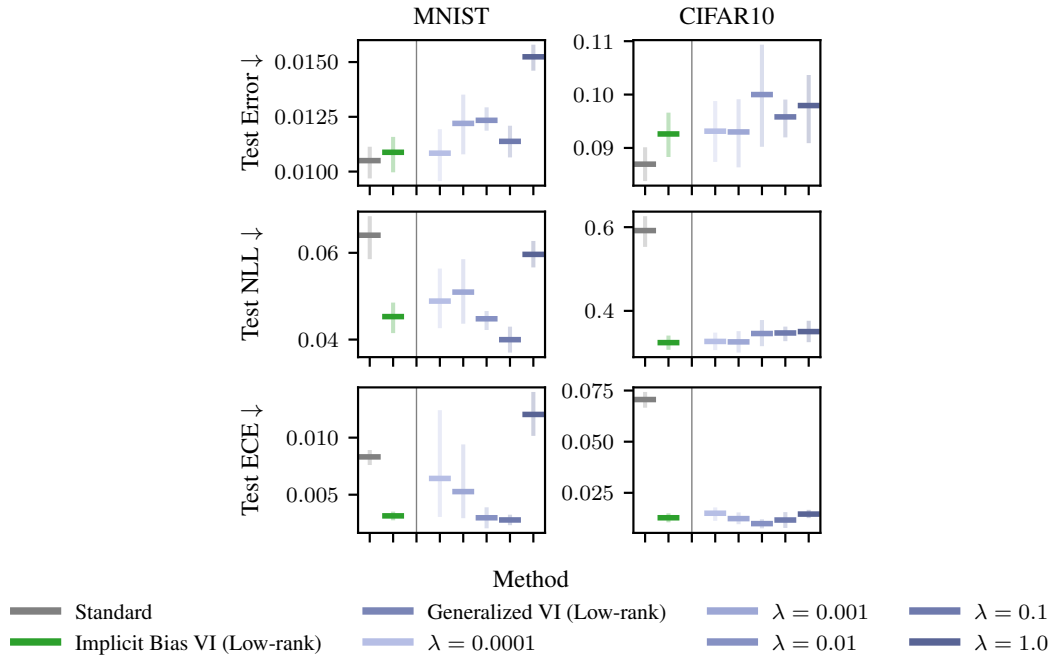


Figure S13: *In-distribution generalization and uncertainty quantification of IBVI and GVI.*

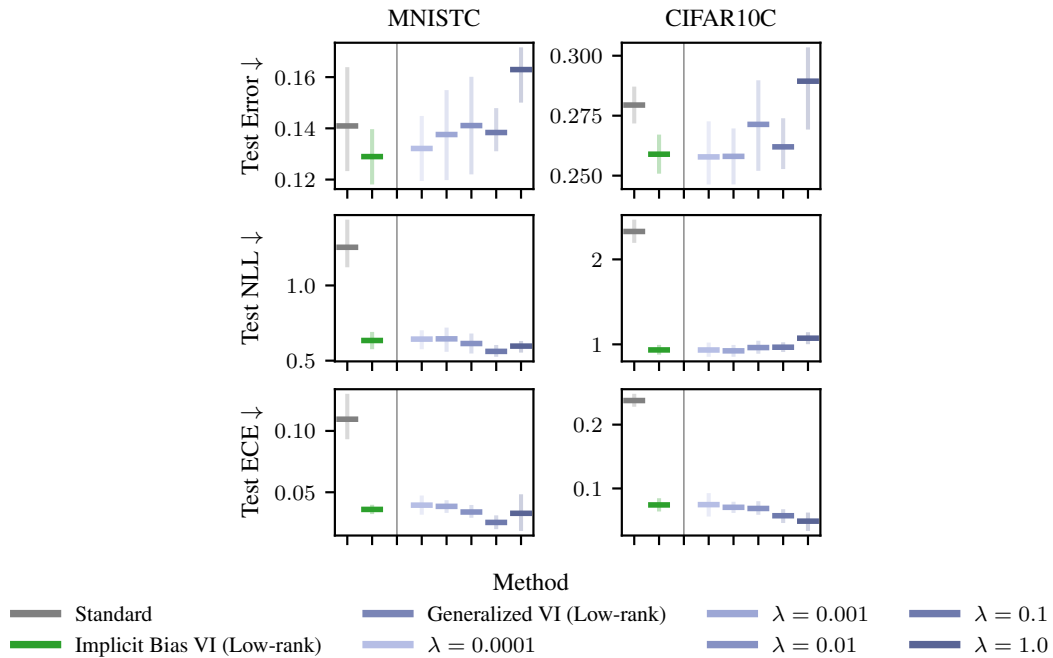


Figure S14: *Out-of-distribution generalization and uncertainty quantification of IBVI and GVI.*