

---

# Epistemic Uncertainty and Observation Noise with the Neural Tangent Kernel

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

Recent work has shown that training wide neural networks with gradient descent is formally equivalent to computing the mean of the posterior distribution in a Gaussian Process (GP) with the Neural Tangent Kernel (NTK) as the prior covariance and zero aleatoric noise [12]. In this paper, we extend this framework in two ways. First, we show how to deal with non-zero aleatoric noise. Second, we derive an estimator for the posterior covariance, giving us a handle on epistemic uncertainty. Our proposed approach integrates seamlessly with standard training pipelines, as it involves training a small number of additional predictors using gradient descent on a mean squared error loss. We demonstrate the proof-of-concept of our method through empirical evaluation on synthetic regression.

## 1 Introduction

Jacot et al. have studied the training of wide neural networks, showing that gradient descent on a standard loss is, in the limit of many iterations, formally equivalent to computing the posterior mean of a Gaussian Process (GP), with the prior covariance specified by the Neural Tangent Kernel (NTK) and with zero aleatoric noise. Crucially, this insight allows us to study complex behaviours of wide networks using Bayesian nonparametrics, which are much better understood.

We extend this analysis by asking two research questions. First, we ask if a similar equivalence exists in cases where we want to do inference for arbitrary values of aleatoric noise. This is crucial in many real-world settings, where measurement accuracy or other data-gathering errors mean that the information in our dataset is only approximate. Second, we ask if it is possible to obtain an estimate of the posterior covariance, not just the mean. Since the posterior covariance measures the epistemic uncertainty about predictions of a model, it is crucial for problems that involve out-of-distribution detection or training with bandit-style feedback.

We answer both of these research questions in the affirmative. Our posterior mean estimator takes the aleatoric noise into account by adding a simple squared norm penalty on the deviation of the network parameters from their initial values, shedding light on regularization in deep learning. Our covariance estimator can be understood as an alternative to existing methods of epistemic uncertainty estimation, such as dropout [7, 20], the Laplace approximation [6, 19], epistemic neural networks [18], deep ensembles [21, 14] and Bayesian Neural Networks [3, 13]. Unlike these approaches, our method has the advantage that it can approximate the NTK-GP posterior arbitrarily well.

**Contributions** We derive estimators for the posterior mean and covariance of an NTK-GP with non-zero aleatoric noise, computable using gradient descent on a standard loss. We evaluate our results empirically on a toy regression problem.

## 2 Preliminaries

**Gaussian Processes** Gaussian Processes (GPs) are a popular non-parametric approach for modeling distributions over functions [22]. Given a dataset of input-output pairs  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ , a GP represents uncertainty about function values by assuming they are jointly Gaussian with a covariance structure

defined by a kernel function  $k(\mathbf{x}, \mathbf{x}')$ . The GP prior is specified as  $f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$ , where  $m(\mathbf{x})$  is the mean function and  $k(\mathbf{x}, \mathbf{x}')$  is the kernel. Assuming  $y_i \sim \mathcal{N}(f(\mathbf{x}), \sigma^2)$  and given new test points  $\mathbf{x}'$ , the posterior mean and covariance are given by:

$$\boldsymbol{\mu}_p(\mathbf{x}') = m(\mathbf{x}') + \mathbf{K}(\mathbf{x}', \mathbf{x})^\top (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1} (\mathbf{y} - m(\mathbf{x})), \quad (1)$$

$$\boldsymbol{\Sigma}_p(\mathbf{x}') = \mathbf{K}(\mathbf{x}', \mathbf{x}') - \mathbf{K}(\mathbf{x}', \mathbf{x})^\top (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1} \mathbf{K}(\mathbf{x}, \mathbf{x}'), \quad (2)$$

where  $\mathbf{K}(\mathbf{x}, \mathbf{x})$  is the covariance matrix computed over the training inputs,  $\mathbf{K}(\mathbf{x}', \mathbf{x})$  is the covariance matrix between the test and training points, and  $\sigma^2$  represents the aleatoric (or observation) noise.

**Neural Tangent Kernel.** The Neural Tangent Kernel (NTK) characterizes the evolution of wide neural network predictions as a linear model in function space. Given a neural network function  $f(\mathbf{x}; \theta)$  parameterized by  $\theta$ , the NTK is defined through the Jacobian  $J(\mathbf{x}) \in \mathbb{R}^{N \times p}$ , where  $J(\mathbf{x}) = \frac{\partial f(\mathbf{x}; \theta)}{\partial \theta}$ ,  $N$  is the number of data points and  $p$  is the number of parameters. The NTK at two sets of inputs  $\mathbf{x}$  and  $\mathbf{x}'$  is given by:

$$\mathbf{K}(\mathbf{x}, \mathbf{x}') = J(\mathbf{x})J(\mathbf{x}')^\top. \quad (3)$$

Interestingly, as shown by [12] the NTK converges to a deterministic kernel and remains constant during training in the infinite-width limit. We call a GP with the kernel (3) the NTK GP.

### 3 Method

We now describe our proposed process of doing inference in the NTK-GP. Our procedure for estimating the posterior mean is given in Algorithm 1, while the procedure for the covariance is given in Algorithm 2. Note that our process is scaleable because both algorithms only use gradient descent, rather than relying on a matrix inverse in equations (1) and (2). While Algorithm 2 relies on the computation of the partial SVD of the Jacobian, we stress that efficient ways of doing so exist and do not require ever storing the full Jacobian. We defer the details of the partial SVD to Appendix E. We describe the theory that justifies our posterior computation in sections 3.1 and 3.2. We defer the discussion of literature to Appendix A.

---

#### Algorithm 1 Algorithm for Computing the Posterior Mean in the NTK-GP

---

```

procedure TRAIN-POSTERIOR-MEAN( $x_i, y_i, \theta_0$ )
   $\hat{y}_i \leftarrow y_i + f(x_i; \theta_0)$  ▷ Shift the targets to get zero prior mean (Lemma 3.2).
   $L \leftarrow \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - f(x_i; \theta))^2 + \beta_N \|\theta - \theta_0\|_2^2$  ▷ Equation (4)
  minimize  $L$  with gradient descent wrt.  $\theta$  until convergence to  $\theta^*$ 
  return  $\theta^*$  ▷ Return the trained weights.
end procedure

procedure QUERY-POSTERIOR-MEAN( $x'_j, \theta^*, \theta^0$ ) ▷  $j = 1, \dots, J$ 
  return  $f(x'_1; \theta^*) - f(x'_1; \theta^0), \dots, f(x'_J; \theta^*) - f(x'_1; \theta^0)$ 
end procedure

```

---

#### 3.1 Aleatoric Noise

**Gradient Descent Converges to the NTK-GP Posterior Mean** We build on the work of [12] by focusing on the computation of the mean posterior in the presence of **non-zero aleatoric noise**. We show that optimizing a regularized mean squared error loss in a neural network is equivalent to computing the mean posterior of an NTK-GP with non-zero aleatoric noise. In the following Lemma, we prove that for a sufficiently long training process, the predictions of the trained network converge to those of an NTK-GP with aleatoric noise characterized by  $\sigma^2 = N\beta_N$ . This is a similar result to [11], but from a Bayesian perspective rather than a frequentist generalization bound. Furthermore, our proof (see Appendix B) focuses on explicitly solving the gradient flows for test and training data points in function space.

**Lemma 3.1.** Consider a parametric model  $f(x; \theta)$  where  $x \in \mathcal{X} \subset \mathbb{R}^N$  and  $\theta \in \mathbb{R}^p$ , initialized under some assumptions with parameters  $\theta_0$ . Minimizing the regularized mean squared error loss

71 with respect to  $\theta$  to find the optimal set of parameters  $\theta^*$  over a dataset  $(\mathbf{x}, \mathbf{y})$  of size  $N$ , and with  
 72 sufficient training time ( $t \rightarrow \infty$ ):

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^p} \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i; \theta))^2 + \beta_N \|\theta - \theta_0\|_2^2, \quad (4)$$

73 is equivalent to computing the mean posterior of a Gaussian process with non-zero aleatoric noise,  
 74  $\sigma^2 = N\beta_N$ , and the NTK as its kernel:

$$f(\mathbf{x}'; \theta_\infty) = f(\mathbf{x}'; \theta_0) + \mathbf{K}(\mathbf{x}', \mathbf{x})(\mathbf{K}(\mathbf{x}, \mathbf{x}) + N\beta_N \mathbf{I})^{-1}(\mathbf{y} - f(\mathbf{x}; \theta_0)). \quad (5)$$

75 **Zero Prior Mean** In many practical scenarios, it is desirable to start with zero prior mean rather  
 76 than with a prior mean that corresponds to random network initialization. To accommodate this, we  
 77 introduce a simple yet effective transformation of the data and the network outputs, to be applied  
 78 together with 3.1. We summarize it into the following lemma (see Appendix B for proof):

79 **Lemma 3.2.** Consider the computational process derived in Lemma 3.1. Define shifted labels  $\tilde{\mathbf{y}}$  and  
 80 predictions  $\tilde{f}(\mathbf{x}; \theta_\infty)$  as follows::

$$\tilde{\mathbf{y}} = \mathbf{y} + f(\mathbf{x}; \theta_0), \quad \tilde{f}(\mathbf{x}; \theta_\infty) = f(\mathbf{x}; \theta_\infty) - f(\mathbf{x}'; \theta_0).$$

81 Using these definitions, the posterior mean of a zero-mean Gaussian process can be computed as:

$$\tilde{f}(\mathbf{x}', \theta_\infty) = \mathbf{K}(\mathbf{x}', \mathbf{x})(\mathbf{K}(\mathbf{x}, \mathbf{x}) + N\beta_N \mathbf{I})^{-1} \tilde{\mathbf{y}}. \quad (6)$$

---

**Algorithm 2** Algorithm for Computing the Posterior Covariance in the NTK-GP

---

```

procedure TRAIN-POSTERIOR-COVARIANCE( $x_i, K, \theta_0$ ) ▷  $K$  is the number of predictors
   $U, \Sigma \leftarrow \text{PARTIAL-SVD}(J_{\theta_0}(\mathbf{x}), K)$  ▷ Partial SVD of the Jacobian - see appendix E.
  for  $i = 1, \dots, K$  do
     $\theta_i^* \leftarrow \text{TRAIN-POSTERIOR-MEAN}(x_i, U_i)$  ▷  $U_i$  is the  $i$ -th column of  $U$ .
  end for
  for  $i = 1, \dots, K'$  do ▷ Setting  $K' = 0$  often works well (see Appendix D).
     $\theta_i'^* \leftarrow \text{TRAIN-POSTERIOR-MEAN}(x_i, \epsilon_i)$  ▷  $\epsilon_i \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ 
  end for
  return  $\Sigma, \theta_1^*, \dots, \theta_K^*, \theta_1'^*, \dots, \theta_{K'}'^*$ 
end procedure

procedure QUERY-POSTERIOR-COVARIANCE( $x_j', \Sigma, \theta_i^*, \theta_i'^*, \theta_0$ ) ▷  $j = 1, \dots, J$ 
   $P \leftarrow \begin{bmatrix} f(x_1'; \theta_1^*) - f(x_1'; \theta_0) & \dots & f(x_1'; \theta_K^*) - f(x_1'; \theta_0) \\ \vdots & \ddots & \vdots \\ f(x_j'; \theta_1^*) - f(x_j'; \theta_0) & \dots & f(x_j'; \theta_K^*) - f(x_j'; \theta_0) \end{bmatrix}, P' \leftarrow \begin{bmatrix} f(x_1'; \theta_1'^*) - f(x_1'; \theta_0) & \dots & f(x_1'; \theta_{K'}'^*) - f(x_1'; \theta_0) \\ \vdots & \ddots & \vdots \\ f(x_j'; \theta_1'^*) - f(x_j'; \theta_0) & \dots & f(x_j'; \theta_{K'}'^*) - f(x_j'; \theta_0) \end{bmatrix}$ 
  return  $J(\mathbf{x}')J(\mathbf{x}')^\top - P\Sigma^2P^\top - P'(P')^\top / K'$  ▷ The last term vanishes for  $K' = 0$ 
end procedure

```

---

### 82 3.2 Estimating the Covariance

83 We now justify Algorithm 2 for estimating the posterior covariance. The main observation that allows  
 84 us to derive our estimator comes from examining the term  $\mathbf{K}(\mathbf{x}', \mathbf{x})^\top (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1} \mathbf{K}(\mathbf{x}', \mathbf{x})$  in  
 85 the posterior covariance formula (2). This is summarized in the following Proposition.

86 **Proposition 3.1.** Diagonalize  $\mathbf{K}(\mathbf{x}, \mathbf{x})$  so that  $\mathbf{K}(\mathbf{x}, \mathbf{x}) = U\Lambda U^\top$ . We have

$$\mathbf{K}(\mathbf{x}', \mathbf{x})^\top (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1} \mathbf{K}(\mathbf{x}', \mathbf{x}) = (MU)\Lambda(MU)^\top + \sigma^2 MM^\top.$$

87 Here,  $M = \mathbf{K}(\mathbf{x}', \mathbf{x})^\top (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1}$ .

88 *Proof.* We can rewrite it as:

$$\mathbf{K}(\mathbf{x}', \mathbf{x})^\top (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1} \mathbf{K}(\mathbf{x}', \mathbf{x}) = \underbrace{\mathbf{K}(\mathbf{x}', \mathbf{x})^\top (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1}}_{M^\top} (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I}) \underbrace{(\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1} \mathbf{K}(\mathbf{x}', \mathbf{x})}_{M}$$

89 Denoting the term  $\mathbf{K}(\mathbf{x}', \mathbf{x})^\top (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1}$  with  $M$ , this can be written as:

$$\mathbf{K}(\mathbf{x}', \mathbf{x})^\top (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1} \mathbf{K}(\mathbf{x}', \mathbf{x}) = (MU)\Lambda(MU)^\top + \sigma^2 MM^\top.$$

90

□

91 The proposition is useful because the matrix  $M$  appears in equation (1). Hence the matrix multi-  
 92 plication  $MU$  is equivalent to estimating the posterior mean using algorithm 1 where targets are  
 93 given by the columns of the matrix  $U$ . Hence the term  $(MU)\Lambda(MU)^\top$  can be computed by gradient  
 94 descent. In order to derive a complete estimator of the covariance, we still need to deal with the term  
 95  $\sigma^2 MM^\top$ . We can either estimate this term by fitting random targets (which corresponds to setting  
 96  $K' > 0$  in algorithm 2) or accept an upper bound on the covariance, setting  $K' = 0$ . We describe this  
 97 in detail in Appendix D.

## 98 4 Experiment

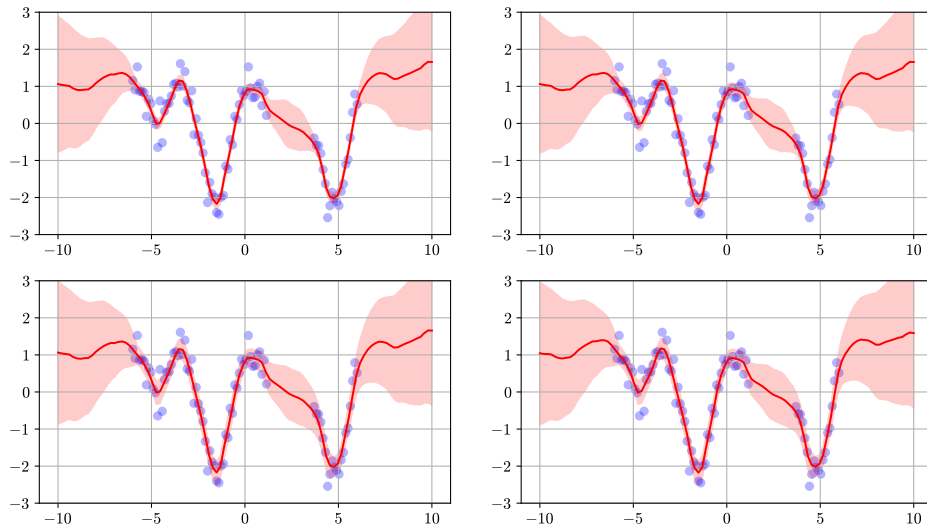


Figure 1: The NTK-GP posterior and its approximations: (top-left) Analytic Posterior, (top-right) Analytic upper bound on posterior (all eigenvectors), (bottom-left) Analytic upper bound on posterior (5 eigenvectors), (bottom-right) Posterior obtained with gradient descent ( $K = 5$  predictors,  $K' = 0$ ).

99 We applied the method to a toy regression problem shown in Figure 1. The problem is a standard  
 100 non-linear 1d regression task which requires both interpolation and extrapolation. The top-left figure  
 101 was obtained by computing the kernel of the NTK-GP using formula (3) and computing the posterior  
 102 mean and covariance using equations (1) and (2). The top-right figure was obtained by analytically  
 103 computing the upper bound defined in appendix D. The bottom-left figure was obtained by taking  
 104 the first 5 eigenvectors of the kernel. Finally, the bottom-right figure was obtained by fitting a mean  
 105 prediction network and 5 predictor networks using the gradient-descent method described in algorithm  
 106 2. The similarity of the figures shows that the method works. Details of network architecture are  
 107 deferred to Appendix C.

## 108 5 Conclusions

109 This paper introduces a method for computing the posterior mean and covariance of NTK-Gaussian  
 110 Processes with non-zero aleatoric noise. Our approach integrates seamlessly with standard training  
 111 procedures using gradient descent, providing a practical tool for uncertainty estimation in contexts  
 112 such as Bayesian optimization. The method has been validated empirically on a toy task, demon-  
 113 strating its effectiveness in capturing uncertainty while maintaining computational efficiency. This  
 114 work opens up opportunities for further research in applying NTK-GP frameworks to more complex  
 115 scenarios and datasets.

## References

- [1] R. Bhatia. *Positive Definite Matrices*. Princeton Series in Applied Mathematics. Princeton University Press, 2015. ISBN: 9780691168258. URL: <https://books.google.co.uk/books?id=Y22YDwAAQBAJ>.
- [2] Avrim Blum, John Hopcroft, and Ravindran Kannan. *Foundations of data science*. Cambridge University Press, 2020.
- [3] Charles Blundell et al. “Weight uncertainty in neural network”. In: *International conference on machine learning*. PMLR. 2015, pp. 1613–1622.
- [4] Yuri Burda et al. “Exploration by random network distillation”. In: *arXiv preprint arXiv:1810.12894* (2018).
- [5] Kamil Ciosek et al. “Conservative uncertainty estimation by fitting prior networks”. In: *International Conference on Learning Representations*. 2019.
- [6] Erik Daxberger et al. “Laplace redux-effortless bayesian deep learning”. In: *Advances in Neural Information Processing Systems* 34 (2021), pp. 20089–20103.
- [7] Yarin Gal and Zoubin Ghahramani. “Dropout as a bayesian approximation: Representing model uncertainty in deep learning”. In: *international conference on machine learning*. PMLR. 2016, pp. 1050–1059.
- [8] Eugene Golikov, Eduard Pokonechnyy, and Vladimir Korviakov. *Neural Tangent Kernel: A Survey*. 2022. arXiv: 2208.13614 [cs.LG]. URL: <https://arxiv.org/abs/2208.13614>.
- [9] Nathan Halko, Per-Gunnar Martinsson, and Joel A Tropp. “Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions”. In: *SIAM review* 53.2 (2011), pp. 217–288.
- [10] Bobby He, Balaji Lakshminarayanan, and Yee Whye Teh. *Bayesian Deep Ensembles via the Neural Tangent Kernel*. 2020. arXiv: 2007.05864 [stat.ML]. URL: <https://arxiv.org/abs/2007.05864>.
- [11] Wei Hu, Zhiyuan Li, and Dingli Yu. *Simple and Effective Regularization Methods for Training on Noisily Labeled Data with Generalization Guarantee*. 2020. arXiv: 1905.11368 [cs.LG]. URL: <https://arxiv.org/abs/1905.11368>.
- [12] Arthur Jacot, Franck Gabriel, and Clément Hongler. “Neural tangent kernel: Convergence and generalization in neural networks”. In: *Advances in neural information processing systems* 31 (2018).
- [13] Diederik P Kingma. “Auto-encoding variational bayes”. In: *arXiv preprint arXiv:1312.6114* (2013).
- [14] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. “Simple and scalable predictive uncertainty estimation using deep ensembles”. In: *Advances in neural information processing systems* 30 (2017).
- [15] Jaehoon Lee et al. “Wide neural networks of any depth evolve as linear models under gradient descent”. In: *Advances in neural information processing systems* 32 (2019).
- [16] Radford M Neal. *Bayesian learning for neural networks*. Vol. 118. Springer Science & Business Media, 2012.
- [17] Ian Osband et al. “Deep exploration via randomized value functions”. In: *Journal of Machine Learning Research* 20.124 (2019), pp. 1–62.
- [18] Ian Osband et al. “Epistemic neural networks”. In: *Advances in Neural Information Processing Systems* 36 (2023), pp. 2795–2823.
- [19] Hippolyt Ritter, Aleksandar Botev, and David Barber. “A scalable laplace approximation for neural networks”. In: *6th international conference on learning representations, ICLR 2018-conference track proceedings*. Vol. 6. International Conference on Representation Learning. 2018.
- [20] Nitish Srivastava et al. “Dropout: a simple way to prevent neural networks from overfitting”. In: *The journal of machine learning research* 15.1 (2014), pp. 1929–1958.
- [21] Robert J Tibshirani and Bradley Efron. “An introduction to the bootstrap”. In: *Monographs on statistics and applied probability* 57.1 (1993), pp. 1–436.
- [22] Christopher KI Williams and Carl Edward Rasmussen. *Gaussian processes for machine learning*. Vol. 2. 3. MIT press Cambridge, MA, 2006.

## 170 A Related Work

171 **Neural Tangent Kernel** The definition of the Neural Tangent Kernel (3), the proof of the fact that  
 172 it stays constant during training and doesn't depend on initialization as well as the link to Gaussian  
 173 Processes with no aleatoric noise are all due to the seminal paper [12]. The work of Lee et al. builds on  
 174 that, showing that wide neural networks can be understood as linear models for purposes of studying  
 175 their training dynamics, a fact we crucially rely on in the proof of our Lemma 3.1. Hu et al. describe  
 176 a regularizer for networks trained in the NTK regime which leads to the same optimization problem  
 177 used in our Lemma 3.1. The difference lies in the fact that we rely on the Bayesian interpretation of  
 178 the network obtained at the end of training, while they focus on a frequentist generalization bound.

179 **Predictor Networks** Prior work [17, 4, 5] has considered epistemic uncertainty estimation by  
 180 fitting functions generated using a process that includes some kind of randomness. Burda et al.  
 181 have applied a similar idea to reinforcement learning, obtaining exceptional results on Montezuma's  
 182 Revenge, a problem where it is known that exploration very is hard. Ciosek et al. provided a link to  
 183 Gaussian Processes, but did not leverage the NTK, instead describing an upper bound on a posterior  
 184 relative to the kernel [16] where sampling corresponds to sampling from the network initialization.  
 185 Osband et al. proposed<sup>1</sup> a way of sampling from a Bayesian linear regression posterior by solving  
 186 an optimization problem with a similar structure to ours. However, this approach is different in two  
 187 crucial ways. First, Osband et al. is interested in obtaining samples from the posterior, while we are  
 188 interested in computing the posterior moments. Second, the sampling process in the paper by Osband  
 189 et al. depends on the true regression targets in a way that our posterior covariance estimate does not.  
 190 Also, our method is framed differently, as we intend it to be used in the context of the NTK regime,  
 191 while Osband et al. discusses vanilla linear regression.

192 **Epistemic Uncertainty** Our method of fitting the posterior covariance about network outputs can  
 193 be thought of as quantifying epistemic uncertainty. There are several established methods in this  
 194 space. Dropout [7, 20], works by randomly disabling neurons in a network and has a Bayesian  
 195 interpretation. The Laplace approximation [6, 19] works by replacing an arbitrary likelihood with  
 196 a Gaussian one. Epistemic neural networks [18] are based on the idea of using an additional input  
 197 (the epistemic index) when training the network. Deep ensembles [21, 14] work by training several  
 198 copies of a network with different initializations and sometimes training sets that are only partially  
 199 overlapping. While classic deep ensembles do not have a Bayesian interpretation, He et al. have  
 200 recently proposed a modification that approximates the posterior in the NTK-GP. Bayesian Neural  
 201 Networks [3, 13] attempt to apply Bayes rule in the space of neural network parameters, applying  
 202 various approximations. A full survey of methods of epistemic uncertainty estimation is beyond the  
 203 scope of this paper.

## 204 B Proofs

205 **Lemma 3.1.** Consider a parametric model  $f(x; \theta)$  where  $x \in \mathcal{X} \subset \mathbb{R}^N$  and  $\theta \in \mathbb{R}^p$ , initialized  
 206 under some assumptions with parameters  $\theta_0$ . Minimizing the regularized mean squared error loss  
 207 with respect to  $\theta$  to find the optimal set of parameters  $\theta^*$  over a dataset  $(\mathbf{x}, \mathbf{y})$  of size  $N$ , and with  
 208 sufficient training time ( $t \rightarrow \infty$ ):

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^p} \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i; \theta))^2 + \beta_N \|\theta - \theta_0\|_2^2, \quad (4)$$

209 is equivalent to computing the mean posterior of a Gaussian process with non-zero aleatoric noise,  
 210  $\sigma^2 = N\beta_N$ , and the NTK as its kernel:

$$f(\mathbf{x}'; \theta_\infty) = f(\mathbf{x}'; \theta_0) + \mathbf{K}(\mathbf{x}', \mathbf{x})(\mathbf{K}(\mathbf{x}, \mathbf{x}) + N\beta_N \mathbf{I})^{-1}(\mathbf{y} - f(\mathbf{x}; \theta_0)). \quad (5)$$

211 *Proof.* Consider a regression problem with the following regularized empirical loss:

$$\mathcal{L}(\mathbf{y}, f(\mathbf{x}; \theta)) = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i; \theta))^2 + \beta_N \|\theta - \theta_0\|_2^2. \quad (7)$$

---

<sup>1</sup>See Section 5.3.1 in the paper by Osband et al.

212 Let us use  $\theta_t$  to represent the parameters of the network evolving in time  $t$  and let  $\alpha$  be the learning  
 213 rate. Assuming we train the network via continuous-time gradient flow, then the evolution of the  
 214 parameters  $\theta_t$  can be expressed as:

$$\frac{d\theta_t}{dt} = -\alpha \left[ \frac{2}{N} \nabla_{\theta} f(\mathbf{x}; \theta_t) (f(\mathbf{x}; \theta_t) - \mathbf{y}) + 2\beta_N (\theta_t - \theta_0) \right]. \quad (8)$$

215 Assuming that our neural network architecture operates in a sufficiently wide regime [15], where the  
 216 first-order approximation remains valid throughout gradient descent, we obtain:

$$f(\mathbf{x}'; \theta_t) = f(\mathbf{x}'; \theta_0) + J_t(\mathbf{x}')(\theta_t - \theta_0) \rightarrow \nabla_{\theta} f(\mathbf{x}'; \theta_t)^{\top} = J_t(\mathbf{x}'). \quad (9)$$

217 The dynamics of the neural network on the training data:

$$\begin{aligned} \frac{df(\mathbf{x}; \theta_t)}{dt} &= J_t(\mathbf{x}) \frac{d\theta_t}{dt} = -\frac{2\alpha}{N} J_t(\mathbf{x}) [J_t(\mathbf{x})^{\top} (f(\mathbf{x}; \theta_t) - \mathbf{y}) + \beta_N (\theta_t - \theta_0)] \\ &= -\frac{2\alpha}{N} (\mathbf{K}(\mathbf{x}, \mathbf{x}) (f(\mathbf{x}; \theta_t) - \mathbf{y}) + \beta_N J_t(\mathbf{x}) (\theta_t - \theta_0)) \\ &= -\frac{2\alpha}{N} (\mathbf{K}(\mathbf{x}, \mathbf{x}) (f(\mathbf{x}; \theta_t) - \mathbf{y}) + \beta_N (f(\mathbf{x}; \theta_t) - f(\mathbf{x}; \theta_0))) \\ &= -\frac{2\alpha}{N} (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N \mathbf{I}) f(\mathbf{x}; \theta_t) + \frac{2\alpha}{N} (\mathbf{K}(\mathbf{x}, \mathbf{x}) \mathbf{y} + \beta_N f(\mathbf{x}; \theta_0)) \end{aligned}$$

218 This is a linear ODE, we can solve this:

$$\begin{aligned} f(\mathbf{x}; \theta_t) &= \exp \left( -\frac{2\alpha}{N} t (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N \mathbf{I}) \right) f(\mathbf{x}; \theta_0) \\ &\quad - \frac{N}{2\alpha} (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N \mathbf{I})^{-1} \left[ \exp \left( -\frac{2\alpha}{N} t (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N \mathbf{I}) \right) - \mathbf{I} \right] \\ &\quad \times \frac{2\alpha}{N} (\mathbf{K}(\mathbf{x}, \mathbf{x}) \mathbf{y} + \beta_N f(\mathbf{x}; \theta_0)) \end{aligned}$$

219 Using  $A^{-1}e^A = e^A A^{-1}$ , and writing  $\mathbf{K}(\mathbf{x}, \mathbf{x}) \mathbf{y} + \beta_N f(\mathbf{x}; \theta_0) = (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N \mathbf{I}) f(\mathbf{x}; \theta_0) +$   
 220  $\mathbf{K}(\mathbf{x}, \mathbf{x}) (\mathbf{y} - f(\mathbf{x}; \theta_0))$ , we get:

$$\begin{aligned} f(\mathbf{x}; \theta_t) &= \exp \left( -\frac{2\alpha}{N} t (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N \mathbf{I}) \right) f(\mathbf{x}; \theta_0) \\ &\quad + \left[ \mathbf{I} - \exp \left( -\frac{2\alpha}{N} t (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N \mathbf{I}) \right) \right] (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N \mathbf{I})^{-1} (\mathbf{K}(\mathbf{x}, \mathbf{x}) \mathbf{y} + \beta_N f(\mathbf{x}; \theta_0)) \\ &= \exp \left( -\frac{2\alpha}{N} t (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N \mathbf{I}) \right) f(\mathbf{x}; \theta_0) + \left[ \mathbf{I} - \exp \left( -\frac{2\alpha}{N} t (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N \mathbf{I}) \right) \right] f(\mathbf{x}; \theta_0) \\ &\quad + \left[ \mathbf{I} - \exp \left( -\frac{2\alpha}{N} t (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N \mathbf{I}) \right) \right] (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N \mathbf{I})^{-1} \mathbf{K}(\mathbf{x}, \mathbf{x}) (\mathbf{y} - f(\mathbf{x}; \theta_0)) \\ &= f(\mathbf{x}; \theta_0) + \left[ \mathbf{I} - \exp \left( -\frac{2\alpha}{N} t (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N \mathbf{I}) \right) \right] (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N \mathbf{I})^{-1} \mathbf{K}(\mathbf{x}, \mathbf{x}) (\mathbf{y} - f(\mathbf{x}; \theta_0)). \end{aligned}$$

221 Now, we consider the dynamics for the neural network of an arbitrary set of test points  $\mathbf{x}'$ :

$$\frac{df(\mathbf{x}'; \theta_t)}{dt} = -\frac{2\alpha}{N} \beta_N f(\mathbf{x}'; \theta_t) - \frac{2\alpha}{N} (\mathbf{K}(\mathbf{x}', \mathbf{x}) (f(\mathbf{x}; \theta_t) - \mathbf{y}) - \beta_N f(\mathbf{x}'; \theta_0)). \quad (10)$$

222 This is a linear ODE with a time-dependent inhomogeneous term, we can solve it as follows:

$$\begin{aligned}
f(x', \theta_t) &= e^{-\frac{2\alpha}{N}\beta_N t} f(x', \theta_0) - \frac{2\alpha}{N} e^{-\frac{2\alpha}{N}\beta_N t} \int_0^t e^{\frac{2\alpha}{N}\beta_N u} (\mathbf{K}(\mathbf{x}', \mathbf{x})(f(x, \theta_u) - y) - \beta_N f(x', \theta_0)) du \\
&= e^{-\frac{2\alpha}{N}\beta_N t} f(x', \theta_0) + \frac{2\alpha}{N} e^{-\frac{2\alpha}{N}\beta_N t} \int_0^t e^{\frac{2\alpha}{N}\beta_N u} du (\mathbf{K}(\mathbf{x}', \mathbf{x})y + \beta_N f(x', \theta_0)) \\
&\quad - \frac{2\alpha}{N} e^{-\frac{2\alpha}{N}\beta_N t} \mathbf{K}(\mathbf{x}', \mathbf{x}) \int_0^t e^{\frac{2\alpha}{N}\beta_N u} f(x, \theta_u) du. \\
&= e^{-\frac{2\alpha}{N}\beta_N t} f(x', \theta_0) + e^{-\frac{2\alpha}{N}\beta_N t} \frac{1}{\beta_N} \left( e^{\frac{2\alpha}{N}\beta_N t} - 1 \right) (\mathbf{K}(\mathbf{x}', \mathbf{x})y + \beta_N f(x', \theta_0)) \\
&\quad - \frac{2\alpha}{N} e^{-\frac{2\alpha}{N}\beta_N t} \mathbf{K}(\mathbf{x}', \mathbf{x}) \int_0^t e^{\frac{2\alpha}{N}\beta_N u} f(x, \theta_u) du \\
&\quad - \frac{2\alpha}{N} e^{-\frac{2\alpha}{N}\beta_N t} \mathbf{K}(\mathbf{x}', \mathbf{x}) \int_0^t e^{\frac{2\alpha}{N}\beta_N u} \left[ I - \exp\left(-\frac{2\alpha}{N}u(\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N I)\right) \right] du \\
&\quad \times (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N I)^{-1} \mathbf{K}(\mathbf{x}, \mathbf{x})(y - f(x, \theta_0)). \\
&= f(x', \theta_0) + \frac{1}{\beta_N} (1 - e^{-\frac{2\alpha}{N}\beta_N t}) \mathbf{K}(\mathbf{x}', \mathbf{x})y - \frac{1}{\beta_N} (1 - e^{-\frac{2\alpha}{N}\beta_N t}) \mathbf{K}(\mathbf{x}', \mathbf{x})f(x, \theta_0) \\
&\quad - \frac{2\alpha}{N} e^{-\frac{2\alpha}{N}\beta_N t} \mathbf{K}(\mathbf{x}', \mathbf{x}) \left[ \frac{N}{2\alpha\beta} (e^{\frac{2\alpha}{N}\beta_N t} - 1)I - \frac{N}{2\alpha\beta} \mathbf{K}(\mathbf{x}, \mathbf{x})^{-1} \left( \exp\left(-\frac{2\alpha}{N}t\mathbf{K}(\mathbf{x}, \mathbf{x})\right) - I \right) \right] \\
&\quad \times (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N I)^{-1} \mathbf{K}(\mathbf{x}, \mathbf{x})(y - f(x, \theta_0)) \\
&= f(x', \theta_0) + \frac{1}{\beta_N} (1 - e^{-\frac{2\alpha}{N}\beta_N t}) \mathbf{K}(\mathbf{x}', \mathbf{x})(y - f(x, \theta_0)) \\
&\quad - \frac{1}{\beta} \mathbf{K}(\mathbf{x}', \mathbf{x}) \left[ (1 - e^{-\frac{2\alpha}{N}\beta_N t})I - \mathbf{K}(\mathbf{x}, \mathbf{x})^{-1} \left( \exp\left(-\frac{2\alpha}{N}t(\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N I)\right) - e^{-\frac{2\alpha}{N}\beta_N t}I \right) \right] \\
&\quad \times (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N I)^{-1} \mathbf{K}(\mathbf{x}, \mathbf{x})(y - f(x, \theta_0)).
\end{aligned}$$

223 Lastly, taking  $t \rightarrow \infty$ , we get

$$\begin{aligned}
f(x', \theta_\infty) &= f(x', \theta_0) + \frac{1}{\beta_N} \mathbf{K}(\mathbf{x}', \mathbf{x})(y - f(x, \theta_0)) - \frac{1}{\beta_N} \mathbf{K}(\mathbf{x}', \mathbf{x})(\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N I)^{-1} \mathbf{K}(\mathbf{x}, \mathbf{x})(y - f(x, \theta_0)) \\
&= f(x', \theta_0) + \frac{1}{\beta_N} \mathbf{K}(\mathbf{x}', \mathbf{x}) (I - (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N I)^{-1} \mathbf{K}(\mathbf{x}, \mathbf{x})) (y - f(x, \theta_0)) \\
&= f(x', \theta_0) + \mathbf{K}(\mathbf{x}', \mathbf{x})(\mathbf{K}(\mathbf{x}, \mathbf{x}) + \beta_N I)^{-1} (y - f(x, \theta_0)),
\end{aligned}$$

224 we achieve the desired result and hence having a regularized gradient flow in the infinite-width limit  
225 is equivalent to inferring the mean posterior of a non-zero aleatoric noise NTK-GP.  $\square$

226 **Lemma 3.2.** Consider the computational process derived in Lemma 3.1. Define shifted labels  $\tilde{\mathbf{y}}$  and  
227 predictions  $\tilde{f}(\mathbf{x}; \theta_\infty)$  as follows::

$$\tilde{\mathbf{y}} = \mathbf{y} + f(\mathbf{x}; \theta_0), \quad \tilde{f}(\mathbf{x}; \theta_\infty) = f(\mathbf{x}; \theta_\infty) - f(\mathbf{x}'; \theta_0).$$

228 Using these definitions, the posterior mean of a zero-mean Gaussian process can be computed as:

$$\tilde{f}(\mathbf{x}', \theta_\infty) = \mathbf{K}(\mathbf{x}', \mathbf{x})(\mathbf{K}(\mathbf{x}, \mathbf{x}) + N\beta_N \mathbf{I})^{-1} \mathbf{y}. \quad (6)$$

229 *Proof.* Firstly, substituting  $\tilde{\mathbf{y}}$  into  $\mathbf{y}$ :

$$\begin{aligned}
f(\mathbf{x}'; \theta_\infty) &= f(\mathbf{x}'; \theta_0) + \mathbf{K}(\mathbf{x}', \mathbf{x}) (\mathbf{K}(\mathbf{x}, \mathbf{x}) + N\beta_N \mathbf{I})^{-1} (\tilde{\mathbf{y}} - f(\mathbf{x}; \theta_0)) \\
&= f(\mathbf{x}'; \theta_0) + \mathbf{K}(\mathbf{x}', \mathbf{x}) (\mathbf{K}(\mathbf{x}, \mathbf{x}) + N\beta_N \mathbf{I})^{-1} \mathbf{y}
\end{aligned}$$

230 Now, using this new computational process, scaling it as  $\tilde{f}(\mathbf{x}; \theta_\infty)$ :

$$\tilde{f}(\mathbf{x}; \theta_\infty) = f(\mathbf{x}; \theta_\infty) - f(\mathbf{x}'; \theta_0) = \mathbf{K}(\mathbf{x}', \mathbf{x}) (\mathbf{K}(\mathbf{x}, \mathbf{x}) + N\beta_N \mathbf{I})^{-1} \mathbf{y},$$

231 achieving the desired zero-mean Gaussian process.  $\square$



## C Details of the Experimental Setup

The Adam optimizer was used whenever our experiments needed gradient descent. A patience-based stopping rule was used where training was stopped if there was no improvement in the loss for 500 epochs. The other hyperparameters are given in the table below.

	hyperparameter	value
	no of hidden layers	2
	size of hidden layer	512
	non-linearity	softplus
	softplus beta	87.09
	scaling multiplier in the output	3.5
	learning rate for network predicting mean	1e-4
	learning rate for covariance predictor networks	5e-5

Moreover, we used trigonometric normalization, where an input point  $x$  is first scaled and shifted to lie between 0 and  $\pi$ , obtaining a normalized point  $x'$ . The point  $x'$  is then represented with a vector  $[\sin(x'), \cos(x')]$ .

## D Details on Estimating The Covariance

We now describe two of dealing with the term  $\sigma^2 MM^\top$  in the covariance formula. Upper bounding the covariance is described in Section D.1, while estimating the exact covariance by fitting noisy targets is described in Section D.2.

### D.1 Upper Bounding the Covariance

First, we can simply ignore the term in our estimator, obtaining an upper bound on the covariance. We now characterize the tightness of the upper bound, i.e. the magnitude of the term

$$\sigma^2 MM^\top = \sigma^2 \mathbf{K}(\mathbf{x}', \mathbf{x})(\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1}(\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1} \mathbf{K}(\mathbf{x}', \mathbf{x})^\top.$$

We do this is the following two lemmas.

**Lemma D.1.** *When  $\mathbf{x} = \mathbf{x}'$ , i.e. on the training set, we have*

$$\sigma^2 \mathbf{K}(\mathbf{x}', \mathbf{x})(\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1}(\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1} \mathbf{K}(\mathbf{x}', \mathbf{x})^\top \preceq \sigma^2 \mathbf{I}.$$

*Proof.* By assumption,  $\mathbf{K}(\mathbf{x}', \mathbf{x}) = \mathbf{K}(\mathbf{x}, \mathbf{x}) = \mathbf{K}$ . Denote the diagonalization of  $\mathbf{K}$  with  $\mathbf{K} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$ . We have

$$\begin{aligned} & \sigma^2 \mathbf{K}(\mathbf{x}', \mathbf{x})(\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1}(\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1} \mathbf{K}(\mathbf{x}', \mathbf{x})^\top \\ &= \sigma^2 \mathbf{K}(\mathbf{K} + \sigma^2 \mathbf{I})^{-2} \mathbf{K}^\top \\ &= \sigma^2 \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top (\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top + \sigma^2 \mathbf{I})^{-2} \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top \\ &= \sigma^2 \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top \mathbf{U}(\mathbf{\Lambda} + \sigma^2 \mathbf{I})^{-2} \mathbf{U}^\top \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top \\ &= \sigma^2 \mathbf{U}\mathbf{\Lambda}(\mathbf{\Lambda} + \sigma^2 \mathbf{I})^{-2} \mathbf{\Lambda}\mathbf{U}^\top. \end{aligned}$$

It can be seen that the diagonal entries of  $\mathbf{\Lambda}(\mathbf{\Lambda} + \sigma^2 \mathbf{I})^{-2} \mathbf{\Lambda}$  are less than or equal one.  $\square$

The Lemma above, stated in words, implies that, on the training set, the variance estimates that come from using the upper bound (which doesn't require us to fit noisy targets as in Section D.2) are off by at most  $\sigma^2$ .

We now give another Lemma, which characterizes the upper bound on arbitrary test points, not just the training set.

**Lemma D.2.** *Denote by  $\lambda_{\max}$  the maximum singular value of  $\mathbf{K}(\mathbf{x}', \mathbf{x}')$ . Then we have*

$$\left\| \sigma^2 \mathbf{K}(\mathbf{x}', \mathbf{x})(\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1}(\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1} \mathbf{K}(\mathbf{x}', \mathbf{x})^\top \right\|_2 \leq \frac{1}{4} \lambda_{\max}.$$

258 *Proof.* By Proposition 1.3.2 from the book by Bhatia, we have that

$$\mathbf{K}(\mathbf{x}', \mathbf{x})^\top = \mathbf{K}(\mathbf{x}, \mathbf{x})^{1/2} \mathbf{C} \mathbf{K}(\mathbf{x}', \mathbf{x}')^{1/2},$$

259 where  $\mathbf{C}$  is a contraction. Denote the diagonalization of  $\mathbf{K}(\mathbf{x}, \mathbf{x})$  with  $\mathbf{K}(\mathbf{x}, \mathbf{x}) = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$ . We have

$$\begin{aligned} & \left\| \sigma^2 \mathbf{K}(\mathbf{x}', \mathbf{x}) (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-2} \mathbf{K}(\mathbf{x}', \mathbf{x})^\top \right\|_2 \\ &= \left\| \sigma^2 \mathbf{K}(\mathbf{x}', \mathbf{x}')^{1/2} \mathbf{C}^\top \mathbf{K}(\mathbf{x}, \mathbf{x})^{1/2} (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-2} \mathbf{K}(\mathbf{x}, \mathbf{x})^{1/2} \mathbf{C} \mathbf{K}(\mathbf{x}', \mathbf{x}')^{1/2} \right\|_2 \\ &\leq \sigma^2 \lambda_{\max} \left\| \mathbf{K}(\mathbf{x}, \mathbf{x})^{1/2} (\mathbf{K}(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-2} \mathbf{K}(\mathbf{x}, \mathbf{x})^{1/2} \right\|_2 \\ &= \sigma^2 \lambda_{\max} \left\| \mathbf{U} \mathbf{\Lambda}^{1/2} \mathbf{U}^\top \mathbf{U} (\mathbf{\Lambda} + \sigma^2 \mathbf{I})^{-2} \mathbf{U}^\top \mathbf{U} \mathbf{\Lambda}^{1/2} \mathbf{U}^\top \right\|_2 \\ &= \sigma^2 \lambda_{\max} \left\| \mathbf{\Lambda}^{1/2} (\mathbf{\Lambda} + \sigma^2 \mathbf{I})^{-2} \mathbf{\Lambda}^{1/2} \right\|_2. \end{aligned}$$

260 We can expand  $\left\| \mathbf{\Lambda}^{1/2} (\mathbf{\Lambda} + \sigma^2 \mathbf{I})^{-2} \mathbf{\Lambda}^{1/2} \right\|_2$  as  $\max_i \left\{ \frac{\lambda_i}{(\lambda_i + \sigma^2)^2} \right\} \leq \frac{1}{4\sigma^2}$ , which gives the desired  
261 result.  $\square$

## 262 D.2 Exact Covariance by Fitting Noisy Targets

263 In certain cases, we might not be satisfied with having an upper bound on the posterior covariance,  
264 even if it is reasonably tight. We can address these scenario by fitting additional predictor networks,  
265 trained on targets sampled from the spherical normal. Formally, we have

$$\sigma^2 M M^\top = M \mathbb{E}_\epsilon [\epsilon \epsilon^\top] M^\top,$$

266 where  $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ . We can take  $K'$  samples  $\epsilon_1, \dots, \epsilon_{K'}$ , obtaining

$$M \mathbb{E}_\epsilon [\epsilon \epsilon^\top] M^\top \approx \frac{1}{K'} \sum_i M \epsilon_i \epsilon_i^\top M^\top = \frac{1}{K'} \sum_i (M \epsilon_i)(M \epsilon_i)^\top, \quad (11)$$

267 where the approximation becomes exact by the law of large numbers as  $K' \rightarrow \infty$ . Since the  
268 multiplication  $M \epsilon_i$  is equivalent to estimating the posterior mean with algorithm 1, we can perform  
269 the computation in equation (11) by gradient descent.

## 270 E Computing The Partial SVD

271 Our Algorithm 2 includes the computation of the partial SVD of the Jacobian:

$$272 \quad U, \Sigma \leftarrow \text{PARTIAL-SVD}(J_{\theta_0}(\mathbf{x}), K).$$

273 We require an SVD which is partial in the sense that we only want to compute the first  $K$  singular  
274 values. For the regression experiment in this submission, we simply called the full SVD on the  
275 Jacobian and took the first  $K$  columns of  $U$  and the first  $K$  diagonal entries of  $\Sigma$ . This process is  
276 infeasible for larger problem instances.

277 This can be addressed by observing that the power method for SVD computation [2] only requires  
278 computing Jacobian-vector products and vector-Jacobian products, which can be efficiently computed  
279 in deep learning frameworks without access to the full Jacobian. Another approach that avoids  
280 constructing the full Jacobian is the use of randomized SVD [9]. We leave the implementation of  
281 these ideas to further work.

## 282 F Network Initialization

283 We consider a neural network model  $f(x; \theta)$ , where  $\theta \in \mathbb{R}^p$  denotes the set of parameters. The model  
284 consists of  $L$  layers with dimensions  $\{n_0, n_1, \dots, n_L\}$ , where  $n_0$  is the input dimension and  $n_L$  is  
285 the output dimension. Note that, as we want to leverage the theory of wide networks, the number of  
286 neurons in the hidden layers,  $\{n_2, \dots, n_{L-1}\}$ , is large.

287 For each fully connected layer  $l$ , the weight matrix  $W^{(l)} \in \mathbb{R}^{n_l \times n_{l-1}}$  and the bias vector  $b^{(l)} \in \mathbb{R}^{n_l}$   
 288 are initialized from a Gaussian distribution with mean zero and standard deviations  $\sigma_w$  and  $\sigma_b$ ,  
 289 respectively:

$$W_{ij}^{(l)} \sim \mathcal{N}(0, \sigma_w^2), \quad b_j^{(l)} \sim \mathcal{N}(0, \sigma_b^2),$$

290 where  $\sigma_w$  and  $\sigma_b$  are fixed values set as hyperparameters during initialization (we use  $\sigma_w = \sigma_b = 1$ ).

291 The network uses a non-linear activation function  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  with bounded second derivative,  
 292 ensuring Lipschitz continuity. The output of each layer  $l$  is scaled by  $1/\sqrt{n_l}$  to maintain the  
 293 appropriate magnitude, particularly when considering the infinite-width limit:

$$a^{(l)} = \sigma \left( \frac{1}{\sqrt{n_l}} W^{(l)} a^{(l-1)} + b^{(l)} \right),$$

294 where  $a^{(l)}$  is the output of layer  $l$ , and  $a^{(0)} = x$  is the input to the network.

295 The final layer output is further scaled by a constant factor  $c_{\text{out}}$  to ensure that the overall network  
 296 output remains within the desired range. Specifically, the output  $f(x; \theta)$  is given by:

$$f(x; \theta) = \frac{c_{\text{out}}}{\sqrt{n_L}} W^{(L)} a^{(L-1)},$$

297 where  $c_{\text{out}}$  is a predefined constant that ensures the final output is of the appropriate scale. In our  
 298 model,  $c_{\text{out}}$  is set to 3.5. For the hidden layers, we choose  $\sigma(\cdot)$  to be Softplus – a smoothed version  
 299 of ReLU. In this case, an additional scaling factor  $\beta$  is introduced to modulate the sharpness of the  
 300 non-linearity:

$$a^{(l)} = \sigma \left( \frac{1}{\sqrt{n_l}} W^{(l)} a^{(l-1)} + b^{(l)}; \beta \right).$$

301 In our model, we set  $\beta = 87.09$  for the Softplus activation to ensure the appropriate range of  
 302 activation values. The process described above is standard. We followed closely the methodology  
 303 provided in several works in the literature [12][15][8]. This initialization strategy ensures that the  
 304 network’s activations and gradients do not explode or vanish as the number of neurons  $n_l$  increases.