

---

# Bayesian Treatment of the Spectrum of the Empirical Kernel in (Sub)Linear-Width Neural Networks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We study Bayesian neural networks (BNNs) in the theoretical limits of infinitely  
2 increasing number of training examples, network width and input space dimen-  
3 sion. Our findings establish new bridges between kernel-theoretic approaches and  
4 techniques derived from statistical mechanics through the correspondence between  
5 Mercer’s eigenvalues and limiting spectral distributions of covariance matrices  
6 studied in random matrix theory. Our theoretical contributions first consist in  
7 novel integral formulas that accurately describe the predictors of BNNs in the  
8 asymptotic linear-width and sublinear-width regimes. Moreover, we extend the  
9 recently developed renormalisation theory of deep linear neural networks, enabling  
10 a rigorous explanation of the mounting empirical evidence that hints at the theory’s  
11 applicability to nonlinear BNNs with ReLU activations in the linear-width regime.  
12 From a practical standpoint, our results introduce a novel technique for estimating  
13 the predictor statistics of a trained BNN that is applicable to the sublinear-width  
14 regime where the predictions of the renormalisation theory are inaccurate.

## 15 1 Introduction

16 Bayesian Neural Networks (BNNs) are a variant of neural networks that incorporate Bayesian infer-  
17 ence techniques to mitigate overfitting, enable learning from small datasets, and capture uncertainty  
18 in predictions [Neal, 2012, Gal, 2016]. In a BNN, prior probability distributions are specified for  
19 weights and biases. During training, the posterior distribution, which represents the updated knowl-  
20 edge about the parameters after observing the data, is updated using Bayes’ rule. A trained BNN  
21 can be interpreted as an infinite ensemble of neural networks where each individual contribution  
22 in the ensemble is weighted by the posterior probability of its parameters given the training data.  
23 Although computing the posterior distribution is intractable and difficult to approximate, BNNs have  
24 gained significant traction with the development of effective estimation techniques [Gal, 2016, Blei  
25 et al., 2017]. BNNs demonstrate generalisation performance on par with deep neural networks trained  
26 using gradient descent [Lee et al., 2020, Magris and Iosifidis, 2023]. BNNs also showcase improved  
27 sensitivity to out-of-distribution examples [Gal, 2016] and the ability to estimate uncertainty.

28 In an effort to analyse the generalisation properties of BNNs, researchers study idealised views  
29 of fully-connected neural architectures defined by the input dimension, the layer widths, and the  
30 activation function. As the width approaches infinity in each layer (the *NNGP limit*), the functions  
31 generated by random weight selection converge in distribution to a Gaussian process (GP) [Rasmussen  
32 and Williams, 2006]. The covariance function of such GP, called the *NNGP kernel*, can be recursively  
33 defined by proceeding on a layer by layer basis [Lee et al., 2018b]. This perspective based on kernel  
34 and GP theory has led to the development of analytical formulas to estimate the generalisation error  
35 of related kernel and random features models [Canatar et al., 2021, Simon et al., 2023]. These  
36 formulas often rely on the *spectral universality assumption* (SUA), which simplifies the derivations

37 by approximating the eigenfunctions of the kernel with independent Gaussian entries [Karoui, 2010,  
 38 Cheng and Singer, 2013, Fan and Montanari, 2015]. Extensive research is being devoted to study the  
 39 accuracy of the SUA [Liu et al., 2021, Lu and Yau, 2023, Schröder et al., 2023].

40 In addition to the NNGP limit, BNNs have also been studied under the *linear-width limit* (also  
 41 referred to as *thermodynamic limit* or *proportional limit*) where the network’s width, the number of  
 42 training examples and the dimension of the input space are taken simultaneously to infinity while  
 43 keeping constant and bounded ratios between them [Engel et al., 2012]. By employing techniques  
 44 from statistical mechanics, such as saddle point approximations [Seung and Sompolinsky, 1992], the  
 45 replica method [Barbier et al., 2018, Canatar et al., 2021], and random matrix theory [Wigner, 1955,  
 46 Livan et al., 2018], researchers have studied the mean and variance of the output generated by trained  
 47 BNNs in this setting. One of the most prominent results in this area is the *renormalisation theory* [Li  
 48 and Sompolinsky, 2021] of linear BNNs (i.e., those without non-linear activations), which establishes  
 49 that the mean predictor and the predictor variance of the BNN coincide with that of Bayesian linear  
 50 regression, but surprisingly the variance must be renormalised by a factor dependent on the training  
 51 data and problem dimensions. Subsequent developments have provided more detailed analysis on  
 52 the linear setting including non-asymptotic results [Hanin and Zlokapa, 2023], and comparison with  
 53 deep random feature models [Zavatone-Veth et al., 2022]. It remains an open question, however,  
 54 whether the insights from the renormalisation theory for linear BNNs can be extended to non-linear  
 55 networks, as suggested by empirical evidence [Li and Sompolinsky, 2021, Ariosto et al., 2023]. A  
 56 recent theoretical work [Cui et al., 2023] further substantiates these observations by deriving the  
 57 predictor learned by non-linear BNNs in the case of Gaussian data.

58 More recently, *sublinear-width regimes*, where the width is small compared to the number of data  
 59 points [Maillard et al., 2024], and related scalings [van Meegen and Sompolinsky, 2024] have been  
 60 studied, and the emergence of strong feature learning has been demonstrated in these scenarios.

61 **Our Contributions** In this paper, we establish new connections between the kernel-theoretic  
 62 perspective associated with the NNGP limit and the statistical mechanics viewpoint associated  
 63 with the linear-width and sublinear-width limits, and contribute new insights to the generalisation  
 64 properties of BNNs. First, we demonstrate that training a (non-linear) BNN in the linear-width and  
 65 sublinear-width limits result in a predictor with identical mean and variance to that of GP regression  
 66 with a modified NNGP kernel, and we observe that the Mercer spectrum [Mercer, 1909, Minh  
 67 et al., 2006] of this kernel is known in the linear-width regime. Second, we prove necessary and  
 68 sufficient conditions (on the data and the architecture) for the application of renormalisation theory to  
 69 non-linear BNNs in the linear-width limit. These conditions also provide a criterion for determining  
 70 the applicability of the spectral universality assumption (SUA) from kernel theory in the context of  
 71 BNNs. Third, we present initial findings on a sublinear-width regime where the relevant quantities  
 72 are simultaneously taken to infinity while the number of training examples remains proportional to  
 73 the product of the network width and the dimension of the input space. In particular, we provide a  
 74 novel mechanism for estimating the mean and variance of the predictions of non-linear BNNs in this  
 75 setting, for which renormalisation theory is not applicable.

## 76 2 Preliminaries

77 We use standard notation for real-valued vectors  $\mathbf{v} \in \mathbb{R}^n$ , matrices  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , and their transposes  
 78  $\mathbf{v}^T$  and  $\mathbf{A}^T$ . We use  $\mathbf{a}_i$  to denote the vector in the  $i$ -th row of  $\mathbf{A}$ . The Moore-Penrose pseudo-inverse  
 79 of a matrix  $\mathbf{A}$  is denoted as  $\mathbf{A}^\dagger$  [Moore, 1920].

80 **Neural networks.** A fully-connected neural (FCN) architecture with  $L$  layers is a tuple  $f =$   
 81  $\langle \{\mathbf{W}^\ell\}_{1 \leq \ell \leq L}, \{\mathbf{b}^\ell\}_{1 \leq \ell \leq L}, \{\sigma^\ell\}_{1 \leq \ell \leq L} \rangle$ . Each layer  $\ell \in \{1, \dots, L\}$  of width  $N_\ell$  is given by a weight  
 82 matrix  $\mathbf{W}^\ell \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$ , a bias  $\mathbf{b}^\ell \in \mathbb{R}^{N_\ell}$  and an activation function  $\sigma^\ell : \mathbb{R} \mapsto \mathbb{R}$ . On input  
 83  $\mathbf{x} \in \mathbb{R}^{N_0}$ , network  $f$  sets  $\mathbf{x}^0 = \mathbf{x}$  and then computes recursively on the depth the sequence of  
 84 pre-activations  $\mathbf{h}^\ell$  and activations  $\mathbf{x}^\ell$  as follows, where the network’s output  $f(\mathbf{x})$  is given by  $\mathbf{x}^L$ :

$$\mathbf{h}^\ell = \mathbf{W}^\ell \cdot \mathbf{x}^{\ell-1} + \mathbf{b}^\ell \quad \mathbf{x}^\ell = \sigma^\ell(\mathbf{h}^\ell) \quad (1)$$

85 We assume that all but the last layer have the same width  $N$ . For the last layer, we assume width  
 86  $N_L = 1$  (ensuring a real-valued output),  $b^L = 0$  and  $\sigma^L = \text{Id}_{\mathbb{R}}$  (ensuring linearity). In this setting,  
 87 the weights  $\mathbf{W}^L$  are referred to as the readout weights [Li and Sompolinsky, 2021].

88 **Kernels.** A kernel on  $\mathbb{R}^{N_0}$  is a positive semi-definite symmetric function  $K : \mathbb{R}^{N_0} \times \mathbb{R}^{N_0} \mapsto \mathbb{R}$ . By  
89 Mercer’s theorem [Minh et al., 2006], given a distribution  $\mathbf{x} \sim p(\mathbf{x})$  with compact support on  $\mathbb{R}^{N_0}$ ,  
90 there exist unique countable collections of Mercer’s eigenvalues  $(\lambda_i)_{i \in \mathbb{N}}$  and eigenfunctions  $(\varphi_i)_{i \in \mathbb{N}}$   
91 such that  $K(\mathbf{x}, \mathbf{x}') = \sum_i \lambda_i \varphi_i(\mathbf{x}) \varphi_i(\mathbf{x}')$  and  $(\varphi_i)$  are orthonormal w.r.t. the data distribution:  
92  $\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} (\varphi_i(\mathbf{x}) \varphi_j(\mathbf{x})) = \delta_{i,j}$  for all  $i, j$ . By Riesz’s theorem, there exists a Hilbert space  $\mathcal{H}$  and a  
93 feature map  $\phi : \mathbb{R}^{N_0} \mapsto \mathcal{H}$  such that  $K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{H}}$ . Kernel regression amounts to linear  
94 regression in the corresponding Hilbert space: when trained on data  $\mathbf{X}, \mathbf{y}$ , the prediction on a new  
95 point  $\mathbf{x}^*$  is given by  $\mathbf{k}_{\mathbf{x}^*, \mathbf{X}}^T \mathbf{K}_{\mathbf{X}, \mathbf{X}}^{-1} \mathbf{y}$  where the vector  $\mathbf{k}_{\mathbf{x}^*, \mathbf{X}}$  is given by  $(\mathbf{k}_{\mathbf{x}^*, \mathbf{X}})_i = K(\mathbf{x}^*, \mathbf{x}_i)$  and  
96 the kernel matrix  $\mathbf{K}_{\mathbf{X}, \mathbf{X}}$  is given by  $(\mathbf{K}_{\mathbf{X}, \mathbf{X}})_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$ . Although the kernel’s eigenfunctions  
97 exhibit the described structure, the *spectral universality assumption* (SUA) is commonly adopted. The  
98 SUA posits that, as  $P$  increases, the eigenfunctions can be approximated by independent Gaussian  
99 entries:  $\varphi_i(\mathbf{x}_j) \sim \mathcal{N}(\mu_K, \sigma_K^2)$ , where  $\mu_K$  and  $\sigma_K^2$  depend on the kernel  $K$  and the data distribution  
100  $p(\mathbf{x})$ , but not on specific instances  $i$  and  $j$ . The SUA works well in practice [Karoui, 2010, Cheng  
101 and Singer, 2013, Fan and Montanari, 2015, Liu et al., 2021, Simon et al., 2023, Lu and Yau, 2023,  
102 Schröder et al., 2023], and research focuses on identifying conditions under which it holds.

103 **Random feature maps.** Let  $\Theta$  represent all parameters of  $f$  up to layer  $L - 1$ . The *random*  
104 *feature map*  $\phi(\Theta, \cdot) : \mathbb{R}^{N_0} \mapsto \mathbb{R}^N$  is a nonlinear transformation (random in  $\Theta$ ) mapping the input  
105 and the activation  $\mathbf{x}^{L-1}$ . By definition,  $f(\mathbf{x}) = (\mathbf{W}^L)^T \phi(\Theta, \mathbf{x})$ , and to highlight the parameter  
106 dependency we denote it as  $f_{\Theta, \mathbf{W}^L}$ . The random feature map is associated to a *random kernel*  
107  $K_{\Theta}^{N, N_0} : (\mathbf{x}, \mathbf{x}') \mapsto \frac{1}{N} \langle \phi(\Theta, \mathbf{x}), \phi(\Theta, \mathbf{x}') \rangle$  expressed as the inner product between the corresponding  
108 random feature map evaluations. For this kernel, the Hilbert space  $\mathcal{H} = \mathbb{R}^N$  is thus known.

109 **Training set.** The training set  $(\mathbf{X}, \mathbf{y})$  consists of  $P$  examples sampled i.i.d. from an unknown  
110 distribution  $\mathbb{P}_{N_0}$  with compact support on  $\mathbb{R}^{N_0} \times \mathbb{R}$ . We assume that in the limit  $N_0 \rightarrow \infty$ ,  $\mathbb{P}_{N_0}$   
111 converges to a well-defined distribution over  $\mathbb{R}^N$  noted  $\lim_{N_0 \rightarrow \infty} \mathbb{P}_{N_0}$ . We denote each example  
112 by  $(\mathbf{x}_i, y_i)$ , so that  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_P)^T \in \mathbb{R}^{P \times N_0}$  and  $\mathbf{y} = (y_1, \dots, y_P)^T \in \mathbb{R}^P$ . We denote the  
113 evaluation of the random feature map on the training set by  $\phi(\Theta, \mathbf{X}) = (\phi(\Theta, \mathbf{x}_1), \dots, \phi(\Theta, \mathbf{x}_P))^T \in$   
114  $\mathbb{R}^{N \times P}$ ; this induces an empirical kernel matrix  $\mathbf{K}_{\Theta}^{P, N, N_0}(\mathbf{X}, \mathbf{X})$  given by  $\frac{1}{N} [\phi(\mathbf{X}, \Theta)]^T \phi(\mathbf{X}, \Theta) \in$   
115  $\mathbb{R}^{P \times P}$ . The training data  $\mathbf{X}$  also induces an empirical distribution  $p_{\mathbf{X}}(\mathbf{x}) = \frac{1}{P} \left( \sum_{i=1}^P \delta_{\mathbf{x}_i}(\mathbf{x}) \right)$  with  
116  $\delta_{\mathbf{x}_i}$  the Dirac measure.

117 **BNNs.** We assume a *prior distribution* over parameters  $(\Theta, \mathbf{W}^L)$  with weights sampled i.i.d. from  
118  $\mathcal{N}(0, \frac{1}{N})$  and biases sampled i.i.d. from  $\mathcal{N}(0, 1)$ ; this yields a density  $p(\Theta, \mathbf{W}^L)$  that is a product of  
119 Gaussian densities. The *posterior distribution* given the training data is given by Bayes’ rule:

$$p(\Theta, \mathbf{W}^L | \mathbf{X}, \mathbf{y}) = p(\Theta, \mathbf{W}^L) \frac{p(\mathbf{y} | \mathbf{X}, \Theta, \mathbf{W}^L)}{p(\mathbf{y} | \mathbf{X})}$$

120 where  $p(\mathbf{y} | \mathbf{X}, \Theta, \mathbf{W}^L)$  is the *likelihood* of the data given a set of parameters, and  $p(\mathbf{y} | \mathbf{X}) =$   
121  $\int p(\mathbf{y} | \mathbf{X}, \Theta, \mathbf{W}^L) p(\Theta, \mathbf{W}^L) d\Theta d\mathbf{W}^L$  is the *marginal likelihood* (or *evidence*). We assume Gaussian  
122 likelihoods, i.e.  $p(\mathbf{y} | \mathbf{X}, \Theta, \mathbf{W}^L) \sim \mathcal{N}(\mathbf{y}, \phi(\Theta, \mathbf{X})^T \mathbf{W}^L \mathbf{W}^{L T} \phi(\Theta, \mathbf{X}))$ . Calculating the posterior  
123 distribution, which is the essence of BNN training, is analytically intractable and remains a core  
124 challenge [Gal, 2016]. In practice, the posterior distribution is estimated via variational inference  
125 [Blei et al., 2017] or Monte-Carlo simulation methods [Rasmussen, 1995].

126 Given the posterior distribution, the predictor defines a distribution over functions  $f_{\Theta, \mathbf{W}^L}$  with  
127  $(\Theta, \mathbf{W}^L) \sim p(\Theta, \mathbf{W}^L | \mathbf{X}, \mathbf{y})$ . The mean-squared generalisation error is defined for any new point  
128  $(\mathbf{x}^*, y^*)$  as the expectation over the predictor error:  $\mathbb{E}_{(\Theta, \mathbf{W}^L) \sim p(\Theta, \mathbf{W}^L | \mathbf{X}, \mathbf{y})} \left( (y^* - f_{\Theta, \mathbf{W}^L}(\mathbf{x}^*))^2 \right)$ .  
129 Only the mean and variance of the predictor are needed to calculate it.

130 **Gaussian processes and NNGPs.** A GP  $g$  over a space  $\mathbb{R}^{N_0}$  is a random scalar field such that its  
131 evaluation at any collection of finitely many points  $(g(x_1), \dots, g(x_P))$  follows a multivariate Gaussian  
132 distribution. A GP is determined by a mean function  $\mu : \mathbb{R}^{N_0} \mapsto \mathbb{R}$ , and a covariance function  
133  $K : \mathbb{R}^{N_0} \times \mathbb{R}^{N_0} \mapsto \mathbb{R}$ , which describe respectively the mean of the Gaussian distribution at each  
134 point and the covariance between the Gaussians at any two points. The covariance function of a GP

135 is a kernel [Rasmussen and Williams, 2006]. We note  $g \sim \mathcal{GP}(\mu, K)$ . Gaussian process regression  
 136 consists in performing Bayesian inference using a Gaussian process as the prior distribution over  
 137 functions. The prediction distribution of GP regression with prior  $\mathcal{GP}(0, K)$  trained on the data  
 138  $\mathbf{X}, \mathbf{y}$  is given, on a new point  $\mathbf{x}^*$ , by  $\mathcal{N}(\mathbf{k}_{\mathbf{x}^*, \mathbf{X}}^T \mathbf{K}_{\mathbf{X}, \mathbf{X}}^{-1} \mathbf{y}, K(\mathbf{x}^*, \mathbf{x}^*) - \mathbf{k}_{\mathbf{x}^*, \mathbf{X}}^T \mathbf{K}_{\mathbf{X}, \mathbf{X}}^{-1} \mathbf{k}_{\mathbf{X}, \mathbf{x}^*})$ . The mean  
 139 prediction of GP regression coincides with the prediction of kernel regression with the same kernel.

140 Applying successively the central limit theorem to each layer, the infinite-width limit of (1) yields  
 141 a GP, called the *Neural Network Gaussian Process (NNGP)*. If we let the width  $N \rightarrow \infty$ , the  
 142  $\mathbf{h}_i^L \sim \mathcal{GP}(\mu^L, K^L)$  are independent and defined inductively by layers as follows for all  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{N_0}$   
 143 and each  $\ell \in 1, \dots, L$ . First,  $\mu^\ell(\mathbf{x}) = 0$  and  $K^0(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$ . Then,  $\mathbf{h}_i^{\ell-1} \sim \mathcal{GP}(\mu^{\ell-1}, K^{\ell-1})$  and  
 144 the covariance functions  $K^\ell(\mathbf{x}, \mathbf{x}')$  are given by  $\mathbb{E}_{\mathbf{h}_i^{\ell-1} \sim \mathcal{GP}(\mu^{\ell-1}, K^{\ell-1})} (\sigma^\ell(\mathbf{h}_i^{\ell-1}(\mathbf{x})) \sigma^\ell(\mathbf{h}_i^{\ell-1}(\mathbf{x}')))$ .

145 The covariance function  $K^L$  is the *NNGP kernel* [Daniely et al., 2016], denoted as  $K^L = K_{\text{NNGP}}$ .  
 146 Infinite-width limits involve various subtleties [Matthews et al., 2018], and we follow the approach in  
 147 Lee et al. [2018a] where infinite limits are taken sequentially. In this limit, the number of examples  
 148  $P$  and the input dimension  $N_0$  remain fixed. Furthermore, we will investigate more comprehensive  
 149 limits where  $P, N$ , and  $N_0$  all tend to infinity simultaneously, first while maintaining constant and  
 150 bounded ratios  $\alpha = \frac{P}{N}$  and  $\alpha_0 = \frac{P}{N_0}$  (linear-width regime), then while  $P \propto N \cdot N_0$ , thus  $\alpha \rightarrow \infty$   
 151 and  $\alpha_0 \rightarrow \infty$  (sublinear-width regime).

152 **Random matrix theory.** Random matrix theory [Wigner, 1955, Livan et al., 2018] is the study  
 153 of the spectral distributions of large matrices of random variables. The spectral measure  $F_P$  for  
 154 a given matrix, with eigenvalues  $\lambda_i$ , is given, for  $x \in \mathbb{R}$ , by  $F_P(x) := \frac{1}{P} \sum_{i=1}^P \delta_{\lambda_i}(x)$ , where  
 155  $\delta_{\lambda_i}(x)$  represents the Dirac measure centered at the eigenvalue  $\lambda_i$ . When the matrix is random, the  
 156 spectral measure becomes a random measure, called the empirical spectral distribution. Our focus  
 157 lies in studying weak convergences (convergences in distribution) of the spectral measures towards  
 158 nonrandom measures [Geronimo and Hill, 2002]. A sufficient condition for weak convergence of  
 159 measures is to have pointwise convergence in their Stieltjes transforms [Geronimo and Hill, 2002].  
 160 We rely on a famous result in random matrix theory. Consider  $\mathbf{W} \in \mathbb{R}^{N \times P}$ , a random matrix with  
 161 i.i.d. entries drawn from  $\mathcal{N}(0, \frac{1}{N})$  and  $\Psi$  a nonrandom positive semi-definite matrix. Suppose that  $\Psi$   
 162 has a limiting spectral measure  $\rho$ , and let  $P, N \rightarrow \infty$  with fixed ratio  $\alpha := \frac{P}{N}$ , then the random matrix  
 163  $\Psi^{1/2} \mathbf{W}^T \mathbf{W} \Psi^{1/2}$  has a limiting nonrandom spectral measure  $\rho_{MP}^\alpha \boxtimes \rho$ . The Marchenko-Pastur  
 164 map of  $\rho$ , denoted  $\rho_\alpha^{MP} \boxtimes \rho$ , is defined by the Stieltjes transform solving the Marchenko-Pastur  
 165 equation [Marchenko and Pastur, 1967]. It also appears in the free probability literature as the free  
 166 multiplicative convolution between the probability measures  $\rho_{MP}^\alpha$  and  $\rho$  [Mingo and Speicher, 2017].  
 167 When considering the specific case where  $\Psi = \mathbf{I}_P$  (identity matrix of size  $P$ ), then  $\rho$  represents the  
 168 Dirac measure at 1 and we recover the well-known Marchenko-Pastur distribution, denoted as  $\rho_{MP}^\alpha$ .  
 169 Furthermore, we denote as  $\rho_{MP} \boxtimes^\ell \rho := \rho_{MP} \boxtimes (\dots(\rho_{MP} \boxtimes \rho))$  the composition of  $\ell$  successive  
 170 Marchenko-Pastur maps.

### 171 3 BNNs as Modified GP Regression

172 First we state our definitions of linear-width and sublinear-width regimes.

173 **Assumption 3.1** (Linear-width regime). Assume that  $\frac{P}{N} \rightarrow \alpha$  and  $\frac{P}{N_0} \rightarrow \alpha_0$  as  $P, N, N_0 \rightarrow \infty$  with  
 174 the ratios  $\alpha, \alpha_0 \in (0, +\infty)$ .

175 **Assumption 3.2** (Sublinear-width regime). Assume that  $\frac{P}{N \cdot N_0} \rightarrow \gamma$  as  $P, N, N_0 \rightarrow \infty$  with the ratio  
 176  $\gamma \in (0, +\infty)$ .

177 Our first aim in this section is to showcase the emergence of a modified NNGP kernel during the  
 178 training of BNNs in the linear-width and sublinear-width limits. We then study the Mercer's spectrum  
 179 of the modified NNGP kernel and exploit it to extend the renormalisation theory to encompass  
 180 nonlinear networks in the linear-width regime. Finally, we outline the fundamental arguments  
 181 supporting the expansion of this theory to the sublinear-width regime.

182 **3.1 The Modified NNGP Kernel**

183 Mercer’s theorem applied to the random kernel  $K_{\Theta}^{N,N_0}$  and the data distribution  $p_{\mathbf{X}}$  decom-  
 184 poses the kernel into terms of eigenvalues and eigenfunctions as follows:  $K_{\Theta}^{N,N_0}(\mathbf{x}, \mathbf{x}') =$   
 185  $\sum_k \lambda_k^{P,N,N_0} \varphi_k^{P,N,N_0}(\mathbf{x}) \varphi_k^{P,N,N_0}(\mathbf{x}')$ . This defines a random spectral measure  $\rho_{\Theta}^{P,N,N_0}$  with spec-  
 186 trum given by the eigenvalues and random eigenfunctions  $\varphi_k^{P,N,N_0}(\cdot)$ . Here, the dependency in  $P$   
 187 stems from the empirical training data distribution. We use the correspondence between Mercer’s  
 188 eigenvalues and the limiting spectral measure of the corresponding empirical kernel matrix to show  
 189 that, to estimate the infinite random matrix  $\mathbf{K}_{\Theta}^{P,N,N_0}(\mathbf{X}, \mathbf{X})$ , there is no need to examine the joint  
 190 distribution of its eigenvalues, as Mercer’s eigenvalues can be sampled independently. This is a crucial  
 191 observation because the correlations between kernel matrix eigenvalues in the classical decomposition  
 192 is a notorious obstacle in the computation of the posterior distributions.

193 **Theorem 3.3.** *Assume that Assumption 3.1 (respectively, Assumption 3.2) holds. Assume that for*  
 194 *each  $k \in \mathbb{N}$  there is a random function  $\varphi_k^{\alpha,\alpha_0} : \mathbb{R}^N \mapsto \mathbb{R}$  (respectively,  $\varphi_k^{\gamma}$ ) such that  $\varphi_k^{P,N,N_0}(\mathbf{x}_i)$*   
 195 *converges in distribution to  $\varphi_k^{\alpha,\alpha_0}(\mathbf{x})$  (respectively,  $\varphi_k^{\gamma}(\mathbf{x})$ ), where  $\mathbf{x} \sim \lim_{N_0 \rightarrow \infty} \mathbb{P}_{N_0}$ . Assume that*  
 196 *the spectrum of  $\mathbf{K}_{\Theta}^{P,N,N_0}(\mathbf{X}, \mathbf{X})$  (respectively, the strictly positive support of the spectrum) admits*  
 197 *a limiting nonrandom measure  $\rho^{\alpha,\alpha_0}$  (respectively,  $\rho^{\gamma}$ ). Consider the random matrix  $\Phi \Lambda \Phi^T$ , with*  
 198  *$\Phi \in \mathbb{R}^{P \times M}$ ,  $\Phi_{i,k} := \varphi_k^{\alpha,\alpha_0}(\tilde{\mathbf{x}}_i)$  (respectively,  $\varphi_k^{\gamma}$ ) and  $\Lambda \in \mathbb{R}^{M \times M}$ ,  $\Lambda_{k,l} := \delta_{k,l} \lambda_k$  with each  $\lambda_k$*   
 199 *follows independently  $\rho^{\alpha,\alpha_0}$  (respectively,  $\rho^{\gamma}$ ) and each  $\tilde{\mathbf{x}}_i$  follows independently  $\lim_{N_0 \rightarrow \infty} \mathbb{P}_{N_0}$ <sup>1</sup>.*  
 200 *Then, the random matrices  $\mathbf{K}_{\Theta}^{P,N,N_0}(\mathbf{X}, \mathbf{X})$  and  $\Phi \Lambda \Phi^T$  converge (in distribution) to the same*  
 201 *distribution over  $\mathbb{R}^{N \times N}$  in the limit  $\frac{M}{P} \rightarrow \infty$ .*

202 This result is non-trivial and only holds if the spectrum admits a *nonrandom* limit: this is the key  
 203 argument that allows us to disregard, in the limit, the correlations between eigenvalues when using  
 204 the Mercer decomposition. Note that the distribution of eigenfunctions  $\varphi_k^{\alpha,\alpha_0}$  and  $\varphi_k^{\gamma}$  are not known  
 205 in general. We will justify in the next section the SUA as a means for alleviating this limitation.  
 206 Similarly, we will denote with  $\Phi^*$  evaluations of the eigenfunctions on an unseen data point  $\mathbf{x}^*$ .

207 The *modified NNGP kernel* is the random kernel  $K_{\Theta}^{\alpha,\alpha_0}$  (respectively,  $K_{\Theta}^{\gamma}$ ) defined over  $\mathbb{R}^N$  in the  
 208 linear-width regime (respectively, the sublinear-width regime). In the limit, the feature map is not  
 209 known explicitly, but it must exist by Riesz’s representation theorem.

210 **The nonrandom spectral measure is known in the linear-width regime.** Observe that, in the  
 211 linear-width regime, for many cases of interest (including ReLU activations),  $\rho_{\Theta}^{\alpha,\alpha_0}$  indeed no longer  
 212 depends on  $\Theta$  and hence becomes a nonrandom measure. To this end, let us first consider the  
 213 kernel random matrix  $\mathbf{K}_{\text{NNGP}}(\mathbf{X}, \mathbf{X})$  associated with the NNGP kernel  $K_{\text{NNGP}}$ . El Harzli et al.  
 214 [2024] have shown that, under mild assumptions on the activation functions  $\sigma^{\ell}$  (namely measurability  
 215 and Lipschitz continuity),  $\mathbf{K}_{\text{NNGP}}(\mathbf{X}, \mathbf{X})$  admits a limiting nonrandom spectral measure  $\rho_{\text{NNGP}}^{\alpha_0}$   
 216 as  $P, N_0 \rightarrow \infty$  with constant ratio  $\alpha_0$ ; and furthermore, in the linear-width limit, the limiting  
 217 spectral distribution of the same random matrix as  $\mathbf{K}_{\Theta}^{P,N,N_0}(\mathbf{X}, \mathbf{X})$  but where the interior widths  
 218 have already been taken to infinity (i.e. when the linear-width limit only pertains to the last-layer  
 219 width) is  $\rho_{MP}^{\alpha} \boxtimes \rho_{\text{NNGP}}^{\alpha_0}$ . By immediate induction, successively applying the linear-width limit to  
 220 the hidden-layer widths and keeping the remaining interior widths infinite until reaching the input  
 221 layer, it follows as a direct corollary of Theorem 2 in El Harzli et al. [2024] that, in the linear-width  
 222 limit,  $\mathbf{K}_{\Theta}^{P,N,N_0}(\mathbf{X}, \mathbf{X})$  also admits a limiting nonrandom spectral measure given by the composed  
 223 Marchenko-Pastur maps  $\rho_{MP}^{\alpha} \boxtimes^L \rho_{\text{NNGP}}^{\alpha_0}$ .<sup>2</sup>

224 **3.2 Training BNNs with the Modified NNGP Kernel**

225 We can now study the predictor statistics of trained BNNs in the linear-width limit and the sublinear-  
 226 width limit. In particular, the following theorem provides integral formulae to estimate, under the  
 227 SUA, the first and second moments of the trained BNN using only the limiting spectral measure.

<sup>1</sup>Expression  $\Phi \Lambda \Phi^T$  is not the usual eigendecomposition of a square matrix: the evaluations of eigenfunctions yield rectangular (infinite) matrices. This decomposition is enabled by Mercer’s theorem and applies to kernels.

<sup>2</sup>This result first appeared in the context of neural networks in Fan and Wang [2020]. The result by El Harzli et al. [2024] extends it to a more general setting.



228 In this section, the results hold indistinctively of the linear-width or the sublinear-width limit, so to  
 229 simplify notations, we will note  $\rho$  for both  $\rho^{\alpha, \alpha_0}$  and  $\rho^\gamma$  and  $K$  for both  $K_\Theta^{\alpha, \alpha_0}$  and  $K_\Theta^\gamma$ .

230 **Theorem 3.4.** Assume that Assumption 3.1 or Assumption 3.2 holds. Let  $\rho$  be the nonrandom spectral  
 231 measure characterising the modified NNGP kernel  $K$ , and assume that the SUA holds.

232 The mean  $\langle f \rangle(\mathbf{x}^*, \mathbf{X}, \mathbf{y})$  and variance  $\langle \delta f \rangle(\mathbf{x}^*, \mathbf{X}, \mathbf{y})$  of the predictor associated to a BNN with  
 233 training data  $(\mathbf{X}, \mathbf{y})$  is given by expressions (2) and (3) respectively:<sup>3</sup>

$$\langle f \rangle = \int \left( \Phi^{*T} \Lambda \Phi^T \Phi^{T\dagger} \Lambda^{-1} \Phi^\dagger \mathbf{y} \right) \cdot \frac{p(\mathbf{y}, \Phi | \Lambda, \mathbf{X})}{p(\mathbf{y} | \mathbf{X})} d(\rho)(\Lambda) \mathcal{D}\Phi \mathcal{D}\Phi^* \quad (2)$$

234

$$\langle \delta f \rangle = \int \left( \Phi^{*T} \Lambda \Phi^* - \Phi^{*T} \Lambda \Phi^T \Phi^{T\dagger} \Lambda^{-1} \Phi^\dagger \Phi \Lambda \Phi^* \right) \cdot \frac{p(\mathbf{y}, \Phi | \Lambda, \mathbf{X})}{p(\mathbf{y} | \mathbf{X})} d(\rho)(\Lambda) \mathcal{D}\Phi \mathcal{D}\Phi^* \quad (3)$$

235 where  $\mathcal{D}\Phi$  is a standard Gaussian matrix measure,  $\Phi_{i,j} \sim_{\text{iid}} \mathcal{N}(\mu_K, \sigma_K^2)$  obtained from the SUA; the  
 236 likelihood is given by  $p(\mathbf{y}, \Phi | \Lambda, \mathbf{X}) \sim \mathcal{N}(\Phi^T \mathbf{y}, \Lambda)$ ; the marginal likelihood is given by  $p(\mathbf{y} | \mathbf{X}) =$   
 237  $\int p(\mathbf{y}, \Phi | \Lambda, \mathbf{X}) d(\rho)(\Lambda) \mathcal{D}\Phi$ .

238 The integral forms (2) and (3) provide a new estimation of the predictor statistics of a trained BNN.  
 239 While these expressions are exact only in the limit, we will present empirical evidence suggesting that  
 240 they constitute a reasonable approximation. A practical challenge arises from the need to estimate the  
 241 spectral distribution  $p(\Lambda | \mathbf{X}, \mathbf{x}^*) = (\rho)(\Lambda)$ , which often involves diagonalising a kernel matrix. This  
 242 process can be computationally intensive, especially for large training datasets; however, it is worth  
 243 mentioning that in many applications of BNNs, where the training datasets are relatively small, this  
 244 computational difficulty becomes less significant.

245 Theorem 3.4 and its proof also offer valuable new perspectives on the applicability of the SUA. In  
 246 particular, in the last steps of the proof, the probability density of  $\Phi, \Phi^*$  no longer appears directly  
 247 in the integral. For given  $\Lambda, \mathbf{X}, \mathbf{x}^*$ , if each orthogonal matrix  $\Phi, \Phi^*$  has a non-zero probability of  
 248 occurring, the integral spans uniformly over the entire space of orthogonal matrices of size  $P \times M$ .  
 249 This is useful because in the limit of infinite dimensions, this space coincides with that of Gaussian  
 250 matrices with independent entries [Haar, 1933]. Remarkably, this property precisely corresponds  
 251 to the SUA in kernel theory [Karoui, 2010, Jacot et al., 2020, Simon et al., 2023], which posits  
 252 that, in terms of the generalisation error statistics in kernel regression, the eigenfunctions can be  
 253 approximated by Gaussian matrices with independent entries, denoted  $\Phi_{i,k} \sim \mathcal{N}(\mu_K, \sigma_K^2)$ . This is  
 254 reminiscent of the Gaussian equivalence assumption [Schröder et al., 2023, Cui et al., 2023]; however,  
 255 to the best of our knowledge, this marks a first connection between BNNs and the SUA from kernel  
 256 theory (i.e. applied to Mercer’s eigenfunctions).

257 We can now reframe the question concerning the correctness of the SUA approximation as follows:  
 258 given  $\Lambda$  and  $\mathbf{X}$ , is it the case that all orthogonal matrices  $\Phi$  have non-zero probabilities (according  
 259 to  $\Theta$ ) to satisfy  $\mathbf{K}(\mathbf{X}, \mathbf{X}) = \Phi \Lambda \Phi^T$ ? If this condition holds, the SUA is applicable and Gaussian  
 260 eigenfunctions can be used for estimating (2) and (3). Since the prior is a Gaussian matrix, any matrix  
 261 has a non-zero probability of occurrence, thus it suffices to show that for any orthogonal  $\Phi$ , there  
 262 exists a  $\Theta$  such that  $\mathbf{K}(\mathbf{X}, \mathbf{X}) = \Phi \Lambda \Phi^T$ . In particular, it is easy to show that the SUA always holds  
 263 in the linear case: for any orthogonal  $\Phi$ , there exists  $\Theta$  such that  $\mathbf{X}^T \Theta^T \Theta \mathbf{X} = \Phi \Lambda \Phi^T$ . With a  
 264 non-linearity, the problem is less obvious: is there a  $\Theta$  such that  $\phi(\Theta, \mathbf{X})^T \phi(\Theta, \mathbf{X}) = \Phi \Lambda \Phi^T$  for  
 265 any orthogonal  $\Phi$ ? In the next section, we show that, in the linear-width limit, this criterion about  
 266  $\phi(\Theta, \cdot)$  and  $\mathbf{X}$  is a necessary and sufficient assumption for the renormalisation theory to hold.

### 267 3.3 An Extended Renormalisation Theory

268 This section only concerns the linear-width regime. In this limit, we can explicitly study the results  
 269 on our integral estimators because the corresponding limiting nonrandom spectral measure is known  
 270 El Harzli et al. [2024] (see Paragraph 3.1).

271 The renormalisation theory for linear BNNs establishes that, in the linear-width limit, the marginal  
 272 likelihood  $p(\mathbf{y} | \mathbf{X})$  follows a multivariate Gaussian with mean vector  $\mathbf{y}$  and covariance matrix

<sup>3</sup>Here,  $\mathbf{X}$  is the infinite matrix representing the linear-width limit of the training data. To be completely rigorous, we should write  $f(\mathbf{x}^*, \mathbf{X}, \mathbf{y}) = \lim_{P, N, N_0 \rightarrow \infty} f_{P, N, N_0}(\mathbf{x}^*_{N_0}, \mathbf{X}_{P, N_0}, \mathbf{y}_P)$ , but by slight abuse of notation the same notation is used for both. In practice, one would use finite (but large) objects in calculations.

273  $u_0^L \mathbf{K}_0$ , with  $\mathbf{K}_0 = \frac{1}{N_0} \mathbf{X} \mathbf{X}^T$  and  $u_0$  the renormalisation factor fulfilling the fixed-point equation  
 274  $1 - u_0 = \alpha(1 - \frac{r_0}{u_0^L})$  with  $r_0 = \frac{1}{P} \mathbf{y}^T \mathbf{K}_0^{-1} \mathbf{y}$ . This result was obtained in Li and Sompolinsky [2021]  
 275 by successively applying the saddle point method when integrating out the weights  $\Theta$ ,  $\mathbf{W}^L$ .

276 The following theorem shows that this result generalises to BNNs with nonlinear activations if and  
 277 only if the SUA is correct (i.e., it gives the correct estimate for the marginal likelihood). Here, we  
 278 exploit the characterisation of the correctness of the SUA developed as a corollary of Theorem 3.4.

279 **Theorem 3.5.** *Assume that Assumption 3.1 holds. Let  $u_{\text{NNGP}}$  fulfil the fixed-point equation  $1 -$   
 280  $u_{\text{NNGP}} = \alpha(1 - \frac{r_{\text{NNGP}}}{u_{\text{NNGP}}^L})$  with  $r_{\text{NNGP}} = \frac{1}{P} \mathbf{y}^T \mathbf{K}_{\text{NNGP}}^{-1} \mathbf{y}$ . The marginal likelihood for a nonlinear  
 281 BNN verifies  $p(\mathbf{y} | \mathbf{X}) \sim \mathcal{N}(\mathbf{y}, u_{\text{NNGP}}^L \mathbf{K}_{\text{NNGP}})$  if and only if, for given  $\Lambda$ ,  $\mathbf{X}$  and orthogonal  $\Phi$ , there  
 282 exists  $\Theta$  such that  $\phi(\Theta, \mathbf{X})^T \phi(\Theta, \mathbf{X}) = \Phi \Lambda \Phi^T$ .*

283 This result characterises the renormalisation theory in the nonlinear case and describes a continuous  
 284 transition between an accurate and a poor approximation. Specifically, if the SUA significantly  
 285 deviates (the feature map spans a small fraction of the space of orthogonal matrices) then the  
 286 equivalence (7) also deviates substantially from the correct value. Conversely, if the SUA is nearly  
 287 accurate (meaning that the feature map encompasses a large portion of the space of orthogonal  
 288 matrices) then (7) closely approximates the true marginal likelihood. Thanks to these insights, future  
 289 research on BNNs can benefit from research advances on the accuracy of the SUA [Liu et al., 2021].

### 290 3.4 Applications to the Sublinear-Width Regime

291 In this section, we consider the application of our integral estimators to the sublinear-width regime.

292 In this regime, the ratios  $\alpha$  and  $\alpha_0$  from the linear-width regime tend to infinity and hence are  
 293 no longer bounded. Here, the renormalisation theory breaks even in the linear case, because the  
 294 random matrix  $\mathbf{K}_{\Theta}^{P,N,N_0}(\mathbf{X}, \mathbf{X})$  becomes degenerate and its limiting spectral distribution is the Dirac  
 295 distribution at 0 and (6) no longer holds. A mismatch with the predictions of the renormalisation  
 296 theory has indeed been observed empirically for high values of  $\alpha$  and  $\alpha_0$  [Li and Sompolinsky, 2021],  
 297 hence the need for a new theory.

298 Remarkably, our kernel-theoretic description of BNNs (Theorem 3.4) still holds as its validity relies  
 299 only on the dot product of random feature maps  $\phi(\Theta, \cdot)$  defining a random kernel. This remains true  
 300 for the sublinear-width regime (as well as for other regimes of interest). Additionally, zero eigenvalues  
 301 in Mercer’s decomposition can be disregarded since they do not contribute to the kernel evaluation.  
 302 An alternative perspective is that, when calculating the first two moments, one takes into account the  
 303 Moore-Penrose pseudo-inverse  $\mathbf{K}_{\Theta}^{P,N,N_0}(\mathbf{X}, \mathbf{X})^{\dagger}$  of the kernel random matrix. Consequently, only  
 304 the contributions from the *strictly positive support* of the limiting spectral distribution are considered  
 305 (see Theorem 3.3). As a result, our integral estimators for the mean and variance of the predictor  
 306 remain applicable under the SUA. The only missing element in the argument is whether  $p(\Lambda | \mathbf{X}, \mathbf{x}^*)$ ,  
 307 which is now the strictly positive support of the limiting spectral distribution of  $\mathbf{K}_{\Theta}^{P,N,N_0}(\mathbf{X}, \mathbf{X})$   
 308 (rescaled to integrate to 1), also converges to a nonrandom spectral measure (in order to apply  
 309 Theorem 3.3). Although further work is needed to precisely characterise the behavior of this Mercer’s  
 310 random spectral measure, it is still possible to numerically compute the strictly positive support of  
 311 the random matrix  $\mathbf{K}_{\Theta}^{P,N,N_0}(\mathbf{X}, \mathbf{X})$  and use our integral forms (4) and (5) to estimate the predictor  
 312 of a trained BNNs in this new regime. Note however that our approach only concerns the predictor  
 313 statistics and is not derived in weight space, thus one limitation of our approach that we anticipate is  
 314 that it might be difficult to characterize (strong) feature learning from this standpoint.

## 315 4 Experiments

316 We consider a synthetic dataset generated by a multivariate Gaussian  $\mathbf{x} \sim \mathcal{N}(0, \frac{1}{N_0} \mathbf{I}_{N_0})$  to which  
 317 we apply a linear teacher and noise  $y = \beta^T \mathbf{x} + \epsilon$  with  $\epsilon \sim \mathcal{N}(0, \sigma_{\epsilon}^2)$ . We also consider a subset of  
 318 MNIST restricted to classes "0" and "1" of size  $P = 105$  and with  $N_0 = 784$  pixels per image.

319 Our first experiment verifies that our estimators coincide with the predictions of the renormalisation  
 320 theory in the linear-width limit both for a single hidden-layer network with ReLU activations and a  
 321 linear network with a hidden layer. We computed the renormalisation factors using the fixed-point

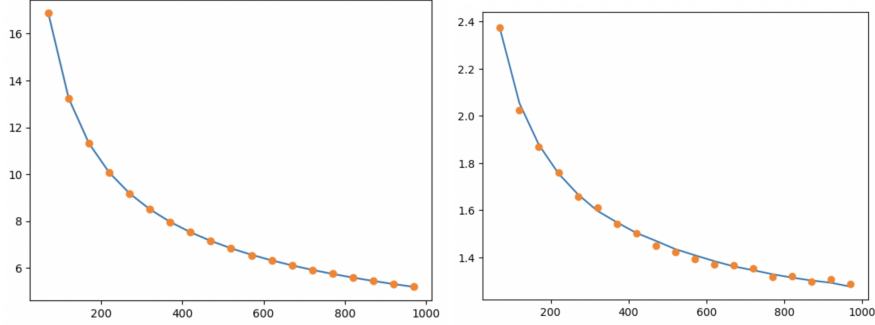


Figure 1: Comparison with Li and Sompolinsky [2021]. X-axis are indexed by the width and Y-axis by the renormalisation factor, which is computed in the linear setting on our synthetic dataset with  $N_0 = 500$  and  $P = 200$  (respectively, in the nonlinear setting on the subset of MNIST) on the left (respectively, on the right). The blue line is computed using the fixed-point equation, and the orange dots are the ratio between the result of our integral estimator (3) and the variance of Bayesian linear regression (respectively, NNGP regression) on the left (respectively, on the right). In the nonlinear case, we use a large width  $\hat{N} = 10000$  to estimate the NNGP kernel matrix for ReLU.

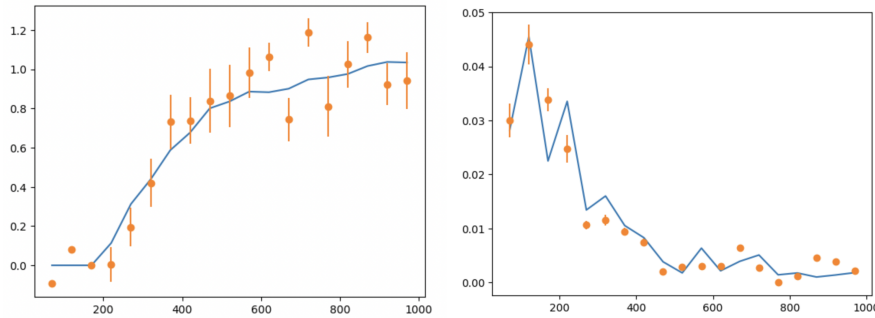


Figure 2: Mean and variance of the predictor against the width  $N$  of the single ReLU hidden-layer on our synthetic dataset with  $P = 200$  and  $N_0 = 40$ . X-axis are indexed by the width and Y-axis are indexed respectively by the mean of the predictor (on the left) and by the variance of the predictor (on the right). In both cases, the blue line is computed using the probabilistic predictions of a BNN trained with variational inference on the synthetic data, and the orange dots correspond to our estimates.

322 equation and used (2) and (3) to estimate the mean and the variance of the predictor in our approach.  
 323 To compute (2) and (3) we first computed the Marchenko-Pastur maps of the empirical spectral  
 324 distributions (of the NNGP kernel) by solving numerically the Marchenko-Pastur fixed-point equation  
 325 in the Stieltjes transform space; then, we relied on the SUA to estimate the integral forms. In a second  
 326 experiment, we simulated the regime  $P \propto N \cdot N_0$  (for which the renormalisation theory breaks) using  
 327 a small value of  $N_0$  (thus making  $\alpha_0$  high). We compared our estimators for the regime as described  
 328 in the previous section with the predictions of BNNs trained with variational inference using the  
 329 library Pyro [Bingham et al., 2019]. For the spectral distribution, we computed the strictly positive  
 330 support of the empirical spectral distributions by sampling and diagonalising the empirical kernel  
 331 matrices several times and shuffling the eigenvalues; we continued to use the SUA for eigenfunctions.

332 As shown in Figure 1, our estimates align with the renormalisation theory in the thermodynamic limit.  
 333 As shown in Figure 2, for a regime where the renormalisation theory is inaccurate, our estimators  
 334 provide reasonable matches to the actual predictions. These results suggest that our estimates are  
 335 better suited to regimes where key assumptions of the renormalisation theory do not hold.



## 336 5 Conclusion

337 In this paper, we have explored bridges between BNNs trained under interesting idealised limits and  
338 kernel theory, which enable an extension of the renormalisation theory to non-linear networks. From  
339 a practical standpoint, our theory offers a new way to estimate the prediction of BNNs with better  
340 accuracy in the sublinear-width regime. Finally, we hope that the theory developed here will motivate  
341 further research on the application of existing kernel-theoretic results in the context of BNNs.

## 342 References

343 S. Ariosto, R. Pacelli, M. Pastore, F. Ginelli, M. Gherardi, and P. Rotondo. Statistical mechanics of  
344 deep learning beyond the infinite-width limit, 2023.

345 C. T. H. Baker. *The Numerical Treatment of Integral Equations*. Oxford University Press, 1977.

346 Jean Barbier, Florent Krzakala, Nicolas Macris, Léo Miolane, and Lenka Zdeborová. Optimal errors  
347 and phase transitions in high-dimensional generalized linear models. In Sébastien Bubeck, Vianney  
348 Perchet, and Philippe Rigollet, editors, *Proceedings of the 31st Conference On Learning Theory*,  
349 volume 75 of *Proceedings of Machine Learning Research*, pages 728–731. PMLR, 06–09 Jul 2018.  
350 URL <https://proceedings.mlr.press/v75/barbier18a.html>.

351 Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis  
352 Karaletsos, Rohit Singh, Paul Szerlip, Paul Horsfall, and Noah D. Goodman. Pyro: Deep universal  
353 probabilistic programming. *Journal of Machine Learning Research*, 20(28):1–6, 2019. URL  
354 <http://jmlr.org/papers/v20/18-403.html>.

355 David M. Blei, Alp Kucukelbir, and Jon D. McAuliffe. Variational inference: A review for statisticians.  
356 *Journal of the American Statistical Association*, 112(518):859–877, apr 2017. doi: 10.1080/  
357 01621459.2017.1285773. URL <https://doi.org/10.1080%2F01621459.2017.1285773>.

358 Abdulkadir Canatar, Blake Bordelon, and Cengiz Pehlevan. Spectral bias and task-model align-  
359 ment explain generalization in kernel regression and infinitely wide neural networks. *Nature*  
360 *Communications*, 12(1), may 2021. doi: 10.1038/s41467-021-23103-1. URL <https://doi.org/10.1038%2Fs41467-021-23103-1>.

362 Xiuyan Cheng and Amit Singer. The spectrum of random inner-product kernel matrices. *Random*  
363 *Matrices: Theory and Applications*, 02(04):1350010, 2013. doi: 10.1142/S201032631350010X.  
364 URL <https://doi.org/10.1142/S201032631350010X>.

365 Hugo Cui, Florent Krzakala, and Lenka Zdeborova. Bayes-optimal learning of deep random networks  
366 of extensive-width. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt,  
367 Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on*  
368 *Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 6468–6521.  
369 PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/cui23b.html>.

370 A. Daniely, R. Frostig, and Y. Singer. Toward deeper understanding of neural networks: The power of  
371 initialization and a dual view on expressivity. *Advances in Neural Information Processing Systems*,  
372 29, 2016.

373 Ouns El Harzli, Bernardo Cuenca Grau, Guillermo Valle-Pérez, and Ard A. Louis. Double-descent  
374 curves in neural networks: A new perspective using gaussian processes. *Proceedings of the AAAI*  
375 *Conference on Artificial Intelligence*, 38(10):11856–11864, Mar. 2024. doi: 10.1609/aaai.v38i10.  
376 29071. URL <https://ojs.aaai.org/index.php/AAAI/article/view/29071>.

377 A. Engel, Otto von Guericke, and C. Van den Broeck. *Statistical mechanics of learning*. Cambridge  
378 University Press, 2012.

379 Z. Fan and Z. Wang. Spectra of the conjugate kernel and neural tangent kernel for linear-width neural  
380 networks. *Advances in Neural Information Processing Systems*, 33, 2020.

381 Zhou Fan and Andrea Montanari. The spectral norm of random inner-product kernel matrices.  
382 *Probability Theory and Related Fields*, 173:27–85, 2015.

- 383 Yarín Gal. *Uncertainty in Deep Learning*. PhD thesis, University of Cambridge, 2016.
- 384 J.S. Geronimo and T.P. Hill. Necessary and sufficient condition that the limit of stieltjes transforms is  
385 a stieltjes transform. *Journal of Approximation Theory*, 2002.
- 386 Alfred Haar. Der massbegriff in der theorie der kontinuierlichen gruppen. *Annals of Mathematics*, 34  
387 (1):147–169, 1933. ISSN 0003486X. URL <http://www.jstor.org/stable/1968346>.
- 388 Boris Hanin and Alexander Zlokapa. Bayesian interpolation with deep linear networks. *Proceedings  
389 of the National Academy of Sciences*, 120(23):e2301345120, 2023. doi: 10.1073/pnas.2301345120.  
390 URL <https://www.pnas.org/doi/abs/10.1073/pnas.2301345120>.
- 391 Arthur Jacot, Berfin Simsek, Francesco Spadaro, Clement Hongler, and Franck Gabriel. Kernel align-  
392 ment risk estimator: Risk prediction from training data. In H. Larochelle, M. Ranzato, R. Hadsell,  
393 M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33,  
394 pages 15568–15578. Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/  
395 paper\\_files/paper/2020/file/b367e525a7e574817c19ad24b7b35607-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/b367e525a7e574817c19ad24b7b35607-Paper.pdf).
- 396 Nouredine El Karoui. The spectrum of kernel random matrices. *The Annals of Statistics*, 38(1), feb  
397 2010. doi: 10.1214/08-aos648. URL <https://doi.org/10.1214/08-aos648>.
- 398 J. Lee, Y. Bahri, R. Novak, S. S. Schoenholz, J. Pennington, and J. Sohl-Dickstein. Deep neural  
399 networks as gaussian processes. *International Conference on Learning Representations*, 2018a.
- 400 J. Lee, S. Schoenholz, J. Pennington, B. Adlam, L. Xiao, R. Novak, and J. Sohl-Dickstein. Finite  
401 versus infinite neural networks: an empirical study. *Advances in Neural Information Processing  
402 Systems*, 33, 2020.
- 403 Jaehoon Lee, Jascha Sohl-dickstein, Jeffrey Pennington, Roman Novak, Sam Schoenholz, and  
404 Yasaman Bahri. Deep neural networks as gaussian processes. In *International Conference on  
405 Learning Representations*, 2018b. URL <https://openreview.net/forum?id=B1EA-M-OZ>.
- 406 Qianyi Li and Haim Sompolinsky. Statistical mechanics of deep linear neural networks: The  
407 backpropagating kernel renormalization. *Physical review X*, 11(3), 2021.
- 408 Fanghui Liu, Zhenyu Liao, and Johan Suykens. Kernel regression in high dimensions: Refined  
409 analysis beyond double descent. In Arindam Banerjee and Kenji Fukumizu, editors, *Proceedings  
410 of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of  
411 *Proceedings of Machine Learning Research*, pages 649–657. PMLR, 13–15 Apr 2021. URL  
412 <https://proceedings.mlr.press/v130/liu21b.html>.
- 413 Giacomo Livan, Marcel Novaes, and Pierpaolo Vivo. *Introduction to Random Matrices*. Springer  
414 International Publishing, 2018. doi: 10.1007/978-3-319-70885-0. URL [https://doi.org/10.  
415 1007/978-3-319-70885-0](https://doi.org/10.1007/978-3-319-70885-0).
- 416 Yue M. Lu and Horng-Tzer Yau. An equivalence principle for the spectrum of random inner-product  
417 kernel matrices with polynomial scalings, 2023.
- 418 Martin Magris and Alexandros Iosifidis. Bayesian learning for neural networks: an algorithmic  
419 survey. *Artificial Intelligence Survey*, 2023. doi: <https://doi.org/10.1007/s10462-023-10443-1>.
- 420 Antoine Maillard, Emanuele Troiani, Simon Martin, Florent Krzakala, and Lenka Zdeborová. Bayes-  
421 optimal learning of an extensive-width neural network from quadratically many samples, 2024.  
422 URL <https://arxiv.org/abs/2408.03733>.
- 423 V.A. Marchenko and L.A. Pastur. Distribution of eigenvalues for some sets of random matrices.  
424 *Matematicheskii Sbornik*, 72, 1967.
- 425 A. G. de G. Matthews, M. Rowland, J. Hron, R. E. Turner, and Z. Ghahramani. Gaussian process  
426 behaviour in wide deep neural networks. *arXiv preprint arXiv:1804.11271*, 2018.
- 427 J. Mercer. Functions of positive and negative type, and their connection the theory of integral  
428 equations. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers  
429 of a Mathematical or Physical Character*, page 209415–446, 1909. doi: [http://doi.org/10.1098/  
430 rsta.1909.0016](http://doi.org/10.1098/rsta.1909.0016).

- 431 James A. Mingo and Roland Speicher. *Free Probability and Random Matrices*. Springer, 2017. doi:  
432 <https://doi.org/10.1007/978-1-4939-6942-5>.
- 433 H.Q. Minh, P. Niyogi, and Y. Yao. Mercer’s theorem, feature maps, and smoothing. In *Learning*  
434 *Theory*, pages 154–168. Springer, 2006.
- 435 E.H. Moore. On the reciprocal of the general algebraic matrix. *Bulletin of the American Mathematical*  
436 *Society*, 1920.
- 437 Radford M Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business  
438 Media, 2012.
- 439 Carl Rasmussen. A practical monte carlo implementation of bayesian learning. In D. Touretzky, M.C.  
440 Mozer, and M. Hasselmo, editors, *Advances in Neural Information Processing Systems*, volume 8.  
441 MIT Press, 1995. URL [https://proceedings.neurips.cc/paper\\_files/paper/1995/  
442 file/84d2004bf28a2095230e8e14993d398d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/1995/file/84d2004bf28a2095230e8e14993d398d-Paper.pdf).
- 443 C.E. Rasmussen and C.K.I. Williams. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- 444 Saburo Saitoh and Yoshihiro Sawano. Springer, 2016. doi: [https://doi.org/10.1007/  
445 978-981-10-0530-5](https://doi.org/10.1007/978-981-10-0530-5).
- 446 Dominik Schröder, Hugo Cui, Daniil Dmitriev, and Bruno Loureiro. Deterministic equivalent  
447 and error universality of deep random features learning. In Andreas Krause, Emma Brunskill,  
448 Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of*  
449 *the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine*  
450 *Learning Research*, pages 30285–30320. PMLR, 23–29 Jul 2023. URL [https://proceedings.  
451 mlr.press/v202/schroder23a.html](https://proceedings.mlr.press/v202/schroder23a.html).
- 452 H. S. Seung and H. Sompolinsky. Statistical mechanics of learning from examples. *Physical Review*  
453 *A*, 45, 1992.
- 454 James B Simon, Madeline Dickens, Dhruva Karkada, and Michael Deweese. The eigenlearn-  
455 ing framework: A conservation law perspective on kernel ridge regression and wide neural  
456 networks. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL  
457 <https://openreview.net/forum?id=FDbQGCAViI>.
- 458 Alexander van Meegen and Haim Sompolinsky. Coding schemes in neural networks learning  
459 classification tasks, 2024. URL <https://arxiv.org/abs/2406.16689>.
- 460 E. Wigner. Characteristic vectors of bordered matrices with infinite dimensions. *Annals of Mathemat-*  
461 *ics*, 62, 1955.
- 462 Jacob A. Zavatone-Veth, William L. Tong, and Cengiz Pehlevan. Contrasting random and learned  
463 features in deep bayesian linear regression. *Phys. Rev. E*, 105:064118, Jun 2022. doi: 10.  
464 1103/PhysRevE.105.064118. URL [https://link.aps.org/doi/10.1103/PhysRevE.105.  
465 064118](https://link.aps.org/doi/10.1103/PhysRevE.105.064118).

## 466 A Proof of Theorem 3.3

467 In the linear-width (respectively, the sub-linear width limit), the positive semi-definiteness of any  
 468 matrix extracted from  $K_{\Theta}^{N, N_0}$  and  $p_{\mathbf{X}}$  is maintained (the limit of a positive sequence remains positive),  
 469 and this suffices to characterise the kernel property over a compact subset of an infinite-dimensional  
 470 space [Saitoh and Sawano, 2016]. Thus, there is a random kernel  $K_{\Theta}^{\alpha, \alpha_0}$  (respectively,  $K_{\Theta}^{\gamma}$ ) defined  
 471 over  $\mathbb{R}^N$  which characterises the convergence in distribution of  $\mathbf{K}_{\Theta}^{P, N, N_0}(\mathbf{X}, \mathbf{X})$ . As per Mercer's  
 472 theorem,  $K_{\Theta}^{\alpha, \alpha_0}$  (respectively,  $K_{\Theta}^{\gamma}$ ) also defines a random spectral measure  $\rho_{\Theta}^{\alpha, \alpha_0}$  (respectively,  $\rho_{\Theta}^{\gamma}$ )  
 473 associated with its Mercer's eigenvalues. By Baker's result [Baker, 1977] stating the convergence of  
 474 eigenvalues in a kernel matrix to the Mercer eigenvalues of the respective kernel, it follows that  $\rho_{\Theta}^{\alpha, \alpha_0}$   
 475 (respectively,  $\rho_{\Theta}^{\gamma}$ ) is the limiting spectral distribution of the random matrices  $\mathbf{K}_{\Theta}^{P, N, N_0}(\mathbf{X}, \mathbf{X})$  in the  
 476 linear-width limit (respectively, the sublinear-width regime). By assumption, this spectral measure  
 477 (respectively, the strictly positive support of this spectral measure) is nonrandom  $\rho_{\Theta}^{\alpha, \alpha_0} = \rho^{\alpha, \alpha_0}$   
 478 (respectively,  $\rho_{\Theta}^{\gamma} = \rho^{\gamma}$ ). Thus, we can reformulate the empirical kernel matrix corresponding to the  
 479 random kernel  $K_{\Theta}^{\alpha, \alpha_0}$  (respectively,  $K_{\Theta}^{\gamma}$ ) as  $\Phi \Lambda \Phi^T$ , where  $\lambda_k$  are drawn independently according  
 480 to  $\rho^{\alpha, \alpha_0}$  (respectively,  $\rho^{\gamma}$ ). Since the spectral measure no longer depends on  $\Theta$ , the eigenvalues  
 481 can be sampled independently from the eigenfunctions. It follows that  $\mathbf{K}_{\Theta}^{P, N, N_0}(\mathbf{X}, \mathbf{X})$  and  $\Phi \Lambda \Phi^T$   
 482 converge to the same distribution over  $\mathbb{R}^{N \times N}$ .

## 483 B Proof of Theorem 3.4

484 We calculate the conditional expectation  $\langle f \rangle(\mathbf{x}^*, \mathbf{X}, \mathbf{y}, \Theta)$  and variance  $\langle \delta f \rangle(\mathbf{x}^*, \mathbf{X}, \mathbf{y}, \Theta)$  of the  
 485 predictor by marginalising over the readout weights  $\mathbf{W}^L$ :

$$\langle f \rangle(\mathbf{x}^*, \mathbf{X}, \mathbf{y}, \Theta) = \int \mathbf{W}^{L T} \phi(\Theta, \mathbf{X}) p(\mathbf{W}^L | \mathbf{X}, \mathbf{y}, \Theta) d\mathbf{W}^L$$

486

$$\langle \delta f \rangle(\mathbf{x}^*, \mathbf{X}, \mathbf{y}, \Theta) = \int \left[ \mathbf{W}^{L T} \phi(\Theta, \mathbf{X}) \right]^2 p(\mathbf{W}^L | \mathbf{X}, \mathbf{y}, \Theta) d\mathbf{W}^L - [\langle f \rangle(\mathbf{x}^*, \mathbf{X}, \mathbf{y}, \Theta)]^2$$

487 where  $p(\mathbf{W}^L | \mathbf{X}, \mathbf{y}, \Theta)$  can be expressed by Bayes rule using Gaussian likelihoods. The result can be  
 488 expressed analytically and yields the same prediction as GP regression with prior  $\mathcal{GP}(0, K_{\Theta}^{N, N_0})$ :

$$\langle f \rangle(\mathbf{x}^*, \mathbf{X}, \mathbf{y}, \Theta) = [\mathbf{k}_{\Theta}^{P, N, N_0}(\mathbf{x}^*, \mathbf{X})]^T [\mathbf{K}_{\Theta}^{P, N, N_0}(\mathbf{X}, \mathbf{X})]^{-1} \mathbf{y}$$

489

$$\langle \delta f \rangle(\mathbf{x}^*, \mathbf{X}, \mathbf{y}, \Theta) = K_{\Theta}^{P, N, N_0}(\mathbf{x}^*, \mathbf{x}^*) - [\mathbf{k}_{\Theta}^{P, N, N_0}(\mathbf{x}^*, \mathbf{X})]^T [\mathbf{K}_{\Theta}^{P, N, N_0}(\mathbf{X}, \mathbf{X})]^{-1} \mathbf{k}_{\Theta}^{P, N, N_0}(\mathbf{x}^*, \mathbf{X}).$$

490 To marginalise over  $\Theta \sim p(\Theta | \mathbf{X}, \mathbf{y})$ , we perform the change of variables  $\Theta \mapsto (\Phi^*, \Phi, \Lambda)$ , relying on  
 491 the fact that all quantities of interest involving  $\Theta$  can be expressed in the limit solely using eigenvalues  
 492 and eigenfunctions, namely  $\mathbf{K}(\mathbf{X}, \mathbf{X}) = \Phi \Lambda \Phi^T$ ,  $\mathbf{k}(\mathbf{x}^*, \mathbf{X}) = \Phi \Lambda \Phi^*$ , and  $K(\mathbf{x}^*, \mathbf{x}^*) = \Phi^{* T} \Lambda \Phi^*$ .  
 493 Since  $\Phi \in \mathbb{R}^{P \times M}$  has orthogonal rows,  $\Phi^\dagger = \Phi^T (\Phi \Phi^T)^{-1}$ , and  $\Phi^{T \dagger} = (\Phi \Phi^T)^{-1} \Phi$ . This  
 494 allows us to express the mean and variance of the predictor as follows:

$$\langle f \rangle = \int \left( \Phi^{* T} \Lambda \Phi^T \Phi^{T \dagger} \Lambda^{-1} \Phi^\dagger \mathbf{y} \right) \cdot p(\Lambda, \Phi, \Phi^* | \mathbf{X}, \mathbf{x}^*) \cdot \frac{p(\mathbf{y} | \Lambda, \Phi, \Phi^*, \mathbf{X}, \mathbf{x}^*)}{p(\mathbf{y} | \mathbf{X})} d\Lambda d\Phi d\Phi^* \quad (4)$$

495

$$\langle \delta f \rangle = \int \left( \Phi^{* T} \Lambda \Phi^* - \Phi^{* T} \Lambda \Phi^T \Phi^{T \dagger} \Lambda^{-1} \Phi^\dagger \Phi \Lambda \Phi^* \right) \cdot p(\Lambda, \Phi, \Phi^* | \mathbf{X}, \mathbf{x}^*) \frac{p(\mathbf{y} | \Lambda, \Phi, \Phi^*, \mathbf{X}, \mathbf{x}^*)}{p(\mathbf{y} | \mathbf{X})} d\Lambda d\Phi d\Phi^* \quad (5)$$

496 Furthermore, it holds that  $p(\Lambda, \Phi, \Phi^* | \mathbf{X}, \mathbf{x}^*) = p(\Lambda | \mathbf{X}, \mathbf{x}^*)p(\Phi, \Phi^* | \Lambda, \mathbf{X}, \mathbf{x}^*)$  if  $p(\Lambda | \mathbf{X}, \mathbf{x}^*) \neq 0$   
 497 and also  $p(\mathbf{y} | \Lambda, \Phi, \Phi^*, \mathbf{X}, \mathbf{x}^*) = \frac{p(\mathbf{y}, \Phi, \Phi^* | \Lambda, \mathbf{X}, \mathbf{x}^*)}{p(\Phi, \Phi^* | \Lambda, \mathbf{X}, \mathbf{x}^*)}$  if  $p(\Phi, \Phi^* | \Lambda, \mathbf{X}, \mathbf{x}^*) \neq 0$ , which yields:

$$\langle f \rangle = \int \left( \Phi^{*T} \Lambda \Phi^T \Phi^{T\dagger} \Lambda^{-1} \Phi^\dagger \mathbf{y} \right) \cdot p(\Lambda | \mathbf{X}, \mathbf{x}^*) \frac{p(\mathbf{y}, \Phi, \Phi^* | \Lambda, \mathbf{X}, \mathbf{x}^*)}{p(\mathbf{y} | \mathbf{X})} d\Lambda d\Phi d\Phi^*$$

498

$$\langle \delta f \rangle = \int \left( \Phi^{*T} \Lambda \Phi^* - \Phi^{*T} \Lambda \Phi^T \Phi^{T\dagger} \Lambda^{-1} \Phi^\dagger \Phi \Lambda \Phi^* \right) \cdot p(\Lambda | \mathbf{X}, \mathbf{x}^*) \frac{p(\mathbf{y}, \Phi, \Phi^* | \Lambda, \mathbf{X}, \mathbf{x}^*)}{p(\mathbf{y} | \mathbf{X})} d\Lambda d\Phi d\Phi^*$$

499 where the integral over  $\Lambda$  is restricted to segments where  $p(\Lambda | \mathbf{X}, \mathbf{x}^*) \neq 0$  and the integrals over  $\Phi$   
 500 and  $\Phi^*$  are restricted to where  $p(\Phi, \Phi^* | \Lambda, \mathbf{X}, \mathbf{x}^*) \neq 0$ . We obtain equations (2) and (3) by replacing  
 501  $d\Phi$  and  $d\Phi^*$  by standard Gaussian matrix measures and the density of  $\Lambda$  by the spectral measure.

### 502 C Proof of Theorem 3.5

503 In the linear case, the true NNGP kernel is simply  $K_{\text{NNGP}}(\mathbf{x}, \mathbf{x}') = \frac{1}{N_0} \mathbf{x}^T \mathbf{x}'$ ; hence,  $\rho_{\text{NNGP}}^{\alpha_0}$   
 504 is the limiting spectral distribution of the kernel random matrix  $\mathbf{K}_0$ , which we denote  $\rho_0$ . The  
 505 renormalisation theory of linear networks thus implies that:

$$\int p(\mathbf{y}, \Phi | \Lambda, \mathbf{X}) d(\rho_{MP}^\alpha \boxtimes^L \rho_0)(\Lambda) \mathcal{D}\Phi \sim \mathcal{N}(\mathbf{y}, u_0^L \mathbf{K}_0) \quad (6)$$

506 This identity is exact in the linear-width limit and holds in general without assumption on  $\mathbf{X}, \mathbf{y}$ , as  
 507 long as the integral  $\mathcal{D}\Phi$  is uniform on the space of orthogonal matrices.

508 Assume the SUA holds in the nonlinear case. Thus, we can express the marginal likelihood as  
 509  $p(\mathbf{y} | \mathbf{X}) = \int p(\mathbf{y}, \Phi | \Lambda, \mathbf{X}) d(\rho_{MP}^\alpha \boxtimes^L \rho_{\text{NNGP}}^{\alpha_0})(\Lambda) \mathcal{D}\Phi$ . Furthermore, we can freely interchange  
 510 the role of  $\mathbf{K}_0$  and  $\mathbf{K}_{\text{NNGP}}(\mathbf{X}, \mathbf{X})$  in (6). Indeed, it suffices to consider the linear case and a new  
 511 training dataset  $\tilde{\mathbf{X}}$  which exhibits the same covariance structure  $\frac{1}{N_0} \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T$  as that of  $\mathbf{K}_{\text{NNGP}}(\mathbf{X}, \mathbf{X})$ .  
 512 As a result, a similar equation to (6) applies to nonlinear networks by replacing the linear kernel  
 513  $(\mathbf{x}, \mathbf{x}') \mapsto \frac{1}{N_0} \mathbf{x}^T \mathbf{x}'$  with the true NNGP kernel in the equations, provided that the SUA holds:

$$\int p(\mathbf{y}, \Phi | \Lambda, \mathbf{X}) d(\rho_{MP}^\alpha \boxtimes^L \rho_{\text{NNGP}}^{\alpha_0})(\Lambda) \mathcal{D}\Phi \sim \mathcal{N}(\mathbf{y}, u_{\text{NNGP}}^L \mathbf{K}_{\text{NNGP}}) \quad (7)$$

514 Conversely, if the SUA does not hold, the integral with respect to  $\Phi$  does not span the space of  
 515 orthogonal matrices, the identity (7) is no longer exact (all integrands are strictly positive), nor is the  
 516 renormalisation. Thus, the SUA is necessary and sufficient for the renormalisation to hold.