

Bayesian Interpolation with Deep Linear Networks

Boris Hanin*, Alexander Zlokapa†

* Princeton ORFE

† MIT Center for Theoretical Physics, Google Quantum AI

May 16, 2023

Abstract

Characterizing how neural network depth, width, and dataset size jointly impact model quality is a central problem in deep learning theory. We give here a complete solution in the special case of linear networks with output dimension one trained using zero noise Bayesian inference with Gaussian weight priors and mean squared error as a negative log-likelihood. For any training dataset, network depth, and hidden layer widths, we find non-asymptotic expressions for the predictive posterior and Bayesian model evidence in terms of Meijer-G functions, a class of meromorphic special functions of a single complex variable. Through novel asymptotic expansions of these Meijer-G functions, a rich new picture of the joint role of depth, width, and dataset size emerges. We show that linear networks make provably optimal predictions at infinite depth: the posterior of infinitely deep linear networks with data-agnostic priors is the same as that of shallow networks with evidence-maximizing data-dependent priors. This yields a principled reason to prefer deeper networks when priors are forced to be data-agnostic. Moreover, we show that with data-agnostic priors, Bayesian model evidence in wide linear networks is maximized at infinite depth, elucidating the salutary role of increased depth for model selection. Underpinning our results is a novel emergent notion of effective depth, given by the number of hidden layers times the number of data points divided by the network width; this determines the structure of the posterior in the large-data limit.

1 Introduction

1.1 Background

A central aim of deep learning theory is to understand the properties of overparameterized networks trained to fit large datasets. Key questions include: how do learned networks use training data to make predictions on test points? Which neural network architectures lead to more parsimonious models? What are the joint scaling laws connecting the quality of learned models to the number of training data points, network depth and network width [GJS⁺20, KMH⁺20, BDK⁺21, RBC⁺21]?

The present article gives the first exact answers to such questions for a class of neural networks in which one can simultaneously vary input dimension, number of training data

*BH is supported by NSF grants DMS-2143754, DMS-1855684, and DMS-2133806

†AZ is supported by the Hertz Foundation, and by the Department of Defense through the National Defense Science and Engineering Graduate Fellowship Program

points, network width, and network depth. This is significant because the limits where these four structural parameters tend to infinity do not commute, causing all prior work to miss important aspects of how they jointly influence learning. Our results pertain specifically to *deep linear networks*

$$f(x) = W^{(L+1)} \dots W^{(1)} x \tag{1.1}$$

with input dimension N_0 , L hidden layers of widths N_ℓ , and output dimension $N_{L+1} = 1$. As a form of learning we take zero noise Bayesian interpolation starting from a Gaussian prior on network weights $W^{(\ell)} \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$ and empirical mean squared error over a training dataset of P examples as the negative log-likelihood. Deep linear networks, while linear in x , are non-linear in their parameters and have been extensively studied as models for learning with neural networks using both gradient-based methods [ACGH19, ACH18, SMG14, SSL22, Kaw16] and Bayesian inference [LS21, ZVP21, ZVTP22].

Since we are considering an output dimension of 1, we may write $f(x) = \theta^T x$ for a vector $\theta \in \mathbb{R}^{N_0}$. What differentiates our work from a classical Bayesian analysis of Gaussian linear regression is that as soon as $L \geq 1$ the components of θ are correlated and non-Gaussian under the prior. Predictions $f(x)$ on inputs x orthogonal to inputs from the training data therefore differ under the prior and posterior. Specifically, as shown in Figure 1, we may decompose $\theta = \theta_{\parallel} + \theta_{\perp}$ into its projections onto directions spanned by the training data and their complement. By our Bayesian construction, θ_{\parallel} is responsible for fitting the training data. Due to the correlations under the prior between θ_{\parallel} and θ_{\perp} , however, information from the training data will influence the posterior distribution of θ_{\perp} . It is precisely this data-dependent extrapolation displayed by deep linear networks with $L \geq 1$ that differentiates them from the linear models obtained by taking $L = 0$.

1.2 Relation to Prior Work

To put our work in context, we briefly summarize prior approaches to understanding learning with neural networks. The neural tangent kernel (NTK) and other kernel-based models [DLT⁺18, DZPS19, AZLL19, JGH18, LZB22] reduce neural networks to linear models. Training by gradient descent or Bayesian inference consequently does not affect predictions $f(x_{\perp})$ of test inputs orthogonal to the training dataset. Moreover, the NTK regime only considers the limit of infinite width with finite depth and dataset size. More recent work [Han18, HN20b, HR18, HN20a, LNR22, RYH22, Yai20] shows that feature learning emerges when taking depth and width to infinity simultaneously with the effective prior depth $\lambda_{\text{prior}} = L/N$ (see (2.7)) determining both the behavior under both gradient descent and Bayesian inference. Still, such analyses are restricted to finite dataset size P (or more precisely dataset sizes that are much smaller than both network depth and width).

Phenomena such as double descent [BHMM19, BHX20] and benign overfitting [BLLT20] appear in linear models when taking width, dataset size, and dataset dimension to infinity simultaneously [HMRT22, MZ22, ALP22, AP20, ASS20, MM19, MMM21]. However, as mentioned previously, these linear models do not learn to make data-adaptive predictions in directions not already present in the training data. Moreover, such approaches are restricted to studying fixed depth. Other approaches consider neural networks in the mean-field limit [MMN18, RVE18, CB18, SS20, SS21, YS21, PN21, YH21]. In this regime, networks do make data-adaptive predictions $f(x_{\perp})$. However, mean-field limits have only been considered at

fixed depth and recast optimization in terms of complex non-linear evolution equations, whose dependence on the training data is typically difficult to access.

The literature most directly related to the present article studies Bayesian inference with either non-linear [LBN⁺18, APP⁺22, HNPSD22, CKZ23, NR21, SR21] or linear networks [ZVP21, NBR⁺21, ZVTP22, LS21]. These works consider only the regime where depth is either fixed or much smaller than both width and size of the training dataset. We find, in contrast, that the full role of depth in model selection and extrapolation can only be understood in the regime where depth, width, and dataset size are simultaneously large. Finally, our results are the first to characterize the behavior of deep neural networks at any joint scaling of depth, width, dataset size, and dataset dimension. In this sense, they can be viewed as giving exact expressions for the predictive posterior over deep Gaussian processes [DL13] with Euclidean covariance in every layer.

1.3 Overview of Results

Our first result, Theorem 3.1 below, gives exact non-asymptotic formulas for both the predictive posterior (i.e., the distribution of $f(x)$ jointly for all inputs x when $W^{(\ell)}$ are drawn from the posterior) and the Bayesian model evidence in terms of a class of meromorphic special functions in one complex variable called Meijer-G functions [Mei36]. These results hold for arbitrary training datasets, input dimensions, hidden layer widths, and network depth. They represent a novel enlargement of the class of priors over θ for which posteriors can be computed in closed form. In particular, they show that zero noise Bayesian inference is exactly solvable for deep Gaussian processes [DL13] with Euclidean covariances.

To glean insights from the non-asymptotic results in Theorem 3.1, we provide in Theorem 3.2 new asymptotic expansions of Meijer-G functions that allow us to compute expressions for the Bayesian model evidence and the predictive posterior under essentially any joint scalings of P, N_ℓ, L in the large-data limit where $P \rightarrow \infty$ with $P/N_0 \rightarrow \alpha_0 \in (0, 1)$. To understand the role of depth, we consider regimes in which L either stays finite or grows together with P, N_ℓ .

We focus in this article on zero-noise, or interpolating, posteriors that fit the training data exactly (see (2.4) and the discussion in Optimal Extrapolation). For parametric models interpolation often causes overfitting at large sample sizes. An important empirical [ZBH⁺17] and theoretical [BLLT20, BHMM19, HMRT22] observation, however, is that in many non-parametric overparameterized models, such as the deep linear neural networks we study, this does not occur. Our results therefore give new information about the joint effects of depth, width, and sample size on the nature of interpolating models.

What emerges from our analysis is a rich new picture of the role of depth in linear networks in determining the nature of extrapolation and Bayesian model selection, given by maximizing the Bayesian model evidence, i.e., the likelihood of the training data under the posterior (cf §4, §5 in [Mac92]). We present here an informal explanation of our main results, starting with Theorem 3.3 which implies the following:

Takeaway: Evidence-maximizing priors give the same Gaussian predictive posterior for any architecture in the large-data limit.

This distinguished posterior represents, from a Bayesian point of view, a notion of optimal extrapolation. Indeed, because $f(x)$ is linear in x , it is natural to decompose

$$x = x_{\parallel} + x_{\perp}, \quad x \in \mathbb{R}^{N_0},$$

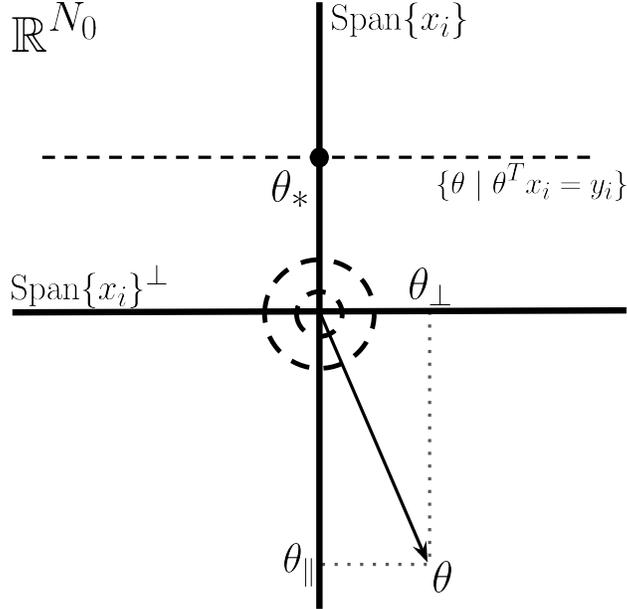


Figure 1: The input space \mathbb{R}^{N_0} is decomposed into directions $\text{Span}\{x_i\}$ spanned by inputs from the training data and its orthogonal complement. The minimal ℓ_2 norm interpolant is the intersection between the space of interpolants (i.e., the dashed line showing all θ satisfying $\theta^T x_i = y_i$ for all i) and $\text{Span}\{x_i\}$. When generating predictions $f(x) = \theta^T x$, the parameter vector θ can be decomposed into its projections θ_\parallel and θ_\perp onto $\text{Span}\{x_i\}$ and $\text{Span}\{x_i\}^\perp$ respectively. When fully trained, θ_\parallel will equal θ_* . Circles centered at the origin represent equi-probable lines for the prior over θ , which is always radial but non-Gaussian for any $L \geq 1$.

where x_\parallel is the projection of x onto directions spanned by inputs from the training data, and x_\perp is the projection of x onto the orthogonal complement (like the decomposition of θ in Figure 1).

Predictions under the evidence-maximizing posterior at test input x are Gaussian and take the form

$$f(x) \sim \mathcal{N}\left(\theta_*^T x_\parallel, \frac{\|\theta_*\|^2}{\alpha_0} \|x_\perp\|^2\right), \quad (1.2)$$

where θ_* is the minimal norm interpolant of the training data (cf. Figure 1 and (3.5)). As we explain in Optimal Extrapolation below, the deterministic mean $\theta_*^T x_\parallel$ appears because, by construction, our posteriors are concentrated on θ that interpolate the training data (see (2.4)) and thus we must have $\theta_\parallel = \theta_*$. The scalar $\|\theta_*\|^2/\alpha_0$, in contrast, sets a data-dependent variance for making predictions in directions orthogonal to inputs in the training data. This particular value for the variance is the most likely given that our posteriors are necessarily isotropic in directions orthogonal to the training data. Moreover, the approximate normality of the evidence-maximizing posterior at large sample sizes is due to the fact that the coordinates of a spherically symmetric random vector in high dimensions are approximately Gaussian (see the discussion just below (2.10)).

In general, a linear network that maximizes Bayesian evidence, and hence produces the

posterior (1.2), may require the prior distribution over network weights to be data-dependent. In machine learning contexts, we hope instead to find optimal but *data-agnostic* priors. Theorems 3.6 and 3.8 (in combination with Theorem 3.3) show this is possible in linear networks with large depth. Informally they give:

Takeaway: Wide linear networks (at comparable depth, width, number of data points) with *data-agnostic* priors give the same predictive posterior as shallow networks with optimal *data-dependent* priors.

This result highlights the remarkable role of depth in shaping the posterior over predictions in directions of feature space orthogonal to those present in the training data. This can only happen in non-linear models such as deep linear networks. Quantifying how large network depth must be to ensure optimal extrapolation is explained in Theorem 3.8, which provides universal scaling laws for Bayesian posteriors in terms of a single parameter that couples depth, width, and dataset size. Informally, we have the following:

Takeaway: Consider linear networks in the regime $1 \ll \text{depth}$, $\text{dataset size} \ll \text{width}$. With data-agnostic priors, the posterior depends only on the effective posterior depth

$$\lambda_{\text{post}} := \frac{(\text{network depth}) \times (\text{dataset size})}{\text{network width}}.$$

As $\lambda_{\text{post}} \rightarrow \infty$, evidence grows and the posterior converges to the evidence-maximizing posterior (1.2).

Since λ_{post} determines both the bias and the variance of the posterior, the preceding takeaway can be viewed as a scaling law relating depth, width, and training set size [RBC⁺21, GJS⁺20, KMH⁺20, BDK⁺21]. In particular, it shows that for large linear networks it is λ_{post} , rather than depth, width, and dataset size separately, that determines the quality of the learned model.

As the preceding statement suggests, at least with data-agnostic priors and wide linear networks, maximizing Bayesian evidence requires large depth, as measured by λ_{post} . Moreover, evidence maximization is not possible at finite λ_{post} . The final result we present (see Theorem 3.6) concerns maximization of Bayesian evidence — a principled method of comparing different architectures [Mac92] — and is summarized as follows:

Takeaway: With data-agnostic priors and width that is proportional to dataset size, Bayesian evidence is maximal in networks with depth equal to a data-dependent constant times width.

Mis-specification of this constant only results in an order one decrease in evidence and does not affect the posterior. In comparison, a network with smaller depth has exponentially smaller evidence and a suboptimal posterior. The preceding takeaways give perhaps the first principled Bayesian justification for preferring neural networks with large depth, albeit in the restricted setting of linear networks.

We then state our first result, Theorem 3.1, in §3.1. This gives an exact non-asymptotic formula for the characteristic function of predictive posterior and the model evidence in terms of Meijer-G functions. We complement this in §3.2 by giving in Theorem 3.2 asymptotic expansions for Meijer-G functions. The next collection of results, provided in §3.3, details the

model evidence and posterior as number of datapoints, input dimension, width, and depth tend to infinity in various regimes. The results in §3.3.1 pertain specifically to the analysis of networks with a finite number of hidden layers in the regime where number of training datapoints, input dimension, and width tend to infinity. The main result is Theorem 3.4. In contrast, §3.3.2 and §3.3.3 consider regimes in which depth also tends to infinity. The main result is Theorem 3.8. Finally, §3.4 contains simple corollaries connecting our results to scaling laws for the generalization error (Theorem 3.10) and double descent (Theorem 3.11).

2 Preliminaries

2.1 Setup

We fix a training set with P examples

$$X_{N_0} = (x_{1,N_0}, \dots, x_{P,N_0}) \in \mathbb{R}^{N_0 \times P}, \quad Y_{N_0} = (y_1, \dots, y_P) \in \mathbb{R}^P.$$

We will assume X_{N_0} has full rank. Since we study zero noise posteriors supported on the models that minimize the likelihood (2.2), we assume also that $1 \leq P \leq N_0$. Otherwise, the set of minima of the likelihood consists of a single θ and our posteriors would have zero variance. A key role in our results will be played by the minimal ℓ_2 -norm solution to ordinary linear least squares regression of Y_{N_0} onto X_{N_0} :

$$\theta_{*,N_0} := \arg \min_{\theta \in \mathbb{R}^{N_0}} \|\theta\|_2 \quad \text{s.t.} \quad \theta^T X_{N_0} = Y_{N_0}. \quad (2.1)$$

Further, we fix $N_1, \dots, N_L \geq 1$ and consider fitting the training data (X_{N_0}, Y_{N_0}) by a linear model

$$f(x) = \theta^T x \in \mathbb{R}, \quad \theta, x \in \mathbb{R}^{N_0}$$

equipped with quadratic negative log-likelihood

$$\mathcal{L}(\theta \mid X_{N_0}, Y_{N_0}) := \frac{1}{2} \|\theta^T X_{N_0} - Y_{N_0}\|_2^2 \quad (2.2)$$

and a *deep linear prior*

$$\theta \sim \mathbb{P}_{\text{prior}} \iff \theta = W^{(L+1)} \dots W^{(1)} \quad (2.3)$$

in which

$$W^{(\ell)} \in \mathbb{R}^{N_\ell \times N_{\ell-1}}, \quad W_{ij}^{(\ell)} \sim \mathcal{N}\left(0, \frac{\sigma^2}{N_{\ell-1}}\right) \quad \text{independent,}$$

where $\sigma > 0$. Our goal is to study the posterior distribution over the set of $\theta \in \mathbb{R}^{N_0}$ that exactly fit the training data. Explicitly, writing $d\mathbb{P}_{\text{prior}}(\theta)$ for the prior density, we study zero noise posteriors

$$\begin{aligned} & d\mathbb{P}_{\text{post}}(\theta \mid L, N_\ell, \sigma^2, X_{N_0}, Y_{N_0}) \\ & := \lim_{\beta \rightarrow \infty} \frac{d\mathbb{P}_{\text{prior}}(\theta \mid N_0, L, N_\ell, \sigma^2) \exp[-\beta \mathcal{L}(\theta \mid X_{N_0}, Y_{N_0})]}{Z_\beta(X_{N_0}, Y_{N_0} \mid L, N_\ell, \sigma^2)}. \end{aligned} \quad (2.4)$$

Writing $\mathbb{E}_{\text{post}}[\cdot]$ for the expectation with respect to the posterior (2.4), we describe the posterior by giving exact formulas for its characteristic function

$$\mathbb{E}_{\text{post}}[\exp\{-i\mathbf{t} \cdot \theta\}] = \frac{Z_{\infty}(\mathbf{t})}{Z_{\infty}(\mathbf{0})}, \quad \mathbf{t}, \theta \in \mathbb{R}^{N_0}, \quad (2.5)$$

where $Z_{\infty}(\mathbf{t}) = Z_{\infty}(\mathbf{t} \mid L, N_{\ell}, \sigma^2, X_{N_0}, Y_{N_0})$ is the zero-temperature partition function given by taking $\beta \rightarrow \infty$ in

$$Z_{\beta}(\mathbf{t}) := A_{\beta} \int \exp \left[-\frac{\beta}{2} \left\| Y - \prod_{\ell=1}^{L+1} W^{(\ell)} X \right\|_2^2 - i\theta \cdot \mathbf{t} - \sum_{\ell=1}^{L+1} \frac{N_{\ell-1}}{2\sigma^2} \left\| W^{(\ell)} \right\|_F^2 \right] \prod_{\ell=1}^{L+1} dW^{(\ell)}. \quad (2.6)$$

The normalizing constant A_{β} cancels in the ratio (2.4) and in any computations involving maximizing ratios of model evidence (see §4.3).

The denominator $Z_{\infty}(\mathbf{0})$ is often called the Bayesian model evidence and represents the probability of the data (X_{N_0}, Y_{N_0}) given the model (i.e., the depth L , layer widths N_1, \dots, N_L and prior variance σ^2). As detailed in §4, §5 of [Mac92], maximizing the Bayesian model evidence is therefore equivalent to maximum likelihood estimation over the space of models and gives a principled way to select among different models, all of which interpolate the training data. Before stating our technical results we briefly explain how to compute effective depth and how to reason about optimal extrapolation in linear networks.

2.2 Effective Depth

The number of layers L does not provide a useful measure of complexity for the prior distribution over network outputs $f(x) = \theta^T x$ when θ is drawn from the deep linear prior (2.3). This is true in both linear and non-linear networks at large width (see e.g., [Han18, HN20a, Han22, RYH22] for a treatment of deep non-linear networks). A more useful notion of depth is

$$\lambda_{\text{prior}} = \text{effective depth of prior} := \sum_{\ell=1}^L \frac{1}{N_{\ell}}, \quad (2.7)$$

and it is indeed λ_{prior} that plays an important role in our results. Let us provide a brief justification for why λ_{prior} is a natural measure of complexity for the prior. With $N_1 = \dots = N_L = N$, Theorem 1.2 in [HP21] shows that when $\sigma^2 = 1$, under the prior, the squared singular values of $(W^{(L)} \dots W^{(1)})^{1/L}$ converge to the uniform distribution on $[0, 1]$. Hence, only the squared singular values of $(W^{(L)} \dots W^{(1)})^{1/L}$ lying in intervals of the form $[1 - CL^{-1}, 1]$ correspond to singular values of $W^{(L)} \dots W^{(1)}$ that remain uniformly bounded away from 0 at large L . At least heuristically, this implies that $W^{(L)} \dots W^{(1)}$ is supported on matrices of rank approximately $\lambda_{\text{prior}}^{-1}$ ¹.

Viewing $\lambda_{\text{prior}}^{-1}$ as a natural measure of the number of degrees of freedom in the prior motivates the introduction of a *posterior effective depth*

$$\lambda_{\text{post}} := \frac{P}{\lambda_{\text{prior}}^{-1}} = \sum_{\ell=1}^L \frac{P}{N_{\ell}}. \quad (2.8)$$

¹We do not know this for sure since uniformity for the distribution of singular values at the right edge of the spectrum does not follow from a result only about the global density of singular values.

λ_{post} is a ratio between the number of degrees of freedom in the training data (given by the number of training data points) and in the prior. We'll see in Theorem 3.8 that it is precisely λ_{post} that controls the structure of the posterior.

2.3 Optimal Extrapolation

This section explains key structural properties of Bayesian posteriors in linear networks. Consider the model $f(x) = \theta^T x$ and decompose the parameters $\theta = \theta_{\parallel} + \theta_{\perp}$ as in Figure 1. Since zero noise posteriors fit the training data, we have

$$\theta \sim \mathbb{P}_{\text{post}} \quad \implies \quad \theta_{\parallel} = \theta_{*,N_0},$$

where θ_{*,N_0} is the minimum-norm interpolant (2.1). Moreover, the prior over θ is invariant under all orthogonal transformations and the likelihood is invariant under arbitrary transformations of θ_{\perp} . Hence, in distribution

$$\theta \sim \mathbb{P}_{\text{post}} \quad \implies \quad \theta_{\perp} \stackrel{d}{=} u \cdot \|\theta_{\perp}\|, \quad (2.9)$$

where u is independent of $\|\theta_{\perp}\|$ and is uniformly distributed on the unit sphere in $\text{col}(X_{N_0})^{\perp} \subseteq \mathbb{R}^{N_0}$. The only degree of freedom in the posterior is therefore the distribution of the radial part $\|\theta_{\perp}\|$ of the vector θ_{\perp} . Given a test data point $x = x_{\parallel} + x_{\perp}$, we find

$$\theta \sim \mathbb{P}_{\text{post}} \quad \implies \quad f(x) = (\theta_{*,N_0})^T x_{\parallel} + \|\theta_{\perp}\| \cdot u^T x_{\perp}.$$

The distribution of $\|\theta_{\perp}\|$ controls the scale of predictions for data not spanned by the training set, i.e., for the task of extrapolation. For example if $x = x_{\parallel} \in \text{col}(X_{N_0})$, then $f(x)$ has zero variance since it is determined completely by the training data. More generally, by the Poincare-Borel Theorem (see [DF87, Bor14]) we have for $N_0, P \gg 1$ that

$$f(x) \approx \mathcal{N} \left((\theta_{*,N_0})^T x_{\parallel}, \|\theta_{\perp}\|^2 \frac{\|x_{\perp}\|^2}{N_0 - P} \right). \quad (2.10)$$

Indeed, since $\hat{\theta}_{\perp} = \theta_{\perp} / \|\theta_{\perp}\| \in \mathbb{R}^{N_0 - P}$ is rotationally invariant under the posterior, the Poincare-Borel Theorem (e.g. Theorem 2 in [21] together with the fact that the mixing measure μ is precisely a point-mass at 1 by (1) in [DF87]) shows that the joint distribution of any fixed (or even slowly growing) number of marginals $\{\hat{\theta}_{\perp}^T x_1, \dots, \hat{\theta}_{\perp}^T x_k\}$ is approximately Gaussian when $N_0 - P$ is large. In particular, since predictions under the posterior are of the form $f(x; \theta) = \theta_{\parallel}^T x_{\parallel} + \|\theta_{\perp}\| \hat{\theta}_{\perp}^T x_{\perp}$ and $\|\theta_{\perp}\|$ is independent of $\hat{\theta}_{\perp}$, we see that in the high-dimensional regime $N_0, P \gg 1$ the finite-dimensional distributions of the posterior are approximately normal.

At $L = 0$, the prior and posterior distributions over $\|\theta_{\perp}\|^2$ are identical, preventing any feature learning from occurring. For $L \geq 1$, all components of θ are correlated under the prior, allowing information from the training data to be encoded into θ_{\perp} . We shall see (Theorem 3.8) that λ_{post} quantifies how much information the model learns about θ_{\perp} . In particular, increasing λ_{post} causes the posterior distribution of $\|\theta_{\perp}\|^2$ to be more and more concentrated around a particular value:

$$\|\theta_{\perp}\|^2 \approx \|\theta_{*,N_0}\|^2 (1 - \alpha_0) / \alpha_0.$$

This special choice of scale maximizes Bayesian evidence, in accordance with the first takeaway described in the Introduction. It corresponds to the natural estimate for the true signal strength $\|v\|^2$ under a zero noise generative process $y_i = v^T x_i$ in which x_i are isotropic, given that one observes only the projection $v_{\parallel} = \theta_{*,N_0}$ of v onto directions in the training data.

3 Main Results

3.1 Non-Asymptotic Results

We are now ready to formulate our first main result (Theorem 3.1), which expresses the partition function $Z_{\infty}(\mathbf{t})$ defined in (2.6) in terms of the Meijer-G function; this allows the Bayesian evidence and predictive posterior to be written in exact closed form for any choice of network depth, hidden layer widths, dataset size, and input dimension. Compared to prior work using either iterative saddle point approximations of integrals encoding the Bayesian posterior [LS21] or more involved methods such as the replica trick [ZVTP22], our method provides a direct representation of the network posterior via the partition function in terms of a single contour integral without specializing to limiting cases. Additional quantities, such as the variance of the posterior, are simply expressed as a ratio of Meijer-G functions. We shall later recover known limiting cases and uncover new asymptotic results from expansions of the Meijer-G function (Theorem 3.2).

Theorem 3.1 (Predictive Posterior and Evidence). *Fix $P, L, N_0, \dots, N_L \geq 1, \sigma^2 > 0$ as well as training data X_{N_0}, Y_{N_0} . Fix $\mathbf{t} \in \mathbb{R}^{N_0}$ and write*

$$\mathbf{t} = \mathbf{t}_{\parallel} + \mathbf{t}_{\perp}, \quad \mathbf{t}_{\parallel} \in \text{col}(X_{N_0}), \quad \mathbf{t}_{\perp} \in \text{col}(X_{N_0})^{\perp}.$$

Define $4M = \prod_{\ell=0}^L 2\sigma^2/N_{\ell}$ and introduce the following shorthand for the Meijer-G functions (see §4.2) with parameters given by layer widths:

$$G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_{*,N_0}\|^2}{4M} \mid \frac{P}{2}, \frac{-}{\frac{N}{2} + k} \right) := G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_{*,N_0}\|^2}{4M} \mid \frac{P}{2}, \frac{N_1}{2} + k, \dots, \frac{N_L}{2} + k \right).$$

The partition function $Z_{\infty}(\mathbf{t})$ of the predictive posterior defined in (2.6) is

$$\begin{aligned} Z_{\infty}(\mathbf{t}) &= \left(\frac{4\pi}{\|\theta_{*}\|^2} \right)^{\frac{P}{2}} \exp[-i \langle \theta_{*,N_0}, \mathbf{t} \rangle] \prod_{\ell=1}^L \Gamma\left(\frac{N_{\ell}}{2}\right)^{-1} \\ &\quad \times \sum_{k=0}^{\infty} \frac{(-1)^k}{k!} \|\mathbf{t}_{\perp}\|^{2k} M^k G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_{*,N_0}\|^2}{4M} \mid \frac{P}{2}, \frac{-}{\frac{N}{2} + k} \right). \end{aligned} \quad (3.1)$$

In particular, the Bayesian model evidence $Z_{\infty}(\mathbf{0})$ equals

$$\frac{(4\pi)^{P/2}}{\|\theta_{*}\|^P \prod_{\ell=1}^L \Gamma\left(\frac{N_{\ell}}{2}\right)} G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_{*,N_0}\|^2}{4M} \mid \frac{P}{2}, \frac{-}{\frac{N}{2}} \right). \quad (3.2)$$

Further, given $x \in \mathbb{R}^{N_0}$, the mean of the predictive posterior is

$$\mathbb{E}_{\text{post}}[f(x)] = (\theta_{*,N_0})^T x, \quad (3.3)$$

while the posterior variance $\text{Var}_{\text{post}} [f(x)]$ is

$$2M \|x_{\perp}\|^2 \frac{G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_{*,N_0}\|^2}{4M} \mid \frac{P}{2}, \frac{\mathbf{N}}{2} + 1 \right)}{G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_{*,N_0}\|^2}{4M} \mid \frac{P}{2}, \frac{\mathbf{N}}{2} \right)}, \quad (3.4)$$

where x_{\perp} is the projection of x onto the orthogonal complement of the span $\text{col}(X_{N_0})$ of the training data.

See §4.4 for a proof.

3.2 Asymptotic Results

To evaluate the predictive posterior and evidence in Theorem 3.1 in the limits where N_0, P, N_{ℓ} (and potentially L) tend to infinity, we require novel expansions of the Meijer-G function obtained by the Laplace method. We are interested in regimes where N_0, P grow and will assume a mild compatibility condition on the training data: for all $\alpha_0 \in (0, 1)$, we assume there exists constant $\|\theta_*\|$ such that

$$\lim_{\substack{P, N_0 \rightarrow \infty \\ P/N_0 \rightarrow \alpha_0 \in (0, 1)}} \|\theta_{*,N_0}\| = \|\theta_*\|, \quad (3.5)$$

where the convergence is in distribution. This assumption is very generic and is (for example) satisfied for a Gaussian data model where inputs are Gaussian

$$X_{N_0} = (x_{i,N_0}, i = 1, \dots, P), \quad x_{i,N_0} \sim \mathcal{N}(0, \Sigma_{N_0}) \quad \text{iid},$$

outputs are linear plus noise

$$Y = V_{N_0} X_{N_0} + \epsilon_{N_0}, \quad V_{N_0} \sim \text{Unif}(S^{N_0-1}), \quad \epsilon_{N_0} \sim \mathcal{N}(0, \sigma_{\epsilon}^2 I_{N_0}),$$

and the spectral density of the design matrices Σ_{N_0} converges weakly as $N_0 \rightarrow \infty$ to a fixed probability measure on \mathbb{R}_+ with finite moments.

To minimize notation, we report here the expansions in terms of a single layer width $N = N_1 = \dots = N_L$, but expansions with distinct N_{ℓ} (and to higher order) are provided in the proof (§4.5).

Theorem 3.2 (Asymptotic Expansions of Meijer-G). *Set $N_1, \dots, N_L = N$ and define $\mathbf{N} = (N_1, \dots, N_L)$. Suppose that the training data satisfies (3.5). In different limiting cases such that $\{P, N\} \rightarrow \infty$ with fixed $P/N_0 = \alpha_0$, we evaluate the quantities*

$$\begin{aligned} \log G &:= \log G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_{*,N_0}\|^2}{4M} \mid \frac{P}{2}, \frac{\mathbf{N}}{2} + k \right) \\ \Delta(\log G)[k] &:= \log G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_{*,N_0}\|^2}{4M} \mid \frac{P}{2}, \frac{\mathbf{N}}{2} + k \right) - \log G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_{*,N_0}\|^2}{4M} \mid \frac{P}{2}, \frac{\mathbf{N}}{2} \right). \end{aligned}$$

We will see in each case that to leading order $\log G$ does not depend on k while $\Delta(\log G)[k]$ does.

(a) Fix $L < \infty$, $\alpha, \sigma^2 > 0$. Suppose $P, N \rightarrow \infty$ with $P/N \rightarrow \alpha$. Then,

$$\begin{aligned} \log G &= \frac{N\alpha}{2} \left[\log \left(\frac{N\alpha}{2} \right) + \log \left(1 + \frac{z_*}{\alpha} \right) - \left(1 + \frac{z_*}{\alpha} \right) \right] \\ &\quad + \frac{NL}{2} \left[\log \left(\frac{N}{2} \right) + \log(1 + z_*) - (1 + z_*) \right] \end{aligned} \quad (3.6)$$

plus an error of size $\tilde{O}(1)$ and

$$\Delta(\log G)[k] = kL \left[\log \left(\frac{N}{2} \right) + \log(1 + z_*) \right], \quad (3.7)$$

plus an error of size $\tilde{O}(1/N)$, where $z_* > \min\{-\alpha, -1\}$ is the unique solution to

$$\left(1 + \frac{z_*}{\alpha} \right) (1 + z_*)^L = \frac{\|\theta_*\|^2}{\sigma^{2(L+1)}\alpha_0}. \quad (3.8)$$

(b) Fix $\lambda_{\text{prior}}, \alpha > 0$ and suppose $L = \lambda_{\text{prior}}N$, $P, N \rightarrow \infty$, $P/N \rightarrow \alpha$, $\sigma^2 = 1$. Then,

$$\begin{aligned} \log G &= \frac{\lambda_{\text{prior}}N^2}{2} \left[\log \left(\frac{N}{2} \right) - 1 \right] + \frac{N\alpha}{2} \left[\log \left(\frac{N\alpha}{2} \right) - 1 \right] \\ &\quad + \frac{\lambda_{\text{prior}}N}{2} \left[-\log \left(\frac{N}{2} \right) + \log(2\pi) \right] + \tilde{O}(1) \end{aligned} \quad (3.9)$$

and, up to an error of size $\tilde{O}(1/N)$,

$$\Delta(\log G)[k] = k \left[\lambda_{\text{prior}}N \log \left(\frac{N}{2} \right) + \log \left(\frac{\|\theta_*\|^2}{\alpha_0} \right) \right]. \quad (3.10)$$

(c) Fix $\lambda_{\text{post}} > 0$. Suppose $N, P, L \rightarrow \infty$ with $LP/N \rightarrow \lambda_{\text{post}}$, $\sigma^2 = 1$ and $L/N \rightarrow 0$. Then,

$$\begin{aligned} \log G &= \frac{P}{2} \left[\log \left(\frac{P}{2} \right) - 1 \right] + \frac{\lambda_{\text{post}}N}{2} \frac{N}{P} \left[\log \left(\frac{N}{2} \right) - 1 \right] \\ &\quad + \frac{P}{2} \left[\log(1 + t_*) - t_* \left(1 + \frac{\lambda_{\text{post}}t_*}{2} \right) \right] \\ &\quad + \frac{\lambda_{\text{post}}N}{2} \left[\frac{N}{P} \left(1 + \frac{P}{N}t_* \right) \log \left(1 + \frac{P}{N}t_* \right) - t_* \right] \\ &\quad + \tilde{O}(1), \end{aligned} \quad (3.11)$$

and

$$\Delta(\log G)[k] = k \left[L \log \left(\frac{N}{2} \right) + \lambda_{\text{post}}t_* \right] + \tilde{O} \left(\frac{1}{N} \right), \quad (3.12)$$

where t_* is the unique solution to

$$e^{\lambda_{\text{post}}t_*} (1 + t_*) = \|\theta_*\|^2 / \alpha_0. \quad (3.13)$$

In all the estimates above, $\tilde{O}(1) = O(\max\{\log P, \log N\})$ suppresses lower-order terms of order 1, up to logarithmic factors; similarly, $\tilde{O}(1/N) = O(\max\{\log P, \log N\}/N)$. Suppressed terms are included in §4.5

3.3 Model Selection and Extrapolation

We combine our non-asymptotic formulas for the posterior and evidence from Theorem 3.1 with the Meijer-G function expansions from Theorem 3.2 to investigate two fundamental questions about Bayesian interpolation with linear networks. Together they give, at least in the restricted setting of deep linear networks with output dimension 1, a principled reason to prefer deeper networks.

- **Model Selection.** How should we choose the prior weight variance σ^2 , model depth L and layer widths N_ℓ ? Recall that the Bayesian model evidence $Z_\infty(\mathbf{0})$ from (2.5) represents the likelihood of observing the data (X_{N_0}, Y_{N_0}) given the architecture L, N_ℓ and the prior weight variance σ^2 . Maximizing the Bayesian evidence is therefore maximum likelihood estimation over the space of models and gives a principled method for model selection. We shall see that data-agnostic priors maximize the evidence for networks with infinite λ_{post} , while shallower networks require data-dependent priors.
- **Extrapolation.** How optimal is the posterior of a linear network? Predictions on inputs from directions orthogonal to the training data are determined by the distribution θ_\perp . At $\lambda_{\text{post}} = 0$, its prior and posterior coincide. As λ_{post} increases, however, we shall find that information from the training set mixes into θ_\perp and ultimately maximizes evidence, producing optimal extrapolation.

For simplicity, we report our results here in terms of a single hidden layer width $N = N_1 = \dots = N_L$. The general form with generic N_ℓ , as well as related results not stated here, are available by direction application of the general expansions in §4.6. For convenience, we define

$$\nu := \|\theta_*\|^2 / \alpha_0. \quad (3.14)$$

As elsewhere, we emphasize that we focus on the regime

$$P/N_0 \rightarrow \alpha_0 \in (0, 1)$$

as $P, N_0 \rightarrow \infty$. When $\alpha_0 \geq 1$, Theorem 3.1 still holds but immediately yields that $\theta = \theta_{*, N_0}$ almost surely since $\theta_\perp = 0$. As described by (2.10) and stated rigorously in §4.3, the predictive posterior in the regime $P < N_0$ is Gaussian with variance determined by $\|\theta_\perp\|^2$, which converges to a constant. We rewrite this posterior in the form

$$f(x) \rightarrow \mathcal{N}(\mu_*, \nu c \Sigma_\perp), \quad x = (x_{i, N_0}, i = 1, \dots, k)$$

for free scalar c , scalar μ_* , and $k \times k$ PSD matrix Σ_\perp given by

$$\mu_* := \langle \theta_{*, N_0}, x \rangle, \quad \Sigma_\perp = \frac{\langle x_{i, N_0}^\perp, x_{j, N_0}^\perp \rangle}{N_0 - P}$$

in the limit $P, N_0 \rightarrow \infty$ and $P/N_0 \rightarrow \alpha_0 \in (0, 1)$. The following results report the Bayesian evidence and value of c under different joint scalings of P, N and L . First, we observe that maximizing Bayesian evidence always results in the posterior corresponding to $c = 1$ (proven §4.6).

Theorem 3.3 (Universal Maximal Evidence Posterior). *Fix $L, N_1, \dots, N_L \geq 1$, $\sigma^2 > 0$, $\alpha_0 \in (0, 1)$, and consider sequences of training data sets $X_{N_0} \in \mathbb{R}^{N_0 \times P}$, $Y_{N_0} \in \mathbb{R}^{1 \times P}$ such that*

$$P, N_0 \rightarrow \infty, \quad P/N_0 \rightarrow \alpha_0$$

and (3.5) holds. Let $Z_\infty(\mathbf{0})$ denote the limiting Bayesian model evidence. The following statements are equivalent:

- (a) σ^2 maximizes $Z_\infty(\mathbf{0})$, and
- (b) the posterior over predictions $f(x)$ converges weakly to a Gaussian $\mathcal{N}(\mu_*, \nu \Sigma_\perp)$.

We emphasize that Theorem 3.3 holds for any choice of depth and hidden layer widths. At finite depth, we shall see that identifying the evidence-maximizing prior σ_*^2 requires knowledge of the data-dependent parameter ν . In contrast, infinite-depth networks have evidence maximized by $\sigma_*^2 = 1$, allowing them to successfully infer ν from training data despite a data-agnostic prior. We proceed to consider different joint scalings of P, N and L in Theorems 3.4, 3.6 and 3.8, which follow directly from writing the partition function in terms of the Meijer-G function (Theorem 3.1) and applying the suitable asymptotic expansion (Theorem 3.2).

3.3.1 Finite Depth

We present here a characterization of the infinite-width posterior and model evidence for networks with a fixed number of hidden layers.

Theorem 3.4 (Posterior and Evidence at Finite L). *For each $N_0, P \geq 1$, consider training data X_{N_0}, Y_{N_0} satisfying (3.5). Fix constants $L \geq 0$, $\sigma^2 > 0$ and suppose that*

$$P, N_0, N \rightarrow \infty, \quad P/N_0 \rightarrow \alpha_0 \in (0, 1), \quad P/N \rightarrow \alpha \in (0, \infty).$$

In this limit, the posterior over predictions $f(x)$ converges weakly to a Gaussian

$$\mathcal{N}\left(\mu_*, \nu \Sigma_\perp \left(1 + \frac{z_*}{\alpha}\right)^{-1}\right),$$

where z_ is the unique solution to*

$$\frac{\nu}{\sigma^{2(L+1)}} = \left(1 + \frac{z_*}{\alpha}\right) (1 + z_*)^L, \quad z_* > \max\{-1, -\alpha\}. \quad (3.15)$$

Additionally, we have the following large- P expansion of the Bayesian model evidence

$$\begin{aligned} \log Z_\infty(\mathbf{0}) &\rightarrow \frac{P}{2} \left[\log\left(\frac{P}{2}\right) + \log\left(1 + \frac{z_*}{\alpha}\right) - \left(1 + \frac{z_*}{\alpha}\right) \right. \\ &\quad \left. - \log\left(\frac{\|\theta_*\|^2}{4\pi}\right) \right] + \frac{NL}{2} [\log(1 + z_*) - z_*] \\ &\quad + O(\max\{\log P, \log N\}). \end{aligned}$$

In regimes where z_* tends to zero, the posterior will converge to the evidence-maximizing posterior independent of the architecture and prior. By (3.15), this occurs for instance in the limit of infinite depth with $\sigma = 1$ or when $\nu = \sigma^{2(L+1)}$. This last condition corresponds to maximizing the Bayesian evidence $Z_\infty(\mathbf{0})$ at finite L .

Corollary 3.5 (Bayesian Model Selection at Finite L). *In the setting of Theorem 3.4 the Bayesian evidence $Z_\infty(\mathbf{0})$ satisfies:*

$$\sigma_*^2 = \arg \max_{\sigma} Z_\infty(\mathbf{0}) = \nu^{\frac{1}{L+1}} \quad (3.16)$$

$$L_* = \arg \max_L Z_\infty(\mathbf{0}) = \frac{\log(\nu)}{\log(\sigma_*^2)} - 1. \quad (3.17)$$

In particular, given a prior variance with $\text{sgn}(\nu - 1) = \text{sgn}(\sigma^2 - 1)$ satisfying $|\sigma^2 - 1| \leq \epsilon$, the optimal depth network satisfies $L_ \geq |\log(\nu)|/\epsilon$.*

To put this Corollary into context, note that in the large-width limit of Theorem 3.4 there are only two remaining model parameters, L and σ^2 . Model selection can therefore be done in two ways. The first is an empirical Bayes approach in which one uses the training data to determine a data-dependent prior σ_*^2 given by (3.16). The other approach, which more closely follows the use of neural networks in practice, is to seek a universal, data-agnostic value of σ^2 and optimize instead the network architecture. The expression (3.16) shows that the only way to choose σ^2, L to (approximately) maximize model evidence for any fixed ν is to take $\sigma^2 \approx 1$ with $\text{sgn}(\sigma^2 - 1) = \text{sgn}(\nu - 1)$ and $L \rightarrow \infty$. Hence, restricting to the data-agnostic prior $\sigma^2 = 1$ naturally leads to a Bayesian preference for infinite-depth networks, regardless of the training data. This motivates us to consider large- N limits in which L tends to infinity, which we take up in the next two sections.

3.3.2 Infinite Depth

We fix $\sigma^2 = 1$ and investigate extrapolation and model selection in regimes where $N, L, P \rightarrow \infty$ simultaneously.

Theorem 3.6 (Posterior and Evidence at Fixed λ_{prior}). *For each $N_0, P \geq 1$, consider training data X_{N_0}, Y_{N_0} satisfying (3.5). Moreover, fix $\lambda_{\text{prior}}, \alpha \in (0, \infty), \alpha_0 \in (0, 1)$. Suppose that $N_1 = \dots = N_L = N$ and that*

$$P, N_\ell, L \rightarrow \infty, \quad P/N_0 \rightarrow \alpha_0, \quad P/N \rightarrow \alpha, \quad L/N \rightarrow \lambda_{\text{prior}}.$$

In this limit, the posterior over predictions $f(x)$ converges weakly to a Gaussian

$$\mathcal{N}(\mu_*, \nu \Sigma_\perp),$$

which is independent of $\alpha, \lambda_{\text{prior}}$. Additionally, the evidence admits the following asymptotic expansion:

$$\log Z_\infty(\mathbf{0}) = \frac{P}{2} \left[\log \left(\frac{P}{2} \right) - 1 - \log \left(\frac{\|\theta_*\|^2}{4\pi} \right) \right] + O(\max \{\log P, \log N\}).$$

This result highlights the remarkable nature of data-driven extrapolation in deep networks.

Corollary 3.7 (Optimal Learning by Deep Networks). *In the setting of Theorem 3.6, the posterior distribution over regression weights θ is the same in the following two settings:*

- *We fix $L \geq 0$, take the data-dependent prior variance σ_* that maximizes the Bayesian model evidence as a function of the training data and network depth as in (3.16), and send the network width N to infinity.*

- We fix $\lambda_{\text{prior}} > 0$, take a data-agnostic prior $\sigma^2 = 1$, and send both the network depth $L := \lambda_{\text{prior}} \cdot N$ and network width N to infinity together.

This corollary makes precise the statement that infinitely deep networks with data-agnostic priors performs as optimally finite depth networks with fine-tuned data-dependent priors.

In the §4.7, we find that the ratio of model evidence at fixed depth to the model evidence at infinite depth vanishes like $\exp[-O(N)]$. In comparison, mis-specifying the value of constant λ_{prior} only results in an $O(1)$ ratio of model evidences, and it does not affect the posterior to leading order. We conclude that, at $\sigma^2 = 1$, wide networks with depth comparable to width are robustly preferred to shallow networks.

3.3.3 Scaling Laws for Optimal Learning

To emphasize the similarity between dataset size and depth, we take the limit of $L, P, N \rightarrow \infty$ while holding LP/N constant. This results in a scaling law for feature learning that only depends on λ_{post} , as seen in the following theorem.

Theorem 3.8 (Posterior and Evidence at Fixed λ_{post}). *For each $N_0, P \geq 1$, consider training data X_{N_0}, Y_{N_0} satisfying (3.5). Moreover, fix constants $\lambda_{\text{post}} > 0$, $\alpha_0 \in (0, 1)$, $\sigma^2 = 1$. Suppose that $N_1 = \dots = N_L = N$ and suppose that*

$$P, N_\ell, L \rightarrow \infty, \quad \frac{P}{N_0} \rightarrow \alpha_0, \quad \frac{P}{N} \rightarrow 0, \quad \frac{L}{N} \rightarrow 0, \quad \frac{LP}{N} \rightarrow \lambda_{\text{post}}.$$

In this limit, the posterior over predictions $f(x)$ converges weakly to a Gaussian $\mathcal{N}(\mu_, \nu \Sigma_\perp (1 + t_*)^{-1})$, where t_* is the unique solution to*

$$\nu = (1 + t_*)e^{\lambda_{\text{post}} t_*}, \quad t_* > -1.$$

Additionally, we have the following asymptotic expansion for the Bayesian model evidence

$$\begin{aligned} \log Z_\infty(\mathbf{0}) &= \frac{P}{2} \left(\log \left(\frac{P}{2} \right) - 1 \right) - \frac{P}{2} \log \left(\frac{\|\theta_*\|^2}{4\pi} \right) \\ &\quad + \frac{P}{2} \left[\log(1 + t_*) - t_* - \frac{1}{2} \lambda_{\text{post}} t_*^2 \right] + \tilde{O}(\max \{ \log P, \log N, \log L \}). \end{aligned}$$

In the limit $\lambda_{\text{post}} \rightarrow \infty$, Theorem 3.8 implies that optimal feature learning is rapidly approached — specifically, like the harmonic mean of 1 and $\lambda_{\text{post}}/\log \nu$. Bayesian model selection also drives λ_{post} to infinity.

Corollary 3.9 (Bayesian Model Selection at Fixed λ_{post}). *In the setting of Theorem 3.8, the Bayesian evidence $Z_\infty(\mathbf{0})$ is monotonically increasing in λ_{post} . Specifically,*

$$\frac{\partial \log Z_\infty(\mathbf{0})}{\partial \lambda_{\text{post}}} = \frac{P t_*^2}{4} \geq 0. \quad (3.18)$$

These results provide simple scaling laws that apply independently of the choice of dataset, demonstrating the behavior of the predictor’s posterior distribution and model evidence in terms of λ_{post} . The coupling of depth and dataset size in λ_{post} provides a novel interpretation of depth as a mechanism to improve learning in a manner similar to additional data: larger datasets and larger depths contribute equally towards aligning the prior $\sigma^2 = 1$ towards the correct posterior.

3.4 Properties of Deep Linear Networks

We relate our work to prior results in the literature: variance-limited scaling laws [BDK⁺21] and sample-wise double descent [BHMM19, BHX20, ZVTP22]. To make the discussion concrete, we shall focus on the architecture introduced in Theorem 3.6, where depth, width, and dataset size scale linearly with each other to produce a posterior that exhibits optimal feature learning given prior $\sigma^2 = 1$.

3.4.1 Variance-Limited Scaling Laws

We examine the scaling behavior of model error in the infinite-width or infinite-dataset limits. The work of [BDK⁺21] shows that the difference between the finite-size loss and the infinite-size loss scales like $1/x$ for $x = N$ or P while the other parameter is held fixed (P or N , respectively). Here, we demonstrate an analogous scaling law when $N \propto P \propto L$. Similarly to the results of [BDK⁺21], this scaling law is independent of the choice of dataset and consequently provides a universal insight into how performance improves with larger models or more data. The proof is found in §4.8.

Theorem 3.10 (Variance-Limited Scaling Law). *For each $N_0 \geq 1$ consider training data X_{N_0}, Y_{N_0} satisfying (3.5). Fix constants $\lambda_{\text{prior}} > 0, \alpha > 0, \sigma^2 = 1$. Suppose that $N_1 = \dots = N_L = N$, that $L = \lambda_{\text{prior}} N$ and number of data points satisfies $P = N\alpha$. Then, as $N \rightarrow \infty$,*

$$\text{Var}_{\text{post}} [f(x)] = \lim_{N \rightarrow \infty} \text{Var}_{\text{post}} [f(x)] + \frac{C}{N} + O\left(\frac{\log N}{N^2}\right),$$

where $C \in \mathbb{R}$ is a universal constant.

3.4.2 Double Descent

We demonstrate double descent in $\alpha_0 = P/N_0$ consistent with previous literature [ZVTP22]. As a concrete example, we shall consider a Gaussian data model and evaluate double descent for the posterior of optimal feature learning; this posterior is achieved by, for example, deep networks with $\sigma^2 = 1$ (Theorem 3.6), or finite-depth networks with data-tuned priors σ_*^2 (Theorem 3.4). The proof is found §4.9.

Theorem 3.11 (Double Descent in α_0). *Consider generative data model*

$$x_i \in \mathbb{R}^{N_0} \sim \mathcal{N}(0, \mathbf{I}), y_i = V_0 x_i + \epsilon_i, \epsilon_i \sim \mathcal{N}(0, \sigma_\epsilon^2), \|V_0\|^2 = 1$$

and posterior distribution $\mathcal{N}(\mu_*, \nu \Sigma_\perp)$. We have error

$$\mathbb{E}_{x, X, \epsilon} \left[\langle f(x) - V_0 x \rangle^2 \right] = \begin{cases} \frac{1}{\alpha_0} - \alpha_0 + \frac{1}{1-\alpha_0} \sigma_\epsilon^2, & \alpha_0 < 1 \\ \frac{1}{\alpha_0 - 1} \sigma_\epsilon^2, & \alpha_0 \geq 1 \end{cases},$$

which diverges at $\alpha_0 = 1$.

4 Discussion

Neural networks are non-linear functions of their parameters, making an analytic understanding of their properties difficult. Here, we adopted the simplification of studying linear neural

networks, which remain non-linear in their parameters but are more tractable. To conclude, we emphasize three limitations of our work. First, we consider only *linear networks*, which are linear as a function of their inputs. However, they are not linear as a function of their parameters, making learning by Bayesian inference non-trivial. Second, we study learning by Bayesian inference and leave extensions to learning by gradient descent to future work. Finally, our results characterize the *predictive posterior*, i.e., the distribution over model predictions after inference. It would be interesting to derive the form of the posterior over network weights as well.

4.1 Acknowledgements

This work started at the 2022 Summer School on the Statistical Physics of Machine Learning held at École de Physique des Houches. We are grateful for the wonderful atmosphere at the school and would like to express our appreciation to the session organizers Florent Krzakala and Lenka Zdeborová as well as to Haim Sompolinsky for his series of lectures on Bayesian analysis of deep linear networks. We further thank Edward George for pointing out the connection between our work and the deep Gaussian process literature. Finally, we thank Matias Cattaneo, Isaac Chuang, David Dunson, Jianqing Fan, Aram Harrow, Jason Klusowski, Cengiz Pehlevan, Veronika Rockova, and Jacob Zavatone-Veth for their feedback and suggestions. BH is supported by NSF grants DMS-2143754, DMS-1855684, and DMS-2133806. AZ is supported by the Hertz Foundation, and by the DoD NDSEG. We are also thank two anonymous reviewers for improving aspects of the exposition and for pointing out a range of typos in the original manuscript.

4.2 Background

In this section, we introduce background and results needed for our proofs. Specifically, §4.2.1 recalls basic definitions and asymptotic expansions for Gamma and Digamma functions (§4.2.1), §4.2.2 recalls the moments of Gamma random variables, and §4.2.3 defines Meijer-G functions and introduces their basic properties.

4.2.1 Gamma and Digamma Functions

We will need the following asymptotic expansions for Euler’s Gamma function

$$\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt, \quad \Re(z) > 0$$

and Digamma function $\phi^{(1)}(z) = \frac{d}{dz} \log \Gamma(z)$ (see Equations 6.1.7, 6.3.13 in [AS64]):

Proposition 4.1. *We have the following analytic asymptotic expansion for the Gamma function*

$$\log \Gamma(z) \sim \left(z - \frac{1}{2}\right) \log(z) - z + \text{const} + O(|z|^{-1}), \quad \text{as } |z| \rightarrow \infty. \quad (4.1)$$

In particular, we also have

$$\phi^{(1)}(z) = \frac{d}{dz} \log \Gamma(z) \sim \log(z) + O(|z|^{-1}), \quad \text{as } |z| \rightarrow \infty. \quad (4.2)$$

Both expansions hold uniformly on sets of the form $|z| < \pi - \delta$ for a fixed $\delta > 0$.

4.2.2 Gamma Random Variables

We will need the following well-known exact formulas for the moments of Gamma random variables, which follow directly for the formula for their density:

$$\text{Den}_X(x) = \begin{cases} \frac{x^{k-1}}{\Gamma(k)\theta^k} e^{-x/\theta}, & x > 0 \\ 0, & x \leq 0 \end{cases}, \quad X \sim \Gamma(k, \theta), \quad k, \theta > 0$$

and the definition of the Gamma function.

Proposition 4.2. *Let $k, \theta > 0$ and suppose $\phi \sim \Gamma(k, \theta)$. Then for any $t \in \mathbb{R}$ we have*

$$\mathbb{E}[\phi^t] = \theta^t \frac{\Gamma(k+t)}{\Gamma(k)}. \quad (4.3)$$

In particular

$$\mathbb{E}[e^{-it \log \phi}] = \exp[-it \log \theta + \log \Gamma(k-it) - \log \Gamma(k)] \quad (4.4)$$

is a meromorphic function of t with poles on the negative imaginary axis:

$$t = -i(\nu + k), \quad \nu = 0, 1, 2, \dots$$

4.2.3 Meijer G-Functions

The Meijer-G function is defined as the contour integral

$$G_{m,n}^{p,q} \left(z \left| \begin{matrix} \mathbf{a} \\ \mathbf{b} \end{matrix} \right. \right) = \frac{1}{2\pi i} \int_{\mathcal{C}} z^s \chi(s) ds, \quad (4.5)$$

where

$$\chi(s) := \frac{\prod_{j=1}^m \Gamma(b_j - s) \prod_{k=1}^n \Gamma(1 - a_k + s)}{\prod_{j=m+1}^q \Gamma(1 - b_j + s) \prod_{k=n+1}^p \Gamma(a_k + s)}$$

and \mathcal{C} is a Mellin-Barnes contour in the complex plane that separates the poles of $\Gamma(b_j - s)$ from those of $\Gamma(1 - a_k + s)$ (see Figure 2).

We will need several properties of Meijer-G functions, which we now recall.

For (4.13) we require

$$|\arg(\sigma)|, |\arg(\omega)| < \pi$$

and

$$-\min\{\Re(b_1), \dots, \Re(b_m)\} < \Re(\alpha) < 2 - \max\{\Re(a_1), \dots, \Re(a_n)\} - \max\{\Re(c_1), \dots, \Re(c_t)\}.$$

Finally, fix $L \geq 1$ and

$$k_\ell, \theta_\ell > 0, \quad \ell = 1, \dots, L$$

and let

$$\phi_\ell \sim \Gamma(k_\ell, \theta_\ell) \quad \text{independent.}$$

We have for every $y > 0$

$$\prod_{\ell=1}^L \Gamma(k_\ell)^{-1} \frac{1}{A} G_{0,L}^{L,0} \left(\frac{y}{A} \mid k_1 - 1, \dots, k_L - 1 \right) = \text{Den}_{\prod_{\ell=1}^L \phi_\ell}(y), \quad A := \prod_{\ell=1}^L \theta_\ell. \quad (4.14)$$

4.3 Setup

As described in the main text, rather than working with the posterior distribution

$$d\mathbb{P}_{\text{post}}(\theta \mid L, N_\ell, \sigma^2, X_{N_0}, Y_{N_0}) := \lim_{\beta \rightarrow \infty} \frac{d\mathbb{P}_{\text{prior}}(\theta \mid N_0, L, N_\ell, \sigma^2) \exp[-\beta \mathcal{L}(\theta \mid X_{N_0}, Y_{N_0})]}{Z_\beta(X_{N_0}, Y_{N_0} \mid L, N_\ell, \sigma^2)}, \quad (4.15)$$

we work with the characteristic function of the posterior:

$$\mathbb{E}_{\text{post}}[\exp\{-i\mathbf{t} \cdot \theta\}] = \frac{Z_\infty(\mathbf{t} \mid L, N_\ell, \sigma^2, X_{N_0}, Y_{N_0})}{Z_\infty(\mathbf{0} \mid L, N_\ell, \sigma^2, X_{N_0}, Y_{N_0})}, \quad \mathbf{t} = (t_1, \dots, t_{N_0}) \in \mathbb{R}^{N_0}. \quad (4.16)$$

Here, $\mathbb{E}_{\text{post}}[\cdot]$ is the expectation with respect to the posterior (4.15). We have defined for $\beta > 0$ the partition function $Z_\beta(\mathbf{t}) = Z_\beta(\mathbf{t} \mid L, N_\ell, \sigma^2, X_{N_0}, Y_{N_0})$ by

$$Z_\beta(\mathbf{t}) := A_\beta \int \exp \left[-\sum_{\ell=1}^{L+1} \frac{N_{\ell-1}}{2\sigma^2} \|W^{(\ell)}\|_F^2 - \frac{\beta}{2} \|Y - \prod_{\ell=1}^{L+1} W^{(\ell)} X\|_2^2 - i\theta \cdot \mathbf{t} \right] \prod_{\ell=1}^{L+1} dW^{(\ell)} \quad (4.17)$$

and set

$$Z_\infty(\mathbf{t} \mid L, N_\ell, \sigma^2, X_{N_0}, Y_{N_0}) := \lim_{\beta \rightarrow \infty} Z_\beta(\mathbf{t} \mid L, N_\ell, \sigma^2, X_{N_0}, Y_{N_0}),$$

where

$$A_\beta = \det(X_{N_0}^T X_{N_0})^{1/2} (2\pi\beta)^{P/2}$$

and $Y_{N_0, \perp}$ is the projection of Y_{N_0} onto the orthogonal complement to the row span of X_{N_0} . The denominator $Z_\infty(\mathbf{0})$ is often called the Bayesian model evidence and represents the probability of the data (X_{N_0}, Y_{N_0}) given the model (i.e. depth L , layer widths N_1, \dots, N_L and prior scale σ^2). Note that the constant A_β cancels in the ratio (4.15) and in any computations involving maximizing ratios of model evidence.

The effective depth of the prior is measured by L/N ; more precisely, by

$$\lambda_{\text{prior}} = \lambda_{\text{prior}}(N_1, \dots, N_L) = \text{effective depth of prior} := \sum_{\ell=1}^L \frac{1}{N_\ell}. \quad (4.18)$$

Beyond the justification of effective rank provided in the main text, we offer related intuition here to justify calling λ_{prior} the effective depth. Aside from the simple multiplicative dependence on σ^2 , the norm of θ under the prior depends at large N, L only on λ_{prior} . For instance, a simple computation (Equation (9) in [HN20b]) shows that with $\sigma^2 = 1$

$$\lim_{\substack{N_1 + \dots + N_L \rightarrow \infty \\ \lambda_{\text{prior}}(N_1, \dots, N_L) \rightarrow \lambda}} \log \|\theta\|^2 = \mathcal{N}\left(-\frac{\lambda}{2}, \lambda\right),$$

where the convergence is in distribution. Due to the rotational invariance θ , we thus see that with $\sigma^2 = 1$, it is λ_{prior} , as opposed to L , that gives a full description of the prior. Moreover, taking $\lambda_{\text{prior}} \rightarrow 0$ (even if $L \rightarrow \infty$) gives the same (Gaussian) prior over predictions $\theta^T x$ as one would obtain by simply starting with $L = 0$.

As argued in the main text by the Poincare-Borel lemma, the posterior in the large-data limit will always have the form of a normal distribution. We formalize this statement here, taking θ_{*, N_0} to be the minimum-norm interpolant

$$\theta_{*, N_0} := \arg \min_{\theta \in \mathbb{R}^{N_0}} \|\theta\|_2 \quad \text{s.t.} \quad \theta^T X_{N_0} = Y_{N_0}. \quad (4.19)$$

Lemma 4.4 (Asymptotic Normality of Posterior). *For each N_0 , consider a collection of $k \geq 1$ test points*

$$\mathbf{x}_{N_0} = (x_{j, N_0}, j = 1, \dots, k).$$

Suppose that

- *For each $\alpha_0 \in (0, 1)$ there exists a vector $\mu_* = (\mu_{*, j}, j = 1, \dots, k) \in \mathbb{R}^k$ such that*

$$\lim_{\substack{P, N_0 \rightarrow \infty \\ P/N_0 \rightarrow \alpha_0 \in (0, 1)}} \langle \theta_{*, N_0}, x_{j, N_0} \rangle = \mu_{*, j} \quad \text{almost surely.} \quad (4.20)$$

- *For each $\alpha_0 \in (0, 1)$ there exists a positive semi-definite $k \times k$ matrix Σ_{\perp} such that*

$$\lim_{\substack{P, N_0 \rightarrow \infty \\ P/N_0 \rightarrow \alpha_0 \in (0, 1)}} \left(\frac{1}{N_0 - P} \left\langle x_{i, N_0}^{\perp}, x_{j, N_0}^{\perp} \right\rangle \right)_{1 \leq i, j \leq k} = \Sigma_{\perp} \quad \text{almost surely,} \quad (4.21)$$

where x_{j, N_0}^{\perp} denotes the projection of x_j onto the orthogonal complement of the column space of X_{N_0} .

Assume that the convergence of the posterior for $\|\theta_{\perp}\|^2$ in the large- N limit satisfies

$$\|\theta_{\perp}\|^2 \rightarrow \nu(1 - \alpha_0)c$$

for some constant c . Then the distribution over posterior predictions evaluated on \mathbf{x}_{N_0}

$$f(\mathbf{x}_{N_0}) = (a^T x_{j, N_0}, j = 1, \dots, k), \quad \theta \sim \mathbb{P}_{\text{post}}$$

converges weakly to a k -dimensional Gaussian:

$$f(\mathbf{x}_{N_0}) \rightarrow \mathcal{N}(\mu_*, \nu c \Sigma_{\perp}). \quad (4.22)$$

4.4 Proof of Theorem 3.1

Recall that, by definition,

$$Z_\beta(\mathbf{x}_{N_0}, \mathbf{t}) = A_\beta \mathbb{E}_{\text{prior}} \left[\exp \left[-\frac{\beta}{2} \left\| Y - \prod_{\ell=1}^{L+1} W^{(\ell)} X \right\|_2^2 - i \langle f(\mathbf{x}_{N_0}), \mathbf{t} \rangle \right] \right],$$

where

$$\mathbf{x}_{N_0} = (x_{j,N_0}, j = 1, \dots, k) \subseteq \mathbb{R}^{N_0}, \quad x_{j,N_0} \in \mathbb{R}^{N_0},$$

the expectation is over $W_{ij}^{(\ell)} \sim \mathcal{N}(0, \sigma^2/N_{\ell-1})$ and

$$A_\beta = \det(X^T X)^{1/2} (2\pi\beta)^{P/2}.$$

The first step in proving Theorem 3.1 is to write

$$Y = \theta_*^T X,$$

where here and throughout the proof we suppress the subscripts in $X_{N_0}, Y_{N_0}, \theta_{*,N_0}$. This yields

$$Z_\beta(\mathbf{x}_{N_0}, \mathbf{t}) = A_\beta \mathbb{E}_{\text{prior}} \left[\exp \left[-\frac{\beta}{2} \left\| \theta_*^T X - \prod_{\ell=1}^{L+1} W^{(\ell)} X \right\|_2^2 - i \langle f(\mathbf{x}_{N_0}), \mathbf{t} \rangle \right] \right]. \quad (4.23)$$

Next, since $f(x)$ is a linear function of x , we have

$$\langle f(\mathbf{x}_{N_0}), \mathbf{t} \rangle = f(\mathbf{x}_{N_0} \cdot \mathbf{t}),$$

and hence, suppressing the dependence on N_0 , we will write

$$Z_\beta(\mathbf{x}_{N_0}, \mathbf{t}) = Z_\beta(x, 1) =: Z_\beta(x), \quad x := \mathbf{x}_{N_0} \cdot \mathbf{t}.$$

To prove 3.1, we first derive the following expression for $Z_\beta(x)$ for general β .

Proposition 4.5. *For any $\beta > 0$, the partition function $Z_\beta(x)$ equals*

$$\begin{aligned} & \prod_{\ell=1}^L \Gamma \left(\frac{N_\ell}{2} \right)^{-1} \exp[-i\theta_*^T x_{\parallel}] \\ & \times \int_{\text{col}(X)} d\zeta \exp \left[-\frac{\|X^\dagger(\tau - x_{\parallel})\|^2}{2\beta} + i\theta_*^T \zeta \right] G_{L,1}^{1,L} \left(M \|\zeta\|^2 + M \|x_{\perp}\|^2 \mid 1 - \frac{N}{2} \right), \end{aligned}$$

where X^\dagger is the pseudo-inverse of X and

$$4M = \prod_{\ell=0}^L \frac{2\sigma^2}{N_\ell}.$$

Proof. We begin the proof of Proposition 4.5 by integrating out the final layer weights $W^{(L+1)}$ and introducing a dual variable $t \in \mathbb{R}^P$, as in the following.

Lemma 4.6. *Write*

$$X^L := W^{(L)} \dots W^{(1)} X, \quad x^L = W^{(L)} \dots W^{(1)} x.$$

We have

$$Z_\beta(x) = \det(X^T X)^{1/2} \int_{\mathbb{R}^P} \mathbb{E}_{\text{prior}} \left[\exp \left[-\frac{\|t\|_2^2}{2\beta} + i\theta_*^T X t - \frac{\sigma^2}{2N_L} \|X^L t + x^L\|^2 \right] \right] dt, \quad (4.24)$$

Proof. From the following identity

$$1 = \int_{\mathbb{R}^P} \frac{dt}{(2\pi\beta)^{P/2}} \exp \left[-\frac{1}{2\beta} \left\| t^T - i\beta \left(\theta_*^T X - W^{(L+1)} X^L \right) \right\|^2 \right]$$

we conclude

$$\begin{aligned} & \exp \left[-\frac{\beta}{2} \left\| \theta_*^T X - W^{(L+1)} X^L \right\|^2 \right] \\ &= \int \frac{dt}{(2\pi\beta)^{P/2}} \exp \left[-\frac{1}{2\beta} \left\| t^T - i\beta \left(\theta_*^T X - W^{(L+1)} X^L \right) \right\|^2 - \frac{\beta}{2} \left\| \theta_*^T X - W^{(L+1)} X^L \right\|^2 \right] \\ &= \int \frac{dt}{(2\pi\beta)^{P/2}} \exp \left[-\frac{1}{2\beta} \|t\|^2 + it \left(\theta_*^T X - W^{(L+1)} X^L \right) \right]. \end{aligned}$$

Substituting this into (4.23) we find

$$Z_\beta(x) = \frac{A_\beta}{(2\pi\beta)^{P/2}} \int_{\mathbb{R}^P} \mathbb{E}_{\text{prior}} \left[\exp \left[i \left(\theta_*^T X t - W^{(L+1)} (X^L t + x^L) \right) - \frac{\|t\|_2^2}{2\beta} \right] \right] dt$$

Now we compute the expectation over the final layer weights $W^{(L+1)}$ by completing the square:

$$\begin{aligned} -\frac{N_L}{2\sigma^2} \left\| W^{(L+1)} \right\|_2^2 - iW^{(L+1)} (X^L t + x^L) &= -\frac{N_L}{2\sigma^2} \left[\left\| W^{(L+1)} - i\frac{\sigma^2}{N_L} (X^L t + x^L) \right\|_2^2 \right] \\ &\quad - \frac{\sigma^2}{2N_L} \|X^L t + x^L\|_2^2. \end{aligned}$$

This yields

$$Z_\beta(x) = \frac{A_\beta}{(2\pi\beta)^{P/2}} \int_{\mathbb{R}^P} \mathbb{E}_{\text{prior}} \left[\exp \left[-\frac{\|t\|_2^2}{2\beta} + i\theta_*^T X t - \frac{\sigma^2}{2N_L} \|X^L t + x^L\|^2 \right] \right] dt,$$

completing the proof. \square

For each fixed t , note that

$$X^L t + x^L = W^{(L)} \dots W^{(1)} (X t + x).$$

Hence, the expectation

$$\mathbb{E}_{\text{prior}} \left[\exp \left[-\frac{\sigma^2}{2N_L} \|X^L t + x^L\|^2 \right] \right]$$

equals the Laplace transform

$$\mathcal{L}_{Q_{N,L}} \left(M' \|Xt + x\|^2 \right)$$

of the random variable

$$Q_{N,L} = \left\| \left\{ \prod_{\ell=1}^L \widehat{W}^{(\ell)} \right\} u \right\|^2, \quad \widehat{W}^{(\ell)} \sim \mathcal{N} \left(0, I_{N_\ell \times N_{\ell-1}} \right) \text{ independent}, \quad (4.25)$$

where u is any unit vector (the distribution is the same for any u since W^1 is rotationally invariant) evaluated at $\|Xt + x\|^2$ times

$$M' := \frac{1}{2} \left(\prod_{\ell=1}^L \frac{\sigma^2}{N_\ell} \right)$$

Thus, we obtain

$$Z_\beta(x) = \det(X^T X)^{1/2} \int dt \exp \left[-\frac{\|t\|^2}{2\beta} + i\theta_*^T X t \right] \mathcal{L}_{Q_{N,L}} \left(M' \|Xt + x\|^2 \right), \quad (4.26)$$

To proceed we rewrite the Laplace transform in the preceding line in terms of a Meijer-G function.

Lemma 4.7. *Let $Q_{N,L}$ be defined as in (4.25). Then,*

$$\mathcal{L}_{Q_{N,L}}(\tau) = \left(\prod_{\ell=1}^L \Gamma \left(\frac{N_\ell}{2} \right) \right)^{-1} G_{L,1}^{1,L} \left(2^L \tau \mid 1 - \frac{N_L}{2}, \dots, 1 - \frac{N_1}{2} \right).$$

Proof. The density of $Q_{N,L}^{1/2}$ is known from [ZVP21]:

$$p_{Q_{N,L}^{1/2}}(\rho) = \frac{2^{1-L/2}}{\Gamma \left(\frac{N_1}{2} \right) \dots \Gamma \left(\frac{N_L}{2} \right)} G_{0,L}^{L,0} \left(\frac{\rho^2}{2L} \mid \frac{N_L-1}{2}, \dots, \frac{N_1-1}{2} \right).$$

Hence, the density of $\widehat{Q}_{N,L}$ is

$$\begin{aligned} p_{Q_{N,L}}(\rho) &= \frac{1}{2\rho^{1/2}} p_{Q_{N,L}}(\rho^{1/2}) \\ &= \frac{2^{-L}}{\Gamma \left(\frac{N_1}{2} \right) \dots \Gamma \left(\frac{N_L}{2} \right)} \left(\frac{2^L}{\rho} \right)^{1/2} G_{0,L}^{L,0} \left(\frac{\rho}{2L} \mid \frac{N_L-1}{2}, \dots, \frac{N_1-1}{2} \right) \\ &= \frac{2^{-L}}{\Gamma \left(\frac{N_1}{2} \right) \dots \Gamma \left(\frac{N_L}{2} \right)} G_{0,L}^{L,0} \left(\frac{\rho}{2L} \mid \frac{N_L}{2} - 1, \dots, \frac{N_1}{2} - 1 \right), \end{aligned}$$

where in the last step we've used (4.9). Hence,

$$\begin{aligned}
\mathcal{L}_{Q_{N,L}}(\tau) &= \int_0^\infty e^{-\tau\rho} G_{0,L}^{L,0} \left(\frac{\rho}{2^L} \mid \frac{N_L}{2} - 1, \dots, \frac{N_1}{2} - 1 \right) d\rho \\
&= \left(\Gamma\left(\frac{N_1}{2}\right) \cdots \Gamma\left(\frac{N_L}{2}\right) \right)^{-1} \frac{1}{2^{L\tau}} G_{1,L}^{L,1} \left(\frac{1}{2^{L\tau}} \mid \frac{N_L}{2} - 1, \dots, \frac{N_1}{2} - 1 \right) \\
&= \left(\Gamma\left(\frac{N_1}{2}\right) \cdots \Gamma\left(\frac{N_L}{2}\right) \right)^{-1} G_{1,L}^{L,1} \left(\frac{1}{2^{L\tau}} \mid \frac{N_L}{2}, \dots, \frac{N_1}{2} \right) \\
&= \left(\Gamma\left(\frac{N_1}{2}\right) \cdots \Gamma\left(\frac{N_L}{2}\right) \right)^{-1} G_{L,1}^{1,L} \left(2^{L\tau} \mid 1 - \frac{N_L}{2}, \dots, 1 - \frac{N_1}{2} \right)
\end{aligned}$$

where we've used (4.7) and (4.8). \square

Combining the preceding Lemma with (4.26) gives

$$Z(x) = \frac{\det(X^T X)^{1/2}}{\prod_{\ell=1}^L \Gamma\left(\frac{N_\ell}{2}\right)} \int_{\mathbb{R}^P} dt \exp \left[-\frac{\|t\|^2}{2\beta} + i\theta_*^T X t \right] G_{L,1}^{1,L} \left(M\|Xt + x\|^2 \mid 1 - \frac{N_L}{2}, \dots, 1 - \frac{N_1}{2} \right), \quad (4.27)$$

where

$$4M = \prod_{\ell=0}^L \frac{2\sigma^2}{N_\ell}.$$

To complete the proof of Proposition 4.5, we write

$$x = x_{\parallel} + x_{\perp}, \quad x_{\parallel} \in \text{col}(X), \quad x_{\perp} \in \text{col}(X)^\perp.$$

Note that X gives an isomorphism from \mathbb{R}^P to $\text{im}(X)$. Thus, we may change variables to write $Z_\beta(x)$ as

$$\prod_{\ell=1}^L \Gamma\left(\frac{N_\ell}{2}\right)^{-1} \int_{\mathbb{R}^P} d\zeta \exp \left[-\frac{\|X^\dagger \zeta\|^2}{2\beta} + i\theta_*^T \zeta \right] G_{L,1}^{1,L} \left(M\|\zeta + x\|^2 \mid 1 - \frac{N}{2} \right),$$

where X^\dagger the pseudo-inverse of X . Finally, note that

$$\|\zeta + x\|^2 = \|\zeta + x_{\parallel}\|^2 + \|x_{\perp}\|^2.$$

Hence, changing variables to $\tau = \zeta + x_{\parallel}$ yields the following expression for $Z_\beta(x)$:

$$\prod_{\ell=1}^L \Gamma\left(\frac{N_\ell}{2}\right)^{-1} \exp[-i\theta_*^T x_{\parallel}] \times \int_{\text{col}(X)} d\zeta \exp \left[-\frac{\|X^\dagger(\tau - x_{\parallel})\|^2}{2\beta} + i\theta_*^T \zeta \right] G_{L,1}^{1,L} \left(M\|\zeta\|^2 + M\|x_{\perp}\|^2 \mid 1 - \frac{N}{2} \right)$$

This is precisely the statement of Proposition 4.5. \square

By taking $\beta \rightarrow \infty$ in Proposition 4.5, we see that

$$Z_\infty(x) = \frac{\exp[-i\theta_*^T x_{\parallel}]}{\prod_{\ell=1}^L \Gamma\left(\frac{N_\ell}{2}\right)} \int_{\mathbb{R}^P} \exp[i\theta_*^T \zeta] G_{L,1}^{1,L} \left(M\|\zeta\|^2 + M\|x_{\perp}\|^2 \mid 1 - \frac{N}{2} \right) d\zeta. \quad (4.28)$$

In order to simplify this expression further, we pass to polar coordinates

$$\begin{aligned} & \int_{\mathbb{R}^P} \exp [i\theta_*^T \zeta] G_{L,1}^{1,L} \left(M \|\zeta\|^2 + M \|x_\perp\|^2 \mid \begin{matrix} 1 - \frac{N}{2} \\ 0 \end{matrix} \right) d\zeta \\ &= \int_0^\infty \rho^{P-1} \left\{ \int_{S^{P-1}} \exp [i\rho\theta_*^T \theta] d\theta \right\} G_{L,1}^{1,L} \left(M\rho^2 + M \|x_\perp\|^2 \mid \begin{matrix} 1 - \frac{N}{2} \\ 0 \end{matrix} \right) d\rho. \end{aligned}$$

By the definition of the Bessel function and the relation (4.10), we have

$$\begin{aligned} \int_{S^{P-1}} \exp [i\rho\theta_*^T \theta] d\theta &= (2\pi)^{P/2} (\rho \|\theta_*\|)^{-\frac{P-2}{2}} J_{\frac{P-2}{2}} (\rho \|\theta_*\|) \\ &= (2\pi)^{P/2} (\rho^2 \|\theta_*\|^2)^{-\frac{P-2}{4}} G_{0,2}^{1,0} \left(\frac{\rho^2 \|\theta_*\|^2}{4} \mid \begin{matrix} - \\ \frac{P-2}{4}, -\frac{P-2}{4} \end{matrix} \right) \\ &= 2\pi^{P/2} G_{0,2}^{1,0} \left(\frac{\rho^2 \|\theta_*\|^2}{4} \mid \begin{matrix} - \\ 0, -\frac{P-2}{2} \end{matrix} \right). \end{aligned}$$

We therefore obtain

$$\begin{aligned} & \int_{\mathbb{R}^P} \exp [i\theta_*^T \zeta] G_{L,1}^{1,L} \left(M \|\zeta\|^2 + M \|x_\perp\|^2 \mid \begin{matrix} 1 - \frac{N}{2} \\ 0 \end{matrix} \right) d\zeta \\ &= \pi^{P/2} \int_0^\infty \rho^{\frac{P-2}{2}} G_{0,2}^{1,0} \left(\frac{\rho \|\theta_*\|^2}{4} \mid \begin{matrix} - \\ 0, -\frac{P-2}{2} \end{matrix} \right) G_{L,1}^{1,L} \left(M\rho + M \|x_\perp\|^2 \mid \begin{matrix} 1 - \frac{N}{2} \\ 0 \end{matrix} \right) d\rho \\ &= \left(\frac{4}{\|\theta_*\|^2} \right)^{\frac{P-2}{2}} \pi^{P/2} \int_0^\infty G_{0,2}^{1,0} \left(\frac{\rho \|\theta_*\|^2}{4} \mid \begin{matrix} - \\ \frac{P-2}{2}, 0 \end{matrix} \right) G_{L,1}^{1,L} \left(M\rho + M \|x_\perp\|^2 \mid \begin{matrix} 1 - \frac{N}{2} \\ 0 \end{matrix} \right) d\rho \\ &= \left(\frac{4}{\|\theta_*\|^2} \right)^{\frac{P}{2}} \pi^{P/2} \frac{\|\theta_*\|^2}{4M} \int_0^\infty G_{0,2}^{1,0} \left(\frac{\rho \|\theta_*\|^2}{4M} \mid \begin{matrix} - \\ \frac{P-2}{2}, 0 \end{matrix} \right) G_{L,1}^{1,L} \left(\rho + M \|x_\perp\|^2 \mid \begin{matrix} 1 - \frac{N}{2} \\ 0 \end{matrix} \right) d\rho. \end{aligned}$$

We now apply (4.13) to find

$$\begin{aligned} & \int_0^\infty G_{0,2}^{1,0} \left(\frac{\rho \|\theta_*\|^2}{4M} \mid \begin{matrix} - \\ \frac{P-2}{2}, 0 \end{matrix} \right) G_{L,1}^{1,L} \left(\rho + M \|x_\perp\|^2 \mid \begin{matrix} 1 - \frac{N}{2} \\ 0 \end{matrix} \right) d\rho \\ &= \sum_{k=0}^\infty \frac{1}{k!} (-M \|x_\perp\|^2)^k G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_*\|^2}{4M} \mid \begin{matrix} - \\ \frac{P-2}{2}, \frac{N}{2} + k - 1 \end{matrix} \right). \end{aligned}$$

Therefore,

$$\begin{aligned} & \int_{\mathbb{R}^P} \exp [i\theta_*^T \zeta] G_{L,1}^{1,L} \left(M \|\zeta\|^2 + M \|x_\perp\|^2 \mid \begin{matrix} 1 - \frac{N}{2} \\ 0 \end{matrix} \right) d\zeta \\ &= \left(\frac{4}{\|\theta_*\|^2} \right)^{\frac{P}{2}} \pi^{P/2} \sum_{k=0}^\infty \frac{1}{k!} (-M \|x_\perp\|^2)^k G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_*\|^2}{4M} \mid \begin{matrix} - \\ \frac{P}{2}, \frac{N}{2} + k \end{matrix} \right). \end{aligned}$$

Putting this all together yields

$$Z_\infty(x) = \left(\frac{4\pi}{\|\theta_*\|^2} \right)^{\frac{P}{2}} \prod_{\ell=1}^L \Gamma \left(\frac{N_\ell}{2} \right)^{-1} \sum_{k=0}^\infty \frac{1}{k!} (-M \|x_\perp\|^2)^k G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_*\|^2}{4M} \mid \begin{matrix} - \\ \frac{P}{2}, \frac{N}{2} + k \end{matrix} \right),$$

completing the proof. \square

4.5 Proof of Theorem 3.2

In this section, we derive several apparently novel asymptotic expansion of the Meijer-G functions of the form

$$G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_*\|^2}{4M} \mid \frac{P}{2}, \frac{N}{2} + k \right) := G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_*\|^2}{4M} \mid \frac{P}{2}, \frac{N_1}{2} + k, \dots, \frac{N_L}{2} + k \right). \quad (4.29)$$

Our first step is to obtain a contour integral representation of the Meijer-G functions we are studying. To state the exact result, consider the following independent Γ random variables:

$$\phi_j \sim \begin{cases} \Gamma \left(\frac{N_j}{2} + k + 1, \frac{2\sigma^2}{N_j} \right), & j = 1, \dots, L \\ \Gamma \left(\frac{P}{2} + 1, \frac{2\sigma^2}{P} \frac{\alpha_0}{\|\theta_*\|^2} \right), & j = 0 \end{cases}. \quad (4.30)$$

As we recalled in §4.2.2, the moments of ϕ_j can be explicitly written in terms of Γ functions. Moreover, the Meijer-G functions (4.29) can be interpreted, up to a scaling factor, as densities of the product of products of ϕ_j 's. This allows us to obtain the following

Lemma 4.8. *Fix $L, N, N_0, \dots, N_L \geq 1$ as well as $\|\theta_*\| > 0$ and define M by*

$$4M := \prod_{\ell=0}^L \frac{2\sigma^2}{N_\ell}. \quad (4.31)$$

For any $N \geq 1$ we have

$$\begin{aligned} \text{Den}_{\phi_0 \dots \phi_L}(1) &= \frac{\|\theta_*\|^2}{4M} G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_*\|^2}{4M} \mid \frac{P}{2}, \frac{N}{2} + k \right) \left[\Gamma \left(\frac{P}{2} + 1 \right) \right]^{-1} \prod_{\ell=1}^L \left[\Gamma \left(\frac{N_\ell}{2} + k + 1 \right) \right]^{-1} \\ &= \frac{1}{2\pi} \int_{\mathcal{C}} \exp[\Phi(z)] dz, \end{aligned} \quad (4.32)$$

where $\mathcal{C} \subseteq \mathbb{C}$ is the contour that runs along the real line from $-\infty$ to ∞ and

$$\begin{aligned} \Phi(z) &= -iz \log \left(\frac{2\sigma^2}{P} \frac{\alpha_0}{\|\theta_*\|^2} \right) + \log \left(\frac{\Gamma \left(\frac{P}{2} + 1 - iz \right)}{\Gamma \left(\frac{P}{2} + 1 \right)} \right) \\ &\quad + \sum_{\ell=1}^L \left\{ -iz \log \left(\frac{2\sigma^2}{N_\ell} \right) + \log \left(\frac{\Gamma \left(\frac{N_\ell}{2} + k + 1 - iz \right)}{\Gamma \left(\frac{N_\ell}{2} + k + 1 \right)} \right) \right\}. \end{aligned} \quad (4.33)$$

Proof. We use the relationship (4.14) between G functions and densities of products of Gamma random variables to write

$$\prod_{\ell=1}^L \left[\Gamma \left(\frac{N_\ell}{2} + k + 1 \right) \right]^{-1} \left[\Gamma \left(\frac{P}{2} + 1 \right) \right]^{-1} \frac{\|\theta_*\|^2}{4M} G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_*\|^2}{4M} \mid \frac{P}{2}, \frac{N}{2} + k \right) = \text{Den}_{\phi_0 \dots \phi_L}(1).$$

For any positive random variable X with density $d\mathbb{P}_X$, we have

$$d\mathbb{P}_X(t) = t^{-1} d\mathbb{P}_{\log(X)}(\log(t)).$$

Hence,

$$\text{Den}_{\phi_0 \dots \phi_L}(1) = \text{Den}_{\sum_{\ell=0}^L \log \phi_\ell}(0)$$

is proportional to the density of a sum of independent random variables evaluated at $\log(1) = 0$. Further, recalling (4.4), we find by Fourier inversion that

$$\frac{1}{2\pi} \int_{-\infty}^{\infty} \exp[\Phi(t)] dt$$

equals

$$\prod_{\ell=1}^L \left[\Gamma\left(\frac{N_\ell}{2} + k + 1\right) \right]^{-1} \left[\Gamma\left(\frac{P}{2} + 1\right) \right]^{-1} \frac{\|\theta_*\|^2}{4M} G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_*\|^2}{4M} \mid \frac{P}{2}, \frac{N}{2} + k \right),$$

where

$$\begin{aligned} \Phi(t) &= -it \log 1 - it \log \left(\frac{2\sigma^2}{P} \frac{\alpha_0}{\|\theta_*\|^2} \right) + \log \left(\frac{\Gamma\left(\frac{P}{2} + 1 - it\right)}{\Gamma\left(\frac{P}{2} + 1\right)} \right) \\ &+ \sum_{\ell=1}^L \left\{ -it \log \left(\frac{2\sigma^2}{N_\ell} \right) + \log \left(\frac{\Gamma\left(\frac{N_\ell}{2} + k + 1 - it\right)}{\Gamma\left(\frac{N_\ell}{2} + k + 1\right)} \right) \right\}. \end{aligned} \quad (4.34)$$

□

Note that (4.43) expresses G as a contour integral of a meromorphic function $\exp(N\Psi)$ with poles at

$$-\frac{m}{2} - k - 1 - \nu, \quad \nu \in \mathbb{N}, m \in \{P, N_1, \dots, N_L\}.$$

Our goal is to evaluate the contour integral representation (4.32) for the Meijer-G function using the Laplace method. To do so, we will use the following standard procedure:

1. Contour deformation: \mathcal{C} into a union of several contours $\mathcal{C}_1 \cup \dots \cup \mathcal{C}_R$. Contours $\mathcal{C}_1, \mathcal{C}_R$ are non-compact, whereas the contours $\mathcal{C}_2, \dots, \mathcal{C}_{R-1}$ do not extend to infinity. We will need to choose the contours $\mathcal{C}_2, \dots, \mathcal{C}_{R-1}$ so that exactly one of them passes through what will turn out to be the dominant critical point ζ_* of Ψ and does so in the direction of steepest descent. Moreover, on the contour \mathcal{C}_2 , the imaginary part of the phase will be constant (in fact equal to 0).
2. Localization to compact domain of integration: Show that the integrand $\exp[N\Psi(z)]$ is integrable and exponentially small in N on the contours $\mathcal{C}_1, \mathcal{C}_R$. In particular, for any $K > 0$, we will find that modulo errors of size $O(e^{-KN})$ we may therefore focus on the integral over $\mathcal{C}_2, \dots, \mathcal{C}_{R-1}$. The ability to choose K will be important since the entire integral is exponentially small in N .
3. Computing derivatives of Ψ at ζ_* : Now that we have reduced the integral (4.32) to a compact domain of integration it remains only to check that the critical point is non-degenerate, to compute $\Psi(\zeta_*)$, $\frac{d}{dz}\Psi(\zeta_*)$, $\frac{d^2}{dz^2}\Psi(\zeta_*)$, and to apply the Laplace method.

We now proceed to give the details in the case when

$$L \text{ is fixed, } N := \min \{N_0, P, N_\ell\} \rightarrow \infty, \quad \frac{P}{N_0} \rightarrow \alpha_0 \in (0, 1), \quad \frac{P}{N} \rightarrow \alpha \in (0, \infty).$$

This is a generalization of case (a) of Theorem 3.2. To proceed, we make the change of variables $t \mapsto Nt$ in (4.32) to get that

$$\frac{\|\theta_*\|^2}{4M} G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_*\|^2}{4M} \mid \frac{P}{2}, \frac{N}{2} + k \right) \left[\Gamma \left(\frac{P}{2} + 1 \right) \right]^{-1} \prod_{\ell=1}^L \left[\Gamma \left(\frac{N_\ell}{2} + k + 1 \right) \right]^{-1} = \frac{N}{2\pi} \int_{\mathcal{C}} \exp[N\Psi(z)] dz, \quad (4.35)$$

where

$$\begin{aligned} \Psi(z) = & -iz \log \left(\frac{2\sigma^2}{P} \frac{\alpha_0}{\|\theta_*\|^2} \right) + \frac{1}{N} \log \left(\frac{\Gamma \left(\frac{P}{2} + 1 - iNz \right)}{\Gamma \left(\frac{P}{2} + 1 \right)} \right) \\ & + \sum_{\ell=1}^L \left\{ -iz \log \left(\frac{2\sigma^2}{N_\ell} \right) + \frac{1}{N} \log \left(\frac{\Gamma \left(\frac{N_\ell}{2} + k + 1 - iNz \right)}{\Gamma \left(\frac{N_\ell}{2} + k + 1 \right)} \right) \right\}. \end{aligned} \quad (4.36)$$

With this rescaling, we now deform the contour of integration as follows by fixing constants $\delta, T > 0$ (the value of T will be determined by Lemma 4.9 and the value of δ will be determined by (4.48) and the sentence directly after it) and deforming the contour \mathcal{C} as follows

$$\mathcal{C} \mapsto \bigcup_{j=0}^4 \mathcal{C}_j,$$

where

$$\begin{aligned} \mathcal{C}_0 &= \mathcal{C}_0(T) = (-\infty, -T] \\ \mathcal{C}_1 &= \mathcal{C}_1(T, C_\delta) = \text{linear interpolation from } -T \in \mathbb{R} \text{ to } -iC_\delta \in i\mathbb{R} \\ \mathcal{C}_2 &= \mathcal{C}_2(C_\delta) = [-iC_\delta, iC_\delta] \\ \mathcal{C}_3 &= \mathcal{C}_3(T, C_\delta) = \text{linear interpolation from } iC_\delta \in i\mathbb{R} \text{ to } T \in \mathbb{R} \\ \mathcal{C}_4 &= \mathcal{C}_4(T) = [T, \infty), \end{aligned}$$

where

$$C_\delta := \delta - \min \left\{ \frac{P}{2N}, \frac{N_1}{2N}, \dots, \frac{N_L}{2N} \right\}.$$

See Figure 3. In order to evaluate the integral over each \mathcal{C}_j , it will be convenient to introduce

$$\Psi(z) = \sum_{\ell=0}^L \Psi_\ell(z)$$

where

$$\Psi_\ell(z) = \begin{cases} -iz \log \left(\frac{2\sigma^2}{P} \frac{\alpha_0}{\|\theta_*\|^2} \right) + \frac{1}{N} \log \left(\frac{\Gamma \left(\frac{P}{2} + 1 - iNz \right)}{\Gamma \left(\frac{P}{2} + 1 \right)} \right), & \ell = 0 \\ -iz \log \left(\frac{2\sigma^2}{N_\ell} \right) + \frac{1}{N} \log \left(\frac{\Gamma \left(\frac{N_\ell}{2} + k + 1 - iNz \right)}{\Gamma \left(\frac{N_\ell}{2} + k + 1 \right)} \right), & \ell = 1, \dots, L \end{cases}. \quad (4.37)$$

The following Lemma allows us to throw away the contribution to (4.43) coming from $\mathcal{C}_0, \mathcal{C}_4$.

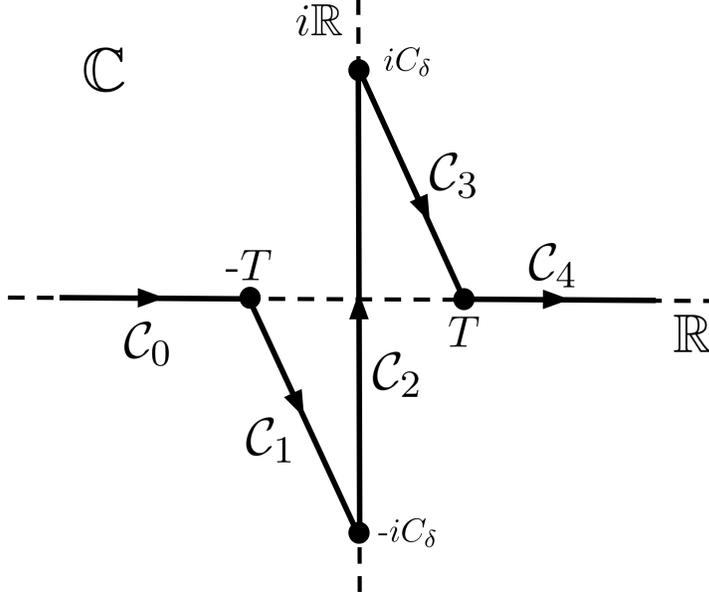


Figure 3: Deformed contour of integration for the proof of Theorem 3.2.

Lemma 4.9. *There exist $c, T_0 > 0$ such that*

$$\sup_{\substack{|z| > T \\ z \in \mathbb{R}}} \frac{\Re \Psi(z)}{1 + |z|} \leq -c, \quad \forall T \geq T_0. \quad (4.38)$$

Moreover, for any $T, \delta > 0$, writing

$$S_{N, \delta, T} := \{z \in \mathbb{C} \mid |z| < T, \Im(z) > C_\delta\}, \quad (4.39)$$

there exists $C > 0$ such that

$$\sup_{z \in S_{N, \delta, T}} \max \left\{ |\Psi(z)|, \left| \frac{d}{dz} \Psi(z) \right| \right\} \leq C. \quad (4.40)$$

Proof. To show both (4.38) and (4.40) we will use the asymptotic expansion

$$\log \Gamma(z) \sim \left(z - \frac{1}{2}\right) \log(z) - z + \text{const} + O(|z|^{-1}), \quad \text{as } |z| \rightarrow \infty, \quad (4.41)$$

which holds uniformly on sets of the form $|z| < \pi - \epsilon$ for a fixed $\epsilon > 0$. With Ψ_ℓ defined in

(4.37), we find for each $\ell = 1, \dots, L$ that uniformly over $z \in \mathbb{R}$

$$\begin{aligned}
\Psi_\ell(z) &= -iz \log\left(\frac{2\sigma^2}{N_\ell}\right) + \frac{1}{N} \left(\log \Gamma\left(\frac{N_\ell}{2} + k + 1 - iNz\right) - \log \Gamma\left(\frac{N_\ell}{2} + k + 1\right) \right) \\
&= -iz \left[-1 + \log\left(\frac{2\sigma^2}{N_\ell}\right) + \log\left(\frac{N_\ell}{2} + k + 1 - iNz\right) \right] \\
&\quad + \frac{1}{N} \left(\frac{N_\ell + 1}{2} + k + 1 \right) \log\left(1 - \frac{iNz}{\frac{N_\ell}{2} + k + 1}\right) + O(N^{-1}) \\
&= -iz \left[-1 + \log(\sigma^2) + \log\left(1 - \frac{2iNz}{N_\ell}\right) \right] + \frac{N_\ell}{2N} \log\left(1 - \frac{2iNz}{N_\ell} \cdot \frac{1}{1 + \frac{2(k+1)}{N_\ell}}\right) + O(N^{-1}) \\
&= -iz \left[-1 + \log(\sigma^2) + \log\left(1 - \frac{2iNz}{N_\ell}\right) \right] + \frac{N_\ell}{2N} \log\left(1 - \frac{2iNz}{N_\ell}\right) + O(N^{-1}).
\end{aligned}$$

Obtaining a similar expression for $\ell = 0$ and summing over ℓ proves (4.40). Moreover, for $\ell = 1, \dots, L$ we obtain uniformly on $z \in \mathbb{R}$

$$\Re \Psi_\ell(z) = z \arg\left(1 - \frac{2iNz}{N_\ell}\right) + \frac{N_\ell}{2N} \log\left|1 - \frac{2iNz}{N_\ell}\right| + O(N^{-1}).$$

Note that for T_0 sufficiently large there we have

$$|z| > T_0 \quad \Rightarrow \quad \arg\left(1 - \frac{2iNz}{N_\ell}\right) \operatorname{sgn}(z) \leq -\frac{\pi}{4} \quad \Rightarrow \quad \Re \Psi_\ell(z) \leq -\frac{\pi}{8} |z|. \quad (4.42)$$

A similar analysis applies to $\ell = 0$ and completes the proof of (4.38). \square

Estimate (4.38) in the previous Lemma shows that, for any $K, \delta > 0$ there exists $T > 0$ such that

$$\frac{N}{2\pi} \int_C \exp[N\Psi(z)] dz = \frac{N}{2\pi} \int_{C_1(T, C_\delta) \cup C_2(C_\delta) \cup C_3(T, C_\delta)} \exp[N\Psi(z)] dz + O(e^{-KN}). \quad (4.43)$$

To evaluate the right hand side of (4.43), we derive the following uniform asymptotic expansion for $\Psi(z)$.

Lemma 4.10. *Fix $\delta, T > 0$ and define $S_{N, \delta, T}$ as in (4.39). Then, uniformly over $S_{N, \delta, T}$, we have*

$$\Psi(z) = \widehat{\Psi}(z) + \frac{1}{N} \widetilde{\Psi}(z) + O(N^{-2}),$$

where

$$\widehat{\Psi}(z) := iz \left[\log\left(\frac{\|\theta_*\|^2}{\sigma^{2(L+1)} \alpha_0}\right) + L + 1 - \log\left(1 - \frac{2iNz}{P}\right) - \sum_{\ell=1}^L \log\left(1 - \frac{2iNz}{N_\ell}\right) \right] \quad (4.44)$$

$$+ \frac{P}{2N} \log\left(1 - \frac{2iNz}{P}\right) + \sum_{\ell=1}^L \frac{N_\ell}{2N} \log\left(1 - \frac{2iNz}{N_\ell}\right)$$

$$\widetilde{\Psi}(z) := \frac{1}{2} \log\left(1 - \frac{2iNz}{P}\right) + \left(k + \frac{1}{2}\right) \sum_{\ell=1}^L \log\left(1 - \frac{2iNz}{N_\ell}\right). \quad (4.45)$$

Proof. This follows directly from expanding $\Psi_\ell(z)$ using (4.1) for each $\ell = 0, \dots, L$. \square

Combining the preceding Lemma with (4.43) shows that for any $K > 0$ there exists $T, \delta > 0$ so that

$$\begin{aligned} & \frac{N}{2\pi} \int_{\mathcal{C}} \exp[N\Psi(z)] dz \\ &= \frac{N}{2\pi} \int_{\mathcal{C}_1(T, C_\delta) \cup \mathcal{C}_2(C_\delta) \cup \mathcal{C}_3(T, C_\delta)} \exp\left[N\widehat{\Psi}(z) + \widetilde{\Psi}(z)\right] dz (1 + O(N^{-1})) + O(e^{-KN}) \end{aligned} \quad (4.46)$$

Moreover, differentiating (4.44) yields

$$\frac{d}{dz} \widehat{\Psi}(z) = i \left[\log \left(\frac{\|\theta_*\|^2}{\sigma^{2(L+1)} \alpha_0} \right) - \log \left(1 - \frac{2iNz}{P} \right) - \sum_{\ell=1}^L \log \left(1 - \frac{2iNz}{N_\ell} \right) \right].$$

Computing the real part of both sides show that

$$\Re \left(\frac{d}{dz} \widehat{\Psi}(z) \right) = \arg \left(1 - \frac{2iNz}{P} \right) + \sum_{\ell=1}^L \arg \left(1 - \frac{2iNz}{N_\ell} \right).$$

Since all the arguments have the same sign, we conclude that the real part of $\frac{d}{dz} \widehat{\Psi}$ vanishes only when $z = i\zeta$ for $\zeta \in \mathbb{R}$. Further,

$$\Im \left(\frac{d}{d\zeta} \widehat{\Psi}(i\zeta) \right) = -\log \left(\frac{\|\theta_*\|^2}{\sigma^{2(L+1)} \alpha_0} \right) + \log \left| 1 + \frac{2N\zeta}{P} \right| + \sum_{\ell=1}^L \log \left| 1 + \frac{2N\zeta}{N_\ell} \right|,$$

which vanishes on \mathcal{C}_2 if and only if $\zeta = \zeta_*$ is the unique solution to

$$\frac{\|\theta_*\|^2}{\sigma^{2(L+1)} \alpha_0} = \left(1 + \frac{2N\zeta_*}{P} \right) \prod_{\ell=1}^L \left(1 + \frac{2N\zeta_*}{N_\ell} \right). \quad (4.47)$$

Indeed, observe that the right hand side of the equation on the preceding line increases monotonically in ζ_* from $-\infty$ to $+\infty$ as ζ_* varies in $(-\min\{\frac{P}{2N}, \frac{N_1}{2N}, \dots, \frac{N_L}{2N}\}, \infty)$. Computing the correction to the saddle point, we differentiate

$$\begin{aligned} \Im \left(\frac{d}{d\zeta} \left(\widehat{\Psi}(i\zeta) + \frac{1}{N} \widetilde{\Psi}(i\zeta) \right) \right) &= -\log \left(\frac{\|\theta_*\|^2}{\sigma^{2(L+1)} \alpha_0} \right) + \log \left| 1 + \frac{2N\zeta}{P} \right| + \sum_{\ell=1}^L \log \left| 1 + \frac{2N\zeta}{N_\ell} \right| \\ &\quad + \frac{1}{N} \left[\frac{1}{2\zeta + P/N} + \sum_{\ell=1}^L \frac{N}{N_\ell} \frac{1 + 2k}{1 + 2\zeta N_\ell/N} \right] \end{aligned}$$

and set the derivative to zero at $\zeta = \zeta_* + \zeta_{**}/N$, giving

$$0 = \frac{1}{N} \left[\frac{2\zeta_{**}}{2\zeta_* + P/N} + \frac{1}{2\zeta_* + P/N} + \sum_{\ell=1}^L \frac{2\zeta_{**}}{2\zeta_* + N_\ell/N} + \frac{N}{N_\ell} \frac{1 + 2k}{1 + 2\zeta_* N_\ell/N} \right].$$

This is solved by

$$\zeta_{**} = -\frac{1}{2} \frac{[2\zeta_* + P/N]^{-1} + (1+2k) \sum_{\ell} [2\zeta_* + N_{\ell}/N]^{-1}}{[2\zeta_* + P/N]^{-1} + \sum_{\ell} [2\zeta_* + N_{\ell}/N]^{-1}}. \quad (4.48)$$

Hence, by choosing δ sufficiently small, we can ensure that \mathcal{C}_2 contains the unique solution to (18) of Theorem 3.2 of the main text. Further, a direct computation shows that

$$\left. \frac{d^2}{d\zeta^2} \widehat{\Psi}(i\zeta) \right|_{\zeta=\zeta_*} = - \left[\frac{1}{\frac{P}{2N} + \zeta_*} + \sum_{\ell=1}^L \frac{1}{\frac{N_{\ell}}{2N} + \zeta_*} \right] < 0,$$

proving that that $i\zeta_*$ is a non-degenerate critical point of $\widehat{\Psi}$. Defining $\zeta_0 = \zeta_* + \zeta_{**}/N$ and

$$\Psi_0(\zeta_0) = \widehat{\Psi}(i\zeta_0) + \frac{1}{N} \widetilde{\Psi}(i\zeta_0),$$

the Laplace method gives

$$\begin{aligned} \log \text{Den}_{\sum_{\ell=0}^L \log \phi_{\ell}}(0) &= \log \left(\frac{N}{2\pi} \int_{\mathcal{C}_1(T, C_{\delta}) \cup \mathcal{C}_2(C_{\delta}) \cup \mathcal{C}_3(T, C_{\delta})} \exp[N\Psi(z)] dz \right) \\ &= N\Psi_0(\zeta_0) + \frac{1}{2} \log \left(\frac{N}{2\pi\Psi_0''(\zeta_0)} \right) \\ &= -N\zeta_0 \left[\log \left(\frac{\|\theta_*\|^2}{\sigma^{2(L+1)}\alpha_0} \right) + L + 1 - \log \left(1 + 2\zeta_0 \frac{N}{P} \right) - \sum_{\ell=1}^L \log \left(1 + 2\zeta_0 \frac{N}{N_{\ell}} \right) \right] \\ &\quad + \frac{N}{2} \left[\frac{P}{N} \log \left(1 + 2\zeta_0 \frac{N}{P} \right) + \frac{N_{\ell}}{N} \sum_{\ell=1}^L \log \left(1 + 2\zeta_0 \frac{N}{N_{\ell}} \right) \right] \\ &\quad + \frac{1}{2} \log \left(1 + 2\zeta_0 \frac{N}{P} \right) + \left(k + \frac{1}{2} \right) \sum_{\ell=1}^L \log \left(1 + 2\zeta_0 \frac{N}{N_{\ell}} \right) \\ &\quad + \frac{1}{2} \log(N) - \frac{1}{2} \log(2\pi) - \frac{1}{2} \log \left[\frac{2}{2\zeta_0 + P/N} + \sum_{\ell=1}^L \frac{2}{2\zeta_0 + N/N_{\ell}} \right]. \end{aligned}$$

Taking

$$N_{\ell} = N, \quad P = \alpha N$$

and observing that

$$\log \left(\frac{\|\theta_*\|^2}{\sigma^{2(L+1)}\alpha_0} \right) = L \log(1 + 2\zeta_*) + \log \left(1 + \frac{2\zeta_*}{\alpha} \right)$$

simplifies the density to

$$\begin{aligned} \log \text{Den}_{\sum_{\ell=0}^L \log \phi_{\ell}}(0) &= \frac{N}{2} \left[\alpha \left(\log \left(1 + \frac{2\zeta_*}{\alpha} \right) - \frac{2\zeta_*}{\alpha} \right) + L \left(\log(1 + 2\zeta_*) - 2\zeta_* \right) \right] + \frac{1}{2} \log(N) - \frac{1}{2} \log(\pi) \\ &\quad + \frac{1}{2} \log \left(1 + \frac{2\zeta_*}{\alpha} \right) + \left(k + \frac{1}{2} \right) L \log(1 + 2\zeta_*) - \frac{1}{2} \log \left(\frac{1}{\alpha + 2\zeta_*} + \frac{L}{1 + 2\zeta_*} \right). \end{aligned}$$

Solving for the G -function given (4.14), we conclude that

$$\begin{aligned}
\log G &= \log \text{Den}_{\sum_{\ell=0}^L \log \phi_\ell}(0) + L \log \left[\Gamma \left(\frac{N}{2} + k + 1 \right) \right] + \log \left[\Gamma \left(\frac{N\alpha}{2} + 1 \right) \right] \\
&\quad - \log \frac{\|\theta_*\|^2}{\alpha_0} - L \log \left(\frac{N}{2\sigma^2} \right) - \log \left(\frac{N\alpha}{2\sigma^2} \right) \\
&= \frac{N}{2} \left\{ \alpha \left[\log \left(\frac{N\alpha}{2} \right) + \log \left(1 + \frac{2\zeta_*}{\alpha} \right) - \left(1 + \frac{2\zeta_*}{\alpha} \right) \right] + L \left[\log \left(\frac{N}{2} \right) + \log(1 + 2\zeta_*) - (1 + 2\zeta_*) \right] \right\} \\
&\quad + \frac{L(2k-1)}{2} \left[\log \left(\frac{N}{2} \right) + \log(1 + 2\zeta_*) \right] + \frac{L}{2} \log(2\pi) - \frac{1}{2} \log \left(1 + \frac{\alpha + 2\zeta_*}{1 + 2\zeta_*} L \right).
\end{aligned}$$

The leading order part of this expression reproduces (16) of Theorem 3.2 of the main text, and taking differences between the value at a fixed k and at $k = 0$ gives (19). Note that at $L = 0$, direct computation gives

$$\log G = -\frac{\|\theta_*\|^2}{\alpha_0} \frac{N\alpha}{2\sigma^2} + \frac{N\alpha}{2} \left(\log \frac{N\alpha}{2\sigma^2} + \log \frac{\|\theta_*\|^2}{\alpha_0} \right),$$

which we see is reproduced by the above. The derivations of the formulas in cases (b) and (c) of Theorem 3.2 are very similar to those of (a). So we indicate only the salient differences, starting with case (b). We will actually consider the following somewhat more general regime:

$$N := \min \{N_1, \dots, N_L\} \rightarrow \infty, \quad P, N_0, N, L \rightarrow \infty, \quad \frac{P}{N_0} \rightarrow \alpha_0, \quad \sum_{\ell=1}^L \frac{1}{N_\ell} \rightarrow \lambda_{\text{prior}},$$

with $\alpha_0 \in (0, 1)$ and $\lambda_{\text{prior}} \in (0, \infty)$. Our starting point is again Lemma 4.8. The first modification in the proof is that we must redefine the contours \mathcal{C}_j , $j = 0, \dots, 4$ by replacing

$$T \mapsto T/L.$$

Next, Lemma 4.9 now reads.

Lemma 4.11. *There exist $c, T_0 > 0$ such that*

$$\sup_{\substack{|z| > T/L \\ z \in \mathbb{R}}} \frac{\Re \Psi(z)}{1 + |z|} \leq -c, \quad \forall T \geq T_0. \tag{4.49}$$

Moreover, for any $T, \delta > 0$, writing

$$S_{N,L,\delta,T} := \{z \in \mathbb{C} \mid |z| < T/L, \Im(z) > C_\delta\}, \tag{4.50}$$

there exists $C > 0$ such that

$$\sup_{z \in S_{N,L,\delta,T}} \max \left\{ |\Psi(z)|, \left| \frac{d}{dz} \Psi(z) \right| \right\} \leq C. \tag{4.51}$$

The proof of Lemma 4.11 is essentially identical to that of Lemma 4.9 with the main difference being that the analogous estimates in (4.42) must now be summed from $\ell = 0$ to $\ell = L$, which involves a growing number of terms. In each term, for T sufficiently large, the

real part of $\Psi_\ell(z)$ is bounded above by $-c * (1 + |z|)$ as soon as $|z|/L > T$. As a result, the previous Lemma shows that, for any $K, \delta > 0$ there exists $T > 0$ such that

$$\frac{N}{2\pi} \int_{\mathcal{C}} \exp[N\Psi(z)] dz = \frac{N}{2\pi} \int_{\mathcal{C}_1(T/L, C_\delta) \cup \mathcal{C}_2(C_\delta) \cup \mathcal{C}_3(T/L, C_\delta)} \exp[N\Psi(z)] dz + O(e^{-KN}). \quad (4.52)$$

Next, exactly as in the derivation (4.47) and (4.48), we find that Ψ has a unique critical point $\zeta_* + \frac{1}{N}\zeta_{**} + O(N^{-2})$ along \mathcal{C}_2 and no critical points along $\mathcal{C}_1(T/L, C_\delta)$ and $\mathcal{C}_3(T/L, C_\delta)$. Moreover, recalling that $\sigma^2 = 1$ in this regime, a direct inspection of (4.47) and (4.48) reveals that there is a constant $C_* > 0$ so that

$$\left| \zeta_* + \frac{1}{N}\zeta_{**} \right| \leq C_* \left(\frac{1}{L} + \frac{1}{N} \right).$$

This allows us to choose δ in the definition of C_δ to be independent of L, N, P, N_0 . The remainder of the derivation is a direct computation of the value, first derivative, and second derivative of Ψ at its critical point followed by a straight-forward application of the Laplace method. Namely, the critical point of Ψ on the contour \mathcal{C}_2 takes the form

$$\frac{d}{dz} \Psi(i\zeta) = 0 \iff \zeta = \zeta_* + \frac{1}{N}\zeta_{**} + O(N^{-2}),$$

where the critical point is given in terms of $\lambda_{\text{prior}} = \sum_{\ell=1}^N N_\ell^{-1}$:

$$\begin{aligned} \zeta_* &= 0 \\ \zeta_{**} &= \frac{1}{2} \left[\frac{1}{\lambda_{\text{prior}}} \log \left(\frac{\|\theta_*\|^2}{\alpha_0} \right) - (2k + 1) \right]. \end{aligned}$$

The corresponding critical value is given by

$$\Psi \left(i \left(\zeta_* + \frac{1}{N}\zeta_{**} \right) \right) = -\frac{\lambda_{\text{prior}}}{4N} \left(2k + 1 - \frac{1}{\lambda_{\text{prior}}} \log \frac{\|\theta_*\|^2}{\alpha_0} \right)^2,$$

and the Hessian is, to leading order,

$$\frac{d^2}{d\zeta^2} \Psi \left(i \left(\zeta_* + \frac{1}{N}\zeta_{**} \right) \right) = -2N\lambda_{\text{prior}} < 0,$$

showing that $i\zeta_{**}/N$ is a non-degenerate critical point. Including the terms from the Gamma function prefactors, we obtain the G function

$$\begin{aligned} \log G_{0, L+1}^{L+1, 0} \left(\frac{\|\theta_*\|^2}{4M} \mid \frac{P}{2}, \frac{N}{2} + k \right) &= \sum_{\ell=1}^L \frac{N_\ell}{2} \left[\log \left(\frac{N_\ell}{2} \right) - 1 \right] + \frac{P}{2} \left[\log \left(\frac{P}{2} \right) - 1 \right] + \frac{L}{2} \log(2\pi) \\ &+ \left(k - \frac{1}{2} \right) \sum_{\ell=1}^L \log \left(\frac{N_\ell}{2} \right) - \frac{1}{2} \log \left(\frac{P}{2} \right) + \left(k - \frac{1}{2} \right) \log \left(\frac{\|\theta_*\|^2}{\alpha_0} \right) \\ &- \frac{1}{12} \lambda_{\text{prior}} - \frac{1}{4\lambda_{\text{prior}}} \left[\log \left(\frac{\|\theta_*\|^2}{\alpha_0} \right) \right]^2 - \frac{1}{2} \log(2\lambda_{\text{prior}}). \end{aligned}$$

When specialized to the case when $N_1 = \dots = N_L = N$, these are the results stated in (19) and (20) of Theorem 3.2. Finally, for case (c), we our results apply to the regime where

$$N := \min \{N_1, \dots, N_L\} \rightarrow \infty, \quad P, N_0, N \rightarrow \infty, \quad \frac{P}{N_0} \rightarrow \alpha_0, \quad P \sum_{\ell=1}^L \frac{1}{N_\ell} \rightarrow \lambda_{\text{post}},$$

with $\alpha_0 \in (0, 1)$ and $\lambda_{\text{post}} \in (0, \infty)$. The analysis in this case mirrors almost exactly case (b), but we use the variable substitution $t \mapsto Pt$ when changing from Φ to Ψ , and use the fact that both P/N and L/N vanish. Here, we record only the result:

$$\begin{aligned} \log G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_*\|^2}{4M} \mid \frac{P}{2}, \frac{N}{2} + k \right) &= \frac{P}{2} \left[\log \left(\frac{P}{2} \right) - 1 + \log(1 + t_*) - t_* \left(1 + \frac{\lambda_{\text{post}} t_*}{2} \right) \right] - \log P \\ &+ \frac{L}{2} \log(2\pi) + \sum_{\ell=1}^L \frac{N_\ell}{2} \left[\log \left(\frac{N_\ell}{2} \right) - 1 \right] + \frac{1}{2} (2k - 1) \log \left(\frac{N_\ell}{2} \right) \\ &- \frac{1}{2} \log \left(\lambda_{\text{post}} + \frac{1}{1 + t_*} \right) + \frac{1}{2} [(2k + 1) \lambda_{\text{post}} t_* + \log(1 + t_*)], \end{aligned}$$

where t_* is the unique solution to

$$e^{\lambda_{\text{post}} t_*} (1 + t_*) = \frac{\|\theta_*\|^2}{\alpha_0}.$$

When specialized to the case when $N_1 = \dots = N_L = N$, this gives the results stated in (21) and (22) of Theorem 3.2. This completes the proof of Theorem 3.2. \square

4.6 Proof of Theorem 3.3

Consider the setup from Theorem 3.1. We adopt the notation

$$z = \frac{\|\theta_*\|^2}{4M} = \frac{\nu}{\sigma^{2(L+1)}} \frac{P}{2} \left(\frac{N}{2} \right)^L$$

so that we have variance

$$\begin{aligned} \text{Var}_{\text{post}}[f(x)] &= \frac{\|x_\perp\|^2 \|\theta_*\|^2}{2} \frac{G_{0,L+1}^{L+1,0} \left(z \mid \frac{P}{2}, \frac{N_1}{2} + 1, \dots, \frac{N_L}{2} + 1 \right)}{z G_{0,L+1}^{L+1,0} \left(z \mid \frac{P}{2}, \frac{N_1}{2}, \dots, \frac{N_L}{2} \right)} \\ &= \frac{\|x_\perp\|^2}{2} \|\theta_*\|^2 \frac{G_{0,L+1}^{L+1,0} \left(z \mid \frac{P}{2} - 1, \frac{N_1}{2}, \dots, \frac{N_L}{2} \right)}{G_{0,L+1}^{L+1,0} \left(z \mid \frac{P}{2}, \frac{N_1}{2}, \dots, \frac{N_L}{2} \right)}. \end{aligned}$$

A Bayes-optimal σ^2 implies

$$\frac{\partial Z_\infty(0)}{\partial \sigma^2} = 0 \implies \frac{d}{dz} G_{0,L+1}^{L+1,0} \left(z \mid \frac{P}{2}, \frac{N_1}{2}, \dots, \frac{N_L}{2} \right) = 0.$$

Using the identity

$$\frac{d}{dz} \left[z^{-b_1} G_{p,q}^{m,n} \left(z \mid \begin{array}{c} \mathbf{a}_p \\ \mathbf{b}_q \end{array} \right) \right] = -z^{-1-b_1} G_{p,q}^{m,n} \left(z \mid b_1 + 1, b_2, \dots, b_q \right)$$

which holds when $m \geq 1$, we find that

$$\frac{d}{dz} G_{p,q}^{m,n} \left(z \mid \begin{array}{c} \mathbf{a}_p \\ \mathbf{b}_q \end{array} \right) = \frac{1}{z} \left[b_1 G_{p,q}^{m,n} \left(z \mid \begin{array}{c} \mathbf{a}_p \\ \mathbf{b}_q \end{array} \right) - G_{p,q}^{m,n} \left(z \mid b_1 + 1, b_2, \dots, b_q \right) \right]$$

and thus

$$\frac{P}{2} G_{0,L+1}^{L+1,0} \left(z \mid \frac{P}{2}, \frac{N_1}{2}, \dots, \frac{N_L}{2} \right) = G_{0,L+1}^{L+1,0} \left(z \mid \frac{P}{2} + 1, \frac{N_1}{2}, \dots, \frac{N_L}{2} \right). \quad (4.53)$$

To show both directions of Theorem 3.3, it hence suffices to show that to leading order in the large- P limit,

$$\frac{G_{0,L+1}^{L+1,0} \left(z \mid \frac{P}{2} - 1, \frac{N_1}{2}, \dots, \frac{N_L}{2} \right)}{G_{0,L+1}^{L+1,0} \left(z \mid \frac{P}{2}, \frac{N_1}{2}, \dots, \frac{N_L}{2} \right)} = \frac{G_{0,L+1}^{L+1,0} \left(z \mid \frac{P}{2}, \frac{N_1}{2}, \dots, \frac{N_L}{2} \right)}{G_{0,L+1}^{L+1,0} \left(z \mid \frac{P}{2} + 1, \frac{N_1}{2}, \dots, \frac{N_L}{2} \right)} = \frac{2}{P}, \quad (4.54)$$

which ensures variance

$$\text{Var}_{\text{post}}[f(x)] = \left(\frac{\|x_\perp\|^2}{2} \|\theta_*\|^2 \right) \left(\frac{2}{P} \right) = \frac{\|x_\perp\|^2 \|\theta_*\|^2}{N_0 \alpha_0}.$$

This simultaneously sets the predictor variance to $\nu \Sigma_\perp$ that we recover above. Conversely, if the variance is given by $\nu \Sigma_\perp$ to leading order, (4.54) implies that $\partial Z_\infty(0)/\partial \sigma^2 = 0$.

To show (4.54), we evaluate a saddle point approximation of the relevant G function

$$G_{0,L+1}^{L+1,0} \left(z \mid \frac{P}{2} + k, \frac{N_1}{2}, \dots, \frac{N_L}{2} \right)$$

with constant L, N_1, \dots, N_L . Similarly to the result of Lemma 4.8, we use the density of independent Γ random variables

$$\phi_j \sim \begin{cases} \Gamma \left(\frac{N_j}{2} + 1, \frac{2\sigma^2}{N_j} \right), & j = 1, \dots, L \\ \Gamma \left(\frac{P}{2} + k + 1, \frac{2\sigma^2}{P} \frac{\alpha_0}{\|\theta_*\|^2} \right), & j = 0 \end{cases}.$$

Note that unlike in Lemma 4.8, these variables offset P by k , not N_ℓ . The density is thus

$$\begin{aligned} \text{Den}_{\phi_0 \dots \phi_L}(1) &= \frac{\|\theta_*\|^2}{4M} G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_*\|^2}{4M} \mid \frac{P}{2} + k, \frac{N}{2} \right) \left[\Gamma \left(\frac{P}{2} + k + 1 \right) \right]^{-1} \prod_{\ell=1}^L \left[\Gamma \left(\frac{N_\ell}{2} + 1 \right) \right]^{-1} \\ &= \frac{1}{2\pi} \int_{\mathcal{C}} \exp[\Phi(z)] dz, \end{aligned}$$

where $\mathcal{C} \subseteq \mathbb{C}$ is the contour that runs along the real line from $-\infty$ to ∞ and

$$\begin{aligned} \Phi(z) &= -iz \log \left(\frac{2\sigma^2}{P} \frac{\alpha_0}{\|\theta_*\|^2} \right) + \log \left(\frac{\Gamma\left(\frac{P}{2} + k + 1 - iz\right)}{\Gamma\left(\frac{P}{2} + k + 1\right)} \right) \\ &\quad + \sum_{\ell=1}^L \left\{ -iz \log \left(\frac{2\sigma^2}{N_\ell} \right) + \log \left(\frac{\Gamma\left(\frac{N_\ell}{2} + 1 - iz\right)}{\Gamma\left(\frac{N_\ell}{2} + 1\right)} \right) \right\}. \end{aligned}$$

The fixed point equation $d\Phi(i\zeta)/d\zeta = 0$ is solved by ζ_* satisfying

$$\sum_{\ell=1}^L -\log \left(\frac{N}{2} \right) + \psi \left(\frac{N_\ell}{2} + \zeta_* + 1 \right) = \log \left(\frac{\|\theta_*\|^2}{\alpha_0 \sigma^{2(L+1)}} \right),$$

i.e., by $\zeta_* = O(1)$. Directly applying Laplace's method

$$\log \frac{1}{2\pi} \int_{\mathcal{C}} \exp[\Phi(z)] dz = \Phi(z_*) - \frac{1}{2} \log(2\pi) - \frac{1}{2} \log \Phi''(z_*)$$

and evaluating the ratios in (4.54), we find

$$\log \frac{G_{0,L+1}^{L+1,0} \left(z \mid \frac{P}{2} - 1, \frac{N_1}{2}, \dots, \frac{N_L}{2} \right)}{G_{0,L+1}^{L+1,0} \left(z \mid \frac{P}{2}, \frac{N_1}{2}, \dots, \frac{N_L}{2} \right)} = \log \frac{G_{0,L+1}^{L+1,0} \left(z \mid \frac{P}{2}, \frac{N_1}{2}, \dots, \frac{N_L}{2} \right)}{G_{0,L+1}^{L+1,0} \left(z \mid \frac{P}{2} + 1, \frac{N_1}{2}, \dots, \frac{N_L}{2} \right)} = \log \left(\frac{2}{P} \right) + O \left(\frac{1}{P} \right),$$

which we note is independent of σ^2 to leading order. That is, the large dataset overwhelms the prior, causing all choices of σ^2 to become optimal. This completes the proof.

4.7 Model Evidence

In the infinite-depth case $L = \lambda_{\text{prior}} N$ as $N, P \rightarrow \infty$, we find that not only does setting $\sigma^2 = 1, \lambda_{\text{prior}} > 0$ give desirable posteriors, but it also enjoys a significant preference with respect to model evidence.

Corollary 4.12 (Bayesian Preference for Infinite Depth). *As in the setting of Theorem 6, fix $\lambda_{\text{prior}} > 0$. For each fixed $L \geq 0$*

$$\lim_{N \rightarrow \infty} e^{cN} \frac{Z_\infty(0 \mid L, N_\ell = N, \sigma^2 = 1, X_{N_0}, Y_{N_0})}{Z_\infty(0 \mid L = N\lambda_{\text{prior}}, N_\ell = N, \sigma^2 = 1, X_{N_0}, Y_{N_0})} = 1,$$

where

$$c = -\frac{\alpha}{2} \left[\log \left(1 + \frac{z_*}{\alpha} \right) - \frac{z_*}{\alpha} \right] - \frac{L}{2} [\log(1 + z_*) - z_*]$$

satisfies $c \leq 0$, and z_* is defined as in (3.8) of Theorem 3.2, i.e.,

$$\frac{\nu}{\sigma^{2(L+1)}} = \left(1 + \frac{z_*}{\alpha} \right) (1 + z_*)^L, \quad z_* > \max\{-1, -\alpha\}. \quad (4.55)$$

This shows that, when $\sigma^2 = 1$, any choice of λ_{prior} results in exponentially greater evidence than a network with finitely many hidden layers in the large- N limit. We omit the proof, since it is a direct manipulation of the reported evidence in the main statement of Theorems 4 and 6.

Moreover, the Bayes-optimal depth is given by an $O(1)$ choice of λ_{prior} . That is, unlike in the finite-depth limit where $L_* \rightarrow \infty$ when using an ML prior, choosing $L = O(N)$ has an attainable optimal depth that maximizes Bayesian evidence.

Corollary 4.13 (Bayes-optimal infinite depth). *In the setting of Theorem 6, we have*

$$\lambda_{\text{prior},*} = \sqrt{1 + \log(\nu)^2} - 1 = \arg \max_{\lambda} \lim_{\substack{P, N_{\ell} \rightarrow \infty \\ P/N_0 \rightarrow \alpha_0 \in (0,1) \\ \lambda_{\text{prior}}(N_1, \dots, N_L) \rightarrow \lambda}} Z_{\infty}(\mathbf{0} \mid L, N_{\ell}, \sigma^2 = 1, X_{N_0}, Y_{N_0}).$$

Moreover, for any $\lambda > 0$, there exists $c > 0$ such that we have

$$c < \lim_{\substack{P, N_{\ell} \rightarrow \infty \\ P/N_0 \rightarrow \alpha_0 \in (0,1)}} \frac{Z_{\infty}(\mathbf{0} \mid L = N\lambda, N_{\ell} = N, \sigma^2 = 1, X_{N_0}, Y_{N_0})}{Z_{\infty}(\mathbf{0} \mid L = N\lambda_{\text{prior},*}, N_{\ell} = N, \sigma^2 = 1, X_{N_0}, Y_{N_0})} \leq 1,$$

i.e., the ratio of evidences is lower-bounded by a constant.

Proof. In Corollary 4.13, we compute the Bayes-optimal neural network depth by maximizing Bayesian evidence [Mac92]. A similar computation to the following is used to find Bayes-optimal parameters throughout the main text; we provide the proof to Corollary 4.13 as an example. To evaluate $\partial \log Z / \partial \lambda_{\text{prior}} = 0$ for the partition function given by Theorem 3.1, we take the partition function with $N_1 = \dots = N_L = N$ obtained from case (b) of Theorem 3.2,

$$\begin{aligned} \log Z_{\infty}(0) &= \frac{P}{2} \log \left(\frac{4\pi}{\|\theta_*\|^2} \right) + \frac{P}{2} \left[\log \left(\frac{P}{2} \right) - 1 \right] - \frac{1}{2} \log \left(\frac{P}{2} \right) \\ &\quad - \frac{1}{2} \log(2\lambda_{\text{prior}}) - \frac{1}{4\lambda_{\text{prior}}} \left(\lambda_{\text{prior}} + \log \left(\frac{\|\theta_*\|^2}{\alpha_0} \right) \right)^2, \end{aligned}$$

and we evaluate the derivative with variable substitution $\nu = \|\theta_*\|^2 / \alpha_0$, yielding

$$\frac{\partial \log Z_{\infty}(0)}{\partial \lambda_{\text{prior}}} = \frac{\log(\nu)^2 - \lambda_{\text{prior}}(2 + \lambda_{\text{prior}})}{4\lambda_{\text{prior}}^2} = 0.$$

Solving gives

$$\lambda_{\text{prior},*} = \sqrt{1 + \log^2 \nu} - 1 \geq 0.$$

Moreover, the second derivative is

$$\frac{\partial^2 \log Z_{\infty}(0)}{\partial \lambda_{\text{prior}}^2} = \frac{\lambda_{\text{prior}} - \log^2 \nu}{2\lambda_{\text{prior}}^3},$$

which at $\lambda_{\text{prior},*}$ is

$$\frac{\partial^2 \log Z_{\infty}(0)}{\partial \lambda_{\text{prior}}^2} = -\frac{\lambda_{\text{prior},*} + 1}{2\lambda_{\text{prior},*}^2} < 0.$$

Hence, Bayesian evidence is maximized at $\lambda_{\text{prior},*} = O(1)$.

Moreover, evaluating the evidence at a different choice of fixed λ_{prior} such that $\lambda_{\text{prior}} \neq \lambda_{\text{prior},*}$, we see that the evidence ratio is

$$\begin{aligned} \log Z_{\infty}(0; \lambda_{\text{prior}}) - \log Z_{\infty}(0; \lambda_{\text{prior},*}) &= -\frac{1}{2} \log \left(\frac{\lambda_{\text{prior}}}{\lambda_{\text{prior},*}} \right) - \frac{1}{4\lambda_{\text{prior}}} \left(\lambda_{\text{prior}} + \log \left(\frac{\|\theta_*\|^2}{\alpha_0} \right) \right)^2 \\ &\quad + \frac{1}{4\lambda_{\text{prior},*}} (\lambda_{\text{prior},*} + \log \nu)^2, \end{aligned}$$

which is a constant. \square

4.8 Proof of Theorem 3.10

In order to prove Theorem 3.10, we expand $\Delta(\log G)[k]$ in the case of $L = \lambda_{\text{prior}}N$, $P/N = \alpha$, $\sigma^2 = 1$ to higher order than reported in Theorem 3.2. Specifically, we find

$$\begin{aligned} \Delta(\log G, k=1) &:= \log G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_*\|^2}{4M} \mid \frac{P}{2}, \frac{N}{2} + 1 \right) - \log G_{0,L+1}^{L+1,0} \left(\frac{\|\theta_*\|^2}{4M} \mid \frac{P}{2}, \frac{N}{2} \right) \\ &= \lambda_{\text{prior}}N \log \left(\frac{N}{2} \right) + \log(\nu) + \frac{c}{N} \end{aligned}$$

for

$$c = -\frac{8\lambda_{\text{prior}}}{3} + 2(1 + \log(\nu)) + \left(\frac{1}{\alpha} + \log(\nu) \right) \left(1 - \frac{\log(\nu)}{\lambda_{\text{prior}}} \right).$$

Applying Theorem 3.1, the difference in variance compared to infinite N is

$$\text{Var}_{\text{post}}[f(x)] - \lim_{N \rightarrow \infty} \text{Var}_{\text{post}}[f(x)] = \frac{c}{N} \nu \Sigma_{\perp} \propto \frac{1}{N} \propto \frac{1}{P} \propto \frac{1}{L}.$$

4.9 Proof of Theorem 3.11

To prove Theorem 3.11 we begin with the bias-variance decomposition:

$$\langle f(x) - V_0x \rangle^2 = (\theta_*x_{\parallel} - V_0x)^2 + \frac{\|x_{\perp}\|^2 \|\theta_*\|^2}{P},$$

where $x_{\perp} = x - x_{\parallel}$ and

$$x_{\parallel} = \text{im}(X)\text{im}(X)^T x.$$

Consider first the case when $\alpha_0 < 1$. Then

$$\theta_* = V_0 + \epsilon(X^T X)^{-1} X^T.$$

Thus, the error introduced by the bias is

$$\begin{aligned} \mathbb{E} \left[(\theta_*x_{\parallel} - V_0x)^2 \right] &= \mathbb{E} [(V_0x_{\perp})^2] + \mathbb{E} [(\epsilon(X^T X)^{-1} X^T x_{\parallel})^2] \\ &= \mathbb{E} [(V_0x)^2 - 2(V_0x)(V_0x_{\parallel}) + (V_0x_{\parallel})^2] + \sigma_{\epsilon}^2 \mathbb{E} [\|(X^T X)^{-1} X^T x_{\parallel}\|^2] \\ &= 1 - \alpha_0 + \sigma_{\epsilon}^2 \mathbb{E} [\|(X^T X)^{-1} X^T x_{\parallel}\|^2]. \end{aligned}$$

We observe that, for $X^\dagger = (X^T X)^{-1} X^T$,

$$\begin{aligned}\mathbb{E} [\| (X^T X)^{-1} X^T x_{\parallel} \|^2] &= \mathbb{E} \left[\text{Tr} \left(x_{\parallel}^T (X^\dagger)^T X^\dagger x_{\parallel} \right) \right] \\ &= \mathbb{E} \left[\text{Tr} \left((X^\dagger)^T X^\dagger x_{\parallel} x_{\parallel}^T \right) \right] \\ &= \frac{\alpha_0}{1 - \alpha_0}.\end{aligned}$$

Hence, the bias is

$$\mathbb{E} \left[(\theta_* x_{\parallel} - V_0 x)^2 \right] = 1 - \alpha_0 + \frac{\alpha_0}{1 - \alpha_0} \sigma_\epsilon^2. \quad (4.56)$$

The error introduced by the variance is

$$\begin{aligned}\mathbb{E} \left[\frac{\|x_{\perp}\|^2 \|\theta_*\|^2}{P} \right] &= \frac{1}{P} \mathbb{E} \left[(\|V_0\|^2 + \|\epsilon (X^T X)^{-1} X^T\|^2) (\|x\|^2 + \|x_{\parallel}\|^2 - 2\|\text{im}(X)^T x\|^2) \right] \\ &= \frac{1 - \alpha_0}{P} \mathbb{E} \left[\|x\|^2 (1 + \sigma_\epsilon^2 \|(X^T X)^{-1} X^T\|^2) \right] \\ &= \left(\frac{1}{\alpha_0} - 1 \right) (1 + \sigma_\epsilon^2 \mathbb{E} [\|(X^T X)^{-1} X^T\|^2]).\end{aligned}$$

Observing that

$$\mathbb{E} \left[\|(X^T X)^{-1} X^T\|^2 \right] = \mathbb{E} \left[\text{Tr} \left((X^T X)^{-1} X^T X (X^T X)^{-1} \right) \right] = \mathbb{E} \left[\text{Tr} \left((X^T X)^{-1} \right) \right] = \frac{\alpha_0}{1 - \alpha_0}$$

from the -1 st moment of the Marchenko-Pastur distribution for $\alpha_0 < 1$, we obtain

$$\begin{aligned}\langle f(x) - V_0 x \rangle^2 &= 1 - \alpha_0 + \frac{\alpha_0}{1 - \alpha_0} \sigma_\epsilon^2 + \left(\frac{1}{\alpha_0} - 1 \right) \left(1 + \frac{\alpha_0}{1 - \alpha_0} \sigma_\epsilon^2 \right) \\ &= \frac{1}{\alpha_0} - \alpha_0 + \frac{\sigma_\epsilon^2}{1 - \alpha_0}.\end{aligned}$$

In the case of $\alpha_0 > 1$, the variance is zero since $\|x_{\perp}\|^2 = 0$. Given

$$\theta_* = V_0 + \epsilon X^T (X X^T)^{-1},$$

the total error originates from the bias, i.e.,

$$\begin{aligned}\langle f(x) - V_0 x \rangle^2 &= \mathbb{E} \left[(\epsilon (X^T X)^{-1} X^T x)^2 \right] \\ &= \frac{\sigma_\epsilon^2}{\alpha_0 - 1}\end{aligned} \quad (4.57)$$

similarly to the bias computation for $\alpha_0 < 1$.

References

- [ACGH19] Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient descent for deep linear neural networks. *ICLR*, 2019.

- [ACH18] Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit acceleration by overparameterization. *ICML*, 2018.
- [ALP22] Ben Adlam, Jake Levinson, and Jeffrey Pennington. A random matrix perspective on mixtures of nonlinearities for deep learning. *AISTATS*, 2022.
- [AP20] Ben Adlam and Jeffrey Pennington. The neural tangent kernel in high dimensions: Triple descent and a multi-scale theory of generalization. In *International Conference on Machine Learning*, pages 74–84. PMLR, 2020.
- [APP⁺22] S Ariosto, R Pacelli, M Pastore, F Ginelli, M Gherardi, and P Rotondo. Statistical mechanics of deep learning beyond the infinite-width limit. *arXiv preprint arXiv:2209.04882*, 2022.
- [AS64] Milton Abramowitz and Irene A Stegun. *Handbook of mathematical functions with formulas, graphs, and mathematical tables*, volume 55. US Government printing office, 1964.
- [ASS20] Madhu S Advani, Andrew M Saxe, and Haim Sompolinsky. High-dimensional dynamics of generalization error in neural networks. *Neural Networks*, 132:428–446, 2020.
- [AZLL19] Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameterized neural networks, going beyond two layers. In *Advances in neural information processing systems*, pages 6158–6169, 2019.
- [BDK⁺21] Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining neural scaling laws. *arXiv preprint arXiv:2102.06701*, 2021.
- [BHMM19] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.
- [BHX20] Mikhail Belkin, Daniel Hsu, and Ji Xu. Two models of double descent for weak features. *SIAM Journal on Mathematics of Data Science*, 2(4):1167–1180, 2020.
- [BLLT20] Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 2020.
- [Bor14] Emile Borel. *Introduction géométrique à quelques théories physiques*. Gauthier-Villars, 1914.
- [CB18] Lenaïc Chizat and Francis Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. In *Advances in neural information processing systems*, pages 3036–3046, 2018.
- [CKZ23] Hugo Cui, Florent Krzakala, and Lenka Zdeborová. Optimal learning of deep random networks of extensive-width. *arXiv preprint arXiv:2302.00375*, 2023.
- [DF87] Persi Diaconis and David Freedman. A dozen de finetti-style results in search of a theory. In *Annales de l’IHP Probabilités et statistiques*, volume 23, pages 397–423, 1987.

- [DL13] Andreas Damianou and Neil D Lawrence. Deep gaussian processes. In *Artificial intelligence and statistics*, pages 207–215. PMLR, 2013.
- [DLT⁺18] Simon S Du, Jason D Lee, Yuandong Tian, Barnabas Poczos, and Aarti Singh. Gradient descent learns one-hidden-layer cnn: Don’t be afraid of spurious local minima. *ICML*, 2018.
- [DZPS19] Simon S. Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations*, 2019.
- [GJS⁺20] Mario Geiger, Arthur Jacot, Stefano Spigler, Franck Gabriel, Levent Sagun, Stéphane d’Ascoli, Giulio Biroli, Clément Hongler, and Matthieu Wyart. Scaling description of generalization with number of parameters in deep learning. *Journal of Statistical Mechanics: Theory and Experiment*, 2020(2):023401, 2020.
- [Han18] Boris Hanin. Which neural net architectures give rise to exploding and vanishing gradients? In *Advances in Neural Information Processing Systems*, 2018.
- [Han22] Boris Hanin. Random fully connected neural networks as perturbatively solvable hierarchies. *arXiv preprint arXiv:2204.01058*, 2022.
- [HMRT22] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. *The Annals of Statistics*, 50(2):949–986, 2022.
- [HN20a] Boris Hanin and Mihai Nica. Finite depth and width corrections to the neural tangent kernel. *ICLR 2020*, 2020.
- [HN20b] Boris Hanin and Mihai Nica. Products of many large random matrices and gradients in deep neural networks. *Communications in Mathematical Physics*, 376(1):287–322, 2020.
- [HNPSD22] Jiri Hron, Roman Novak, Jeffrey Pennington, and Jascha Sohl-Dickstein. Wide bayesian neural networks have a simple weight posterior: theory and accelerated sampling. In *International Conference on Machine Learning*, pages 8926–8945. PMLR, 2022.
- [HP21] Boris Hanin and Grigoris Paouris. Non-asymptotic results for singular values of gaussian matrix products. *Geometric and Functional Analysis*, 31(2):268–324, 2021.
- [HR18] Boris Hanin and David Rolnick. How to start training: The effect of initialization and architecture. In *Advances in Neural Information Processing Systems*, pages 571–581, 2018.
- [JGH18] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018.
- [Kaw16] Kenji Kawaguchi. Deep learning without poor local minima. In *Advances in Neural Information Processing Systems*, pages 586–594, 2016.

- [KMH⁺20] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [LBN⁺18] Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein. Deep neural networks as gaussian processes. *ICML 2018 and arXiv:1711.00165*, 2018.
- [LNR22] Mufan Bill Li, Mihai Nica, and Daniel M Roy. The neural covariance sde: Shaped infinite depth-and-width networks at initialization. *NeurIPS 2022*, 2022.
- [LS21] Qianyi Li and Haim Sompolinsky. Statistical mechanics of deep linear neural networks: The backpropagating kernel renormalization. *Physical Review X*, 11(3):031059, 2021.
- [LZB22] Chaoyue Liu, Libin Zhu, and Mikhail Belkin. Loss landscapes and optimization in over-parameterized non-linear systems and neural networks. *Applied and Computational Harmonic Analysis*, 59:85–116, 2022.
- [Mac92] David JC MacKay. Bayesian interpolation. *Neural computation*, 4(3):415–447, 1992.
- [Mei36] CS Meijer. Über whittakersche bzw. besselsche funktionen und deren produkte. *Nieuw Archief voor Wiskunde*, 18(2):10–29, 1936.
- [MM19] Song Mei and Andrea Montanari. The generalization error of random features regression: Precise asymptotics and the double descent curve. *Communications on Pure and Applied Mathematics*, 2019.
- [MMM21] Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Generalization error of random feature and kernel methods: hypercontractivity and kernel matrix concentration. *Applied and Computational Harmonic Analysis*, 2021.
- [MMN18] Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671, 2018.
- [MZ22] Andrea Montanari and Yiqiao Zhong. The interpolation phase transition in neural networks: Memorization and generalization under lazy training. *The Annals of Statistics*, 50(5):2816–2847, 2022.
- [NBR⁺21] Lorenzo Noci, Gregor Bachmann, Kevin Roth, Sebastian Nowozin, and Thomas Hofmann. Precise characterization of the prior predictive distribution of deep relu networks. *Advances in Neural Information Processing Systems*, 34:20851–20862, 2021.
- [NR21] Gadi Naveh and Zohar Ringel. A self consistent theory of gaussian processes captures feature learning effects in finite cnns. *Advances in Neural Information Processing Systems*, 34, 2021.
- [PN21] Huy Tuan Pham and Phan-Minh Nguyen. Global convergence of three-layer neural networks in the mean field regime. *ICLR*, 2021.

- [RBC⁺21] Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*, 2021.
- [RVE18] Grant Rotskoff and Eric Vanden-Eijnden. Parameters as interacting particles: long time convergence and asymptotic error scaling of neural networks. *Advances in neural information processing systems*, 31, 2018.
- [RYH22] Daniel A Roberts, Sho Yaida, and Boris Hanin. *The Principles of Deep Learning Theory: An Effective Theory Approach to Understanding Neural Networks*. Cambridge University Press, 2022.
- [SMG14] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *ICLR*, 2014.
- [SR21] Inbar Seroussi and Zohar Ringel. Separation of scales and a thermodynamic description of feature learning in some cnns. *arXiv preprint arXiv:2112.15383*, 2021.
- [SS20] Justin Sirignano and Konstantinos Spiliopoulos. Mean field analysis of neural networks: A central limit theorem. *Stochastic Processes and their Applications*, 130(3):1820–1852, 2020.
- [SS21] Justin Sirignano and Konstantinos Spiliopoulos. Mean field analysis of deep neural networks. *Mathematics of Operations Research*, 2021.
- [SSL22] Andrew Saxe, Shagun Sodhani, and Sam Jay Lewallen. The neural race reduction: dynamics of abstraction in gated networks. In *International Conference on Machine Learning*, pages 19287–19309. PMLR, 2022.
- [Yai20] Sho Yaida. Non-gaussian processes and neural networks at finite widths. *MSML*, 2020.
- [YH21] Greg Yang and Edward J Hu. Tensor programs iv: Feature learning in infinite-width neural networks. In *International Conference on Machine Learning*, pages 11727–11737. PMLR, 2021.
- [YS21] Jiahui Yu and Konstantinos Spiliopoulos. Normalization effects on shallow neural networks and related asymptotic expansions. *Foundations of Data Science*, 3(2):151–200, 2021.
- [ZBH⁺17] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations, (ICLR)*, 2017.
- [ZVP21] Jacob Zavatore-Veth and Cengiz Pehlevan. Exact marginal prior distributions of finite bayesian neural networks. *Advances in Neural Information Processing Systems*, 34, 2021.

- [ZVTP22] Jacob A. Zavatone-Veth, William L. Tong, and Cengiz Pehlevan. Contrasting random and learned features in deep bayesian linear regression. *Phys. Rev. E*, 105:064118, 2022.