
The Spectrum of the Fisher Information Matrix of a Single-Hidden-Layer Neural Network

Jeffrey Pennington
Google Brain
jpennin@google.com

Pratik Worah
Google Research
pworah@google.com

Abstract

An important factor contributing to the success of deep learning has been the remarkable ability to optimize large neural networks using simple first-order optimization algorithms like stochastic gradient descent. While the efficiency of such methods depends crucially on the local curvature of the loss surface, very little is actually known about how this geometry depends on network architecture and hyperparameters. In this work, we extend a recently-developed framework for studying spectra of nonlinear random matrices to characterize an important measure of curvature, namely the eigenvalues of the Fisher information matrix. We focus on a single-hidden-layer neural network with Gaussian data and weights and provide an exact expression for the spectrum in the limit of infinite width. We find that linear networks suffer worse conditioning than nonlinear networks and that nonlinear networks are generically non-degenerate. We also predict and demonstrate empirically that by adjusting the nonlinearity, the spectrum can be tuned so as to improve the efficiency of first-order optimization methods.

1 Introduction

In recent years, the success of deep learning has spread from classical problems in image recognition [1], audio synthesis [2], translation [3], and speech recognition [4] to more diverse applications in unexpected areas such as protein structure prediction [5], quantum chemistry [5] and drug discovery [6]. These empirical successes continue to outpace the development of a concrete theoretical understanding of *how* and *in what contexts* deep learning works. A central difficulty in analyzing deep learning systems stems from the complexity of neural network loss surfaces, which are highly non-convex functions, often of millions or even billions [7] of parameters.

Optimization in such high-dimensional spaces poses many challenges. For most problems in deep learning, second-order methods are too costly to perform exactly. Despite recent developments on efficient approximations of these methods, such as the Neumann optimizer [8] and K-FAC [9], most practitioners use gradient descent and its variants [10], [11]. Despite their widespread use, it is not obvious why first-order methods are often successful in deep learning since it is known that first-order methods perform poorly in the presence of pathological curvature. An important open question in this direction is to what extent pathological curvature pervades deep learning and how it can be mitigated. More broadly, in order to continue improving neural network models and performance, we aim to better understand the conditions under which first-order methods will work well, and how those conditions depend on model design choices and hyperparameters.

Among the variety of objects that may be used to quantify the geometry of the loss surface, two matrices have elevated importance: the Hessian matrix and the Fisher information matrix. From the perspective of Euclidean coordinate space, the Hessian matrix is the natural object with which to quantify the local geometry of the loss surface. It is also the fundamental object underlying many second-order optimization schemes and its spectrum provides insights as to the nature of critical

points. From the perspective of information geometry, distances are measured in model space rather than in coordinate space, and the Fisher information matrix defines the metric and determines the update directions in natural gradient descent [12]. In contrast to the standard gradient, the natural gradient defines the direction in the parameter space which gives the largest change in the objective per unit change in the model, as measured by Kullback-Leibler divergence. As we discuss in Section 2, the Hessian and the Fisher are related; for the squared error loss functions that we consider in this work, it turns out that the Fisher equals the Gauss-Newton approximation of the Hessian, so the connection is concrete.

A central difficulty in building up a robust understanding of the properties of these curvature matrices stems from the fact that they are high-dimensional and the empirical estimation of their spectra is limited by memory and computational constraints. These limitations typically prevent direct calculations for models with more than a few tens of thousands of parameters and it is difficult to know whether conclusions drawn from such small models would generalize to the mega- or giga-dimensional networks used in practice.

It is therefore important to develop theoretical tools to analyze the spectra of these matrices. In general, the spectra will depend in intimate ways on the specific parameter values of the weights and the distribution of input data to the network. It is not feasible to precisely capture all of these details, and even if a theory were developed that did so, it would not be clear how to derive generalizable conclusions from it. We therefore focus on a simplified configuration in which the weights and inputs are taken to be random variables. The analysis then becomes a well-defined computation in random matrix theory.

The Fisher is a nonlinear function of the weights and data. To compute its spectrum, we extend the framework developed by Pennington and Worah [13] to study random matrices with nonlinear dependencies. As we describe in Section 2.4, the Fisher also has an internal block structure that complicates the resulting combinatorial analysis. The main technical contribution of this work is to extend the nonlinear random matrix theory of [13] to matrices with nontrivial internal structure.

The result of our analysis is an explicit characterization of the spectrum of the Fisher information matrix of a single-hidden-layer neural network with squared loss, random Gaussian weights and random Gaussian input data in the limit of large width. We draw several nontrivial and potentially surprising conclusions about the spectrum. For example, linear networks suffer worse conditioning than any nonlinear network, and although nonlinear networks may have many small eigenvalues they are generically non-degenerate. Our results also suggest precise ways to tune the nonlinearity in order to improve conditioning of the spectrum, and our empirical simulations show improvements in the speed of first-order optimization as a result.

2 Preliminaries

2.1 Notation and problem statement

Consider a single-hidden-layer neural network with weight matrices $W^{(1)}, W^{(2)} \in \mathbb{R}^{n \times n}$ and pointwise activation function $f : \mathbb{R} \rightarrow \mathbb{R}$. For input $X \in \mathbb{R}^n$, the output of the network $\hat{Y}(X) \in \mathbb{R}^n$ is given by $\hat{Y}(X) = W^{(2)} f(W^{(1)} X)$. For concreteness, we focus our analysis on the case of squared loss, in which case,

$$\mathcal{L}(\theta) = \mathbb{E}_{X,Y} \frac{1}{2} \|Y - \hat{Y}(X)\|_2^2, \quad (1)$$

where $Y \in \mathbb{R}^n$ are the regression targets and θ denotes the vector of all parameters $\{W^{(1)}, W^{(2)}\}$. The matrix of second derivatives or *Hessian* of the loss with respect to the parameters can be written as,

$$H = H^{(0)} + H^{(1)}, \quad (2)$$

where,

$$H_{ij}^{(0)} = \mathbb{E}_X \sum_{\alpha} \frac{\partial \hat{Y}_{\alpha}}{\partial \theta_i} \frac{\partial \hat{Y}_{\alpha}}{\partial \theta_j}, \quad \text{and} \quad H_{ij}^{(1)} = \mathbb{E}_X \sum_{\alpha} (\hat{Y}(X) - Y)_{\alpha} \frac{\partial^2 \hat{Y}_{\alpha}}{\partial \theta_i \partial \theta_j}. \quad (3)$$

In this work we focus on the positive-semi-definite matrix $H^{(0)}$, which is known as the Gauss-Newton matrix. It can also be written as $H^{(0)} = J^T J$, where $J \in \mathbb{R}^{n \times 2n^2}$ is the Jacobian matrix of \hat{Y} with

respect to the parameters θ . For models with squared loss, it is known that the Gauss-Newton matrix is equal to the Fisher information matrix of the model distribution with respect to its parameters [14]. As such, by studying $H^{(0)}$ we simultaneously examine the Gauss-Newton matrix and the Fisher information matrix.

The distribution of eigenvalues or *spectrum* of curvature matrices like $H^{(0)}$ plays an important role in optimization, as it characterizes the local geometry of the loss surface and the efficiency of first-order optimization methods. In this work, we seek to build a detailed understanding of this spectrum and how the architectural components of the neural network influence it. In order to isolate these factors from idiosyncratic behavior related to the specifics of the data and weight configurations, we focus on the a vanilla baseline configuration in which the data and the weights are both taken to be iid Gaussian random variables.

Concretely, we take $X \sim \mathcal{N}(0, I_n)$, $W_{ij}^{(l)} \sim \mathcal{N}(0, \frac{1}{n})$, and we will be interested in computing the expected distribution of eigenvalues $H^{(0)}$ for large n . From this perspective, the problem can be framed as a computation in random matrix theory, the principles behind which we now review.

2.2 Spectral density and the Stieltjes transform

The *empirical spectral density* of a matrix M is defined as,

$$\rho_M(\lambda) = \frac{1}{n} \sum_{j=1}^n \delta(\lambda - \lambda_j(M)), \quad (4)$$

where the $\lambda_j(M)$, $j = 1, \dots, n$, denote the n eigenvalues of M , including multiplicity, and δ is the Dirac delta function. The *limiting spectral density* is the limit of eqn. (4) as $n \rightarrow \infty$, if it exists.

For $z \in \mathbb{C} \setminus \text{supp}(\rho_M)$ the *Stieltjes transform* G of ρ_M is defined as,

$$G(z) = \int \frac{\rho_M(t)}{z - t} dt = -\frac{1}{n} \mathbb{E} \text{tr}(M - zI_n)^{-1}, \quad (5)$$

where the expectation is with respect to the random variables W and X . The quantity $(M - zI_n)^{-1}$ is the *resolvent* of M . The spectral density can be recovered from the Stieltjes transform using the inversion formula,

$$\rho_M(\lambda) = -\frac{1}{\pi} \lim_{\epsilon \rightarrow 0^+} \text{Im} G(\lambda + i\epsilon). \quad (6)$$

2.3 Moment method

One of the main tools for computing the limiting spectral distributions of random matrices is the moment method, which, as the name suggests, is based on computations of the moments of ρ_M . The asymptotic expansion of eqn. (5) for large z gives the Laurent series,

$$G(z) = \sum_{k=0}^{\infty} \frac{m_k}{z^{k+1}}, \quad (7)$$

where m_k is the k th moment of the distribution ρ_M ,

$$m_k = \int dt \rho_M(t) t^k = \frac{1}{n} \mathbb{E} \text{tr} M^k. \quad (8)$$

If one can compute m_k , then the density ρ_M can be obtained via eqns. (7) and (6). The idea behind the moment method is to compute m_k by expanding out powers of M inside the trace as,

$$\frac{1}{n} \mathbb{E} \text{tr} M^k = \frac{1}{n} \mathbb{E} \sum_{i_1, \dots, i_k \in [n]} M_{i_1 i_2} M_{i_2 i_3} \cdots M_{i_{k-1} i_k} M_{i_k i_1}, \quad (9)$$

and evaluating the leading contributions to the sum as $n \rightarrow \infty$. We will use the moment method in order to compute the limiting spectral density of the Fisher. As a first step in that direction, we focus on the properties of the layer-wise block structure in the Fisher induced by the neural network architecture.

2.4 Block structure of the Fisher

As described above, in our single-hidden-layer setting with squared loss, the Fisher is given by

$$H^{(0)} = \mathbb{E}_X [J^T J], \quad J_{\alpha i} = \frac{\partial \hat{Y}_\alpha}{\partial \theta_i}. \quad (10)$$

Because the parameters of the model are organized into two layers, it is convenient to partition the Fisher into a 2×2 block matrix,

$$H^{(0)} = \begin{pmatrix} H_{11}^{(0)} & H_{12}^{(0)} \\ H_{12}^{(0)T} & H_{22}^{(0)} \end{pmatrix}.$$

Furthermore, because the parameters of each layer are matrices, it is useful to regard each block of the Fisher as a four-index tensor. In particular,

$$\begin{aligned} [H_{11}^{(0)}]_{a_1 b_1, a_2 b_2} &= \mathbb{E}_X \left[\sum_i J_{i, a_1 b_1}^{(1)} J_{i, a_2 b_2}^{(1)} \right], \\ [H_{12}^{(0)}]_{a_1 b_1, c_1 d_1} &= \mathbb{E}_X \left[\sum_i J_{i, a_1 b_1}^{(1)} J_{i, c_1 d_1}^{(2)} \right], \\ [H_{22}^{(0)}]_{c_1 d_1, c_2 d_2} &= \mathbb{E}_X \left[\sum_i J_{i, c_1 d_1}^{(2)} J_{i, c_2 d_2}^{(2)} \right]. \end{aligned}$$

The Jacobian entries $J_{i, ab}^{(l)}$ equal the derivatives of \hat{Y}_i with respect to the weight variables $W_{ab}^{(l)}$,

$$J_{i, ab}^{(1)} = W_{ia}^{(2)} f' \left(\sum_k W_{ak}^{(1)} X_k \right) X_b, \quad J_{j, cd}^{(2)} = \delta_{cj} f \left(\sum_l W_{dl}^{(1)} X_l \right), \quad (11)$$

where δ_{cj} denotes the Kronecker delta function i.e., it is 1 if $c = j$, and 0 otherwise.

In order to proceed by the method of moments, we will need to compute the normalized trace of powers of the Fisher, i.e. $\frac{1}{n} \text{tr}[H^{(0)}]^d$, for any d . The block structure of the Fisher makes the explicit representation of these traces somewhat complicated. The following proposition helps simplify the expressions.

Proposition 1. *Let $A_1 \in \mathbb{R}^{n \times k_1}$, $A_2 \in \mathbb{R}^{n \times k_2}$ and $A = [A_1, A_2] \in \mathbb{R}^{n \times (k_1 + k_2)}$. Then,*

$$\text{tr}[(A^T A)^d] = \sum_{b \in \{1, 2\}^d} \text{tr} \prod_{i=1}^d A_{b_i} A_{b_i}^T = \sum_{b \in \{1, 2\}^d} \text{tr} A_{b_d}^T A_{b_1} \prod_{i=1}^{d-1} A_{b_i}^T A_{b_{i+1}}. \quad (12)$$

Using Proposition 1 with $A_1 = J^{(1)}$ and $A_2 = J^{(2)}$, we have,

$$\text{tr}[(H^{(0)})^d] = \sum_{b \in \{1, 2\}^d} \text{tr} \mathbb{E}_X [J^{(b_d)T} J^{(b_1)}] \prod_{i=1}^{d-1} \mathbb{E}_X [J^{(b_i)T} J^{(b_{i+1})}], \quad (13)$$

which expresses the traces of the block Fisher entirely in terms of products of its constituent blocks.

In order to carry out the moment method to completion, we need the expected normalized traces m_k ,

$$m_k = \frac{1}{n} \mathbb{E}_W \text{tr}[(H^{(0)})^k], \quad (14)$$

in the limit of large n . Because the nonlinearity significantly complicates the analysis, we first illustrate the basics of the methodology in the linear case before moving on to the general case.

2.5 An Illustrative Example: The Linear Case

Let us assume that f is the identity function i.e., $f(z) = z$. In this case, eqn. (11) can be written as,

$$J^{(1)} = W^{(2)T} \otimes X, \quad J^{(2)} = I \otimes W^{(1)} X. \quad (15)$$

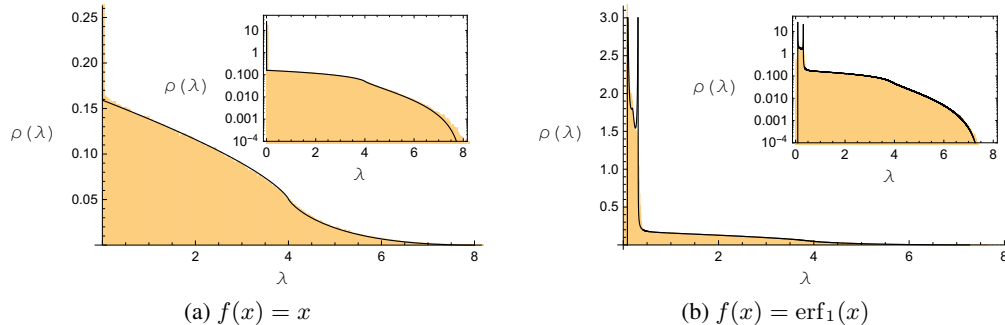


Figure 1: Empirical spectra of Fisher for single-hidden-layer networks of width 128 (orange) and theoretical prediction of spectra (black) for (a) linear and (b) erf_1 (see eqn. 30) networks. Insets show logarithmic scale.

Using the fact that $\mathbb{E}_X [XX^T] = I_n$, eqn. (13) gives,

$$\text{tr}[(H^{(0)})^d] = \mathbb{E}_W \sum_{k=0}^d \binom{d}{k} \text{tr}(W^{(2)}W^{(2)T})^{d-k} \text{tr}(W^{(1)}W^{(1)T})^k = \sum_{k=0}^d \binom{d}{k} C_{d-k} C_k, \quad (16)$$

where C_n is the n th Catalan number. The series can be summed to obtain the Stieltjes transform, whose imaginary part gives the following explicit form for the spectrum,

$$\rho(\lambda) = \frac{1}{2}\delta(\lambda) + \left[\frac{1}{2\pi^2} E\left(\frac{1}{16}(8-\lambda)\lambda\right) + \frac{4-\lambda}{8\pi^2} K\left(\frac{1}{16}(8-\lambda)\lambda\right) \right] \mathbf{1}_{[0,8]}, \quad (17)$$

where K and E are the complete elliptic integrals of the first- and second-kind,

$$K(k) = \int_0^{\frac{\pi}{2}} d\theta \frac{1}{\sqrt{1-k\sin^2\theta}}, \quad E(k) = \int_0^{\frac{\pi}{2}} d\theta \sqrt{1-k\sin^2\theta}. \quad (18)$$

Notice that the spectrum is highly degenerate, with half of the eigenvalues equaling zero. This degeneracy can be attributed to the $GL(n^2)$ symmetry of the product $W^{(2)}W^{(1)}$ under $\{W^{(1)}, W^{(2)}\} \rightarrow \{GW^{(1)}, W^{(2)}G^{-1}\}$. Fig. 1a shows excellent agreement between the predicted spectral density and finite-width empirical simulations.

3 The Stieltjes transform of $H^{(0)}$

3.1 Main Result

If $f : \mathbb{R} \rightarrow \mathbb{R}$ is an activation function with zero Gaussian mean and finite Gaussian moments,

$$\int \frac{dx}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} f(x) = 0, \quad \left| \int \frac{dx}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} f(x)^k \right| < \infty, \text{ for } k > 1, \quad (19)$$

then the Stieltjes transform of the limiting spectral density of $H^{(0)}$ is given by the following theorem.

Theorem 1. *The Stieltjes transform of the spectral density of the Fisher information matrix of a single-hidden-layer neural network with squared loss, activation function f , weight matrices $W^{(1)}, W^{(2)} \in \mathbb{R}^{n \times n}$ with iid entries $W_{ij}^{(l)} \sim \mathcal{N}(0, \frac{1}{n})$, no biases, and iid inputs $X \sim \mathcal{N}(0, I_n)$ is given by the following integral as $n \rightarrow \infty$:*

$$G(z) = \int_{\mathbb{R}} \int_{\mathbb{R}} \frac{\lambda_1 + \lambda_2 - 2z}{2\zeta^2((\eta - \zeta)(\eta' - \zeta) + \lambda_1(z - \eta + \zeta) + \lambda_2(z - \eta' + \zeta) - z^2)} d\mu_1(\lambda_1) d\mu_2(\lambda_2), \quad (20)$$

where the constants η , η' , and ζ are determined by the nonlinearity,

$$\eta = \int_{\mathbb{R}} f(x)^2 \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx, \quad \eta' = \int_{\mathbb{R}} f'(x)^2 \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx, \quad \zeta = \left(\int_{\mathbb{R}} f'(x) \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx \right)^2, \quad (21)$$

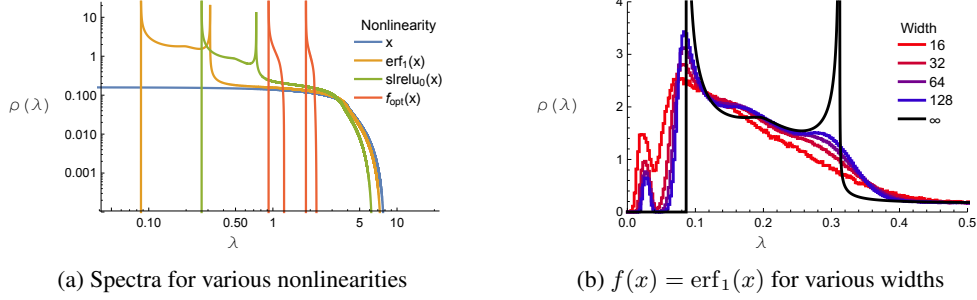


Figure 2: (a) Theoretical predictions for spectra of various nonlinearities; see eqns. (28) and (30). The linear case is degenerate and more poorly conditioned than the nonlinear cases. (b) Theoretical prediction of spectrum for erf_1 compared with empirical simulations. Practical constraints restrict the width to small values, but slow convergence toward the asymptotic prediction can be observed.

and the measures $d\mu_1$ and $d\mu_2$ are given by,

$$d\mu_1(\lambda_1) = \frac{1}{2\pi} \sqrt{\frac{\eta' + 3\zeta - \lambda_1}{\lambda_1 - \eta' + \zeta}} \mathbf{1}_{[\eta' - \zeta, \eta' + 3\zeta]}, \quad d\mu_2(\lambda_2) = \frac{1}{2\pi} \sqrt{\frac{\eta + 3\zeta - \lambda_2}{\lambda_2 - \eta + \zeta}} \mathbf{1}_{[\eta - \zeta, \eta + 3\zeta]}. \quad (22)$$

Remark 1. A straightforward application of Carlson's algorithm [15] can reduce the integral in eqn. (20) to a combination of three standard elliptic integrals.

Remark 2. The spectral density can be recovered from eqn. (20) through the inversion formula, eqn. (6).

The proof of Theorem 1 is quite long and technical, so it's deferred to the Supplementary Material. The basic idea underlying the proof is very similar to that utilized in [13]. The calculation of the moments is divided into two sub-problems, one of enumerating certain connected outer-planar graphs, and another of evaluating certain high-dimensional integrals that correspond to walks in those graphs.

Fig. 1 shows the excellent agreement of the predicted spectrum with empirical simulations of finite-width networks. Fig. 2 highlights the region of the spectrum for which the asymptotic behavior is slow to set in and suggests that empirical simulations with small networks may not provide an accurate portrayal of the behavior of large networks. Fig. 2a shows the predicted spectra for a variety of nonlinearities.

3.2 Features of the spectrum

Owing to eqn. (6), the branch points and poles of $G(z)$ encode information about the delta function peaks, spectral edges, and discontinuities in the derivative of $\rho(\lambda)$. These special points can be determined directly from the integral representation for $G(z)$ in eqn. (20) by examining the zeros of the denominator of the integrand. In particular, the following six values of z are locations of the poles at the integration endpoints and determine the salient features of the spectral density:

$$z_1 = \eta - \zeta, \quad z_2 = \eta + 3\zeta, \quad z_3 = \frac{1}{2}(\eta + \eta' + 6\zeta - \sqrt{(\eta' - \eta)^2 + 64\zeta^2}), \quad (23)$$

$$z_4 = \eta' - \zeta, \quad z_5 = \eta' + 3\zeta, \quad z_6 = \frac{1}{2}(\eta + \eta' + 6\zeta + \sqrt{(\eta' - \eta)^2 + 64\zeta^2}). \quad (24)$$

In the Supplementary Material, we establish the relative ordering of constants $0 \leq \zeta \leq \eta \leq \eta'$, which implies that the minimum and maximum eigenvalues of $H^{(0)}$ are given by,

$$\lambda_{\min} = z_1, \quad \text{and} \quad \lambda_{\max} = z_6. \quad (25)$$

The Supplementary Material also shows that the equality $\eta = \zeta$ only holds for linear networks, which implies that the minimum eigenvalue is nonzero for every nonlinear activation function. There are two delta function peaks in spectrum, which are located at,

$$\lambda_{\text{peak}}^{(1)} = \lambda_{\min} = z_1, \quad \text{and} \quad \lambda_{\text{peak}}^{(2)} = z_4. \quad (26)$$

Table 1: Properties of nonlinearities

	η	η'	ζ	Locations of spectral features					
				z_1	z_2	z_3	z_4	z_5	z_6
x	1	1	1	0	4	0	0	4	8
$\text{erf}_1(x)$	1	1.226	0.914	0.086	3.741	0.198	0.312	3.966	7.51
$\text{srelu}_0(x)$	1	1.467	0.733	0.267	3.200	0.491	0.733	3.667	6.377
f_{opt}	1	1.923	0.077	0.923	1.231	1.138	1.846	2.154	2.247

These peaks indicate specific eigenvalues that have nonvanishing probability of occurrence. These peaks coalesce when $\eta = \eta'$, which can only happen for linear activation functions, in which case $\eta = \eta' = \zeta$, so the peaks occur at $\lambda = 0$, as illustrated in Fig. 2a. That figure also shows that the spectrum may consist of two disconnected components, in which case z_2 is the location of the right edge of the left component. Finally, the derivative of the spectrum is discontinuous at z_3 and z_5 . These predictions can be verified in Fig. 2a by consulting Table 1, which provides numerical values for these special points for the various nonlinearities appearing in the figure.

4 Empirical analysis

4.1 A measure of conditioning

Using the results from Section 4.1, the first two moments can be given explicitly as,

$$\begin{aligned}
 m_1 &= \lim_{n \rightarrow \infty} \frac{1}{n} \text{tr}[H^{(0)}] = \frac{1}{2}(\eta + \eta') \\
 m_2 &= \lim_{n \rightarrow \infty} \frac{1}{n} \text{tr}[H^{(0)^2}] = \frac{1}{2}(\eta^2 + \eta'^2 + 4\zeta^2)
 \end{aligned}
 \tag{27}$$

A scale-invariant measure of conditioning of the Fisher is just m_2/m_1^2 , which is lower-bounded by 1, and which quantifies how tightly concentrated the spectrum is around its mean. Ideally, this quantity should be as small as possible to avoid pathological curvature and to enable fast first-order optimization. One advantage of m_2/m_1^2 compared to other condition numbers such as $\lambda_{\max}/\lambda_{\min}$ or λ_{\max} is that it is scale-invariant and well-defined even in the presence of degeneracy in the spectrum.

By expanding f in a basis of Hermite polynomials, we show in the Supplementary Material that among the functions with zero Gaussian mean that

$$f_{\text{opt}}(x) = \frac{1}{\sqrt{13}}(x + \sqrt{6}(x^2 - 1))
 \tag{28}$$

minimizes the ratio m_2/m_1^2 . Note that we have removed the freedom to rescale f_{opt} by a constant by enforcing $\eta = 1$. Curiously, a linear activation function actually maximizes the ratio, implying that nonlinearity invariably improves conditioning, at least by this measure. The relative conditioning of spectra resulting from various activation functions can be observed in Fig. 2a.

The function $f_{\text{opt}}(x)$ grows quickly for large $|x|$ and may be too unstable to use in actual neural networks. Alternative functions could be found by solving the optimization problem,

$$f_* = \arg \min_f \frac{m_2}{m_1^2},
 \tag{29}$$

subject to some constraints, for example that f be monotone increasing, have zero Gaussian mean, and saturate for large $|x|$. Such a problem could be solved via variational calculus; see the Supplementary Material.

4.2 Efficiency of gradient descent

Another way to investigate the ratio m_2/m_1^2 is to see how well it correlates with the efficiency of first-order optimization. For this purpose, we examine two one-parameter classes of well-behaved

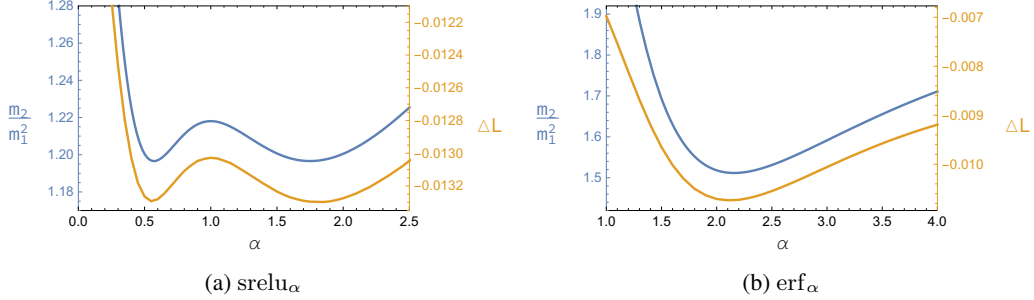


Figure 3: Comparison of the conditioning measure m_2/m_1^2 and single-step loss reduction ΔL (eqn. (33)) as the activation function changes for (a) srelu_α and (b) erf_α (eqn. (30)). The curves are highly correlated, suggesting the possibility of improved first-order optimization performance by tuning the spectrum of the Fisher through the choice of activation function.

activation functions related to ReLU and the error function,

$$\text{srelu}_\alpha(x) = \frac{[x]_+ + \alpha[-x]_+ - \frac{1+\alpha}{\sqrt{2\pi}}}{\sqrt{\frac{1}{2}(1+\alpha^2) - \frac{1}{2\pi}(1+\alpha)^2}}, \quad \text{erf}_\alpha(x) = \frac{\text{erf}(\alpha^2 x)}{\sqrt{\frac{4}{\pi} \tan^{-1} \sqrt{1+4\alpha^4} - 1}}. \quad (30)$$

Here srelu_α is the shifted leaky ReLU function studied in [13]. Both srelu_α and erf_α have zero Gaussian mean and are normalized such that $\eta = 1$ for all α . Changing α does affect η' , ζ and the ratio m_2/m_1^2 , which implies that different functions within these one-parameter families may behave quite differently under gradient descent.

We designed a simple and controlled experiment to explore these differences in the context of neural network training. The setup is a modified student-teacher framework, in which the student is initialized with the teacher’s parameters, but the regression targets are perturbed so that student’s parameters are suboptimal. Then we ask by how much can the student decrease the loss by one optimally-chosen step in the gradient direction. Concretely, we define

$$Y_i = W_t^{(2)} f(W_t^{(1)} X_i) + \epsilon_i, \quad i = 1, \dots, M, \quad (31)$$

for teacher weights $[W_t^{(l)}]_{ij} \sim \mathcal{N}(0, \frac{1}{n})$, $X_i \sim \mathcal{N}(0, I_n)$, and $\epsilon_i \sim \mathcal{N}(0, \varepsilon^2 I_n)$, with width $n = 2^7$, number of samples $M = 2^{17}$, and perturbation size $\varepsilon = 10^{-3}$. The loss is defined as,

$$L(W_s) = \sum_{i=1}^M \frac{1}{2} \|Y_i - W_s^{(2)} f(W_s^{(1)} X_i)\|_2^2. \quad (32)$$

We are interested in the maximal single-step loss decrease when W_s is initialized at W_t , i.e.,

$$\Delta L = \min_{\eta} [L(W_t - \eta \nabla L|_{W_t}) - L(W_t)]. \quad (33)$$

For the two classes of activation functions in eqn. (30), we empirically measured ΔL as a function of α . In Fig. 3 we compare the results with our theoretical predictions for m_2/m_1^2 as a function of α . The agreement is excellent, suggesting that our theory may be able to make practical predictions regarding training efficiency of actual neural networks.

5 Conclusions

In this work, we computed the spectrum of the Fisher information matrix of a single-hidden-layer neural network with squared loss and Gaussian weights and Gaussian data in the limit of large network width. Our explicit results indicate that linear networks suffer worse conditioning than nonlinear networks and that although nonlinear networks may have numerous small eigenvalues they are generically non-degenerate. We also showed that by tuning the nonlinearity it is possible to adjust the spectrum in such a way that the efficiency of first-order optimization methods can be improved. By undertaking this analysis, we demonstrated how to extend the techniques developed in [13] for studying random matrices with nonlinear dependencies to the block-structured curvature matrices that are relevant for optimization in deep learning. The techniques presented here pave the way for future work studying deep learning via random matrix theory.

References

- [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [2] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
- [3] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*, 2016.
- [4] Geoffrey Hinton, Li Deng, Dong Yu, George E. Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N Sainath, et al. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Processing Magazine*, 29(6):82–97, 2012.
- [5] Garrett Goh, Nathan Hodas, and Abhinav Vishnu. Deep Learning for Computational Chemistry. *arXiv preprint arXiv:1701.04503*, 2017.
- [6] Han Altae-Tran, Bharath Ramsundar, Aneesh S. Pappu, and Vijay Pande. Low Data Drug Discovery with One-Shot Learning. *American Chemical Society Central Science*, 2017.
- [7] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean. Outrageously large neural language models using sparsely gated mixtures of experts. *ICLR*, 2017. URL <http://arxiv.org/abs/1701.06538>.
- [8] Shankar Krishnan, Ying Xiao, and Rif A. Saurous. Neumann optimizer: A practical optimization algorithm for deep neural networks. In *International Conference on Learning Representations*, 2018.
- [9] Roger B. Grosse and James Martens. A kronecker-factored approximate fisher matrix for convolution layers. In *Proceedings of the 33rd International Conference on Machine Learning, ICML*, pages 573–582, 2016.
- [10] John Duchi, Elad Hazan, and Yoram Singer. Adaptive Subgradient Methods for Online Learning and Stochastic Optimization. *Journal of Machine Learning Research*, 2011.
- [11] Diedrik Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. *arxiv:1412.6980*, 2014.
- [12] S.I. Amari. Natural gradient works efficiently in learning. *Neural Computation*, 1998.
- [13] Jeffrey Pennington and Pratik Worah. Nonlinear random matrix theory for deep learning. In *Advances in Neural Information Processing Systems*, pages 2634–2643, 2017.
- [14] Tom Heskes. On “natural” learning and pruning in multilayered perceptrons. *Neural Computation*, 12(4):881–901, 2000.
- [15] BC Carlson. A table of elliptic integrals of the third kind. *Mathematics of computation*, 51(183):267–280, 1988.
- [16] Mariano Giaquinta and Stefan Hilderbrandt. Calculus of Variations 1. *Springer*, 1994.
- [17] Richard Stanley. Polygon Dissections and Standard Young Tableaux. *Journal of Combinatorial Theory, Series A*, 1996.

Supplemental Material: The Spectrum of the Fisher Information Matrix of a Single-Hidden-Layer Neural Network

1 Hermite expansion

Any function with finite Gaussian moments can be expanded in a basis of Hermite polynomials. Defining

$$H_n(x) = (-1)^n e^{\frac{x^2}{2}} \frac{\partial^n}{\partial x^n} e^{-\frac{x^2}{2}} \quad (\text{S1})$$

we can write,

$$f(x) = \sum_{n=0}^{\infty} \frac{f_n}{\sqrt{n!}} H_n(x), \quad (\text{S2})$$

for some constants f_n . Owing the orthogonality of the Hermite polynomials, this representation is useful for evaluating Gaussian integrals. In particular, the condition that f be centered is equivalent the vanishing of f_0 ,

$$\begin{aligned} 0 &= \int_{-\infty}^{\infty} dx \frac{e^{-x^2/2}}{\sqrt{2\pi}} f(x) \\ &= f_0. \end{aligned} \quad (\text{S3})$$

The constants η , η' , and ζ are also easily expressed in terms of the coefficients,

$$\begin{aligned} \zeta &= \left[\int_{-\infty}^{\infty} dx \frac{e^{-x^2/2}}{\sqrt{2\pi}} f'(x) \right]^2 = f_1^2 \\ \eta &= \int_{-\infty}^{\infty} dx \frac{e^{-x^2/2}}{\sqrt{2\pi}} f(x)^2 = \sum_{n=0}^{\infty} f_n^2 = \zeta + \sum_{n=2}^{\infty} f_n^2 \equiv \zeta + r_1^2 \\ \eta' &= \int_{-\infty}^{\infty} dx \frac{e^{-x^2/2}}{\sqrt{2\pi}} f'(x)^2 = \sum_{n=0}^{\infty} n f_n^2 \equiv \zeta + 2r_1^2 + r_2^2, \end{aligned} \quad (\text{S4})$$

where,

$$r_1^2 = \sum_{n=2}^{\infty} f_n^2, \quad r_2^2 = \sum_{n=3}^{\infty} (n-2) f_n^2. \quad (\text{S5})$$

From this representation it is easy to see that $0 \leq \zeta \leq \eta \leq \eta'$. The first and second moments of the Fisher are then given by,

$$\begin{aligned} m_1 &= \frac{1}{2}(\eta + \eta') = \frac{1}{2}(2\zeta + 3r_1^2 + r_2^2) \\ m_2 &= \frac{1}{2}(\eta^2 + \eta'^2 + 4\zeta^2) = \frac{1}{2}(6\zeta^2 + 6r_1^2\zeta + 2r_2^2\zeta + 5r_1^4 + 4r_1^2r_2^2 + r_2^4). \end{aligned} \quad (\text{S6})$$

Viewed as a function of ζ , r_1 , and r_2 , the ratio m_2/m_1^2 has three critical points over the positive real numbers,

$$\text{(a) } r_1 = 0, \quad r_2 = 0, \quad \frac{m_2}{m_1^2} = 3 \quad (\text{S7})$$

$$\text{(b) } r_1 = 0, \quad r_2^2 = 4\zeta, \quad \frac{m_2}{m_1^2} = \frac{5}{3} \quad (\text{S8})$$

$$\text{(c) } r_1^2 = 12\zeta, \quad r_2 = 0, \quad \frac{m_2}{m_1^2} = \frac{21}{19}. \quad (\text{S9})$$

Solution (c) is the minimum. Note that $r_2 = 0$ implies $f_k = 0$ for $k \geq 3$. Without loss of generality we can set $\eta = 1$, in which case,

$$f_{\text{opt}}(x) = \frac{1}{\sqrt{13}}(x + \sqrt{6}(x^2 - 1)). \quad (\text{S10})$$

2 Variational Calculus

It is not easy to systematically generalize the calculations in the previous section for arbitrary constraints on f . For example, placing lower bounds on f , requiring that f is monotone, placing bounds on the derivative of f , and any other such conditions that are typically seen for activation functions are all valid additional constraints to compute a good conditioned f . Variational calculus is more suited for such a generalization. However, unlike the above explicit solution, the output is a PDE with boundary conditions. For example, adding constraints on the first derivatives beyond the above, only adds non-holonomic constraints to the problem [16]. Below, we outline the same calculation leading to the PDE necessarily satisfied by any extremal f .

Without loss of generality, let $\eta + \eta' = 1$, then the optimization problem becomes

$$f_* = \arg \min_f \frac{1}{2} \left(1 - 2\eta\eta' + 4\zeta^2 \right), \quad \text{s.t.} \quad \eta + \eta' = 1. \quad (\text{S11})$$

Writing the ζ^2 term as $\int F(x_1, x_2, x_3, x_4) d\vec{x}$, where

$$F(x_1, x_2, x_3, x_4) := \kappa \prod_{i=1}^4 f(x_i) \frac{e^{-x_i^2/2}}{\sqrt{2\pi}}. \quad (\text{S12})$$

Here κ is a function of x_i s that is sharply concentrated around the line $x_1 = x_2 = x_3 = x_4$. The sharper the concentration of κ around the line, the closer the solution is to the optimum value.

We have effectively reduced the polynomial objective function in x to a multi-dimensional linear integral function in x_i s:

$$\min_f \left(1 - 2 \int_{\mathbb{R}^2} f(x_1)^2 f'(x_2)^2 \frac{e^{-(x_1^2+x_2^2)/2}}{2\pi} d\vec{x} + 4 \int_{\mathbb{R}^4} F(x_1, x_2, x_3, x_4) d\vec{x} \right), \quad (\text{S13})$$

where each x_i follows a isoperimetric constraint of the form:

$$\int_{\mathbb{R}} (f(x)^2 + f'(x)^2) \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx = 1. \quad (\text{S14})$$

Moreover, we have the (isoperimetric) ‘‘centering’’ constraint:

$$\int_{\mathbb{R}} f(x) \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx = 0. \quad (\text{S15})$$

Eqns. (S13), (S14) and (S15) lead to a standard Euler-Lagrange PDE, which can be simplified from symmetry considerations into an ODE, in the case above. See for example, Chapter 2 in [16]. The solution of that PDE gives the necessary condition, which is also usually sufficient, for the optimal activation function f .

3 Moments

Our main technical result will be stated in terms of generating functions arising from certain combinatorial settings, which are variants of certain standard problems and interesting by themselves. We introduce the following notation. Given a generating function $c(t) = \sum_k c_k t^k$, we denote $[c(t)]_i := c_i$ and for $c(t_1, t_2) = \sum_{k_1, k_2} c_{k_1, k_2} t_1^{k_1} t_2^{k_2}$, we denote $[c(t_1, t_2)]_{i, j} := c_{i, j}$.

Lemma 1. *The Stieltjes transform $G(z)$ of the spectral density of the Fisher information matrix of a single-hidden-layer neural network with squared loss, activation function f , weight matrices $W^{(1)}, W^{(2)} \in \mathbb{R}^{n \times n}$ with i.i.d. entries $W_{ij}^{(l)} \sim \mathcal{N}(0, \frac{1}{n})$, no biases, and i.i.d. inputs $X \sim \mathcal{N}(0, I_n)$ is given by the following integral as $n \rightarrow \infty$:*

$$G(z) = \frac{1}{z} P\left(\frac{1}{z}\right), \quad (\text{S16})$$

where the function $P(t)$ is given by the following series:

$$P(t) = \frac{P(t; \eta, \zeta)}{2} + \frac{P(t; \eta', \zeta)}{2} + \frac{1}{2} \sum_{n=1}^{\infty} \sum_{\substack{n_1, n_2 \\ n_1+n_2=n}} \sum_{d=1}^{\infty} [H(d, \lambda_1, \lambda_2)]_{n_1, n_2} [P_1(d, t)]_{n_1} [P_2(d, t)]_{n_2} t^{n+2(d-1)}, \quad (\text{S17})$$

The generating functions H , with formal variables λ_1 and λ_2 , is defined as follows:

$$H(d, \lambda_1, \lambda_2) = \frac{2\lambda_1\lambda_2 - \lambda_1^2\lambda_2 - \lambda_1\lambda_2^2}{(-1 + \lambda_1)^{d+1}(-1 + \lambda_2)^{d+1}}. \quad (\text{S18})$$

The generating functions $P_1(d, t)$ and $P_2(d, t)$ can be characterized in terms of the generating function $P(t)$ obtained in the paper [13]:

$$P_1(d, t) := \zeta^{d-1} \sum_{i=0}^{d-1} \frac{\binom{d-i-1}{d+i-1} \binom{d+i-1}{i}}{(1-P(t; \eta', \zeta))^{d-i-1}} \quad (\text{S19})$$

$$P_2(d, t) := \zeta^{d-1} \sum_{i=0}^{d-1} \frac{\binom{d-i-1}{d+i-1} \binom{d+i-1}{i}}{(1-P(t; \eta, \zeta))^{d-i-1}}, \quad (\text{S20})$$

where the generating function $P(t; \theta_1, \theta_2)$ with parameters θ_1 and θ_2 is given by the quadratic recurrence:

$$P(t; \theta_1, \theta_2) = 1 + P(t; \theta_1, \theta_2)(\theta_1 - \theta_2)t + \frac{P(t; \theta_1, \theta_2)\theta_2 t}{1 - P(t; \theta_1, \theta_2)\theta_2 t}. \quad (\text{S21})$$

Remark 3. The significance of $P(t)$ in [13] is that it completely characterizes the (Stieltjes transform of) the singular values of the resolvent of the matrix $f(WX)$ i.e., the output obtained from a single-hidden-layer neural network.

Lemma 2. The coefficients of the series P_1 and P_2 can be obtained by the following 1D integrals.

$$\begin{aligned} [P_1(d, t)]_{n_1} &= \frac{1}{\zeta} \int_{\mathbb{R}} (\lambda_1 - \eta' + \zeta)^d \lambda_1^{n_1-1} d\mu_1(\lambda_1) \\ [P_2(d, t)]_{n_2} &= \frac{1}{\zeta} \int_{\mathbb{R}} (\lambda_2 - \eta + \zeta)^d \lambda_2^{n_2-1} d\mu_2(\lambda_2) \end{aligned} \quad (\text{S22})$$

The proof idea is to plug-in the fractional binomial expansion inside the integral and verify that the corresponding two equations for P_1 and P_2 are indeed equal. The sums over n , n_1 , and n_2 are now trivial,

$$\begin{aligned} 2P(t) &= P(t; \eta, \zeta) + P(t; \eta', \zeta) + \\ &\sum_{d=1}^{\infty} t^{2(d-1)} \int_{\mathbb{R}} \int_{\mathbb{R}} \frac{t^2(2-t\lambda_1-t\lambda_2)(\lambda_1-\eta'+\zeta)^d(\lambda_2-\eta+\zeta)^d}{\zeta^2(-1+t\lambda_1)^{d+1}(-1+t\lambda_2)^{d+1}} d\mu_1(\lambda_1)d\mu_2(\lambda_2) \\ &= \sum_{d=0}^{\infty} t^{2(d-1)} \int_{\mathbb{R}} \int_{\mathbb{R}} \frac{t^2(2-t\lambda_1-t\lambda_2)(\lambda_1-\eta'+\zeta)^d(\lambda_2-\eta+\zeta)^d}{\zeta^2(-1+t\lambda_1)^{d+1}(-1+t\lambda_2)^{d+1}} d\mu_1(\lambda_1)d\mu_2(\lambda_2) \\ &= \int_{\mathbb{R}} \int_{\mathbb{R}} \frac{2-t(\lambda_1+\lambda_2)}{\zeta^2(\eta t^2(\lambda_1-\eta'+\zeta) + \eta' t^2(\zeta+\lambda_2) - (1+\zeta t)(-1+\zeta t+t(\lambda_1+\lambda_2)))} d\mu_1(\lambda_1)d\mu_2(\lambda_2) \end{aligned} \quad (\text{S23})$$

Simplifying this expression and utilizing eqn. (S16) yields the expression in eqn. (20).

4 Proof outline for general case

4.1 General Case

In this supplementary subsection, we remove the assumption that f is required to be linear. In the linear case, in eq. (16), we could directly compute the traces by applying the ‘‘mixed-product property’’ of Kronecker products to the expressions in eq. (15). Moreover, summing the resulting

series to obtain the Stieltjes transform was possible because the individual traces corresponded to Catalan numbers for which a generating function is known. In general, there is no analogous mixed product property to simplify the trace calculations, and we believe that the resulting series does not support a characterization with a single dimensional elliptic integral. With that caveat, we now proceed with the general case.

Our first task is to (asymptotically) evaluate traces of the form:

$$\sum_{\substack{i_1, \dots, i_{2k} \in \{0,1\} \\ i_1 + \dots + i_{2k} = k}} \text{tr} \left[M^{(1)^{i_1}} M^{(2)^{i_2}} \dots M^{(1)^{i_{2k-1}}} M^{(2)^{i_{2k}}} \right], \quad (\text{S24})$$

where $M^{(1)} = J^{(1)T} J^{(1)}$ and $M^{(2)} = J^{(2)T} J^{(2)}$.

Suppose that there are m examples¹, with each example i indexed by μ_i . For a given value of k , any trace as in eq. (S24), eventually consists of a sum over a product of per-example Jacobian matrices $J_{\mu_i}^{(1)}$ and $J_{\mu_i}^{(2)}$. Observe that,

$$\begin{aligned} \text{tr} H^{(0)k} &= \text{tr} \frac{1}{m^k} \sum_{\vec{\mu}} \left[\left([J_{\mu_1}^{(1)} J_{\mu_1}^{(2)}] [J_{\mu_1}^{(1)} J_{\mu_1}^{(2)}]^T \right) \right. \\ &\quad \times \left([J_{\mu_2}^{(1)} J_{\mu_2}^{(2)}] [J_{\mu_2}^{(1)} J_{\mu_2}^{(2)}]^T \right) \\ &\quad \left. \dots \left([J_{\mu_k}^{(1)} J_{\mu_k}^{(2)}] [J_{\mu_k}^{(1)} J_{\mu_k}^{(2)}]^T \right) \right], \end{aligned} \quad (\text{S25})$$

where $\mu_1 = \mu_k$ (by definition of the trace). For $n, m \rightarrow \infty$, where we take the limit over m first and then over n ,² observe that the μ_i s are all pairwise unequal, except that $\mu_1 = \mu_k$ as required. Of course, n and m are equal, by assumption.

Focusing on each term in the previous sum and expanding the “ 2×2 ” (block) matrices in the J s, we get traces over terms of the form:

$$\text{tr} \prod_{i=1}^d J_{\mu_i}^{(b_i)} J_{\mu_i}^{(b_i)T}, \quad (\text{S26})$$

where $b_i \in \{1, 2\}$ and μ_i s are unequal. By the cyclicity of the trace, we can rotate the last Jacobian to the front to re-pair terms and rewrite the above trace as:

$$\text{tr} \prod_{i=1}^d J_{\mu_{i+1}}^{(b_{i+1})T} J_{\mu_i}^{(b_i)}, \quad (\text{S27})$$

where addition in the μ_i subscripts is such that $d+1 \mapsto 1$. Finally, expanding the Jacobians into their constituent entries, using equations:

$$J_{ab, i\mu}^{(1)} = W_{ia}^{(2)} f' \left(\sum_k W_{ak}^{(1)} x_{k\mu} \right) x_{b\mu} \quad (\text{S28})$$

$$J_{cd, j\nu}^{(2)} = \delta_{cj} f \left(\sum_l W_{dl}^{(1)} x_{l\nu} \right). \quad (\text{S29})$$

we can write down each trace as in eq. (S27) as a polynomial in terms of the weights in W , the data points in X and f . The set of subindices occurring within the polynomial have some cyclic symmetries (follows from the cyclic arrangement of the Jacobians, as above). For example, in the trace calculations below:

¹In fact, we will assume $n = m$ i.e., we assume the width of the network and the number of examples both go to infinity at exactly the same rate. This way all our matrices are square, and our calculations are simplified.

²This is indeed the case when computing the limiting spectrum for the Fisher.

$$\begin{aligned}
\text{tr} \left[M^{(1)} M^{(1)} M^{(1)} \right] &= \sum_{abi\mu} \left[W_{i_1 a_1}^{(2)} W_{i_2 a_1}^{(2)} W_{i_2 a_2}^{(2)} W_{i_3 a_2}^{(2)} W_{i_3 a_3}^{(2)} W_{i_1 a_3}^{(2)} \right. \\
&\quad \times f'(z_{a_3 \mu_3}) f'(z_{a_3 \mu_1}) f'(z_{a_1 \mu_1}) f'(z_{a_1 \mu_2}) f'(z_{a_2 \mu_2}) f'(z_{a_2 \mu_3}) \\
&\quad \left. \times x_{b_1 \mu_1} x_{b_1 \mu_2} x_{b_2 \mu_2} x_{b_2 \mu_3} x_{b_3 \mu_3} x_{b_3 \mu_1} \right] \\
\text{tr} \left[M^{(2)} M^{(2)} M^{(2)} \right] &= \sum_{i\mu d} \left[\delta_{i_1 i_1} f(z_{d_1 \mu_1}) f(z_{d_1 \mu_2}) f(z_{d_2 \mu_2}) f(z_{d_2 \mu_3}) f(z_{d_3 \mu_3}) f(z_{d_3 \mu_1}) \right] \\
\text{tr} \left[M^{(1)} M^{(2)} \right] &= \sum_{abdi\mu} \left[W_{i_1 a_1}^{(2)} W_{i_1 a_1}^{(2)} f'(z_{a_1 \mu_1}) f'(z_{a_1 \mu_2}) x_{b_1 \mu_1} x_{b_1 \mu_2} f(z_{d_1 \mu_2}) f(z_{d_1 \mu_1}) \right] \\
\text{tr} \left[M^{(1)} M^{(2)} M^{(1)} M^{(2)} \right] &= \sum_{abdi} \left[W_{i_1 a_1}^{(2)} W_{i_2 a_1}^{(2)} W_{i_2 a_2}^{(2)} W_{i_1 a_2}^{(2)} f'(z_{a_1 \mu_1}) f'(z_{a_1 \mu_2}) f'(z_{a_2 \mu_3}) f'(z_{a_2 \mu_4}) \right. \\
&\quad \left. \times x_{b_1 \mu_1} x_{b_1 \mu_2} x_{b_2 \mu_3} x_{b_2 \mu_4} f(z_{d_1 \mu_2}) f(z_{d_1 \mu_3}) f(z_{d_2 \mu_4}) f(z_{d_2 \mu_1}) \right], \quad (\text{S30})
\end{aligned}$$

note that the subscript indices in i , a , b and μ have some cyclic symmetry.

Note that the polynomial is not necessarily multilinear because there may be identifications between various indices.³ Similarly, the polynomial is not completely symmetric because of restrictions on the indices induced by matrix multiplication and chain rule for taking derivatives. So, the cyclic symmetries are not completely trivial. Still, the structure induced among the indices is key to evaluating the trace. We map this structure to certain outer-planar graphs (as in [13]) and follow their machinery in evaluating the asymptotic expression for the trace. The latter effectively means that certain analytic details, like computing the saddle point asymptotic approximations can be hidden under the carpet.

Recall that the normalized trace, that we need to evaluate for the moment method, is of the form:

$$\mathbb{E} \frac{1}{n_1} \text{tr} M^k, \quad (\text{S31})$$

where the matrix $M := f(WX)$, f applied point-wise, and the weights W and input X are Gaussian distributed i.i.d. variables. The crux of the argument that we need from Section 4 of [13], is that the normalized trace, can be written as the integral:

$$\begin{aligned}
&\int \left[f(\sum_l W_{i_1 l} X_{l \mu_1}) f(\sum_l W_{i_2 l} X_{l \mu_1}) \cdots f(\sum_l W_{i_k l} X_{l \mu_k}) \right. \\
&\quad \left. \times f(\sum_l W_{i_1 l} X_{l \mu_k}) \right] \mathcal{D}W \mathcal{D}X. \quad (\text{S32})
\end{aligned}$$

After introducing auxiliary matrix valued variables Z and Λ , evaluating the X and W integrals, they are able to simplify the last integral to:

$$\begin{aligned}
&\int \left[\exp \left[-\frac{n}{2} \log \det \left| 1 + \frac{1}{n} \Lambda \Lambda^T \right| - i \text{tr} \Lambda Z \right] \right. \\
&\quad \left. \times f(Z_{i_1 \mu_1}) \cdots f(Z_{i_1 \mu_k}) \right] \mathcal{D}\lambda \mathcal{D}z, \quad (\text{S33})
\end{aligned}$$

where $\mathcal{D}\lambda = \prod_{\lambda_{\alpha\beta} \in \Lambda} \frac{d\lambda_{\alpha\beta}}{2\pi}$ and $\mathcal{D}z = \prod_{z_{\alpha\beta} \in Z} dz_{\alpha\beta}$. Finally, using saddle point approximations near the origin (the Gaussians are all mean zero), allows them to evaluate the last integral, and therefore the normalized trace, asymptotically as a polynomial in terms of η , ζ (the same η and ζ in our main result). And, they also compute the Stieltjes transform of M .

However, unlike [13], we have two matrices M_1 and M_2 and in order to evaluate traces of the form:

$$\sum_{\substack{i_1, \dots, i_{2k} \in \{0,1\} \\ i_1 + \dots + i_{2k} = k}} \text{tr} \left[M^{(1)^{i_1}} M^{(2)^{i_2}} \cdots M^{(1)^{i_{2k-1}}} M^{(2)^{i_{2k}}} \right], \quad (\text{S34})$$

³Except in the case of μs as mentioned above.

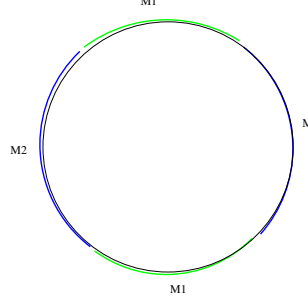


Figure S1: The configuration of a and d indices for $\text{tr}(M_1 M_2 M_1 M_2)$ is shown. In general, each green and blue arc corresponds to a admissible graph [13], connected by a single edge, corresponding to the indices of the common variable.

the resulting integrals as in eq. S32 are replaced by those of the form (cf. eq. S30):

$$\sum_{abdi} \int W_{i_1 a_1}^{(2)} W_{i_2 a_1}^{(2)} W_{i_2 a_2}^{(2)} W_{i_1 a_2}^{(2)} f'(Z_{a_1 \mu_1}) f'(Z_{a_1 \mu_2}) f'(Z_{a_2 \mu_3}) f'(Z_{a_2 \mu_4}) \\ \times X_{b_1 \mu_1} X_{b_1 \mu_2} X_{b_2 \mu_3} X_{b_2 \mu_4} f(Z_{d_1 \mu_2}) f(Z_{d_1 \mu_3}) f(z_{d_2 \mu_4}) f(Z_{d_2 \mu_1}) DWDX, \quad (\text{S35})$$

where the matrix $Z = WX$. A slightly more complicated scenario. On the other hand, we assume that all our matrices, weights as well as inputs are square i.e., dimension $n \times n$, and we assume unit variance throughout, which helps simplify the situation a little.

The crux of the question, when trying to evaluate an integral of the form S35 is what is the relative contribution of terms where certain sets of subscript indices, and therefore coefficients, are identified? What kind of terms dominate the expression when calculating the asymptotic value of the trace? What kind of identifications lead to sub-leading terms?

For example, the following lemma shows that when evaluating $\text{tr} M_1^k$ for $k \geq 3$, all b indices must be equal or else the trace is asymptotically zero. The reason being that the underlying covariances of $X_{b\mu}$ and $X_{b'\mu'}$ are zero when $b \neq b'$.

Lemma 3. *Given an expansion of $\text{tr} M_1^k$ in terms of the entries of W and X , the left indices of the X terms i.e., the b_i s in $\prod_i X_{b_i \mu_i}$, are either all equal, or the contribution to the trace is zero. Furthermore, the statement also holds for each run of M_1 s in eq. S34.*

The above is a structural result for the b indices. Similarly, consider the a and the d indices in eq. S30. We can arrange them as vertices of a cyclic graph to obtain a two-colored cycle corresponding to the a and d indices in $\text{tr}(M_1 M_2 M_1 M_2)$.

The green arcs correspond to the a indices (coming from the W s in M_1) and the blue arcs correspond to the d indices (coming from the Z terms in M_2 s). Note that for a given term, some of the a indices may be equal, in which case those vertices within the green arcs would be identified, and similarly for the vertices / indices in the blue arcs. This identification of vertices results in a complicated graph structure for the blue and green graphs, as opposed to a simple path structure.⁴

The next question can now be framed as follows: Every term arising from the trace corresponds to a graph, so which type of graphs lead to dominant terms i.e., terms that are asymptotically significant?

In [13], it was shown that only terms corresponding to the ‘‘admissible graphs’’, which are graphs consisting of edge disjoint cyclic blocks such that their planar dual forms a tree, contribute asymptotically to the trace in eq. S31.

The asymptotically dominant terms in the trace, corresponding to blue or green graphs still correspond to the *admissible graphs* defined in [13] i.e., the dominant terms will correspond to graphs that can be partitioned into edge-disjoint cyclic blocks, whose planar duals are trees. However, in our case the planar duals of the blue and green graphs, taken separately, may be disconnected (if there are no vertex identifications across arcs), and may therefore form forests. Despite this, the techniques

⁴ So the figure of a circle with arcs is deceptively simple!

of [13] are all still applicable and such “admissible graphs” form the leading terms of the trace in eq. S24. We skip the lengthy proof since the idea is the same as that in [13].

Now, keeping the above in mind, consider the evaluation of traces of the form $\text{tr} \left[M^{(1)^{i_1}} M^{(2)^{i_2}} \dots M^{(1)^{i_{2k-1}}} M^{(2)^{i_{2k}}} \right]$. In particular, consider the trace in eq. S30 and the corresponding integral in eq. S35 as a concrete example. Suppose that there are $2d$ alternations between M_1 and M_2 runs⁵, that the total degree of M_1 is n_1 , and that of M_2 is n_2 . Therefore, $d = 2$, $n_1 = 2$ and $n_2 = 2$ in our concrete example. We then define the following quantities:

- Let $[H(d, \lambda_1, \lambda_2)]_{n_1, n_2}$ denote the number of such monomials terms i.e., those having d alternations, and total degrees n_1 and n_2 in M_1 and M_2 , respectively.
- Let $P_1(d, t)$ denote the contribution of the M_1 terms to the trace. Equivalently, the expected value of the integral corresponding to the trace when any variables with “d” indices i.e., those that belong to the “blue” graph, are dropped.
- Let $P_2(d, t)$ denote the contribution of the M_2 terms to the trace. Equivalently, the expected value of the integral corresponding to the trace when any variables with “a” indices, those that belong to the “green” graph, are dropped.

The following lemmas can then be shown using only elementary methods.

Lemma 4.

$$P_1(d, t) = \zeta^{d-1} \sum_{i=0}^{d-1} \frac{\binom{d-i-1}{d+i-1} \binom{d+i-1}{i}}{(1 - P(t; \eta', \zeta))^{d-i-1}} \quad (\text{S36})$$

Lemma 5.

$$P_2(d, t) = \zeta^{d-1} \sum_{i=0}^{d-1} \frac{\binom{d-i-1}{d+i-1} \binom{d+i-1}{i}}{(1 - P(t; \eta, \zeta))^{d-i-1}} \quad (\text{S37})$$

Lemma 6.

$$H(d, \lambda_1, \lambda_2) = \frac{2\lambda_1\lambda_2 - \lambda_1^2\lambda_2 - \lambda_1\lambda_2^2}{(-1 + \lambda_1)^{d+1}(-1 + \lambda_2)^{d+1}}. \quad (\text{S38})$$

The proof idea is similar to that of the proofs in [13].

For the first two proofs, one assumes that the sequence of blue and green arcs (admissible graphs) on the “circle”, comprising of total n_1 and n_2 vertices is fixed. The proof follows by separating out i (say) out of d “arcs” and holding them to be disconnected i.e., no vertex identifications in-between those arcs. The remaining arcs are assumed to be connected i.e., there are non-crossing vertex identifications between those arcs. Recall that each arc corresponds to an admissible graph, so the proof follows by counting the contribution for each such configuration of admissible graphs. This is just a simple extension of the proof in [13], where the calculation is for a configuration consisting of one admissible graph. Counting the number of ways of selecting the i connected graphs is similar to (but not the same as) counting the number of non-crossing dissections of a $(i + 2)$ -gon. The latter, however, has a bijection to a standard Young’s tableaux (cf. [17]); while in our case, vertex identifications do not formally lead to lines but to “cyclic blocks” (cf. [13]), and that leads to a subtle, but asymptotically significant, difference in the final calculated value. As in [13], each configuration with i connected graphs leads to an integral over admissible graphs, which can be evaluated only in terms of the number of cyclic blocks in the block structure of the admissible graphs. The last property is crucial – it is what allows us to separate and write each term of the generating function as a multiple of three terms! This comprises the essential outline of the proofs of the first two lemmas above.

Finally, one needs to count the number of configurations of blue and green arcs given the total degree n . This is the last lemma. The proof of which consists of counting the number of ways of interlacing d blue and green “arcs”, with blue arcs covering n_1 points and green arcs covering n_2 points on a circle with n points. The above sketches the proof of Lemma 1.

⁵ Assume that the terms $M^{(1)^{i_1}} M^{(2)^{i_2}} \dots M^{(1)^{i_{2k-1}}} M^{(2)^{i_{2k}}}$ are laid out in a circle.