
Hessian Eigenspectra of More Realistic Nonlinear Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Given an optimization problem, the Hessian matrix and its eigenspectrum can be
2 used in many ways, ranging from designing more efficient second-order algorithms
3 to performing model analysis and regression diagnostics. When nonlinear models
4 and non-convex problems are considered, strong simplifying assumptions are often
5 made to make Hessian spectral analysis more tractable. This leads to the question
6 of how relevant the conclusions of such analyses are for realistic nonlinear models.
7 In this paper, we exploit tools from random matrix theory to make a *precise*
8 characterization of the Hessian eigenspectra for a broad family of nonlinear models
9 that generalizes the classical generalized linear models, without relying on strong
10 simplifying assumptions used previously. We show that, depending on the data
11 properties, the nonlinear response model, and the loss function, the Hessian can
12 have *qualitatively* different spectral behaviors: of bounded or unbounded support,
13 with single- or multi-bulk, and with isolated eigenvalues on the left- or right-hand
14 side of the main eigenvalue bulk. By focusing on such a simple but nontrivial
15 model, our analysis takes a step forward to unveil the theoretical origin of many
16 visually striking features observed in more realistic machine learning models.

17 1 Introduction

18 The Hessian is ubiquitous in applied mathematics, statistics, and machine learning (ML). Given
19 a (loss) function $L(\mathbf{w})$ of some parameters $\mathbf{w} \in \mathbb{R}^p$, the Hessian $\mathbf{H}(\mathbf{w}) \in \mathbb{R}^{p \times p}$ is defined as
20 the second derivative of the objective function with respect to the model parameter, i.e., $\mathbf{H}(\mathbf{w}) =$
21 $\partial^2 L(\mathbf{w}) / (\partial \mathbf{w} \partial \mathbf{w}^T)$. When a ML model is being trained, it is common to parameterize that model by
22 \mathbf{w} , and then train that model by minimizing some (smooth) loss function $L(\mathbf{w})$, with the associated
23 Hessian $\mathbf{H}(\mathbf{w})$, e.g., by backpropagating the error to improve \mathbf{w} [25]. Alternatively, once a ML
24 model is trained, the Hessian (and the related Fisher information matrix [60, 62]) can be examined to
25 identify outliers, perform diagnostics, and/or engage in other sorts of model validation [29, 72, 57].

26 For convex problems, the Hessian $\mathbf{H}(\mathbf{w})$ provides detailed information on how to adjust the gradient
27 to achieve improved convergence, e.g., in Newton-like methods. For non-convex problems, the
28 properties of the local loss “landscape” around a given point \mathbf{w} in the parameter space is of central
29 significance [17, 34, 12, 37, 70, 71, 72]. In this case, most obviously, the signs of the smallest
30 and largest Hessian eigenvalue can be used to test whether a given \mathbf{w} is a local maximum, local
31 minimum, or a saddle point. More subtly, the Hessian eigenvalue distribution characterizes the local
32 curvature of the loss function and provides direct access to, for instance, the fraction of negative
33 Hessian eigenvalues that determines the number of (local) descent directions, a quantity that is directly
34 connected to the rates of convergence of various optimization algorithms [31].

35 For theoretical analysis of neural network (NN) models, Hessian eigenspectra are often assumed to
36 follow well-known random matrix distributions such the Marčenko–Pastur law [42] or the Wigner’s

37 semicircle law [64]. This enables one to use Random Matrix Theory (RMT), but it involves (for NNs,
38 at least) making relatively strong simplifying assumptions (e.g., the Hessian can be decomposed as
39 the sum of the two freely independent matrices, the residual error, data feature, and weights are all
40 composed of i.i.d. zero mean normal random variables) [52, 53, 14]. A somewhat more realistic
41 setup involves using a so-called *spiked model* (or a spiked covariance model) [2, 4, 39]. In this
42 case, the matrix follows a *signal-plus-noise* model and consists of *full rank* random noise matrix and
43 *low rank* statistical information structure.¹ The “signal” eigenvalues are generally larger than the
44 noisy “bulk” eigenvalues; and the maximum eigenvalues, when isolated from the bulk, are referred
45 to as the “spikes.” A substantial theory-practice gap exists, however. In both toy examples [26] and
46 state-of-the-art NN models [70, 71, 72, 73, 58, 19], the strong simplifying assumptions are far from
47 satisfactory. (A similar theory-practice gap has been observed for other NN matrices to which RMT
48 has been applied, perhaps most notably weight matrices [43, 44].) A more precise understanding of
49 the Hessian eigenspectra (and its dependence on input data structure, the underlying response model
50 and model parameters, as well as the loss function) for more practical models is needed.

51 1.1 Our approach

52 In this article, we address these issues, in a setting that is simple enough to be analytically tractable
53 but complex enough to shed light on realistic large-scale models. We consider a family of generalized
54 generalized linear models (G-GLMs) that extends the popular generalized linear model (GLM)
55 [18, 29]; and we show that, even for such simple models, the key simplifying assumptions used
56 in previous theoretical analyses of Hessian can be very inexact. In particular, apart from a few
57 special cases (including linear least squares and logistic regression with homogeneous features), most
58 Hessians of G-GLM are *not* close to the Marčenko–Pastur and/or the semicircle law. Instead, the
59 corresponding Hessian depends on the input feature structure, the underlying response model, and the
60 loss function, in a more involved fashion that can be precisely characterized by the proposed analysis.

61 The G-GLM describes a generalized linear relation between the input feature $\mathbf{x}_i \in \mathbb{R}^p$ and the
62 corresponding response y_i , in the sense that there exists some parameters $\mathbf{w}_* \in \mathbb{R}^p$ such that for
63 given $\mathbf{w}_*^\top \mathbf{x}_i$, the response y_i is independently drawn from

$$y_i \sim f(y \mid \mathbf{w}_*^\top \mathbf{x}_i) \quad (1)$$

64 for some conditional density function $f(\cdot \mid \cdot)$. This extends the classical GLM such as

$$\text{logistic model: } \mathbb{P}(y = 1 \mid \mathbf{w}_*^\top \mathbf{x}) = (1 + e^{-\mathbf{w}_*^\top \mathbf{x}})^{-1}, \quad y \in \{-1, 1\}, \quad (2)$$

65 and covers a large family of models in applications in statistics and ML. Other examples include: i)
66 the (noisy) nonlinear factor model [13] where $y \sim \mathcal{N}(g(\mathbf{w}_*^\top \mathbf{x}), \sigma^2)$ for some nonlinear $g: \mathbb{R} \rightarrow \mathbb{R}$
67 and $\sigma > 0$; ii) the (noiseless) phase retrieval model [20] with $y = |\mathbf{w}_*^\top \mathbf{x}|^2$, in which case one wishes
68 to reconstruct \mathbf{w}_* from its (squared) magnitude measurements; and iii) the single-layer NN model
69 $y = \sigma(\mathbf{w}_*^\top \mathbf{x})$ for some nonlinear activation function $\sigma(t)$ such as the tanh-sigmoid $\sigma(t) = \tanh(t)$.

70 For a given training set $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ of size n , the standard approach to obtain/recover the parameter
71 $\mathbf{w}_* \in \mathbb{R}^p$ of a G-GLM is to solve the following optimization problem

$$\min_{\mathbf{w}} L(\mathbf{w}) = \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, \mathbf{w}^\top \mathbf{x}_i), \quad (3)$$

72 for some loss function $\ell(y, h): \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, e.g., the negative log-likelihood of the observation
73 model within the maximum likelihood estimation framework [29] such as the logistic loss $\ell(y, h) =$
74 $\ln(1 + e^{-yh})$ in the case of logistic model (2). In many applications, however, the optimization
75 problem in (3) may *not* be convex, for example to achieve superior robustness and/or accuracy
76 [45, 65, 9], and can be NP-hard in general (the noiseless phase retrieval model $y = (\mathbf{w}_*^\top \mathbf{x})^2$ with the
77 square loss $\ell(y, h) = (y - h^2)^2$ as an example [8]). As we shall see, in such non-convex G-GLMs, the
78 dominant Hessian eigenvector can be shown, in some cases, to positively correlate with the sought-for
79 parameter \mathbf{w}_* and therefore be used as the initialization of gradient descent methods [8, 35, 32]. This
80 particularly motivates our study of the possible isolated Hessian eigenvalue-eigenvector pairs.

81 1.2 Our main contributions

82 The main contribution of this work is the *exact* characterization of Hessian eigenspectra for the family
83 of G-GLMs, in the high-dimensional regime where the feature dimension p and the sample size n are
84 both large and comparable. Precisely, we establish:

¹Since the *same* informative pattern is repeated in each row or column of the matrix.

- 85 1. the limiting eigenvalue distribution of the Hessian matrix (Theorem 1); and
 86 2. the behavior of (possible) isolated eigenvalue-eigenvector pairs (Theorem 2 and 3),
 87 as a function of the dimension ratio $c = \lim p/n$, feature statistics, the loss function ℓ in (3), and the
 88 underlying response model in (1). Our results are based on a technical result of independent interest:
 89 3. a *deterministic equivalent* (Theorem 4) of the random *resolvent* $\mathbf{Q}(z) = (\mathbf{H} - z\mathbf{I}_p)^{-1}$, for
 90 $z \in \mathbb{C}$ not an eigenvalue of \mathbf{H} , of the generalized sample covariance:

$$\mathbf{H} \equiv \mathbf{H}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell''(y_i, \mathbf{w}^\top \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top \equiv \frac{1}{n} \mathbf{X} \mathbf{D} \mathbf{X}^\top, \quad \mathbf{D} \equiv \text{diag}\{\ell''(y_i, \mathbf{w}^\top \mathbf{x}_i)\}_{i=1}^n \quad (4)$$

91 for $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$ and $\ell''(y, h) \equiv \partial^2 \ell(y, h) / \partial h^2$, as $n, p \rightarrow \infty$ with $p/n \rightarrow c \in (0, \infty)$,
 92 under the setting of *generic* Gaussian feature $\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$, for $\boldsymbol{\mu} \in \mathbb{R}^p$ and positive definite
 93 covariance $\mathbf{C} \in \mathbb{R}^{p \times p}$. We also demonstrate our results empirically by showing that:

- 94 4. for a given response model (1), the Hessian eigenvalue distribution depends on the choice
 95 of loss function and the data/feature statistics in an intrinsic manner, e.g., bounded versus
 96 unbounded support and single- versus multi-bulk in Fig 2; and
 97 5. there may exist two *qualitatively* different spikes—one due to data *signal* $\boldsymbol{\mu}$ and the other
 98 due to \mathbf{w}_* or \mathbf{w} and thus the *underlying model*—which may appear on different sides of the
 99 main bulk, and their associated phase transition behaviors are characterized (Fig 4 versus 5).

100 To have a more clear picture of our contribution, we compare, in Fig 1a and 1b, the Hessian eigenvalues
 101 for the logistic model (2) with the (maximum likelihood) logistic loss $\ell(y, h) = \ln(1 + e^{-yh})$, for
 102 different choices of \mathbf{w} in the parameter space. A nontrivial interplay between the response model,
 103 feature statistics and the parameter \mathbf{w} is reflected by the range of the Hessian eigenvalue support
 104 and an additional right-hand spike in Fig 1b, as confirmed by our theory. For phase retrieval model
 105 $y = (\mathbf{w}_*^\top \mathbf{x})^2$ with square loss $\ell(y, h) = (y - h^2)^2/4$, the non-convex nature of the problem is
 106 reflected by a (relatively large) fraction of negative Hessian eigenvalues in Fig 1c. We also note that
 107 the top eigenvector (that corresponds to the largest eigenvalue) contains structural information of the
 108 underlying model, in the sense that it is positively correlated with \mathbf{w}_* , as predicted by our theory. This
 109 is indeed connected to the Hessian-based initialization scheme widely used in non-convex problems.

110 We conclude by emphasizing that, by focusing on the simple yet fundamental G-GLM, we obtain
 111 results that improve upon and are different than previous efforts in the following aspects:

- 112 i) We provide *precise* asymptotic characterizations of the Hessian eigenspectra that goes
 113 beyond, e.g., [6], where only Hessian lower bounds are given in the case of logistic model
 114 with logistic loss: our methodology and theoretical results hold much more generally for
 115 the family of G-GLM with arbitrary loss. As illustrating examples, we discuss linear least
 116 squares in Sec 3.1, logistic model with different choices of loss function in Fig 2, phase
 117 retrieval model in Figure 1c, and more in Sec 4 in the appendix.
 118 ii) We extend the results in [52, 48, 41, 46] to G-GLMs by considering *generic* data statistics
 119 and loss function, whereas in [52, 48, 41, 46] only much more homogeneous models are
 120 discussed, and sometimes under unrealistic assumption, e.g., the Hessian can be decomposed
 121 as the sum of the two freely independent matrices, the residual error, data feature, and weights
 122 are all composed of i.i.d. zero mean normal random variables [52, 53]).
 123 iii) Instead of focusing solely on the main eigenvalue bulk as in [52, 53], our results also shed
 124 novel light on the isolated eigenvalues (above and/or below the bulk) that are empirically
 125 observed in the Hessian of modern NNs [55, 21, 40, 48, 47], as well as on the associated
 126 eigenvectors that are shown closely connected to NN training dynamics [27]. Also, relative
 127 to [52, 53], we show *qualitatively* different behaviors for the Hessian eigenspectra, e.g.,
 128 bounded versus unbounded support, single- versus multi-bulk as in Figure 2. To our
 129 knowledge, these are *not* covered in the existing Hessian literature.

130 1.3 Related work

131 Here, we provide a brief review of related previous efforts, see more discussions in the appendix.

132 **Random matrix theory.** Random matrices of the type (4) are related to the *separable covariance*
 133 *model* [74, 16] in the RMT literature, which is of the form $\mathbf{C}^{\frac{1}{2}} \mathbf{Z} \mathbf{D} \mathbf{Z}^\top \mathbf{C}^{\frac{1}{2}}$, for random \mathbf{Z} and \mathbf{C}, \mathbf{D}

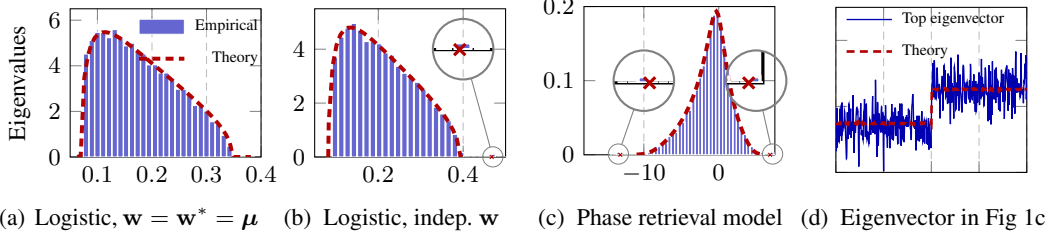


Figure 1: Illustration of our main results: eigenspectral properties of the Hessian of G-GLMs with $p = 800$, $n = 6000$ and $\mathbf{C} = \mathbf{I}_p$. **Fig 1a versus 1b**: absence versus presence of a right-hand side spike for different choices of \mathbf{w} , logistic model (2) with logistic loss, and $\mathbf{w}_* = \boldsymbol{\mu} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p/p)$. **Fig 1c versus 1d**: the Hessian eigenspectra have a rather different shape (as opposed to the Marčenko–Pastur-like in Fig 1a and 1b) for the (non-convex) phase retrieval model (1c) and the top eigenvector is known in this case to be a (noisy) estimate of \mathbf{w}_* (1d), as confirmed by our theory. With square loss $\ell(y, h) = (y - h)^2/4$, $\mathbf{w}_* = [-2 \cdot \mathbf{1}_{p/2}; 2 \cdot \mathbf{1}_{p/2}]/\sqrt{p}$, $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p/p)$ and $\boldsymbol{\mu} = \mathbf{0}$.

134 independent of \mathbf{Z} . Our results generalize this, in the sense that we allow \mathbf{D} to depend on \mathbf{Z} , in a
 135 possibly nonlinear fashion, per (4). This is of direct interest for the Hessian of G-GLMs.

136 **Hessian eigenspectra.** The eigenspectra of Hessian matrices arising in ML models (in particular,
 137 for NNs) have attracted considerable interest recently [55, 56, 10, 21, 66, 23, 30, 19, 58, 72, 73, 24].
 138 However, these investigations are either limited to empirical evaluation [55, 56] or built upon some-
 139 what unrealistic simplifying assumptions and reduce to the “mixed” behavior of Marčenko–Pastur
 140 and semicircle law [52, 14]. In contrast, here we focus on the more tractable example of G-GLM and
 141 provide *precise* results on the Hessian eigenspectra for structural feature on arbitrary loss.

142 **Spectral initialization in non-convex problems.** A popular initialization scheme (of gradient-based
 143 methods) for non-convex problems is the *spectral initialization*, where the top eigenvectors of some
 144 Hessian-type matrices are used as gradient descent initialization [7, 35, 32, 38, 1]. In [41], which
 145 was generalized in [46], the authors evaluated the eigenspectrum asymptotics of $\frac{1}{n} \sum_{i=1}^n f(y_i) \mathbf{x}_i \mathbf{x}_i^T$,
 146 for some $f: \mathbb{R} \rightarrow \mathbb{R}$ and $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$. Their technical approach is, however, limited to the case of
 147 very homogeneous features. Here we generalize the analysis in [41, 46] to the Hessian of G-GLM,
 148 by developing a systematic approach to account for both feature structures and loss functions.

149 **Scalable second-order methods.** Second order methods are among the most powerful optimization
 150 methods that have been designed, and there have been several attempts to use their many advantages
 151 for machine learning applications [68, 63, 54], particularly for training NNs [72, 73, 58, 19, 67]. We
 152 expect that our precise characterization of the Hessian sheds new light on the understanding and
 153 improved design of (e.g., computationally) more efficient second-order methods.

154 2 Main results

155 In the section, we present our main results: on the limiting Hessian eigenspectrum (in Sec 2.1); and
 156 on the behavior of the (possible) isolated eigenvalue-eigenvector(s) (in Sec 2.2). These two main
 157 results depend on a technical deterministic equivalent result for the Hessian resolvent (in Sec 2.3),
 158 which is of independent interest. We position ourselves in the following high-dimensional regime.

159 **Assumption 1** (High-dimensional asymptotics). *As $n, p \rightarrow \infty$ with $p/n \rightarrow c \in (0, \infty)$, we have*
 160 $\max\{\|\mathbf{w}\|, \|\mathbf{w}_*\|\} = O(1)$ and $\mathbf{x}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$ with $\max\{\|\boldsymbol{\mu}\|, \|\mathbf{C}\|\} = O(1)$.

161 2.1 Limiting spectral measure

162 Our first result is the limiting Hessian eigenvalue distribution. This is a direct consequence of our
 163 main technical Theorem 4 and is proven in Sec A.2 of the appendix.

164 **Theorem 1** (Limiting spectral measure). *Let Assumption 1 hold, we have, as $n, p \rightarrow \infty$ with*
 165 $p/n \rightarrow c \in (0, \infty)$, *the empirical spectral measure² $\mu_{\mathbf{H}}$ of the Hessian matrix \mathbf{H} defined in (4)*

²That is, the normalized counting measure of the eigenvalues $\lambda_i(\mathbf{H})$ of \mathbf{H} , i.e., $\mu_{\mathbf{H}} \equiv \frac{1}{p} \sum_{i=1}^p \delta_{\lambda_i(\mathbf{H})}$.

166 converges weakly and almost surely to a probability measure μ , defined through its Stieltjes transform
 167 $m(z) = \int (t - z)^{-1} \mu(dt)$ as the unique solution to³

$$m(z) = \frac{1}{p} \text{tr} \bar{\mathbf{Q}}_b(z), \quad \delta(z) = \frac{1}{n} \text{tr} (\mathbf{C} \bar{\mathbf{Q}}_b(z)), \quad \bar{\mathbf{Q}}_b^{-1}(z) \equiv \mathbb{E} \left[\frac{g \mathbf{C}}{1 + g \delta(z)} \right] - z \mathbf{I}_p, \quad (5)$$

168 where

$$g \equiv \partial^2 \ell(y, h) / \partial h^2, \quad \text{for } h = \mathbf{w}^\top \mathbf{x} \sim \mathcal{N}(\mathbf{w}^\top \boldsymbol{\mu}, \mathbf{w}^\top \mathbf{C} \mathbf{w}), \quad (6)$$

169 and y and ℓ defined respectively in (1) and (3). Moreover, if we denote ν the law of g and assume the
 170 empirical spectral measure of \mathbf{C} converges to $\tilde{\nu}$ as $p \rightarrow \infty$, then (5) can be compactly written as

$$m(z) = \int \left(-z + \tilde{t} \int \frac{t}{1 + t \delta(z)} \nu(dt) \right)^{-1} \tilde{\nu}(d\tilde{t}), \quad \delta(z) = \int \frac{c \tilde{t}}{-z + \tilde{t} \int \frac{t}{1 + t \delta(z)} \nu(dt)} \tilde{\nu}(d\tilde{t}). \quad (7)$$

171 In the form of (7), the (Stieltjes transform of the) limiting Hessian spectral measure μ is determined
 172 by the ratio $c = \lim p/n$ and the two measures ν and $\tilde{\nu}$. This formulation is closely connected to the
 173 separable covariance model [36, 5, 50, 16, 69] in RMT. Moreover, if $\nu(dt) = \delta_1(t)$ is a Dirac mass
 174 at one, this reduces to the popular sample covariance model [59]; taking further $\tilde{\nu}(d\tilde{t}) = \delta_1(\tilde{t})$ gives
 175 the Marčenko-Pastur law. See Sec 3.1 for numerical evaluations of these special cases. In particular,
 176 the support of the (limiting) Hessian spectrum μ is directly linked to that of ν and $\tilde{\nu}$.

177 **Remark 1** (Hessian eigen-support). *Under Assumption 1, the (limiting) spectral measure $\tilde{\nu}$ of \mathbf{C}*
 178 *has bounded support. However, this may not be the case for ν , the law of g defined in (6). Since the*
 179 *Hessian eigenvalue distribution μ is of compact support if and only if both ν and $\tilde{\nu}$ have compact*
 180 *support [16, Proposition 3.4], μ may be of unbounded support, depending on the model and the loss.*

181 An example of unbounded μ is the phase retrieval model with $y = (\mathbf{w}_*^\top \mathbf{x})^2$ and square loss $\ell(y, h) =$
 182 $(y - h^2)^2/4$, for which we have $g = 3(\mathbf{w}^\top \mathbf{x})^2 - (\mathbf{w}_*^\top \mathbf{x})^2$ for $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$. As a consequence,
 183 with say $\mathbf{w}_* = \mathbf{w}$, g follows a chi-square distribution with one degree of freedom and has thus
 184 unbounded support. This corresponds to Fig 1c, where the Hessian spectrum has a “heavier” tail
 185 compared to Fig 1a (logistic model), and the empirically observed “isolated” eigenvalue is due to
 186 a finite-dimensional effect and will be “buried” in the noisy main bulk for larger values of n, p .
 187 Therefore, aiming for an (almost surely) isolated eigenvalue-eigenvector (e.g., to recover the model
 188 parameter \mathbf{w}_* using the top Hessian eigenvector), some preprocessing function f must be applied.
 189 This has been discussed in previous work [41, 46] and corresponds to the so-called trimming strategy
 190 in phase retrieval [11], with for instance the truncation function $f(t) = \delta_{|t| \leq \epsilon}$ for some $\epsilon > 0$.

191 Another example of unbounded μ is when the exponential loss [22] is used. Precisely, consider the
 192 logistic model (2) with $\ell(y, h) = \exp(-yh)$, we have that $g = \exp(-yh)$ for $h \sim \mathcal{N}(\mathbf{w}^\top \boldsymbol{\mu}, \mathbf{w}^\top \mathbf{C} \mathbf{w})$
 193 which follows a log-normal distribution and has unbounded support. As such, the (limiting) Hessian
 194 eigenvalue distribution μ has also unbounded support. On the other hand, with logistic loss $\ell(y, h) =$
 195 $\ln(1 + e^{-yh})$, one has $g \leq 1/4$ and μ is guaranteed to have bounded support. In Fig 2a and 2b, the
 196 empirical Hessian eigenvalues and the limiting distributions are compared for logistic and exponential
 197 losses, with a more “heavy-tailed” behavior observed for the exponential loss.

198 Clearly, depending on the measures ν (of g , which depends on feature statistics, loss and underlying
 199 model) and $\tilde{\nu}$ (of spectrum of feature covariance \mathbf{C}), the Hessian spectrum can have very different
 200 forms. Here we compare the empirical Hessian eigenvalues with their limiting behaviors per Theo-
 201 rem 1 for different feature covariance structures⁴. In particular, one may observe a single main bulk
 202 with more “compact” Hessian spectrum as in Fig 2c or multiple bulks (two in the case of Fig 2d) with
 203 Hessian eigenvalues more “spread-out”, depending on the feature covariance structure $\tilde{\nu}$. In the form
 204 of (7), the condition for the existence of multi-bulk eigenspectrum has been thoroughly discussed in
 205 [16, Sec 3.2–3.4] and can be numerically evaluated with ease.

206 As a side remark, the “multi-bulk” behavior similar to Fig 2d has been empirically observed in
 207 Hessians of NNs in [40, 48] and is believed to be due to the classification structure within data
 208 (i.e., the data vectors are drawn from a mixture of distributions). Here, we provide an alternative
 209 explanation via feature covariance structure that can be observed beyond the classification setting.

³Uniqueness is ensured in such a way that $\Im[m(z)] \cdot \Im[z] > 0$ for $\Im[z] \neq 0$ and $zm(z) < 0$ for $\Im[z] = 0$, so that $(z, m(z))$ is a valid Stieltjes transform couple, see more details in [28].

⁴Covariance describes the joint variability or the “correlation” between entries of the feature vector and it is of particular significance in the analysis of image (with local structure) and time series data.

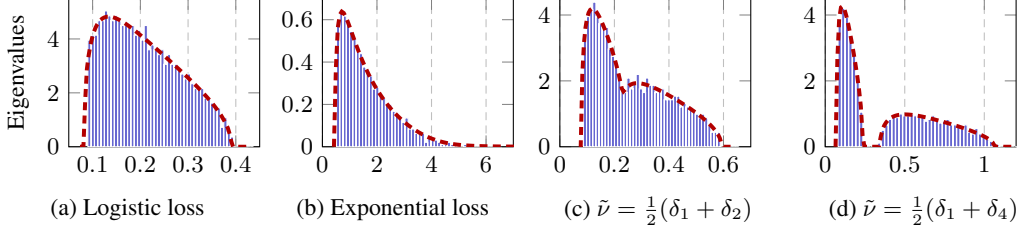


Figure 2: Impact of loss function: bounded (2a, with logistic loss) versus unbounded (2b, with exponential loss) Hessian eigenvalues, with $p = 800$, $n = 6\,000$, logistic model in (2) with $\boldsymbol{\mu} = \mathbf{0}$, $\mathbf{C} = \mathbf{I}_p$, $\mathbf{w}_* = \mathbf{0}$ and $\mathbf{w} = [-\mathbf{1}_{p/2}, \mathbf{1}_{p/2}]/\sqrt{p}$. Impact of feature covariance: Hessian spectrum of single- (2c, with $\mathbf{C} = \text{diag}[\mathbf{1}_{p/2}; 2 \cdot \mathbf{1}_{p/2}]$) versus multi-bulk (2d, with $\mathbf{C} = \text{diag}[\mathbf{1}_{p/2}; 4 \cdot \mathbf{1}_{p/2}]$), with $p = 800$, $n = 6\,000$, logistic model with $\mathbf{w}^* = \mathbf{0}_p$, $\mathbf{w} = \boldsymbol{\mu} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p/p)$.

2.2 Isolated eigenvalues and eigenvectors

As discussed in Remark 1, under Assumption 1, the Hessian has bounded (limiting) eigen-support if and only if ν , the law of g , has bounded support. Under this condition (or, after the application of some function f so that $f(g)$ is bounded), we can then talk about the (possible) isolated Hessian eigenvalues, as in the following result, the proof of which is given in Sec A.3 of the appendix.

Theorem 2 (Isolated eigenvalues). *In the setting of Theorem 1, assume that the law ν of the random variable g defined in (6) is of bounded support, define*

$$\mathbf{G}(z) = \mathbf{I}_3 + \boldsymbol{\Lambda}(z)\mathbf{V}^T\bar{\mathbf{Q}}_b(z)\mathbf{V} \in \mathbb{R}^{3 \times 3}, \quad (8)$$

with $\bar{\mathbf{Q}}_b(z)$, $\delta(z)$ defined in (5), $\mathbf{V} \equiv [\boldsymbol{\mu}, \mathbf{C}\mathbf{w}_*, \mathbf{C}\mathbf{w}] \in \mathbb{R}^{p \times 3}$, $\mathbf{U} \equiv \mathbf{C}^{\frac{1}{2}}[\mathbf{w}_*, \mathbf{w}] \in \mathbb{R}^{p \times 2}$ and

$$\boldsymbol{\Lambda}(z) \equiv \mathbb{E} \frac{g}{1 + g \cdot \delta(z)} \begin{bmatrix} 1 & (\mathbf{U}^+ \mathbf{z})^T \\ \mathbf{U}^+ \mathbf{z} & \mathbf{U}^+ \mathbf{z} (\mathbf{U}^+ \mathbf{z})^T - (\mathbf{U}^T \mathbf{U})^+ \end{bmatrix}, \quad \mathbf{z} = \mathbf{C}^{-\frac{1}{2}}(\mathbf{x} - \boldsymbol{\mu}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p), \quad (9)$$

where we denote \mathbf{U}^+ the Moore–Penrose pseudoinverse of \mathbf{U} . Then, for λ such that $\mathbf{G}(\lambda)$ has a zero eigenvalue (of multiplicity one), there exists an eigenvalue $\hat{\lambda}$ of \mathbf{H} such that $\hat{\lambda} - \lambda \xrightarrow{a.s.} 0$.

Theorem 2 provides an asymptotic characterization of the possible isolated Hessian eigenvalues by computing the determinant of the much smaller (three-by-three) deterministic matrix \mathbf{G} closely related to the key quantity $\delta(z)$ defined in Theorem 1. Note that, Theorem 2 does not provide, at least explicitly, the *phase transition* condition under which these spikes become “isolated” from the main bulk. As we shall see in more details in Sec 3.2, two types of quantitatively different phase transitions can be characterized, due to the data “signal” $\boldsymbol{\mu}$ and the underlying model, respectively.

We can also analyze the associated isolated eigenvectors. First note that, in the infinite data regime (i.e., for $n \rightarrow \infty$ with p fixed), we have, by the strong law of large numbers, that $\mathbf{H} \xrightarrow{a.s.} \mathbb{E}[\mathbf{H}]$, with

$$\mathbb{E}[\mathbf{H}] = \mathbb{E}[\ell''(y, \mathbf{w}^T \mathbf{x}) \mathbf{x} \mathbf{x}^T] = \mathbb{E}[g] \cdot \mathbf{C} + \mathbf{V} \begin{bmatrix} 1 & \mathbb{E}[g \cdot \mathbf{U}^+ \mathbf{z}]^T \\ \mathbb{E}[g \cdot \mathbf{U}^+ \mathbf{z}] & \mathbf{U}^+ \mathbb{E}[g \cdot (\mathbf{z} \mathbf{z}^T - \mathbf{I}_p)] (\mathbf{U}^+)^T \end{bmatrix} \mathbf{V}^T.$$

As a consequence, it is expected that in the large $n, p \rightarrow \infty$ limit, the top eigenvectors of \mathbf{H} could also be related to the columns of \mathbf{V} . This is the case in Fig 1d, where the top eigenvector is observed to be a “noisy” version of the model parameter \mathbf{w}_* . More precisely, for $(\hat{\lambda}, \hat{\mathbf{u}})$ an isolated eigenpair of \mathbf{H} , the projection $\mathbf{V}^T \hat{\mathbf{u}} \mathbf{V} \in \mathbb{R}^{3 \times 3}$ can be shown to be asymptotically close to a deterministic matrix. This measures the “cosine-similarly” between the Hessian isolated eigenvector $\hat{\mathbf{u}}$ with any column of \mathbf{V} and consequently the performance of using $\hat{\mathbf{u}}$ as an estimate of, for instance the model parameter \mathbf{w}_* for $\mathbf{C} = \mathbf{I}_p$. This result is given in the following theorem, which is proven in Appendix A.3.

Theorem 3 (Isolated eigenvectors). *In the setting of Theorem 2, for an isolated eigenvalue-eigenvector pair $(\hat{\lambda}, \hat{\mathbf{u}})$ of \mathbf{H} and λ the asymptotic position (of $\hat{\lambda}$) given in Theorem 2, then*

$$\mathbf{V}^T \hat{\mathbf{u}} \hat{\mathbf{u}}^T \mathbf{V} = -\mathbf{V}^T \bar{\mathbf{Q}}_b(\lambda) \mathbf{V} \cdot \Xi(\lambda) + o(1), \quad \Xi(\lambda) = (\mathbf{v}_{l, \mathbf{G}}^T \mathbf{G}'(\lambda) \mathbf{v}_{r, \mathbf{G}})^{-1} \cdot \mathbf{v}_{r, \mathbf{G}} \mathbf{v}_{l, \mathbf{G}}^T,$$

for $\bar{\mathbf{Q}}_b(z)$ and $\mathbf{G}(z)$ defined in (5) and (8), respectively, $\mathbf{v}_{l, \mathbf{G}}, \mathbf{v}_{r, \mathbf{G}} \in \mathbb{R}^3$ the left and right eigenvectors of $\mathbf{G}(\lambda)$ associated with eigenvalue zero, and $\mathbf{G}'(\lambda)$ the derivative of $\mathbf{G}(z)$ with respect to z evaluated at $z = \lambda$.

240 **2.3 Technical tool: deterministic equivalent**

241 Our main technical tool to derive Theorem 1, 2 and 3 is a so-called deterministic equivalent [28, 15]
 242 result for the Hessian resolvent $\mathbf{Q}(z) = (\mathbf{H} - z\mathbf{I}_p)^{-1}$, that provides simultaneous access to the Hessian
 243 limiting eigenvalue distribution and the behavior of the possible isolated eigenpairs. Precisely, the
 244 normalized trace $\text{tr } \mathbf{Q}(z)/p$ gives the Stieltjes transform $m_{\mathbf{H}}(z) = \int (t-z)^{-1} \mu_{\mathbf{H}}(dt)$ of the empirical
 245 spectral measure $\mu_{\mathbf{H}}$ of \mathbf{H} , from which $\mu_{\mathbf{H}}$ can be recovered. Also, for $(\hat{\lambda}, \hat{\mathbf{u}})$ an eigenpair of interest,
 246 with Cauchy’s integral formula we have $|\mathbf{w}^T \hat{\mathbf{u}}|^2 = -\frac{1}{2\pi i} \oint_{\Gamma_\lambda} \mathbf{w}^T \mathbf{Q}(z) \mathbf{w} dz$, for a deterministic
 247 vector $\mathbf{w} \in \mathbb{R}^p$ and Γ_λ a positively oriented contour surrounding *only* $\hat{\lambda}$. As such, for $\bar{\mathbf{Q}}(z)$ a
 248 deterministic equivalent of $\mathbf{Q}(z)$, that is, $\mathbf{Q}(z) \leftrightarrow \bar{\mathbf{Q}}(z)$ with $\text{tr } \mathbf{A}(\mathbf{Q}(z) - \bar{\mathbf{Q}}(z))/p \rightarrow 0$ and
 249 $\mathbf{a}^T (\mathbf{Q}(z) - \bar{\mathbf{Q}}(z)) \mathbf{b} \rightarrow 0$ almost surely as $n, p \rightarrow \infty$, for $\mathbf{A} \in \mathbb{R}^{p \times p}$ and $\mathbf{a}, \mathbf{b} \in \mathbb{R}^p$ of bounded
 250 (Euclidean and spectral) norms, the limiting spectral measure (via the associated Stieltjes transform)
 251 and the isolated eigenpairs of \mathbf{H} are directly accessible via the study of the deterministic equivalent
 252 $\bar{\mathbf{Q}}(z)$. This result is given as follows, with the proof deferred to Sec A.1 in the appendix.

253 **Theorem 4** (Deterministic equivalent). *Let $\mathbf{Q}(z) \equiv (\mathbf{H} - z\mathbf{I}_p)^{-1}$ be the resolvent of \mathbf{H} defined in*
 254 *(4). Then, under Assumption 1, as $n, p \rightarrow \infty$ with $p/n \rightarrow c \in (0, \infty)$,*

$$\mathbf{Q}(z) \leftrightarrow \bar{\mathbf{Q}}(z), \quad \text{with } \bar{\mathbf{Q}}^{-1}(z) = \mathbb{E} \left[\frac{g}{1+g\delta(z)} (\mathbf{C}^{\frac{1}{2}} (\mathbf{I}_p - \mathbf{P}_U) \mathbf{C}^{\frac{1}{2}} + \boldsymbol{\alpha} \boldsymbol{\alpha}^T) \right] - z\mathbf{I}_p,$$

255 *for random vector $\boldsymbol{\alpha} \equiv \boldsymbol{\mu} + \mathbf{C}^{\frac{1}{2}} \mathbf{P}_U \mathbf{z} \in \mathbb{R}^p$ and $g = \ell''(y, \mathbf{w}^T \boldsymbol{\mu} + \mathbf{w}^T \mathbf{C}^{\frac{1}{2}} \mathbf{z})$ for $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ as*
 256 *defined in (6), y and $\delta(z)$ defined in (1) and (5), respectively, and $\mathbf{P}_U \in \mathbb{R}^{p \times p}$ the projection onto*
 257 *the subspace spanned by the columns of $\mathbf{U} \equiv \mathbf{C}^{\frac{1}{2}} [\mathbf{w}_*, \mathbf{w}]$.*

258 **3 Evaluations and Discussions**

259 In this section, we provide further discussions on the consequences of Theorem 1, 2 and 3, together
 260 with numerical evaluations. Implications of Theorem 1 on the Hessian eigenvalue distribution is
 261 discussed in Sec 3.1. In Sec 3.2, we discuss the consequences of Theorem 2 and 3 on the possible
 262 isolated eigenpairs, for which two fundamentally different phase transitions are characterized.

263 **3.1 Hessian eigenvalues distribution**

264 For a better interpretation of Theorem 1 on the Hessian eigenspectrum, we consider here the special
 265 case of $\mathbf{C} = \mathbf{I}_p$, and start with the simple setting where the random variable g in (6) is constant,
 266 say $g = 1$: this happens, e.g., when the square loss $\ell(y, h) = (y - h)^2/2$ is employed. In this
 267 case, the Hessian does *not* depend on \mathbf{w}, \mathbf{w}_* and the Stieltjes transform $m(z)$ is the solution to
 268 $zcm^2(z) - (1 - c - z)m(z) + 1 = 0$ and corresponds to the Marčenko-Pastur law.

269 As long as g is *not* constant, the limiting Hessian spectrum is, a priori, different from the Marčenko-
 270 Pastur law, even in the $\mathbf{C} = \mathbf{I}_p$ setting, since the associated Stieltjes transform $m(z)$ is different from
 271 the solution to the Marčenko-Pastur equation. However, we see in Fig 3a that, for the logistic model
 272 (2) with logistic loss, the Hessian spectrum is close, at least visually, to a (rescaled) Marčenko-Pastur
 273 law. This can be understood with Theorem 1 and is due to the fact that, the distribution of g is more
 274 “concentrated” (around some constant, see Fig 3b versus 3d for a comparison between different cases).
 275 This is in sharp contrast to Fig 3c where with the exponential loss, the law of g has a much larger
 276 spread and the Hessian is therefore away from a Marčenko-Pastur-shape.

277 This “empirical fit” has been observed in [51, Fig 5], where acceleration methods proposed for
 278 a Marčenko-Pastur distributed Hessian (in linear least squares) work reasonably well on logistic
 279 regression models. Our theory proposes a convincing theoretical explanation of this empirical
 280 observation on logistic regression, and possibly for others more involved ML models. Nonetheless, it
 281 must be pointed out that this “visual approximation” by Marčenko-Pastur law is *not robust*, in the
 282 sense that it “visually” holds only for, yet formally different from, the case of (i) logistic model with
 283 (ii) logistic loss and (iii) identity covariance $\mathbf{C} = \mathbf{I}_p$: any change in the response model (e.g., the
 284 phase retrieval model in Fig 1c), in the choice of loss function (e.g., the exponential loss in Fig 2b),
 285 or beyond the identity covariance setting (as in Fig 2c and 2d) would induce a Hessian spectrum
 286 that is very different from the Marčenko-Pastur law. In this vein, our Theorem 1 goes beyond such
 287 “loose” Marčenko-Pastur approximation and acts as a more accurate first example in the understanding

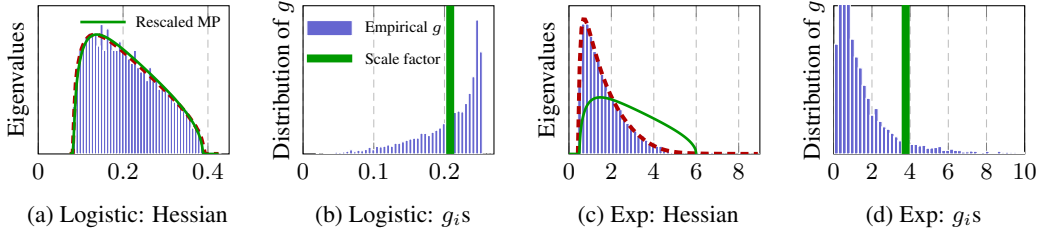


Figure 3: Comparison of Hessian eigenspectra with (rescaled and shifted) Marčenko-Pastur laws in the setting of Fig 2. **Fig 3a versus 3b**: Marčenko-Pastur-like Hessian with logistic loss, Hessian eigenvalues (3a) and empirical distribution of the g_i s versus the scaling factor (3b, empirically obtained by matching the minimal and maximal empirical Hessian eigenvalues to the Marčenko-Pastur law). **Fig 3c versus 3d**: an example of non-Marčenko-Pastur-like Hessian with exponential loss and the associated g_i s. Note that the scales of the axes are different in different subfigures.

288 of Hessian in more involved ML models beyond linear least squares that accounts for nonlinear
 289 transformations (such as activation function in NNs) and feature statistics.

290 While Theorem 1 is proven here only for Gaussian features, we conjecture, as is the case for many
 291 random matrix asymptotics, that it holds more generally beyond Gaussian distribution, see Fig 5 in
 292 the appendix for more discussions on this point.

293 3.2 Isolated eigenvalues and their phase transitions

294 In this section, we discuss the implications of Theorem 2 and 3 on the possible isolated eigenvalue-
 295 eigenvector pairs. More precisely, we show that, different from the classical spiked models extensively
 296 studied in RMT literature [2, 4, 39], for which (i) the isolated spike appears due to the presence
 297 of some statistical “signal” in the data and (ii) a “monotonic” phase transition behavior can be
 298 characterized as a function of the signal strength; here another type of Hessian spike arises due to the
 299 underlying G-GLM model (i.e., \mathbf{w}_* and \mathbf{w}) and exhibits a rather different behavior.

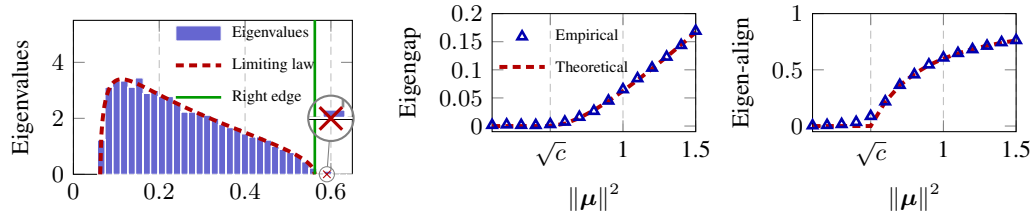


Figure 4: Spike due to data signal in Corollary 1: theory versus practice of (left) Hessian eigen-spectrum with $\|\mu\|^2 = 0.8$, (middle) eigengap $\text{dist}(\lambda_\mu, \text{supp}(\mu))$, and (right) top eigenvector alignment α in (10), as a function of the signal strength $\|\mu\|^2$, on logistic model with logistic loss, for $\mu \propto [-\mathbf{1}_{p/2}, \mathbf{1}_{p/2}]$, $\mathbf{w} = \mathbf{w}_* = \mathbf{0}$, $\mathbf{C} = \mathbf{I}_p$, $p = 512$ and $n = 2048$. Results averaged over 50 runs.

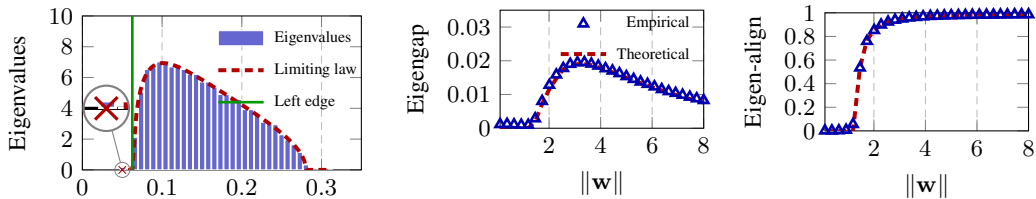


Figure 5: Left-hand side spike due to response model in Corollary 2 in the absence of data signal: (left) Hessian spectrum for $\|\mathbf{w}\| = 2$, with a left isolated eigenvalue $\hat{\lambda}_w$, (middle) eigengap $\text{dist}(\lambda_w, \text{supp}(\mu))$, and (right) dominant eigenvector alignment (with \mathbf{w}), as a function $\|\mathbf{w}\|$ with $\mathbf{w} \propto [-\mathbf{1}_{p/2}, \mathbf{1}_{p/2}]$, $\mathbf{w}_* = \mu = \mathbf{0}$, $\mathbf{C} = \mathbf{I}_p$, $p = 800$ and $n = 8000$. Results averaged over 50 runs.

300 3.2.1 Spike due to data signal

301 To study the spike due to data “signal” $\boldsymbol{\mu}$ and its phase transition behavior, we focus here on the case
 302 $\mathbf{w}_* = \mathbf{w} = \mathbf{0}$. This, in the case of logistic model (2) for example, gives rise to a much simpler form
 303 of limiting spectrum (per Theorem 1) and possible isolated eigenpairs (per Theorem 2 and 3), as
 304 summarized in the following corollary, with detailed derivations given in Sec B.3 of the appendix.

305 **Corollary 1** (Spike due to data signal: logistic model). *Consider the logistic model in (2) with*
 306 *logistic loss, for $\mathbf{w} = \mathbf{w}_* = \mathbf{0}$ and $\mathbf{C} = \mathbf{I}_p$, the limiting Hessian eigenvalue distribution is the*
 307 *Marčenko-Pastur law, but rescaled by a factor of $g = 1/4$. Moreover, there is at most one isolated*
 308 *eigenpair $(\hat{\lambda}_\mu, \hat{\mathbf{u}}_\mu)$ of \mathbf{H} and it satisfies*

$$\hat{\lambda}_\mu \xrightarrow{a.s.} \begin{cases} \lambda_\mu = \frac{1}{4}(1 + \rho + c \cdot \frac{\rho+1}{\rho}) & \rho > \sqrt{c}, \\ \frac{1}{4}(1 + \sqrt{c})^2 & \rho \leq \sqrt{c}; \end{cases}, \quad \frac{|\boldsymbol{\mu}^\top \hat{\mathbf{u}}_\mu|^2}{\|\boldsymbol{\mu}\|^2} \xrightarrow{a.s.} \begin{cases} \alpha = \frac{\rho^2 - c}{\rho^2 + c\rho} & \rho > \sqrt{c}, \\ 0 & \rho \leq \sqrt{c}; \end{cases} \quad (10)$$

309 with the signal strength $\rho = \lim_{p \rightarrow \infty} \|\boldsymbol{\mu}\|^2$ and $c = \lim p/n$.

310 The behavior of the isolated eigen-pairs described in Corollary 1 follows the “classical” phase
 311 transition [3, 2, 49]: (i) the isolated eigenvalue always appears on the right-hand side of the main
 312 (Marčenko-Pastur) bulk and (ii) the eigenvalue amplitude and eigenvector alignment is “monotonic”
 313 with respect to the signal strength $\|\boldsymbol{\mu}\|^2$ in the sense that, for a fixed dimension ratio c , the largest
 314 Hessian eigenvalue is bound to become asymptotically isolated once $\|\boldsymbol{\mu}\|^2$ exceeds \sqrt{c} and its value,
 315 as well as the eigenvector alignment, increase monotonically as $\|\boldsymbol{\mu}\|^2$ grows. This is confirmed in
 316 Fig 4. As we shall see below, this is *not* the case for, e.g., the spike due to model parameter \mathbf{w} .

317 3.2.2 Spike due to model

318 To investigate the spike due to the underlying model (i.e., \mathbf{w}_* and \mathbf{w}), we position ourselves in the
 319 situation where $\boldsymbol{\mu} = \mathbf{0}$, that is, in the absence of data “signal”. This leads to the following corollary,
 320 the proof of which is given in Sec B.4 in the appendix.

321 **Corollary 2** (Spike due to model: logistic model). *Consider the logistic model in (2) with logistic*
 322 *loss, $\boldsymbol{\mu} = \mathbf{0}$ and $\mathbf{C} = \mathbf{I}_p$, then the Stieltjes transform $m(z)$ satisfies $m(z) = 1/(\mathbb{E}[f(r, z)] - z)$*
 323 *for $f(r, z) = 1/(cm(z) + 2 + e^{-r} + e^r)$ and $r \sim \mathcal{N}(0, \|\mathbf{w}\|^2)$ that depends on \mathbf{w} but not on*
 324 *\mathbf{w}_* . Moreover, there is at most one isolated eigenvalue $\hat{\lambda}_\mathbf{w}$ of \mathbf{H} that is due to \mathbf{w} and satisfies*
 325 *$\hat{\lambda}_\mathbf{w} - \lambda_\mathbf{w} \xrightarrow{a.s.} 0$ with $\lambda_\mathbf{w}$ solution to $0 = \det \mathbf{G}(\lambda_\mathbf{w}) = 1 + m(\lambda_\mathbf{w}) \frac{\mathbb{E}[f(r, \lambda_\mathbf{w})(r^2 - \|\mathbf{w}\|^2)]}{\|\mathbf{w}\|^2}$.*

326 The situation here is more subtle (than the spike due to data signal discussed in Sec 3.2.1): as the
 327 model parameter \mathbf{w} changes (e.g., as the “energy” $\|\mathbf{w}\|$ grows), both the Hessian (limiting) eigenvalue
 328 distribution and the possible spike location are impacted. Fig 5 illustrates the behavior of the spike due
 329 to \mathbf{w} in the setting of Corollary 2. Note first that, different from the case of spike due to data signal
 330 $\boldsymbol{\mu}$, the spike in Fig 5-(left) appears on the *left-hand side* of the main bulk: this particularly means
 331 that the Hessian may admit an eigenvalue that is *significantly smaller* than all the other eigenvalues.⁵
 332 Also, note from Fig 5-(middle) that, different from the spike due to $\boldsymbol{\mu}$, the spike due to \mathbf{w} exhibits
 333 here a “non-monotonic” behavior in the sense that, it is absent for small values of $\|\mathbf{w}\|$ (as for small
 334 $\|\boldsymbol{\mu}\|$ in Fig 4-middle) and becomes “isolated” as $\|\mathbf{w}\|$ increases, but then again “merges into” the
 335 main bulk as $\|\mathbf{w}\|$ continues to increase, resulting an eigengap that falls back to zero.

336 It is perhaps even more surprising to observe in Fig 5-(right) that, the alignment between the associated
 337 isolated eigenvector and the parameter \mathbf{w} is, unlike the eigengap in Fig 5-(middle), monotonically
 338 increasing as $\|\mathbf{w}\|$ grows large, as in the case of Fig 4-(right). This suggests that, in the case of spike
 339 due to model, a *smaller* eigengap may not always imply *less* statistical “information” contained in
 340 the associated eigenvector, which somehow goes against the conventional eigengap heuristic [61, 33].
 341 It is worthy mentioning that, while, technically speaking, the proposed analysis is not capable of
 342 charactering the behavior as $\|\mathbf{w}\| \rightarrow \infty$ under Assumption 1, empirical results suggest that for
 343 extremely large $\|\mathbf{w}\|$, the eigengap tends to vanish, the associated “dominant” eigenvector can still be
 344 used to recover \mathbf{w} almost perfectly, see Fig 8 in the appendix as an example.

⁵Depending on the response model and loss, the spike due to model may also appear on the right-hand side of the bulk or even establish a left-to-right transition. The eigenvector alignment can behave differently from Fig 4-5 and establish a non-monotonic behavior as a function of $\|\mathbf{w}\|$ or $\|\mathbf{w}_*\|$; see Sec 4.3 in the appendix.

345 **References**

- 346 [1] Sanjeev Arora, Rong Ge, Tengyu Ma, and Ankur Moitra. Simple, efficient, and neural algorithms for
347 sparse coding. *Journal of Machine Learning Research*, 40(2015), 2015.
- 348 [2] Jinho Baik, Gérard Ben Arous, and Sandrine Péché. Phase transition of the largest eigenvalue for nonnull
349 complex sample covariance matrices. *The Annals of Probability*, 33(5):1643–1697, 2005.
- 350 [3] Jinho Baik and Jack W Silverstein. Eigenvalues of large sample covariance matrices of spiked population
351 models. *Journal of multivariate analysis*, 97(6):1382–1408, 2006.
- 352 [4] Florent Benaych-Georges and Raj Rao Nadakuditi. The eigenvalues and eigenvectors of finite, low rank
353 perturbations of large random matrices. *Advances in Mathematics*, 227(1):494–521, 2011.
- 354 [5] Zdzisław Burda, Jerzy Jurkiewicz, and Bartłomiej Waclaw. Spectral moments of correlated wishart
355 matrices. *Physical Review E*, 71(2):026111, 2005.
- 356 [6] Dankmar Böhning. Multinomial logistic regression algorithm. *Annals of the Institute of Statistical
357 Mathematics*, 44(1):197–200, 1992.
- 358 [7] Emmanuel J Candes, Thomas Strohmer, and Vladislav Voroninski. Phaselift: Exact and stable signal
359 recovery from magnitude measurements via convex programming. *Communications on Pure and Applied
360 Mathematics*, 66(8):1241–1274, 2013.
- 361 [8] Emmanuel J. Candès, Xiaodong Li, and Mahdi Soltanolkotabi. Phase Retrieval via Wirtinger Flow: Theory
362 and Algorithms. *IEEE Transactions on Information Theory*, 61(4):1985–2007, 2015.
- 363 [9] Olivier Chapelle, Choon Teo, Quoc Le, and Alex Smola. Tighter bounds for structured estimation.
364 *Advances in neural information processing systems*, 21:281–288, 2008.
- 365 [10] Pratik Chaudhari, Anna Choromanska, Stefano Soatto, Yann LeCun, Carlo Baldassi, Christian Borgs,
366 Jennifer Chayes, Levent Sagun, and Riccardo Zecchina. Entropy-SGD: Biasing gradient descent into wide
367 valleys. *Journal of Statistical Mechanics: Theory and Experiment*, 2019(12):124018, 2019.
- 368 [11] Yuxin Chen and Emmanuel Candes. Solving random quadratic systems of equations is nearly as easy as
369 solving linear systems. In *Advances in Neural Information Processing Systems*, pages 739–747, 2015.
- 370 [12] Yuejie Chi, Yue M. Lu, and Yuxin Chen. Nonconvex Optimization Meets Low-Rank Matrix Factorization:
371 An Overview. *IEEE Transactions on Signal Processing*, 67(20):5239–5269, 2018.
- 372 [13] Dennis Child. *The essentials of factor analysis*. Cassell Educational, 1990.
- 373 [14] Anna Choromanska, Mikael Henaff, Michael Mathieu, Gérard Ben Arous, and Yann LeCun. The loss
374 surfaces of multilayer networks. In *Artificial intelligence and statistics*, pages 192–204, 2015.
- 375 [15] Romain Couillet and Merouane Debbah. *Random matrix methods for wireless communications*. Cambridge
376 University Press, 2011.
- 377 [16] Romain Couillet and Walid Hachem. Analysis of the limiting spectral measure of large random matrices
378 of the separable covariance type. *Random Matrices: Theory and Applications*, 3(04):1450016, 2014.
- 379 [17] Yann N Dauphin, Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, Surya Ganguli, and Yoshua Bengio.
380 Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. In
381 *Advances in Neural Information Processing Systems*, volume 27 of *NIPS’14*, pages 2933–2941. Curran
382 Associates, Inc., 2014.
- 383 [18] Annette J Dobson and Adrian G Barnett. *An introduction to generalized linear models*. CRC press, 2018.
- 384 [19] Zhen Dong, Zhewei Yao, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. HAWQ: Hessian
385 aware quantization of neural networks with mixed-precision. In *Proceedings of the IEEE International
386 Conference on Computer Vision*, pages 293–302, 2019.
- 387 [20] James R. Fienup. Phase retrieval algorithms: a comparison. *Applied Optics*, 21(15):2758, 1982.
- 388 [21] Stanislav Fort and Surya Ganguli. Emergent properties of the local geometry of neural loss landscapes.
389 *arXiv preprint arXiv:1910.05929*, 2019.
- 390 [22] Yoav Freund, Robert Schapire, and Naoki Abe. A short introduction to boosting. *Journal-Japanese Society
391 For Artificial Intelligence*, 14(771-780):1612, 1999.

- 392 [23] Mario Geiger, Stefano Spigler, Stéphane d’Ascoli, Levent Sagun, Marco Baity-Jesi, Giulio Biroli, and
393 Matthieu Wyart. Jamming transition as a paradigm to understand the loss landscape of deep neural
394 networks. *Physical Review E*, 100(1):012115, 2019.
- 395 [24] Behrooz Ghorbani, Shankar Krishnan, and Ying Xiao. An investigation into neural net optimization via
396 hessian eigenvalue density. In *International Conference on Machine Learning*, pages 2232–2241, 2019.
- 397 [25] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016.
- 398 [26] Diego Granzio. Beyond random matrix theory for deep networks. *arXiv preprint arXiv:2006.07721*, 2020.
- 399 [27] Guy Gur-Ari, Daniel A. Roberts, and Ethan Dyer. Gradient descent happens in a tiny subspace, 2018.
- 400 [28] Walid Hachem, Philippe Loubaton, and Jamal Najim. Deterministic equivalents for certain functionals of
401 large random matrices. *The Annals of Applied Probability*, 17(3):875–930, 2007.
- 402 [29] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning: data mining,
403 inference, and prediction*. Springer Science & Business Media, 2009.
- 404 [30] Arthur Jacot, Franck Gabriel, and Clement Hongler. The asymptotic spectrum of the hessian of dnn
405 throughout training. In *International Conference on Learning Representations*, 2019.
- 406 [31] Prateek Jain and Purushottam Kar. Non-convex optimization for machine learning. *Foundations and
407 Trends® in Machine Learning*, 10(3-4):142–336, 2017.
- 408 [32] Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi. Low-rank matrix completion using alternating
409 minimization. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages
410 665–674, 2013.
- 411 [33] Antony Joseph and Bin Yu. Impact of regularization on spectral clustering. *The Annals of Statistics*,
412 44(4):1765–1791, 2016.
- 413 [34] Kenji Kawaguchi. Deep Learning without Poor Local Minima. In *Advances in Neural Information
414 Processing Systems*, volume 29 of *NIPS’16*, pages 586–594. Curran Associates, Inc., 2016.
- 415 [35] Raghunandan H Keshavan, Andrea Montanari, and Sewoong Oh. Matrix completion from a few entries.
416 *IEEE transactions on information theory*, 56(6):2980–2998, 2010.
- 417 [36] A Khorunzhy. Eigenvalue distribution of large random matrices with correlated entries. *Mat. Fiz. Anal.
418 Geom.*, 3(1-2):80–101, 1996.
- 419 [37] Jason D. Lee, Ioannis Panageas, Georgios Piliouras, Max Simchowitz, Michael I. Jordan, and Benjamin
420 Recht. First-order methods almost always avoid strict saddle points. *Mathematical Programming*, 176(1-
421 2):311–337, 2019.
- 422 [38] Kiryung Lee, Yanjun Li, Marius Junge, and Yoram Bresler. Blind recovery of sparse signals from
423 subsampled convolution. *IEEE Transactions on Information Theory*, 63(2):802–821, 2016.
- 424 [39] Thibault Lesieur, Florent Krzakala, and Lenka Zdeborová. Phase transitions in sparse pca. In *2015 IEEE
425 International Symposium on Information Theory (ISIT)*, pages 1635–1639. IEEE, 2015.
- 426 [40] Xinyan Li, Qilong Gu, Yingxue Zhou, Tiancong Chen, and Arindam Banerjee. Hessian based analysis
427 of sSGD for deep nets: Dynamics and generalization. In *Proceedings of the 2020 SIAM International
428 Conference on Data Mining*, pages 190–198. SIAM, 2020.
- 429 [41] Yue M Lu and Gen Li. Phase transitions of spectral initialization for high-dimensional non-convex
430 estimation. *Information and Inference: A Journal of the IMA*, 2019.
- 431 [42] Vladimir Alexandrovich Marchenko and Leonid Andreevich Pastur. Distribution of eigenvalues for some
432 sets of random matrices. *Matematicheskii Sbornik*, 114(4):507–536, 1967.
- 433 [43] C. H. Martin and M. W. Mahoney. Implicit self-regularization in deep neural networks: Evidence from
434 random matrix theory and implications for learning. Technical Report Preprint: arXiv:1810.01075, 2018.
- 435 [44] C. H. Martin, T. S. Peng, and M. W. Mahoney. Predicting trends in the quality of state-of-the-art neural
436 networks without access to training or testing data. Technical Report Preprint: arXiv:2002.06716, 2020.
- 437 [45] Llew Mason, Jonathan Baxter, Peter L Bartlett, and Marcus R Frean. Boosting algorithms as gradient
438 descent. In *Advances in neural information processing systems*, pages 512–518, 2000.

- 439 [46] Marco Mondelli and Andrea Montanari. Fundamental Limits of Weak Recovery with Applications to
440 Phase Retrieval. *Foundations of Computational Mathematics*, 19(3):703–773, 2019.
- 441 [47] Vardan Papyan. Measurements of three-level hierarchical structure in the outliers in the spectrum of
442 deepnet Hessians. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th*
443 *International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*,
444 pages 5012–5021. PMLR, 09–15 Jun 2019.
- 445 [48] Vardan Papyan. Traces of class/cross-class structure pervade deep learning spectra. *Journal of Machine*
446 *Learning Research*, 21(252):1–64, 2020.
- 447 [49] Debashis Paul. Asymptotics of sample eigenstructure for a large dimensional spiked covariance model.
448 *Statistica Sinica*, pages 1617–1642, 2007.
- 449 [50] Debashis Paul and Jack W Silverstein. No eigenvalues outside the support of the limiting empirical spectral
450 distribution of a separable covariance matrix. *Journal of Multivariate Analysis*, 100(1):37–57, 2009.
- 451 [51] Fabian Pedregosa and Damien Scieur. Average-case acceleration through spectral density estimation. *arXiv*
452 *preprint arXiv:2002.04756*, 2020.
- 453 [52] Jeffrey Pennington and Yasaman Bahri. Geometry of neural network loss surfaces via random matrix
454 theory. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages
455 2798–2806. JMLR. org, 2017.
- 456 [53] Jeffrey Pennington and Pratik Worah. The spectrum of the fisher information matrix of a single-hidden-layer
457 neural network. *Advances in Neural Information Processing Systems*, 31:5410–5419, 2018.
- 458 [54] Farbod Roosta-Khorasani and Michael W. Mahoney. Sub-sampled Newton methods. *Mathematical*
459 *Programming*, 174(1-2):293–326, 2019.
- 460 [55] Levent Sagun, Leon Bottou, and Yann LeCun. Eigenvalues of the hessian in deep learning: Singularity and
461 beyond. *arXiv preprint arXiv:1611.07476*, 2016.
- 462 [56] Levent Sagun, Utku Evci, V Ugur Guney, Yann Dauphin, and Leon Bottou. Empirical analysis of the
463 hessian of over-parametrized neural networks. *arXiv preprint arXiv:1706.04454*, 2017.
- 464 [57] Adepu Ravi Sankar, Yash Khasbage, Rahul Vigneswaran, and Vineeth N Balasubramanian. A deeper look
465 at the hessian eigenspectrum of deep neural networks and its applications to regularization. *arXiv preprint*
466 *arXiv:2012.03801*, 2020.
- 467 [58] Sheng Shen, Zhen Dong, Jiayu Ye, Linjian Ma, Zhewei Yao, Amir Gholami, Michael W Mahoney, and
468 Kurt Keutzer. Q-BERT: Hessian Based Ultra Low Precision Quantization of BERT. In *AAAI*, pages
469 8815–8821, 2020.
- 470 [59] Jack W Silverstein and ZD Bai. On the empirical distribution of eigenvalues of a class of large dimensional
471 random matrices. *Journal of Multivariate analysis*, 54(2):175–192, 1995.
- 472 [60] Aad W. Van der Vaart. *Asymptotic Statistics*, volume 3 of *Cambridge Series in Statistical and Probabilistic*
473 *Mathematics*. Cambridge University Press, 2000.
- 474 [61] Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007.
- 475 [62] Martin J. Wainwright. *High-Dimensional Statistics: : A Non-Asymptotic Viewpoint*. Cambridge Series in
476 *Statistical and Probabilistic Mathematics*. Cambridge University Press, 2019.
- 477 [63] Shusen Wang, Farbod Roosta-Khorasani, Peng Xu, and Michael W. Mahoney. GIANT: Globally Improved
478 Approximate Newton Method for Distributed Optimization. In *Advances in Neural Information Processing*
479 *Systems*, volume 31, pages 2332–2342. Curran Associates, Inc., 2018.
- 480 [64] Eugene P. Wigner. Characteristic vectors of bordered matrices with infinite dimensions. *Annals of*
481 *Mathematics*, 62(3):548–564, 1955.
- 482 [65] Yichao Wu and Yufeng Liu. Robust truncated hinge loss support vector machines. *Journal of the American*
483 *Statistical Association*, 102(479):974–983, 2007.
- 484 [66] Yikai Wu, Xingyu Zhu, Chenwei Wu, Annie Wang, and Rong Ge. Dissecting hessian: Understanding
485 common structure of hessian in neural networks. *arXiv preprint arXiv:2010.04261*, 2020.

- 486 [67] Peng Xu and Michael W. Mahoney. Second-order Optimization for Non-convex Machine Learning: an
 487 Empirical Study. In *Proceedings of the 2020 SIAM International Conference on Data Mining*, pages
 488 199–207, 2020.
- 489 [68] Peng Xu, Jiyan Yang, Farbod Roosta-Khorasani, Christopher Ré, and Michael W. Mahoney. Sub-sampled
 490 Newton Methods with Non-uniform Sampling. In *Advances in Neural Information Processing Systems*,
 491 volume 29, pages 3000–3008. Curran Associates, Inc., 2016.
- 492 [69] Fan Yang. Edge universality of separable covariance matrices. *Electronic Journal of Probability*, 24, 2019.
- 493 [70] Z. Yao, A. Gholami, Q. Lei, K. Keutzer, and M. W. Mahoney. Hessian-based analysis of large batch
 494 training and robustness to adversaries. Technical report, 2018. Preprint: arXiv:1802.08241.
- 495 [71] Z. Yao, A. Gholami, P. Xu, K. Keutzer, and M. W. Mahoney. Trust region based adversarial attack on
 496 neural networks. Technical report, 2018. Preprint: arXiv:1812.06371.
- 497 [72] Zhewei Yao, Amir Gholami, Kurt Keutzer, and Michael Mahoney. PyHessian: Neural networks through
 498 the lens of the Hessian. *arXiv preprint arXiv:1912.07145*, 2019.
- 499 [73] Zhewei Yao, Amir Gholami, Sheng Shen, Kurt Keutzer, and Michael W Mahoney. ADAHESIAN: An
 500 Adaptive Second Order Optimizer for Machine Learning. *arXiv preprint arXiv:2006.00719*, 2020.
- 501 [74] Lixin Zhang. Spectral analysis of large dimensional random matrices. *National University of Singapore*
 502 *PHD Thesis*, 2006.

503 Checklist

504 The checklist follows the references. Please read the checklist guidelines carefully for information on
 505 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]** , **[No]** , or
 506 **[N/A]** . You are strongly encouraged to include a **justification to your answer**, either by referencing
 507 the appropriate section of your paper or providing a brief inline description. For example:

- 508 • Did you include the license to the code and datasets? **[Yes]** See Section ??
- 509 • Did you include the license to the code and datasets? **[No]** The code and the data are
 510 proprietary.
- 511 • Did you include the license to the code and datasets? **[N/A]**

512 Please do not modify the questions and only use the provided macros for your answers. Note that the
 513 Checklist section does not count towards the page limit. In your paper, please delete this instructions
 514 block and only keep the Checklist section heading above along with the questions/answers below.

515 1. For all authors...

- 516 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 517 contributions and scope? **[Yes]**
- 518 (b) Did you describe the limitations of your work? **[Yes]** Theorem 1-3 is stated under the
 519 technical Assumption 1, Theorem 2 and 3 are stated under the additional assumption
 520 that the measure ν is of bounded support.
- 521 (c) Did you discuss any potential negative societal impacts of your work? **[No]** This work
 522 is theoretical.
- 523 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 524 them? **[Yes]**

525 2. If you are including theoretical results...

- 526 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
- 527 (b) Did you include complete proofs of all theoretical results? **[Yes]** In the appendix.

528 3. If you ran experiments...

- 529 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 530 mental results (either in the supplemental material or as a URL)? **[No]**
- 531 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 532 were chosen)? **[Yes]**

- 533 (c) Did you report error bars (e.g., with respect to the random seed after running exper-
534 iments multiple times)? [No] For the ease of visualization, we only report one-shot
535 results or the empirical means of multiple trials (versus our theoretical predictions).
- 536 (d) Did you include the total amount of compute and the type of resources used (e.g., type
537 of GPUs, internal cluster, or cloud provider)? [No]
- 538 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 539 (a) If your work uses existing assets, did you cite the creators? [N/A]
540 (b) Did you mention the license of the assets? [N/A]
541 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
542
- 543 (d) Did you discuss whether and how consent was obtained from people whose data you're
544 using/curating? [N/A]
545 (e) Did you discuss whether the data you are using/curating contains personally identifiable
546 information or offensive content? [N/A]
- 547 5. If you used crowdsourcing or conducted research with human subjects...
- 548 (a) Did you include the full text of instructions given to participants and screenshots, if
549 applicable? [N/A]
550 (b) Did you describe any potential participant risks, with links to Institutional Review
551 Board (IRB) approvals, if applicable? [N/A]
552 (c) Did you include the estimated hourly wage paid to participants and the total amount
553 spent on participant compensation? [N/A]