

On progressive sharpening, flat minima and generalisation

Lachlan E. MacDonald*

Australian Institute for Machine Learning
University of Adelaide
Adelaide, SA, Australia, 5000

lachlan.macdonald@adelaide.edu.au

Jack Valmadre

Australian Institute for Machine Learning
University of Adelaide
Adelaide, SA, Australia, 5000

Simon Lucey

Australian Institute for Machine Learning
University of Adelaide
Adelaide, SA, Australia, 5000

Abstract

We present a new approach to understanding the relationship between loss curvature and generalisation in deep learning. Specifically, we use existing empirical analyses of the spectrum of deep network loss Hessians to ground an ansatz tying together the loss Hessian and the input-output Jacobian of a deep neural network. We then prove a series of theoretical results which quantify the degree to which the input-output Jacobian of a model approximates its Lipschitz norm over a data distribution, and deduce a novel generalisation bound in terms of the empirical Jacobian. We use our ansatz, together with our theoretical results, to give a new account of the recently observed progressive sharpening phenomenon, as well as the generalisation properties of flat minima. Experimental evidence is provided to validate our claims.

1 Introduction

In this paper, we attempt to clarify how the curvature of the loss landscape of a deep neural network is related to the network’s generalisation properties. The hypothesis that network parameter vectors which occur in flat regions of the loss landscape correspond to more generalisable models has formed the groundwork of a great deal of research for some decades [26, 30, 11, 20], and loss flatness has proven empirically to be one of the most reliable correlates of generalisation [29]. Despite the substantial body of work utilising this hypothesis, however, the question of precisely *why* and *when* loss flatness implies good generalisation remains controversial and understudied [14, 23, 41]. It is this question that we study in this paper.

The mechanism we propose to mediate the relationship between loss curvature and generalisation is the input-output Jacobian of the model, which at least implicitly is known to play a part in determining generalisation [15, 5, 37, 27, 9, 22, 38]. Our proposal is based on empirical observations made of the eigenspectrum of the Hessian of deep neural networks [39, 40, 21]. It has been observed empirically across a variety of deep learning tasks that the loss Hessian of a deep neural network admits a number of “outlier” eigenvalues, which can be attributed to a summand of the Hessian known as the Gauss-Newton matrix. Crucially, the Gauss-Newton matrix is second-order only in the cost function, and solely first-order in the network layers.

*lemacdonald@protonmail.com

The Gauss-Newton matrix is a Gram matrix (i.e. a product $A^T A$ for some matrix A), and its conjugate AA^T , closely related to the tangent kernel identified in [28], has the same nonzero eigenvalues. Expanding AA^T reveals that it is determined in part by composite input-output layer Jacobians. Thus, insofar as the outlying Hessian eigenvalues are determined by those of the Gauss-Newton matrix, and insofar as layer Jacobians determine the spectrum of the Gauss-Newton matrix via its isospectrality to its conjugate, one expects the largest eigenvalue of the loss Hessian (commonly known as the *sharpness*) to be closely related to the magnitude of the model’s input-output Jacobian.

The explanatory power provided by our proposed mechanism extends beyond relating loss curvature to generalisation. Our proposed mechanism also enables us to give an account of the progressive sharpening phenomenon recently observed in the literature [12], which has up until this point remained largely mysterious.

Our contributions in this paper are as follows.

1. Based on previous empirical work identifying outlying loss Hessian eigenvalues with those of the Gauss-Newton matrix, we propose an ansatz: that the loss sharpness of a deep neural network acts as (loosely speaking) an upper bound for the model’s maximum Jacobian norm over training samples.
2. We prove theoretical results which quantify the extent to which the maximum input-output Jacobian norm of a model over a training set will approximate the Lipschitz norm of the model over the underlying data distribution. We extend this result to a novel generalisation bound in terms of the maximum input-output Jacobian norm over training data. Our results hold in particular for data generated by a neural network such as a generative adversarial network or an implicit neural function.
3. Based on our theory and our ansatz, we give a new account of progressive sharpening. Our account enables us to change the severity of sharpening by scaling inputs and outputs. We report on experiments across a variety of learning tasks that validate our account.
4. We report on experiments measuring the effect of hyperparameters such as learning rate and batch size on loss sharpness, Jacobian norm and generalisation, revealing the validity of our proposed mediating link between loss sharpness and generalisation. Finally, we present experiments measuring the loss sharpness, Jacobian norm and generalisation of networks trained with a variety of regularisation measures, revealing a tight correlation between generalisation and Jacobian norm whether or not such measures target loss sharpness.

2 Related work

Generalisation bounds: One of our contributions is a generalisation bound, which bounds the distance between training loss and population loss. The most common approaches to generalisation bounds are PAC-Bayes approaches [18, 20, 19] (which frequently invoke loss sharpness), which bound the gap in terms of the KL divergence between distributions over weight space; and algorithmic stability approaches [8, 24, 10, 31, 6], which bound contributions to the gap coming from each step of training. Our approach is distinct from both: unlike PAC-Bayes approaches, ours does not invoke probability distributions over parameters, and unlike algorithmic stability approaches, ours provides a bound in terms of the *end product* of training rather than the error introduced at each step. Ours is also distinguished from previous work in invoking special hypotheses on the data distribution (effectively, that it be generated by a neural network), which nonetheless hold for many practical datasets [35]. Our bound is most closely related to [5, 9], which recognise the Lipschitz constant as a determining factor in how well a model generalises, but in addition formalises the intuitive idea that Jacobian norm regulates generalisation [15, 27, 38]. Since in this paper we are focused on the connection between loss geometry and generalisation, we leave numerical evaluation of our bound, which would require generation of the data from a GAN, to future work.

Flatness and generalisation: The position that flatter minima generalise better dates back to [26]. It has since become a staple concept in the toolbox of deep learning practitioners [30, 11, 20, 3], with its incorporation into training schemes yielding state-of-the-art performance in many tasks. Motivated by this, many researchers have studied how training algorithms interact with loss geometry, finding that learning hyperparameters such as learning rate and batch size have an effect on the sharpness of minima to which the algorithm will converge [46, 50, 36]. Outside of PAC-Bayes theory,

the relationship between loss geometry and generalisation is studied in [41], where it is shown that, via a link between input- and parameter-perturbations for linear layers, a kind of loss flatness does indeed regulate generalisation. Our approach is distinct in focusing explicitly on model Jacobian norm, enabling us to provide new explanations of both the relationship between loss flatness and generalisation *and* progressive sharpening. It is not yet clear to us if or how the two approaches are related. The flat minima hypothesis is not without controversy. Indeed, it was theoretically shown in [14] that the loss curvature of a ReLU model can be scaled arbitrarily without changing the underlying model or its generalisation capacity. Moreover, in [23] it was noted that networks trained with cross-entropy generalise better with weight-decay than without, despite the former ending up in sharper minima. In [42] it is shown that flatness does not necessarily imply better generalisation for Gaussian-activated networks on regression tasks. The relationship we propose between loss curvature and generalisation is compatible with all of these observations.

Progressive sharpening and edge of stability: The effect of learning rate η on loss sharpness has been understood to some extent for several years [46]. In [12], this relationship was attended to for deep networks with a rigorous empirical study which showed that, after an initial period of sharpness increase during training called *progressive sharpening*, the sharpness increase halted at around $2/\eta$ and oscillated there while loss continued to decrease non-monotonically, a phase called *edge of stability*. These phenomena show that the typical assumptions made in theoretical work, namely that the learning rate is always smaller than twice the reciprocal of the sharpness [2, 17, 16, 34], do not hold in practice. Consequently, significant work has since been conducted to understand these phenomena [33, 32, 4, 45, 49, 1], with [13] showing in particular that edge of stability is a universal phenomenon arising from a third-order stabilising effect that must be taken into account when the sharpness grows above $2/\eta$. In contrast, progressive sharpening is not universal, and has been observed mainly in the context of neural networks. Although it has been correlated to growth in the norm of the output layer [45], the cause of progressive sharpening has so far remained mysterious. The mechanism we propose is the first account of which we are aware for a *cause* of progressive sharpening.

3 Background and paper outline

3.1 Outliers in the spectrum

We follow the formal framework introduced in [34], which allows us to treat all neural network layers on a common footing. Specifically, we consider a multilayer parameterised system $\{f_i : \mathbb{R}^{p_i} \times \mathbb{R}^{d_{i-1}} \rightarrow \mathbb{R}^{d_i}\}_{i=1}^L$ (the layers), with a data matrix $X \in \mathbb{R}^{d_0 \times N}$ consisting of N data vectors in \mathbb{R}^{d_0} . We denote by $F : \mathbb{R}^{p_1 + \dots + p_L} \rightarrow \mathbb{R}^{d_L \times N}$ the associated parameter-function map defined by

$$F_X(\theta_1, \dots, \theta_L)_i := f_L(\theta_L) \circ \dots \circ f_1(\theta_1)(X)_i \in \mathbb{R}^{d_L} \quad i = 1, \dots, N. \quad (1)$$

Given a cost function $c : \mathbb{R}^{d_L} \times \mathbb{R}^{d_L} \rightarrow \mathbb{R}$ and a target matrix $Y \in \mathbb{R}^{d_L \times N}$, we consider the associated loss

$$\ell(\vec{\theta}) := \gamma_Y \circ F_X(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N c(F_X(\vec{\theta})_i, Y_i), \quad (2)$$

where $\gamma_Y : \mathbb{R}^{d_L \times N} \rightarrow \mathbb{R}$ is defined by $\gamma_Y(Z) := N^{-1} \sum_{i=1}^N c(Z_i, Y_i)$.

Using the chain rule and the product rule, one observes that the Hessian $D^2\ell$ of ℓ admits the decomposition

$$D^2\ell = DF_X^T D^2\gamma_Y DF_X + D\gamma_Y D^2F_X. \quad (3)$$

The first of these terms, often called the *Gauss-Newton matrix*, is positive-semidefinite by convexity of γ_Y . It has been demonstrated empirically in a vast number of practical settings that the largest, outlying eigenvalues of the Hessian throughout training are due to the Gauss-Newton matrix [39, 40, 12]. In particular, the *sharpness*, or the largest eigenvalue of the Hessian, is, empirically, almost identical to that of the Gauss-Newton matrix. In attempting to understand the relationship between loss sharpness and generalisation, empirical evidence invites us to devote special attention to the Gauss-Newton matrix.

3.2 The Gauss-Newton matrix and input-output Jacobians

Letting C denote the square root $(D^2\gamma_Y)^{\frac{1}{2}}$, we see that the Gauss-Newton matrix $DF_X^T C^2 DF_X$ has the same nonzero eigenvalues as its conjugate matrix $C DF_X DF_X^T C^T$, which is closely related to the *tangent kernel* $DF_X DF_X^T$ identified in [28]. Letting Jf_l and Df_l denote the input-output and parameter derivatives of a layer $f_l : \mathbb{R}^p \times \mathbb{R}^{d_{l-1}} \rightarrow \mathbb{R}^{d_l}$ respectively, one sees that

$$C DF_X DF_X^T C^T = C \left(\sum_{l=1}^L \left(Jf_L \cdots Jf_{l+1} Df_l Df_l^T Jf_{l+1}^T \cdots Jf_L^T \right) \right) C^T \quad (4)$$

is a sum of positive-semidefinite matrices. Each summand is determined in large part by the composites of input-output layer Jacobians. Note, however, that the input-output Jacobian of the *first* layer does not appear; only its parameter derivative does.

It is clear then that insofar as the largest eigenvalues of the Hessian are determined by the Gauss-Newton matrix, these eigenvalues are determined in part by the input-output Jacobians of all layers following the first. For deep networks in particular, therefore, where there are many more layers than the first, one expects the outlying eigenvalues of the loss Hessian to be closely related to the input-output Jacobian of the network as a whole. Although we cannot yet make this relationship mathematically rigorous, we feel the idea is sufficiently clear and intuitive to act as an ansatz for what follows.

Ansatz 3.1. *Under certain conditions, an increase in the magnitude of the input-output Jacobian of a deep neural network will cause an increase in the loss sharpness. Conversely, a decrease in the sharpness will cause a decrease in the magnitude of the input-output Jacobian.*

We are deliberately vague about the “certain conditions” under which Ansatz 3.1 applies, in recognition of the fact that the relationship proposed in Ansatz 3.1 is *necessarily* mediated by both the cost function used and the parameter derivatives of the layers (see Equation (4)). Moreover, the Jacobian of the very first layer does *not* appear in Equation (4) and therefore is not directly related to the loss Hessian. In Appendix C, we will give explicit examples of cases where Ansatz 3.1 does *not* hold, due to the presence of these mediating factors. We emphasise that ***these mediating factors prevent us from identifying a simple and precise upper bound of the Jacobian norm in terms of the largest Hessian eigenvalue, as such an upper bound is contingent and does not hold in general.*** We leave a mathematically precise formulation of Ansatz 3.1 to future work; in the present paper, we will use the ansatz to provide new explanations for progressive sharpening and the relationship between loss geometry and generalisation.

In preparation for this, in Section 4, we will prove that the maximum input-output Jacobian of a neural net over a training sample approximates the Lipschitz norm of the network over the underlying data distribution exponentially well as a function of sample size (Theorem 4.5). We will then prove a generalisation bound relating this maximum input-output Jacobian to generalisation (Theorem 4.6). Together with the ansatz, these theorems are the groundwork for our explanations and experimental validation to follow.

Specifically, we will describe in Section 5 how Theorem 4.5 couples with Ansatz 3.1 to allow for a new account of progressive sharpening: in order to reduce the loss, a neural net needs to fit data to distinct targets, which requires the network to have a sufficiently large Lipschitz norm. As the Lipschitz norm grows during training to fit the data, so too does the Jacobian (by Theorem 4.5) and hence (by Ansatz 3.1) so too does loss sharpness.

In Section 6, we will argue that insofar as regularising loss sharpness causes good generalisation, it does so *through* Jacobian magnitude. Indeed, by Ansatz 3.1, controlling loss sharpness implicitly controls model Jacobian norm. By Theorem 4.6, this Jacobian control narrows the gap between train loss and test loss. Importantly, we show that many regularisation techniques, such as weight decay, mixup [48] and data augmentation, control Jacobian norm *without* controlling loss sharpness, leading to good test performance even in relatively sharp minima.

4 Theory

Our experiments in this paper are motivated by the following idea: the Lipschitz norm of a differentiable function is upper-bounded by the supremum, over the relevant data distribution, of the spectral

norm of its Jacobian. Intuitively, this supremum will itself be approximated by the maximum Jacobian norm over a finite sample of points from the distribution. Thus, insofar as the model’s Lipschitz norm determines its generalisation, we can therefore expect the sample-maximum Jacobian norm to determine its generalisation too. It is this intuition we seek to formalise in this section.

4.1 Definitions

We begin by defining what we mean by a *good* data distribution. In what follows, we will always use \mathfrak{X} to denote a subset of \mathbb{R}^d , equipped with its inherited Euclidean metric, and use \mathbb{P} to denote a probability measure on \mathfrak{X} . Given a Lipschitz function $g : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$, we use $\|g\|_{Lip}$ to denote the Lipschitz norm of g and $\|g\|_{Lip, \mathfrak{X}}$ to denote the Lipschitz norm of $g|_{\mathfrak{X}}$. We use $\|g\|_{\infty}$ to denote the global supremum of $\|g\|_2$, and $\|g\|_{\infty, \mathfrak{X}}$ to denote its supremum over \mathfrak{X} . Note the dual meaning of $\|g\|_2$: when g is *vector*-valued, it refers to the Euclidean norm, while when g is *operator*-valued (e.g., when $g = Jf$ is a Jacobian), $\|g\|_2$ refers to the spectral norm.

Definition 4.1. We say that \mathfrak{X} is *good* if for any differentiable Lipschitz function $g : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ one has

$$\|g\|_{Lip, \mathfrak{X}} \leq \|Df\|_{\infty, \mathfrak{X}} \quad (5)$$

We will also require our data distribution to satisfy the following inequalities. The first of them, as far as we can find, has not yet been introduced in the literature. It will be used to formalise the intuition that the maximum Jacobian norm of a function over a sample is a good approximation of the supremum Jacobian norm over the underlying distribution.

Definition 4.2. Let $h : \mathbb{R}_{\geq 0} \rightarrow [0, 1]$ be a non-decreasing function. We say that $(\mathfrak{X}, \mathbb{P})$ satisfies the *h-Lipschitz maximum inequality* if for any Lipschitz function $g : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ and $\epsilon > 0$ one has

$$\mathbb{P} \left[\|g(x)\|_2 \leq \|g\|_{\infty, \mathfrak{X}} - \epsilon \right] \leq 1 - h \left(\frac{\epsilon}{\|g\|_{Lip, \mathfrak{X}}} \right) \quad (6)$$

The next inequality is well-known in the machine learning literature [9, 7] as a hypothesis on tail behaviour under Lipschitz functions. We will use it in combination with our previous two definitions to prove our generalisation bound.

Definition 4.3. Let $C > 0$. We say that $(\mathfrak{X}, \mathbb{P})$ satisfies the *C-Lipschitz concentration inequality* if, for any Lipschitz function $g : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ and any $\epsilon > 0$, one has

$$\mathbb{P} \left[\|g(x) - \mathbb{E}g\|_2 \geq \epsilon \right] \leq 2 \exp \left(- \frac{C\epsilon^2}{\|g\|_{Lip, \mathfrak{X}}^2} \right) \quad (7)$$

4.2 Theorems

Our first theorem is that many data distributions of interest satisfy all of Definitions 4.1, 4.2 and 4.3. Thus our theorems below can be expected to hold in situations of practical interest. By an *immersion* we mean a differentiable function whose derivative (thought of as a linear map) is everywhere injective.

Theorem 4.4. *Let $(\mathfrak{X}, \mathbb{P})$ be either the unit hypercube with the uniform distribution, or its pushforward by a Lipschitz immersion. Then $(\mathfrak{X}, \mathbb{P})$ is good and satisfies both the Lipschitz maximum and Lipschitz concentration inequalities.*

Note that any GAN generator whose weight matrices are non-degenerate and whose nonlinearities are Lipschitz diffeomorphisms (e.g. smooth leaky ReLU) is a Lipschitz immersion. Thus all of Definitions 4.1, 4.2 and 4.3 can be expected to hold for data generated by a sufficiently nice GAN whose latent space is the uniform distribution on the unit hypercube. We believe it likely that Theorem 4.4 can be strengthened to include the Gaussian distribution and pushforwards by Lipschitz functions more generally, however we leave this to future work.

Our next theorem formalises our intuition that the maximum Jacobian norm over a sample from some distribution is a good approximation of the supremum of the Jacobian norm over the distribution, and that this approximation gets better as the sample gets larger. Note that, given $(\mathfrak{X}, \mathbb{P})$, we use \mathbb{P}^N to denote the product distribution over \mathfrak{X}^N , so that $(\mathfrak{X}^N, \mathbb{P}^N)$ is the probability space of i.i.d. samples from \mathbb{P} of size N .

Theorem 4.5. Suppose that $(\mathfrak{X}, \mathbb{P})$ satisfies the h -Lipschitz maximum inequality. Then for any twice continuously differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ with Lipschitz first and second derivative, one has:

$$\mathbb{P}^N \left[\max_i \|Jf(x_i)\|_2 \leq \|Jf\|_{\infty, \mathfrak{X}} - \epsilon \right] \leq \left(1 - h\left(\frac{\epsilon}{\|Jf\|_{\text{Lip}, \mathfrak{X}}}\right) \right)^N. \quad (8)$$

Our final theorem is our generalisation bound. It essentially follows from the Lipschitz concentration inequality together with the well-known Hoeffding inequality. Our innovation is to replace the Lipschitz norm appearing in the bound with the maximum Jacobian norm over a sample by using Theorem 4.5 and goodness of the data distribution. In this way, we can be assured of a rigorous relationship between the Jacobian norms we compute in experiments and generalisation.

Theorem 4.6. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ be a twice continuously differentiable Lipschitz function, with Lipschitz first and second derivative, and $c : \mathbb{R}^{d'} \times \mathbb{R}^{d'} \rightarrow \mathbb{R}$ a Lipschitz cost function. Denote by $\ell : \mathbb{R}^d \times \mathbb{R}^{d'} \rightarrow \mathbb{R}$ the composite $\ell(x, y) := c(f(x), y)$.

Let $\mathfrak{Z} \subset \mathbb{R}^d \times \mathbb{R}^{d'}$, and denote by \mathfrak{X} and \mathfrak{Y} its projections onto the first and second factors respectively. Let \mathbb{P} be a probability distribution on \mathfrak{Z} such that $(\mathfrak{Z}, \mathbb{P})$ is good and satisfies the C -Lipschitz concentration inequality. Assume moreover that the projection of $(\mathfrak{Z}, \mathbb{P})$ onto \mathfrak{X} satisfies the h -Lipschitz maximum inequality.

Then, there is a constant C' such that, given an i.i.d. tuple $((x_1, y_1), \dots, (x_N, y_N))$ drawn at random from \mathbb{P}^N , for any $\epsilon, \delta > 0$ one has

$$\left| \frac{1}{N} \sum_{i=1}^N \ell(x_i, y_i) - \mathbb{E} \ell \right| \leq \epsilon \quad (9)$$

with probability at least

$$1 - \left(1 - h\left(\frac{\delta}{\|Jf\|_{\text{Lip}, \mathfrak{X}}}\right) \right)^N - 2 \exp\left(\frac{-NC'\epsilon^2}{(\max_i \|Jf(x_i)\|_2 + \delta)^2}\right). \quad (10)$$

5 Progressive sharpening

By our Theorem 4.5, for many interesting data distributions, the Lipschitz norm of a function over the distribution is approximately upper-bounded by the maximum of its Jacobian norms over a sufficiently large sample. By Ansatz 3.1, this implies that an increase in the Lipschitz norm of a network during training will result in a corresponding increase of the sharpness of the associated loss function. We claim that this is a significant cause of progressive sharpening (when it occurs).

To formalise this claim, consider the simple case of fitting two points x_1 and x_2 to distinct values y_1 and y_2 with a Lipschitz function f . We measure the distance between outputs and targets using a cost function c , whose global minimum we assume to be zero, and for which we assume that there is $\gamma > 0$ such that $c(z_1, z_2) \geq \gamma \|z_1 - z_2\|_2^2$. This is clearly always the case for the square cost; by Pinsker's inequality, it also holds for the Kullback-Leibler divergence on softmax outputs, and hence also for the cross-entropy cost provided one subtracts the label entropy.

Suppose now that $c(f(x_i), y_i) \leq \epsilon^2 \gamma$ for each i . Thus $\|f(x_i) - y_i\|_2 \leq \epsilon$ for each i . Since $f(x_1)$ and $f(x_2)$ are at their closest when on the straight line segment between y_1 and y_2 , it follows that $\|f(x_1) - f(x_2)\|_2 \geq \|y_1 - y_2\|_2 - 2\epsilon$. Invoking the Lipschitz property of f , one has:

$$\|f\|_{\text{Lip}} \geq \frac{\|y_1 - y_2\|_2 - 2\epsilon}{\|x_1 - x_2\|_2}. \quad (11)$$

Clearly, this lower bound gets *higher* the more one lowers ϵ . Thus, provided the targets are sufficiently far apart and the Lipschitz norm of the model starts from a low basepoint, one will expect this Lipschitz norm to increase during any effective training procedure.

Since the Lipschitz norm is upper bounded by the supremum Jacobian norm, Theorem 4.5 tells us that training will thus also increase the sample-maximum Jacobian norm, which, by Ansatz 3.1, can be expected to cause a corresponding increase in the magnitude of the loss Hessian (see Appendix A.1): progressive sharpening. Moreover, Equation (11) tells us that this effect can be made more or

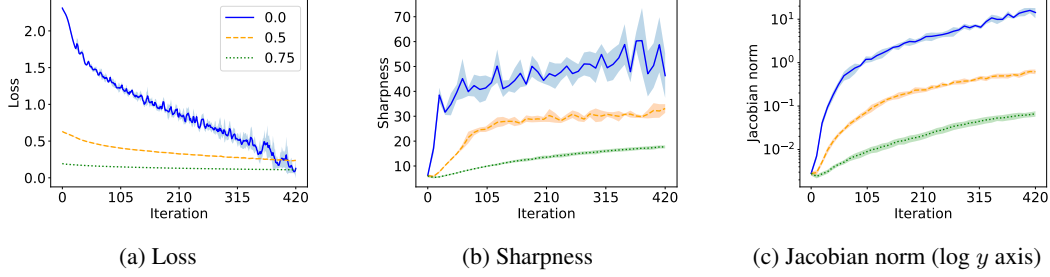


Figure 1: Plots of cross entropy loss (with label entropies subtracted), loss sharpness and softmaxed-model Jacobian norm with varying degrees of label smoothing, for VGG11 trained with gradient descent on CIFAR10, with a learning rate of 0.08. Increasing label smoothing (indicated by line style) brings targets closer together, meaning less growth is necessary in the Jacobian norm to reduce loss. This coincides with less sharpening of the Hessian during training, in line with Ansatz 3.1. Similar plots at different learning rates, and with ResNet18, are in Appendix B.1. Note log y axis on Jacobian norm for ease of distinction.

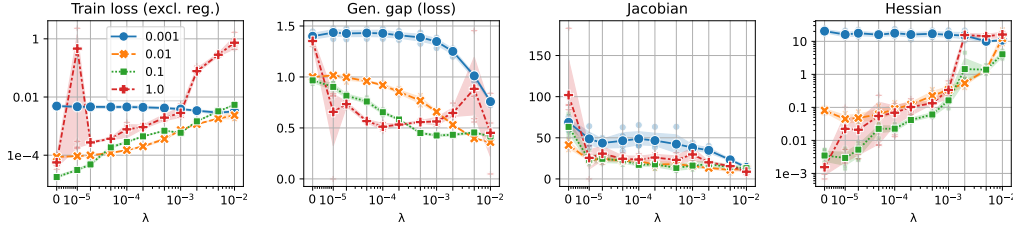


Figure 2: The impact of weight decay (x axis) on Jacobian norm, sharpness and generalisation gap when training with cross-entropy loss at the end of training (90 epochs). Line style indicates learning rate. Without weight decay, SGD finds a solution with near-zero loss and vanishing sharpness, but for which Jacobian norm and generalisation gap are relatively large.

less severe by scaling the distances between target values: indeed this is what we observe (Figure 1, see also Appendix B.1).

Note that from Equation (11), one might expect that one could also reduce progressive sharpening by scaling the input data larger (at least for non-BN networks). Indeed, while it is true that such scaling reduces the *Jacobian norm*, this reduction in Jacobian norm does *not* necessarily reduce sharpening, since it coincides with an *increase* in the magnitude of the parameter derivatives (see Appendix C.2).

Note finally that although the edge of stability mechanism explicated in [13] can be expected to put some downward pressure on the model Jacobian, due to the presence of the mediating factors discussed following Ansatz 3.1 it *need not* cause the Jacobian to *plateau* in the same way that loss sharpness does. Nonetheless, this downward pressure on Jacobian is important: it is to this that we attribute the better generalisation of models trained with larger learning rate, which we validate empirically in the next section.

6 Flat minima and generalisation

In this section we report on experiments measuring the generalisation gap (absolute value of train loss minus test loss), sharpness and Jacobian norm at the end of training with varying degrees of various regularisation techniques. Our results validate our claim that, insofar as Hessian control does yield good generalisation, it does so *because* Hessian control implicitly regularises the Jacobian norm. Moreover, Jacobian norm continues to correlate well with generalisation for all regularisation measures we studied, even where sharpness does not. Note that the Jacobians we plot were all computed in evaluation mode, meaning all BN layers had their statistics fixed: that the relationship with the train mode loss Hessian can still be expected to hold is justified in Appendix A.1.

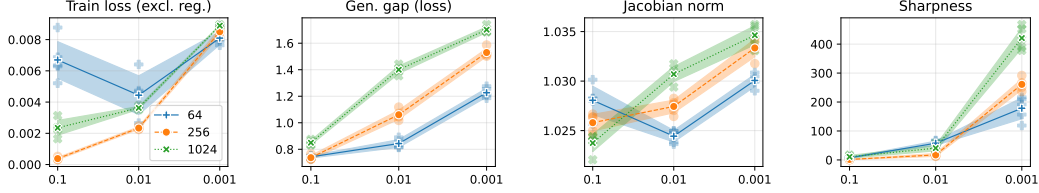


Figure 3: The impact of differing (constant) learning rates (x axis) at end of training with ResNet18 on CIFAR10 trained with SGD (five trials). Line style indicates batch size. In line with conventional wisdom, larger learning rates typically have smaller sharpness and generalise better. The downward pressure on sharpness and generalisation gap correlates with smaller Jacobian norms, as anticipated by Ansatz 3.1 and Theorem 4.6. Models trained with weight decay (0.0001).

6.1 Cross entropy and weight decay

It was observed in [23] that when training with the cross-entropy cost, networks trained with weight decay generalised better than those trained without, despite converging to sharper minima. We confirm this observation in Figure 2. This phenomenon is easily understood through our theory. For the cross-entropy cost, both terms $D^2\gamma_Y$ and $D\gamma_Y$ appearing in the Hessian (Equation (3)) vanish at infinity in parameter space. By preventing convergence of the parameter vector to infinity in parameter space, weight-decay encourages convergence to a sharper minimum than would otherwise be the case. However, since the Frobenius norm of a matrix is an upper bound on the spectral norm, weight decay also implicitly regularises the Jacobian of the model, leading to better generalisation by Theorem 4.6.

Since training with weight decay is the norm in practice, and since any correlations between *vanishing* Hessian and generalisation will likely be unreliable, we used weight decay in all of what follows.

6.2 Learning rate and batch size

It is commonly understood that larger learning rates and smaller batch sizes serve to encourage convergence to flatter minima [30, 46]. By Theorem 4.6 such hyperparameter choice can be expected to lead to better generalisation. We verify this with ResNet18 on CIFAR10. Mostly we see that these expectations are borne out by experiment (Figures 3 and 4), with a notable exception being at the largest learning rate, where the relationship between generalisation gap and Jacobian is the reverse of what is anticipated (Figure 4). The phenomenon appears to be data-agnostic (the same occurs with ResNet18 on CIFAR100, see Appendix B.2), but architecture-dependent (we do not observe this problem with VGG11 on either CIFAR10 or CIFAR100). We speculate that the reason for this is that skip connected architectures trained with large, constant learning rates will tend to produce models whose Jacobian has a larger Lipschitz constant, meaning that the maximum Jacobian norm will be a poorer estimate of the Lipschitz constant of the model (Theorem 4.5) and hence of the generalisation gap (Theorem 4.6). We emphasise, however, that the constant learning rate training used for these experiments is not used in practice, and the decoupling of Jacobian norm from generalisation gap that we see in Figure 4 does not persist with the more realistic training schemes considered in Figure 5.

6.3 Other regularisation techniques

We also test to see the extent that other common regularisation techniques, including label smoothing, data augmentation, mixup [48] and sharpness-aware minimisation (SAM) [20] regularise Hessian and Jacobian to result in better generalisation. We find that while Jacobian norm is regularised across all techniques, only in SAM is sharpness regularised. This validates our proposal that Jacobian norm is a key mediating factor between sharpness and generalisation, as well as the position of [14] that flatness is not necessary for good generalisation. See Appendix B.3 for experimental details and similar results for VGG11 and CIFAR100.

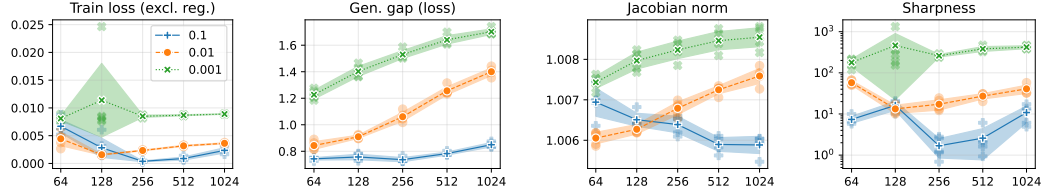


Figure 4: The impact of differing batch sizes (x axis) at the end of training with ResNet18 on CIFAR10 trained with SGD (five trials). Line style indicates (constant) learning rate. At least for the smaller learning rates, increasing batch size typically increases sharpness, Jacobian norm and generalisation gap in line with Theorem 4.6. At the largest learning rate, the relationship between Jacobian norm and generalisation gap is the inverse of what is expected, suggesting that larger learning rates tend to increase the Lipschitz constant of the Jacobian, causing it to underestimate the Lipschitz constant of the model (Theorem 4.6). Note log scale on y axis in Sharpness plot for ease of distinction.

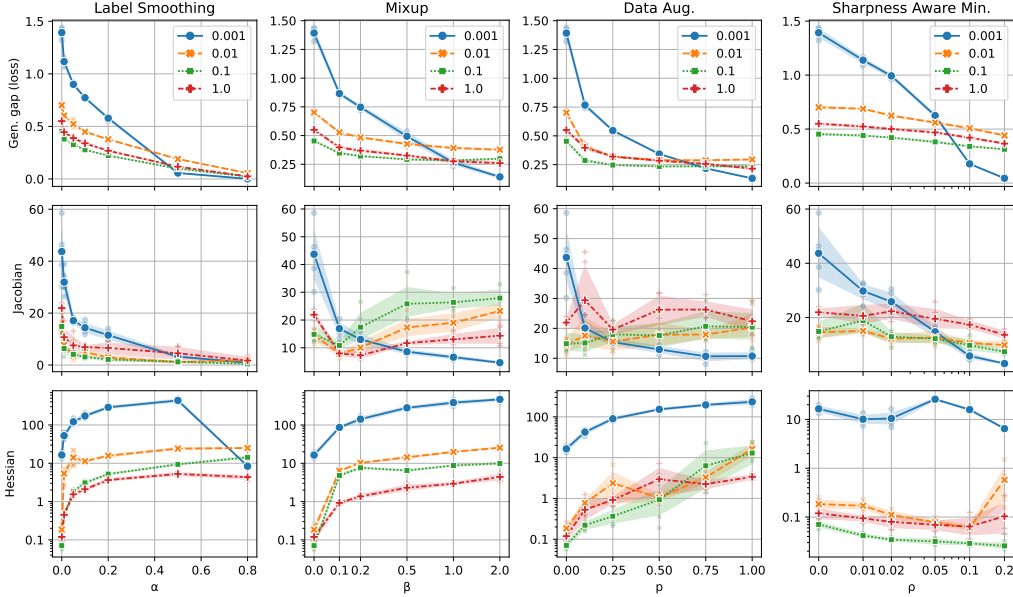


Figure 5: The impact of different regularisation strategies (x axis) when training a ResNet18 model with batch-norm on CIFAR10 (five trials). The norm of the input-output Jacobian is much better correlated with the generalisation gap than that of the Hessian; all regularisation strategies were effective in controlling the Jacobian norm without necessarily controlling the sharpness. Line style indicates initial learning rate. Besides the regularisation strategy being studied, all experiments include weight decay (0.0005). Refer to the appendix for extended results.

7 Discussion

One limitation of our work is that Ansatz 3.1 is not mathematically rigorous. The ansatz calls for a close theoretical analysis, which would ideally provide rigorous conditions under which one can upper bound the model Jacobian norm in terms of loss sharpness. Since the Hessian involves all layer Jacobians following the first, one expects this bound to be *better* for deeper networks, in which a higher percentage of the layer Jacobians are explicitly regularised. Such a bound would help explain why *deeper* networks are frequently observed to generalise better [25].

Another limitation of our work is that we do not numerically evaluate our generalisation bound. Since the bound depends strongly on both the data (via the Lipschitz constant of a network which generates the data) and on the whole training procedure (via the maximum Jacobian norm over the training set *at convergence*), it circumvents the data- and/or algorithm-agnosticism that makes many

existing bounds so ineffective in deep learning [47]. We believe quantitative evaluation of the bound in practical settings is an important future research problem.

8 Conclusion

We proposed a new relationship between the loss Hessian of a deep neural network and the input-output Jacobian of the model. We proved several theorems allowing us to leverage this relationship in providing new explanations for progressive sharpening and the link between loss geometry and generalisation. An experimental study was conducted which validates our proposal.

References

- [1] K. Ahn, J. Zhang, and S. Sra. Understanding the unstable convergence of gradient descent. In *ICML*, 2022.
- [2] Z. Allen-Zhu, Y. Li, and Z. Song. A Convergence Theory for Deep Learning via Over-Parameterization. In *ICML*, pages 242–252, 2019.
- [3] X. C. and C.-J. Hsieh and B. Gong. When Vision Transformers Outperform ResNets without Pre-training or Strong Data Augmentations. In *ICLR*, 2020.
- [4] S. Arora, Z. Li, and A. Panigrahi. Understanding Gradient Descent on Edge of Stability in Deep Learning. In *ICML*, 2022.
- [5] P. Bartlett, D. J. Foster, and M. Telgarsky. Spectrally-normalized margin bounds for neural networks. In *NeurIPS*, 2017.
- [6] R. Bassily, V. Feldman, C. Guzmán, and K. Talwar. Stability of Stochastic Gradient Descent on Nonsmooth Convex Losses. In *NeurIPS*, 2020.
- [7] S. Bombari, M. H. Amani, and M. Mondelli. Memorization and Optimization in Deep Neural Networks with Minimum Over-parameterization. In *NeurIPS*, 2022.
- [8] O. Bousquet and A. Elisseeff. Stability and Generalization. *Journal of Machine Learning Research*, 2:499–526, 2002.
- [9] S. Bubeck and M. Sellke. A Universal Law of Robustness via Isoperimetry . In *NeurIPS*, 2021.
- [10] Z. Charles and D. Papailiopoulos. Stability and Generalization of Learning Algorithms that Converge to Global Optima. In *ICML*, 2018.
- [11] P. Chaudhari, A. Choromanska, S. Soatto, Y. LeCun, C. Baldassi, C. Borgs, J. Chayes, L. Sagun, and R. Zecchina. Entropy-SGD: Biasing Gradient Descent into Wide Valleys. In *ICLR*, 2017.
- [12] J. Cohen, S. Kaur, Y. Li, J. Z. Kolter, and A. Talwalkar. Gradient Descent on Neural Networks Typically Occurs at the Edge of Stability. In *ICLR*, 2021.
- [13] A. Damian, E. Nichani, and J. Lee. Self-Stabilization: The Implicit Bias of Gradient Descent at the Edge of Stability. In *ICLR*, 2023.
- [14] L. Dinh, R. Pascanu, S. Bengio, and Y. Bengio. Sharp minima can generalize for deep nets. In *ICML*, pages 1019–1028, 2017.
- [15] H. Drucker and Y. Le Cun. Improving generalization performance using double backpropagation. *IEEE Transactions on Neural Networks*, 3(6):991–997, 1992.
- [16] S. S. Du, J. Lee, H. Li, L. Wang, and X. Zhai. Gradient Descent Finds Global Minima of Deep Neural Networks. In *ICML*, pages 1675–1685, 2019.
- [17] S. S. Du, X. Zhai, B. Póczos, and A. Singh. Gradient Descent Provably Optimizes Over-parameterized Neural Networks. In *ICLR*, 2019.
- [18] G. K. Dziugaite and D. M. Roy. Computing Nonvacuous Generalization Bounds for Deep (Stochastic) Neural Networks with Many More Parameters than Training Data. In *UAI*, 2017.
- [19] F. He and T. Liu and D. Tao. Control Batch Size and Learning Rate to Generalize Well: Theoretical and Empirical Evidence . In *NeurIPS*, 2019.
- [20] P. Foret, A. Kleiner, H. Mobahi, and B. Neyshabur. Sharpness-Aware Minimization for Efficiently Improving Generalization. In *ICLR*, 2021.

- [21] B. Ghorbani, S. Krishnan, and Y. Xiao. An Investigation into Neural Net Optimization via Hessian Eigenvalue Density. In *ICML*, 2019.
- [22] H. Gouk, E. Frank, B. Pfahringer, and M. J. Cree. Regularisation of neural networks by enforcing Lipschitz continuity. *Machine Learning*, 110:393–416, 2021.
- [23] D. Granzio. Flatness is a False Friend. arXiv:2006.09091, 2020.
- [24] M. Hardt, B. Recht, and Y. Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *ICML*, 2016.
- [25] K. He, X. Zhang, S. Ren, and J. Sun. Identity Mappings in Deep Residual Networks. In *ECCV*, 2016.
- [26] S. Hochreiter and J. Schmidhuber. Flat minima. *Neural Computation*, 9(1):1–42, 1997.
- [27] J. Hoffman, D. A. Roberts, and S. Yaida. Robust Learning with Jacobian Regularization. arXiv:1908.02729, 2019.
- [28] A. Jacot, F. Gabriel, and C. Hongler. Neural Tangent Kernel: Convergence and Generalization in Neural Networks. In *NeurIPS*, pages 8571–8580, 2018.
- [29] Y. Jiang, B. Neyshabur, H. Mobahi, D. Krishnan, and S. Bengio. Fantastic Generalization Measures and Where to Find Them. In *ICLR*, 2020.
- [30] N. S. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, P. Tak, and P. Tang. On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima. In *ICLR*, 2017.
- [31] I. Kuzborskij and C. H. Lampert. Data-Dependent Stability of Stochastic Gradient Descent. In *ICML*, 2018.
- [32] S. Lee and C. Jang. A new characterization of the edge of stability based on a sharpness measure aware of batch gradient distribution. In *ICLR*, 2023.
- [33] K. Lyu, Z. Li, and S. Arora. Understanding the Generalization Benefit of Normalization Layers: Sharpness Reduction. In *NeurIPS*, 2022.
- [34] L. E. MacDonald, H. Saratchandran, J. Valmadre, and S. Lucey. A global analysis of global optimisation. arXiv:2210.05371, 2022.
- [35] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. Spectral Normalization for Generative Adversarial Networks. In *ICML*, 2018.
- [36] R. Mulayoff, T. Michaeli, and D. Soudry. The Implicit Bias of Minima Stability: A View from Function Space. In *NeurIPS*, 2021.
- [37] B. Neyshabur, S. Bhojanapalli, D. Mcallester, and N. Srebro. Exploring Generalization in Deep Learning. In *NeurIPS*, 2017.
- [38] R. Novak, Y. Bahri, D. A. Abolafia, J. Pennington, and J. Sohl-Dickstein. Sensitivity and generalization in neural networks: an empirical study. In *International Conference on Learning Representations*, 2018.
- [39] V. Pappas. The Full Spectrum of Deepnet Hessians at Scale: Dynamics with SGD Training and Sample Size. arXiv:1811.07062, 2018.
- [40] V. Pappas. Measurements of three-level hierarchical structure in the outliers in the spectrum of deepnet hessians. In *ICML*, 2019.
- [41] H. Petzka, M. Kamp, L. Adilova, C. Sminchisescu, and M. Boley. Relative flatness and generalization. *Advances in Neural Information Processing Systems*, 34:18420–18432, 2021.
- [42] S. Ramasinghe, L. E. MacDonald, M. Farazi, H. Saratchandran, and S. Lucey. How much does Initialization Affect Generalization? In *ICML*, 2023.
- [43] M. E. A. Seddik, C. Louart, M. Tamaazousti, and R. Couillet. Random Matrix Theory Proves that Deep Learning Representations of GAN-data Behave as Gaussian Mixtures. In *ICML*, 2020.
- [44] R. Vershynin. *High Dimensional Probability, An Introduction with Applications in Data Science*. Cambridge University Press, 2018.
- [45] Z. Wang, Z. Li, and J. Li. Analyzing Sharpness along GD Trajectory: Progressive Sharpening and Edge of Stability. In *NeurIPS*, 2022.

- [46] L. Wu, C. Ma, and W. E. How SGD Selects the Global Minima in Over-parameterized Learning: A Dynamical Stability Perspective. In *NeurIPS*, 2018.
- [47] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals. Understanding deep learning requires rethinking generalization. In *ICLR*, 2017.
- [48] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz. mixup: Beyond Empirical Risk Minimization . In *ICLR*, 2018.
- [49] X. Zhu, Z. Wang, X. Wang, M. Zhou, and R. Ge. Understanding Edge-of-Stability Training Dynamics with a Minimalist Example. In *ICLR*, 2023.
- [50] Z. Zhu, J. Wu, B. Yu, L. Wu, and J. Ma. The Anisotropic Noise in Stochastic Gradient Descent: Its Behavior of Escaping from Sharp Minima and Regularization Effects. In *ICML*, 2019.

Computational resources: Exploratory experiments were conducted using a desktop machine with two Nvidia RTX A6000 GPUs. The extended experimental results in the appendix were obtained on a shared HPC system, where a total of 3500 GPU hours were used across the lifetime of the project (Nvidia A100 GPUs).

A Proofs

Here we give proofs of the theorems stated in the main body of the paper.

Proof of Theorem 4.4. We first show that $(\mathfrak{X}, \mathbb{P})$ is good. This follows from the fact that in either case, \mathfrak{X} is a connected submanifold of \mathbb{R}^d . Then, given a differentiable Lipschitz function $g : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$, the estimate $\|g\|_{\text{Lip}, \mathfrak{X}} \leq \|Jg\|_{\infty, \mathfrak{X}}$ follows from the mean value theorem.

Indeed, fixing $x, y \in \mathfrak{X}$, let $\gamma : [0, 1] \rightarrow \mathfrak{X}$ be a differentiable path such that $\gamma(0) = x$ and $\gamma(1) = y$. Since \mathfrak{X} is a manifold, we can assume without loss of generality that $\dot{\gamma}(t) \neq 0$ for all t , and hence that $\|\dot{\gamma}(t)\|_2 = \|y - x\|_2$ for all t . Now, define $\varphi(t) := (f(y) - f(x)) \cdot f(\gamma(t))$. Applying the mean value theorem, there exists $t_0 \in [0, 1]$ such that

$$\varphi(1) - \varphi(0) = \varphi'(t_0) = (f(y) - f(x)) \cdot Jf(\gamma(t_0))\dot{\gamma}(t_0). \quad (12)$$

Noting that $\varphi(1) - \varphi(0) = \|f(y) - f(x)\|_2^2$ and applying Cauchy-Schwarz yields:

$$\|f(y) - f(x)\|_2 \leq \|Jf(\gamma(c))\dot{\gamma}(c)\|_2 = \|Jf(\gamma(c))\|_2 \|y - x\|_2 \leq \|Jf\|_{\infty, \mathfrak{X}} \|y - x\|_2 \quad (13)$$

from which the result follows.

We now turn to showing that $(\mathfrak{X}, \mathbb{P})$ satisfies the Lipschitz maximum inequality. We first prove the result when $(\mathfrak{X}, \mathbb{P})$ is the uniform distribution on the unit hypercube. Let g be a Lipschitz function on \mathfrak{X} . Since g is continuous, its supremum is attained at some point x_0 in \mathfrak{X} . By the reverse triangle inequality and the Lipschitz identity, one therefore has

$$\mathbb{P} \left[\|g(x)\|_2 \geq \|g\|_{\infty, \mathfrak{X}} - \epsilon \right] = \mathbb{P} \left[\|g(x)\|_2 \geq \|g(x_0)\|_2 - \epsilon \right] \geq \mathbb{P} \left[\|x - x_0\|_2 \leq \frac{\epsilon}{\|g\|_{\text{Lip}, \mathfrak{X}}} \right]. \quad (14)$$

The quantity on the right is now easily bounded. The volume of an R -ball in \mathbb{R}^n is given by $\Gamma(n/2 + 1)^{-1} \pi^{n/2} R^n =: CR^n$. Centering such a ball at x_0 , its intersection with the unit cube has volume at least $2^{-n} R^n$ (when x_0 is a corner) assuming without loss of generality that $R < 1$ (as we will typically want to take ϵ small). Thus, $(\mathfrak{X}, \mathbb{P})$ satisfies the Lipschitz maximum inequality with $h(t) := 2^{-n} C t^n$.

Let us now assume that $(\mathfrak{X}, \mathbb{P}) = (F(\mathfrak{X}'), F_*\mathbb{P}')$ is the pushforward of the uniform distribution $(\mathfrak{X}', \mathbb{P}')$ on the unit d' -cube by some Lipschitz function $F : \mathbb{R}^{d'} \rightarrow \mathbb{R}^d$. Let $g : \mathbb{R}^d \rightarrow \mathbb{R}^{d''}$ be any other Lipschitz function. Then $g \circ F$ is also Lipschitz, with Lipschitz norm bounded by $\|g\|_{\text{Lip}, \mathfrak{X}} \|F\|_{\text{Lip}, \mathfrak{X}'}$. Since $(\mathfrak{X}', \mathbb{P}')$ satisfies the Lipschitz maximum inequality, one has

$$\mathbb{P} \left[\|g(x)\|_2 \geq \|g\|_{\infty, \mathfrak{X}} - \epsilon \right] = \mathbb{P}' \left[\|(g \circ F)(y)\|_2 \geq \|g \circ F\|_{\infty, \mathfrak{X}'} - \epsilon \right] \quad (15)$$

$$\geq h \left(\frac{\epsilon}{\|g \circ F\|_{\text{Lip}, \mathfrak{X}'}} \right) \quad (16)$$

$$\geq h' \left(\frac{\epsilon}{\|g\|_{\text{Lip}, \mathfrak{X}}} \right), \quad (17)$$

where h is defined in the above paragraph and $h'(t) := h(t/\|F\|_{\text{Lip}, \mathfrak{X}'})$. This concludes the proof.

Finally, we come to the Lipschitz concentration inequality. This, however, is well-known for the distributions under consideration: see [44, Theorem 5.2.10] for the uniform case and [43, Proposition 2.3] for the pushforward case. \square

Proof of Theorem 4.5. Clearly, for any constant $K > 0$,

$$\mathbb{P}^N \left[\max_i \|g(x_i)\|_2 \leq K \right] = \mathbb{P}^N \left[\|g(x_1)\|_2 \leq K, \dots, \|g(x_N)\|_2 \leq K \right]. \quad (18)$$

By independence:

$$\mathbb{P}^N \left[\|g(x_1)\|_2 \leq K, \dots, \|g(x_N)\|_2 \leq K \right] = \prod_{i=1}^N \mathbb{P} \left[\|g(x_i)\|_2 \leq K \right]. \quad (19)$$

Now, assuming g is Lipschitz and that $(\mathfrak{X}, \mathbb{P})$ satisfies the Lipschitz maximum inequality yields, for any $\epsilon > 0$:

$$\mathbb{P}^N \left[\max_i \|g(x_i)\|_2 \leq \|g\|_{\infty, \mathfrak{X}} - \epsilon \right] \leq \left(1 - h \left(\frac{\epsilon}{\|g\|_{Lip, \mathfrak{X}}} \right) \right)^N. \quad (20)$$

The result now follows by setting $g = Jf$. \square

Proof of Theorem 4.6. Let $\mathfrak{X} \subset \mathbb{R}^d$ and $\mathfrak{Y} \subset \mathbb{R}^{d'}$ denote the projections of Z onto the first and second factor respectively. For each $\epsilon, \delta > 0$, there are then two events to consider:

$$A_\delta := \left\{ \max_i \|Jf(x_i)\|_2 \geq \|Jf\|_{\infty, \mathfrak{X}} - \delta \right\} \quad (21)$$

and

$$B_\epsilon := \left\{ \left| \frac{1}{N} \sum_{i=1}^N \ell(x_i, y_i) - \mathbb{E} \ell \right| \leq \epsilon \right\}. \quad (22)$$

We are interested in the intersection $A_\delta \cap B_\epsilon$: the event where the Lipschitz of f is δ -approximated by the maximum Jacobian norm over samples *and* where the generalisation gap is ϵ -small.

By Theorem 4.5 one has

$$\mathbb{P}^N(A_\delta) \geq 1 - \left(1 - h \left(\frac{\delta}{\|Jf\|_{Lip, \mathfrak{X}}} \right) \right)^N. \quad (23)$$

Since $\ell = c \circ (f \times \text{Id}_Y)$, one has $\|\ell\|_{Lip, Z} \leq \|c\|_{Lip} \|f\|_{Lip, \mathfrak{X}}$. Then Hoeffding's inequality together with the C -Lipschitz concentration inequality yield

$$\mathbb{P}^N(B_\epsilon) \geq 1 - 2 \exp \left(\frac{-NC'\epsilon^2}{\|f\|_{Lip, \mathfrak{X}}^2} \right), \quad (24)$$

where $C' = C/\|c\|_{Lip}$. Thus one has

$$\mathbb{P}^N(A_\delta \cap B_\epsilon) \geq 1 - \left(1 - h \left(\frac{\delta}{\|Jf\|_{Lip, \mathfrak{X}}} \right) \right)^N - 2 \exp \left(\frac{-NC'\epsilon^2}{\|f\|_{Lip, \mathfrak{X}}^2} \right). \quad (25)$$

Since the marginal distribution induced by \mathbb{P} on \mathfrak{X} is good by hypothesis, we have $\|f\|_{Lip, \mathfrak{X}} \leq \|Jf\|_{\infty, \mathfrak{X}}$. Moreover, for any $((x_1, y_1), \dots, (x_N, y_N)) \in A_\delta \cap B_\epsilon$ one has $\|Jf\|_{\infty, \mathfrak{X}} \leq \max_i \|Jf(x_i)\|_2 + \delta$. The result follows. \square

A.1 A note on Jacobian computations

In our experiments, the Jacobian norm and sharpness are estimated using the power method and the Hessian/Jacobian-vector product functions in PyTorch. Recalling that we use $\mathbb{R}^{d \times N}$ to denote the set of data matrices, consisting of a batch of data column vectors concatenated together, there are two different kinds of functions $f : \mathbb{R}^{d_0 \times N} \rightarrow \mathbb{R}^{d_1 \times N}$ whose Jacobians we would like to compute. The first of these are functions defined columnwise by some other map $g : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^{d_1}$: that is

$$f(X)_j = g(X_j), \quad (26)$$

where lower indices index columns so that $X_j \in \mathbb{R}^{d_0}$ for all j . In this case, flattening (columns first) shows that $Jf(X)$ is just a block-diagonal matrix, whose j^{th} block is $Jf(X_j)$. Hence $\|Jf(X)\|_2 = \max_j \|Jf(X_j)\|_2$. For a network in evaluation mode, *every* layer Jacobian is of this form.

On the other hand, for a network in *train* mode (as is the case when evaluating the Hessian of the loss), batch normalisation (BN) layers do not have this form. BN layers in train mode are *nonlinear* transformations defined by

$$bn(X) := \frac{X - \mathbb{E}X}{(\epsilon + \sigma^2 X)^{\frac{1}{2}}}, \quad (27)$$

where \mathbb{E} and σ^2 denote mean and variance computed across the column index (again we omit the parameters as they are not essential to the discussion). The skeptical reader would rightly question whether the Jacobian of BN in train mode (to which Equation (4) applies) is sufficiently close to the Jacobian of BN in evaluation mode (to which Theorem 4.6 applies) for Ansatz 3.1 to be valid for BN networks. Our next and final theorem says that for sufficiently large datasets, this is not a problem.

Theorem A.1. *Let bn_T and bn_V denote a batch normalisation layer in train and evaluation mode respectively. Given a data matrix $X \in \mathbb{R}^{d \times N}$, assume the coordinates of X are $O(1)$. Then $Jbn_T = Jbn_V + O(N^{-1})$.*

Proof of Theorem A.1. Given a data matrix X , let its row-wise mean and variance, thought of as constant vectors, be denoted by \mathbb{E} and σ^2 . The evaluation mode BN map $bn_V : \mathbb{R}^{d \times N} \rightarrow \mathbb{R}^{d \times N}$ is simply an affine transformation, with Jacobian coordinates given by

$$\frac{\partial(bn_V)_j^i}{\partial x_l^k} = \frac{\delta_k^i \delta_l^j}{(\epsilon + \sigma^2)^{\frac{1}{2}}}. \quad (28)$$

By [34, Lemma 8.3], the train mode BN map $bn_T : \mathbb{R}^{d \times N} \rightarrow \mathbb{R}^{d \times N}$ can be thought of as the composite $v \circ m$, where $m, v : \mathbb{R}^{d \times N} \rightarrow \mathbb{R}^{d \times N}$ are defined by

$$m(X) := X - \frac{1}{N} X 1_{N \times N} \quad (29)$$

and

$$v(Y) := (\epsilon + N^{-1} \|Y\|_{row}^2)^{-\frac{1}{2}} Y \quad (30)$$

respectively. Here $1_{N \times N}$ denotes the $N \times N$ matrix of 1s, and $\|Y\|_{row}$ denotes the vector of row-wise norms of Y .

As in [34, Lemma 8.3], one has

$$\frac{\partial v_j^i}{\partial y_l^k}(Y) = \frac{\delta_k^i}{(\epsilon + N^{-1} \|Y\|^2)^{\frac{1}{2}}} \left(\delta_l^j - \frac{y_l^i y_j^i}{(\epsilon + N^{-1} \|Y\|^2)} \right), \quad \frac{\partial m_j^i}{\partial x_l^k}(Y) = \delta_k^i (\delta_l^j - N^{-1}) \quad (31)$$

for any data matrix Y . Thus, invoking the chain rule and simplifying, we therefore get

$$\frac{\partial(bn_T)_j^i}{\partial x_l^k}(X) = \frac{\delta_k^i \delta_l^j}{(\epsilon + \sigma^2)^{\frac{1}{2}}} - \frac{1}{N} \frac{\delta_k^i (x_l^i - \mathbb{E}^i)(x_l^j - \mathbb{E}^j)}{(\epsilon + \sigma^2)^{\frac{3}{2}}} + O(N^{-1}) \quad (32)$$

Since the components x are $O(1)$ by hypothesis, the result follows. \square

B Additional experimental results and details

B.1 Progressive sharpening

We trained ResNet18 and VGG11 (superficially modifying <https://github.com/kuangliu/pytorch-cifar/blob/master/models/resnet.py> and <https://github.com/chengyangfu/pytorch-vgg-cifar10/blob/master/vgg.py> respectively, with VGG11 in addition having its dropout layers removed, but BN layers retained) with full batch gradient descent on CIFAR10 with learning rates 0.08, 0.04 and 0.02, using label smoothings of 0.0, 0.5 and 0.75. The models with no label smoothing were trained to 99 percent accuracy, and the number of iterations required to do this were then used as the number of iterations to train the models using nonzero label smoothing. Sharpness and Jacobian norm were computed every 5, 10, or 20 iterations depending on the number of iterations required for convergence of the non-label-smoothed model. Five trials were conducted in each case, with mean and standard deviation plotted. Line style indicates degree of label smoothing.

The dataset was standardised so that each vector component is centered, with unit standard deviation. Full batch gradient descent was approximated by averaging the gradients over 10 ‘ghost batches’ of size 5000 each, as in [12]. This is justified with BN networks by Theorem A.1. The Jacobian norm we plot is for the Jacobian of the softmaxed model, as required by our derivation of Equation (11). The sharpness and Jacobian norms were estimated using a randomly chosen subset of 2500 data points. In all cases, the smoother labels are associated with less severe increase in Jacobian and sharpness

during training, as predicted by Equation (11). Refer to `fullbatch.py` in the supplementary material for the code we used to run these experiments.

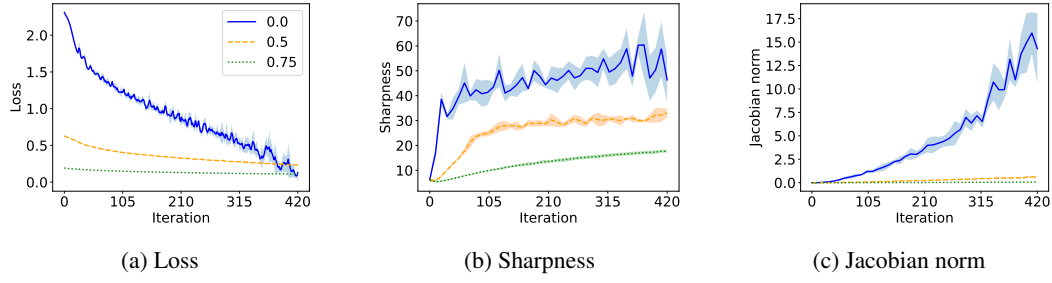


Figure 6: VGG11, learning rate 0.08

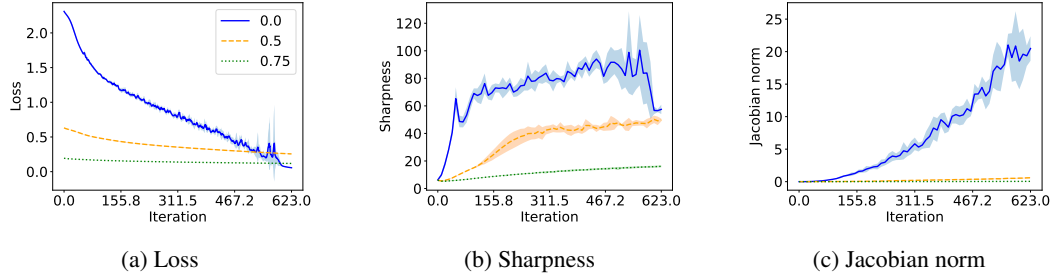


Figure 7: VGG11, learning rate 0.04

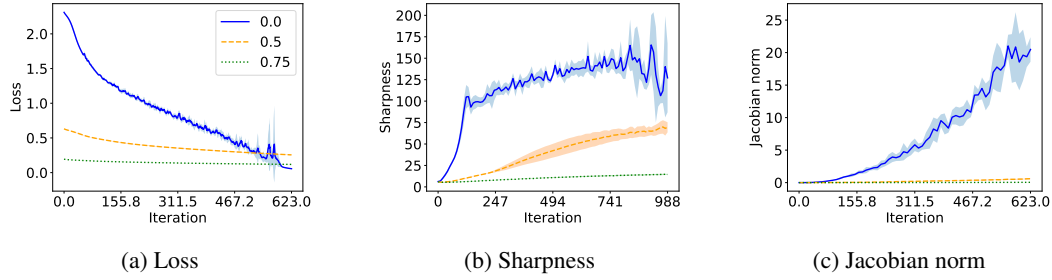


Figure 8: VGG11, learning rate 0.02

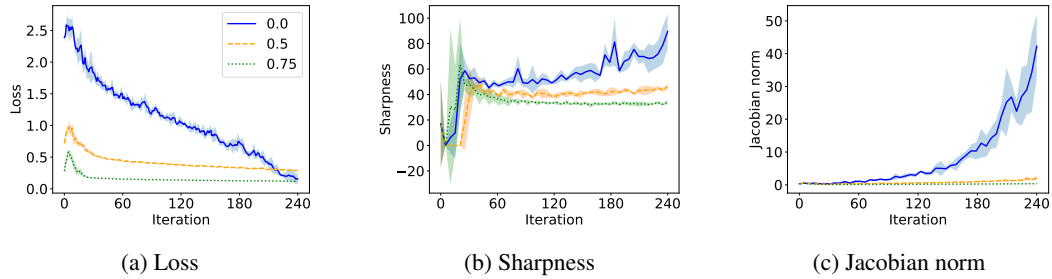


Figure 9: ResNet18, learning rate 0.08

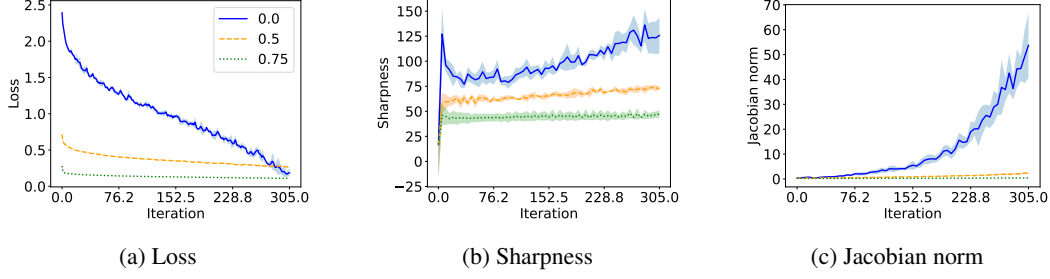


Figure 10: ResNet18, learning rate 0.04

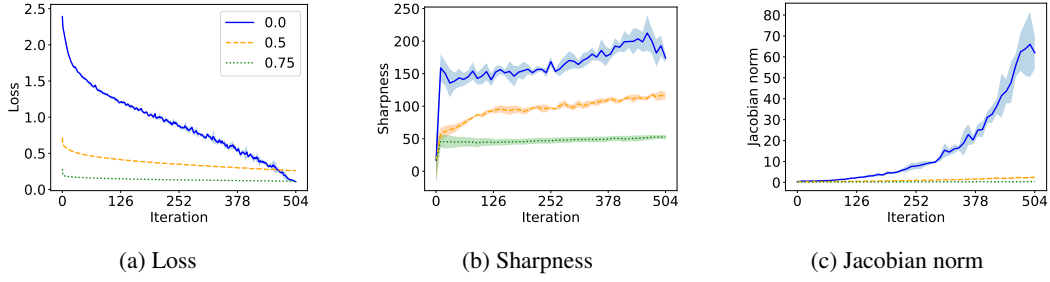


Figure 11: ResNet18, learning rate 0.02

B.2 Batch size and learning rate

We trained a VGG11 and ResNet18 (superficially modifying <https://github.com/chengyangfu/pytorch-vgg-cifar10/blob/master/vgg.py> and <https://github.com/kuangliu/pytorch-cifar/blob/master/models/resnet.py> respectively, additionally removing dropout and retaining BN in the VGG11) on CIFAR10 and CIFAR100 at constant learning rates 0.1, 0.01, and 0.001, with batch sizes of 64, 128, 256, 512 and 1024. The data were normalised by their RGB-channelwise mean and standard deviation, computed across all pixel coordinates and all elements of the training set. Five trials of each were conducted. Training loss on the full data set was computed every 100 iterations, and training was terminated when the average of the most recent 10 such losses was smaller than 0.01 for CIFAR10 and 0.02 for CIFAR100. Weight decay regularisation of 0.0001 was used throughout. Refer to `batch_lr.py` in the supplementary material.

At the termination of training, training loss and test loss were computed, with estimates of the Jacobian norm and sharpness computed using a randomly selected subset of size 2000 from the training set. The same seed was used to generate the parameters for a given trail as was used to generate the random subset of the training set.

We see that, on the whole, the claim that smaller batch sizes and larger learning rates lead to flatter minima and better generalisation is supported by our experiments, with Jacobian norm in particular being well-correlated to generalisation gap as anticipated by Theorem 4.6 and Ansatz 3.1. As noted in the main body of the paper, the notable exception to this is ResNet models trained with the largest learning rate, where Jacobian norm appears to underestimate the Lipschitz constant of the model and hence the generalisation gap. We speculate that this is due to large learning rate training of skip connected models leading to larger Jacobian Lipschitz constant, making the Jacobian norm a poorer estimate of the Lipschitz constant of the model and hence of generalisation (Theorems 4.5 and 4.6).

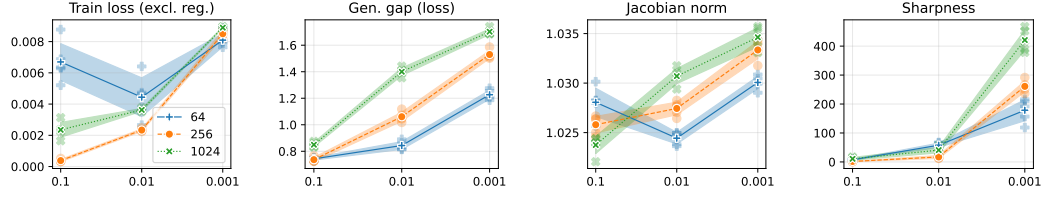


Figure 12: ResNet18 on CIFAR10, learning rate on x axis. Line style indicates batch size.

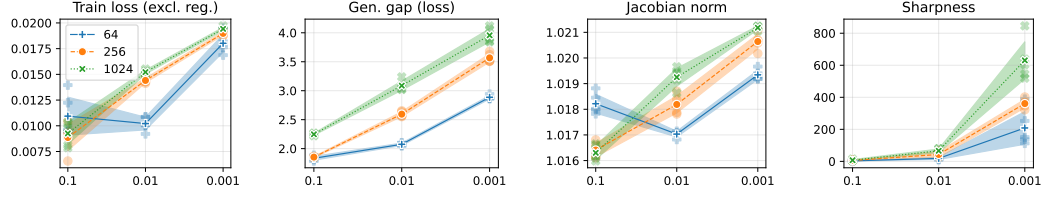


Figure 13: ResNet18 on CIFAR100, learning rate on x axis. Line style indicates batch size.

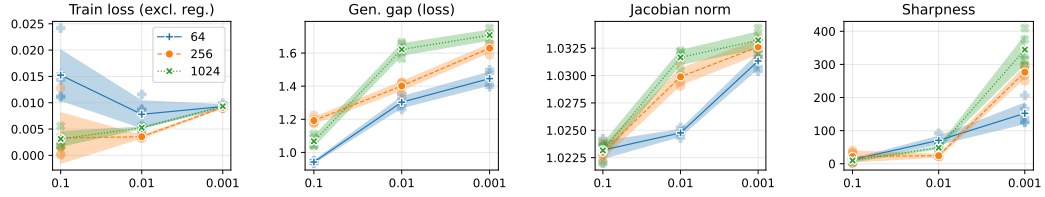


Figure 14: VGG11 on CIFAR10, learning rate on x axis. Line style indicates batch size.

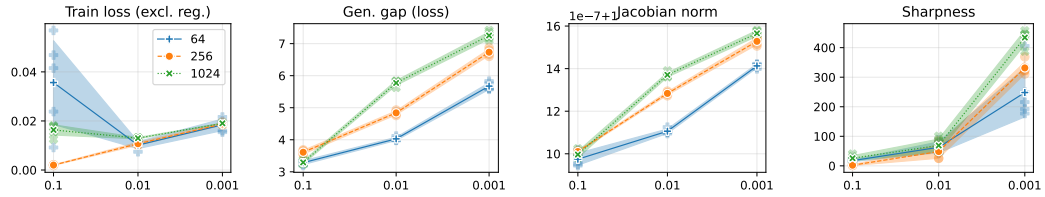


Figure 15: VGG11 on CIFAR100, learning rate on x axis. Line style indicates batch size.

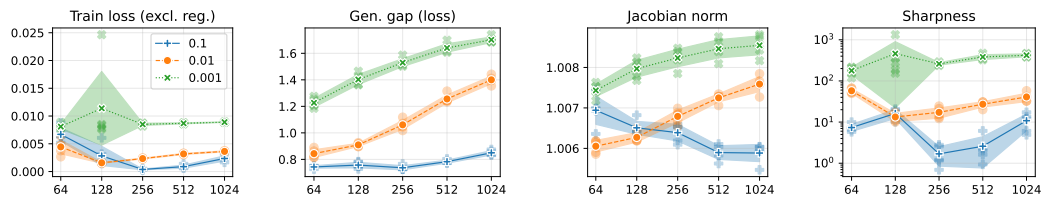


Figure 16: ResNet18 on CIFAR10, batch size on x axis. Line style indicates learning rate.

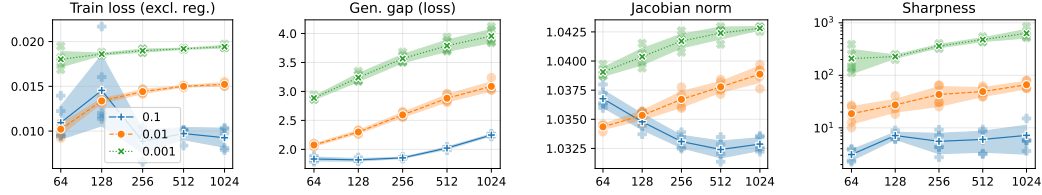


Figure 17: ResNet18 on CIFAR100, batch size on x axis. Line style indicates learning rate.

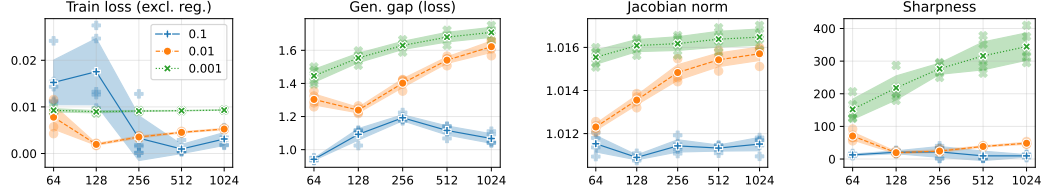


Figure 18: VGG11 on CIFAR10, batch size on x axis. Line style indicates learning rate.

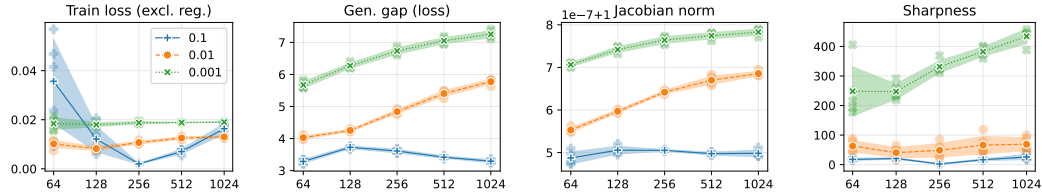


Figure 19: VGG11 on CIFAR100, batch size on x axis. Line style indicates learning rate.

B.3 Extended practical results

Here we provide the experimental details and additional experimental results for Section 6.3. In addition to the generalisation gap, sharpness, and (input-output) Jacobian norm, these plots include the operator norm of the Gauss-Newton matrix, the loss on the training set, and the accuracy on the validation set. Note that the figures in this section present the *final* metrics at the conclusion of training as a function of the regularisation parameter, and hence progressive sharpening (as a function of time) will not be observed directly. Rather, the purpose is to examine the impact of different regularisation strategies on the Hessian and model Jacobian.

In all cases, the model Jacobian is correlated with the generalisation gap, whereas the Hessian is often uncorrelated. Moreover, the norm of the Gauss-Newton matrix is (almost) always very similar to that of Hessian, in line with previous empirical observation [39, 40, 12]. The results support the arguments that (i) the model Jacobian is strongly related to the generalisation gap and (ii) while the Hessian is connected to the model Jacobian, and forcing the Hessian to be sufficiently small likewise appears to force the Jacobian to be smaller (see the SAM plots), the Hessian is also influenced by other factors and can be large even without the Jacobian being large (see other plots). All of these phenomena are consistent with our ansatz.

The figures in this section present results for $\{\text{CIFAR10, CIFAR100}\} \times \{\text{VGG11, ResNet18}\} \times \{\text{batch-norm, no batch-norm}\}$. Line colour and style denote different initial learning rates. The models were trained for 90 epochs with a minibatch size of 128 examples using Polyak momentum of 0.9. The learning rate was decayed by a factor of 10 after 50 and 80 epochs. Models were trained using the softmax cross-entropy loss with weight decay of 0.0005. (If weight decay were disabled, we would often observe the Hessian to collapse to zero as the training loss went to zero, due to the vanishing second gradient of the cost function in this region [12, Appendix C]. This would make the correlation between sharpness and generalisation gap even worse.) Refer to `train_jax.py`

and `slurm/launch_wd_sweep.sh` for the default configuration and experiment configuration, respectively.

For each regularisation strategy, the degree of regularisation is parametrised as follows. Label smoothing considers smoothed labels $(1 - \alpha)y + \alpha(1/n)$ with $\alpha \in [0, 1]$. Mixup takes a convex combination of two examples $(1 - \theta)(x, y) + \theta(x', y')$. The coefficient θ is drawn from a symmetric beta distribution $\theta \sim \text{Beta}(\beta, \beta)$ with $\beta \in [0, \infty)$, which corresponds to a Bernoulli distribution when $\beta = 0$ and a uniform distribution when $\beta = 1$. Data augmentation modifies the inputs x with probability $p \in [0, 1]$; the image transforms which we adopt are the standard choices for CIFAR (four-pixel padding, random crop, random horizontal flip). Sharpness Aware Minimisation (SAM) restricts the distance ρ between the current parameter vector and that which is used to compute the update, which simplifies to SGD when $\rho = 0$.

For all results presented in this section, the loss and generalisation gap were computed using “clean” training examples; that is, without mixup or data augmentation in the respective experiments (despite this requiring an additional evaluation of the model for each step). The Hessian, Gauss-Newton and input-output Jacobian norms were similarly computed using a random batch of 1000 clean training examples. The losses and matrix norms in this section similarly exclude weight decay. Label smoothing was interpreted as a modification of the loss function rather than the training distribution, and therefore *was* included in the calculation of the losses and matrix norms. Batch-norm layers were used in “train mode” (statistical moments computed from batch) for the Hessian and Gauss-Newton matrix, and in “eval mode” (statistical moments are constants) for the Jacobian matrix, although this difference was observed to have negligible effect on the Jacobian norm.

B.3.1 CIFAR10

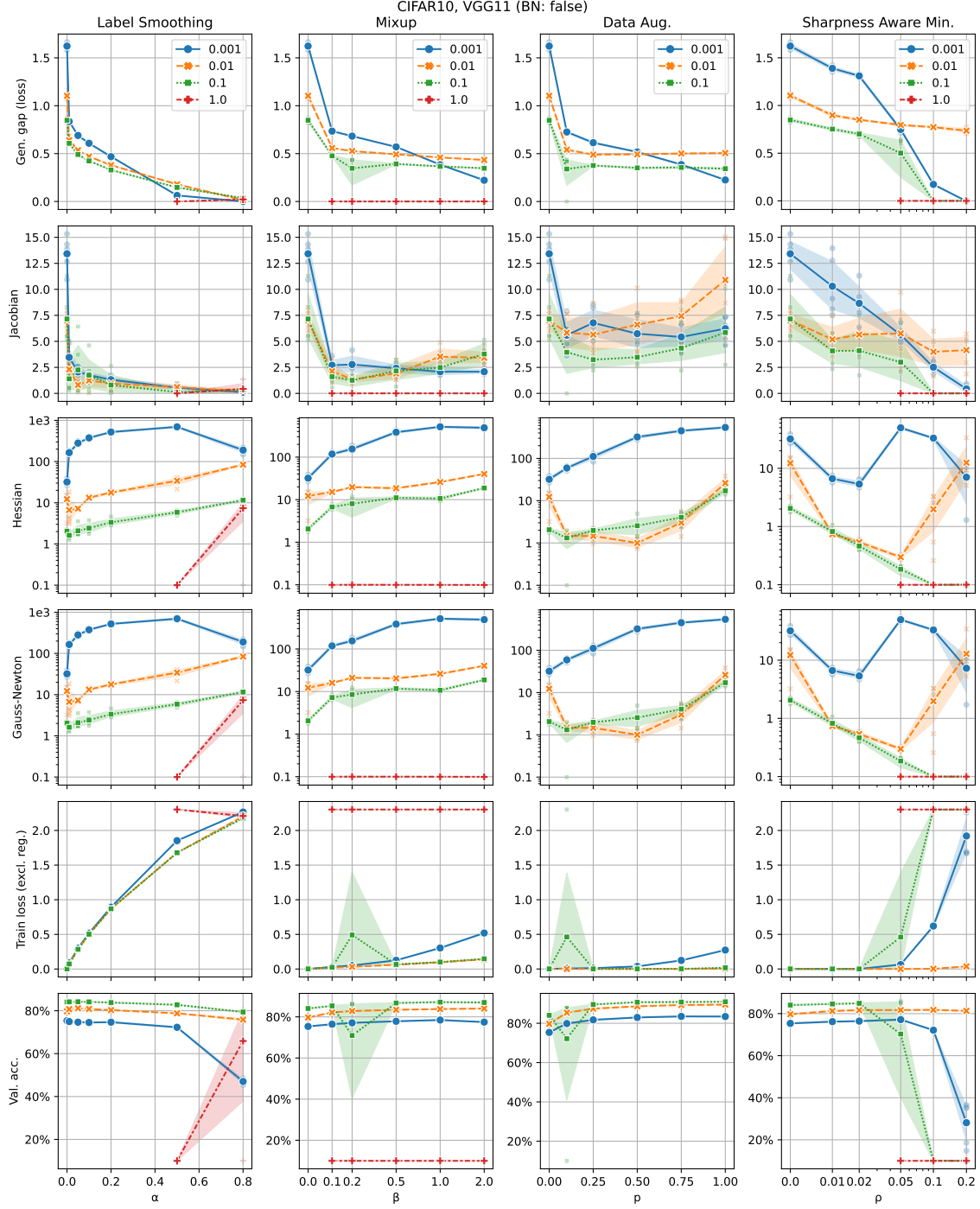


Figure 20: VGG11 without batch-norm on CIFAR10.

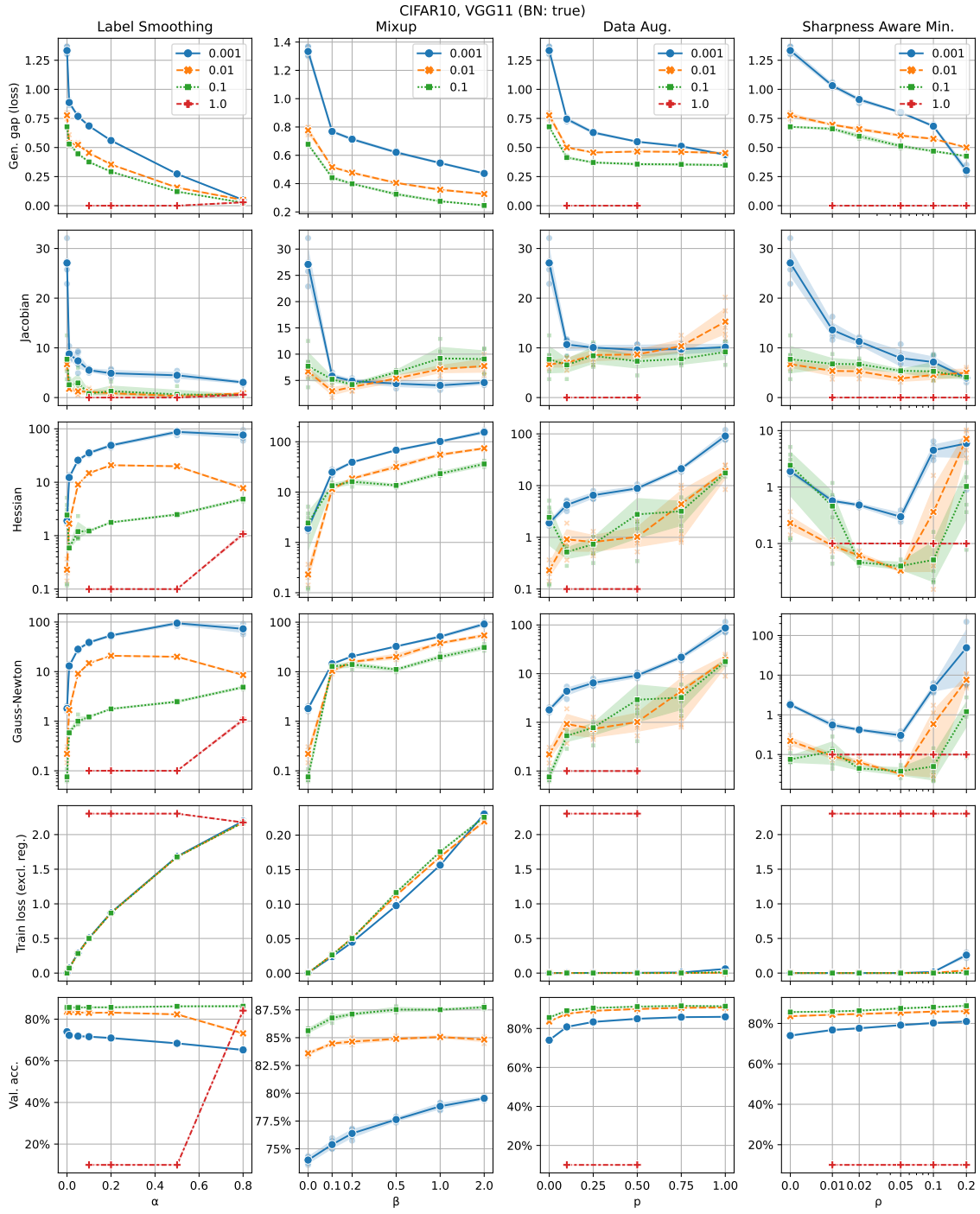


Figure 21: VGG11 with batch-norm on CIFAR10.

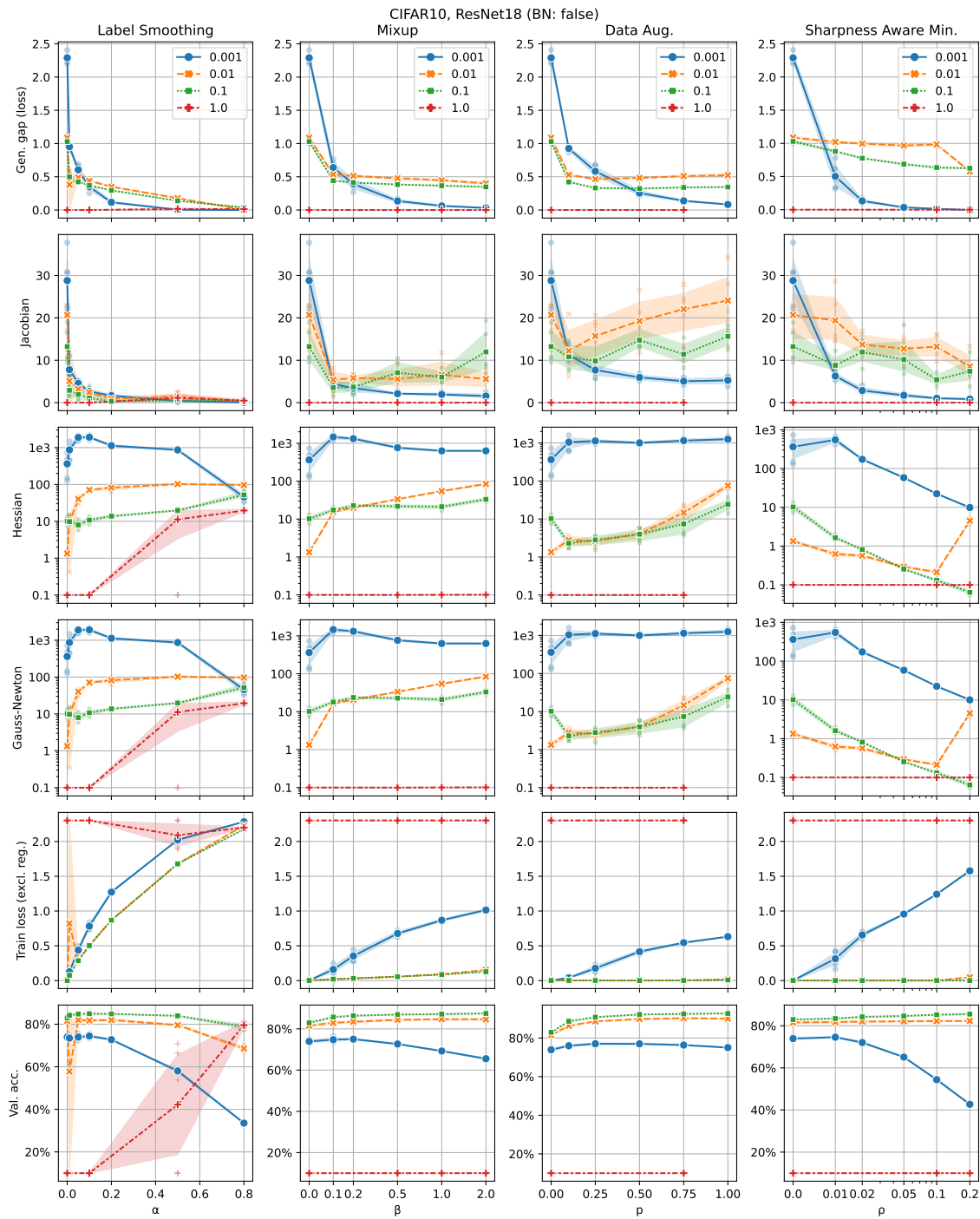


Figure 22: ResNet18 without batch-norm on CIFAR10.

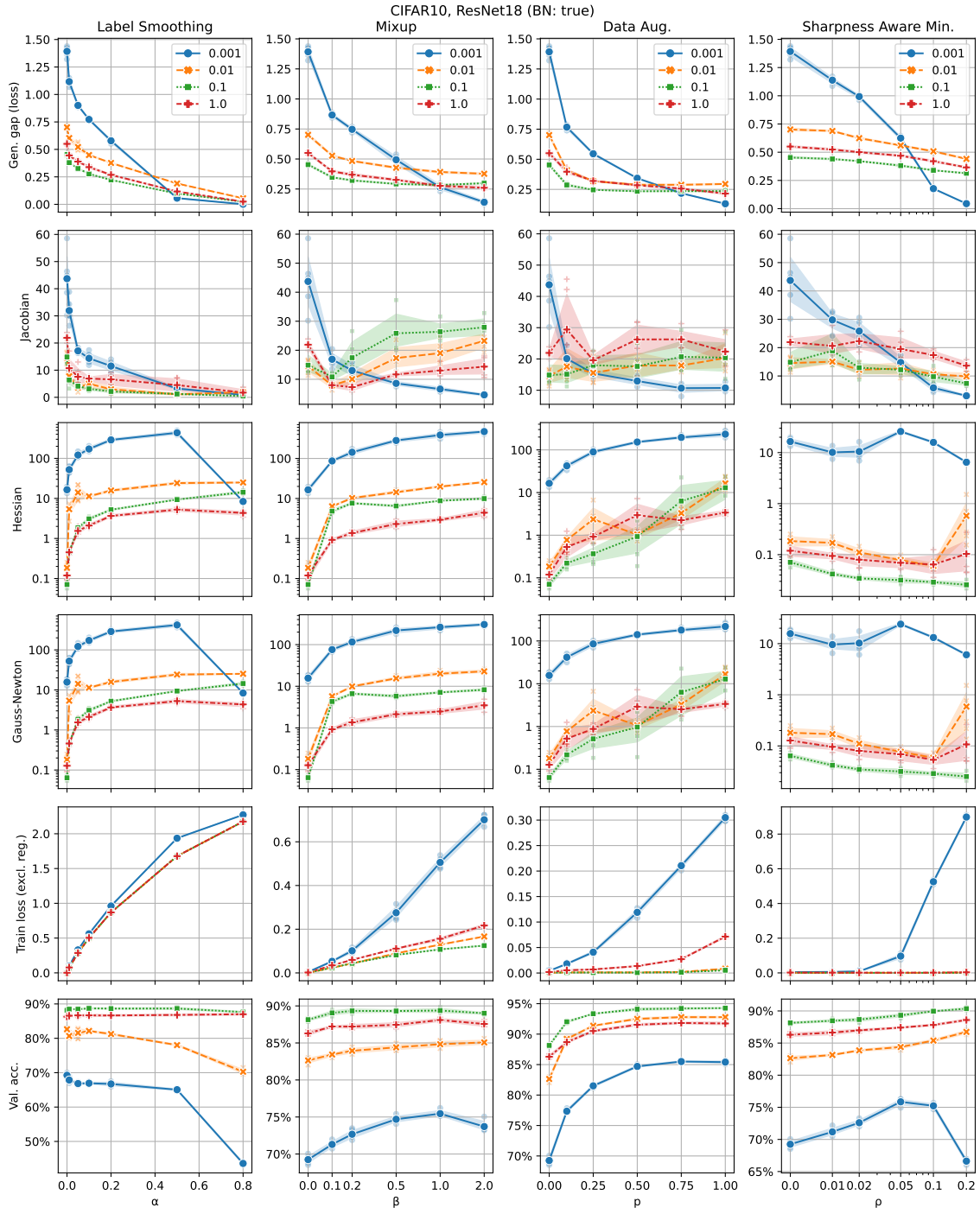


Figure 23: ResNet18 with batch-norm on CIFAR10.

B.3.2 CIFAR100

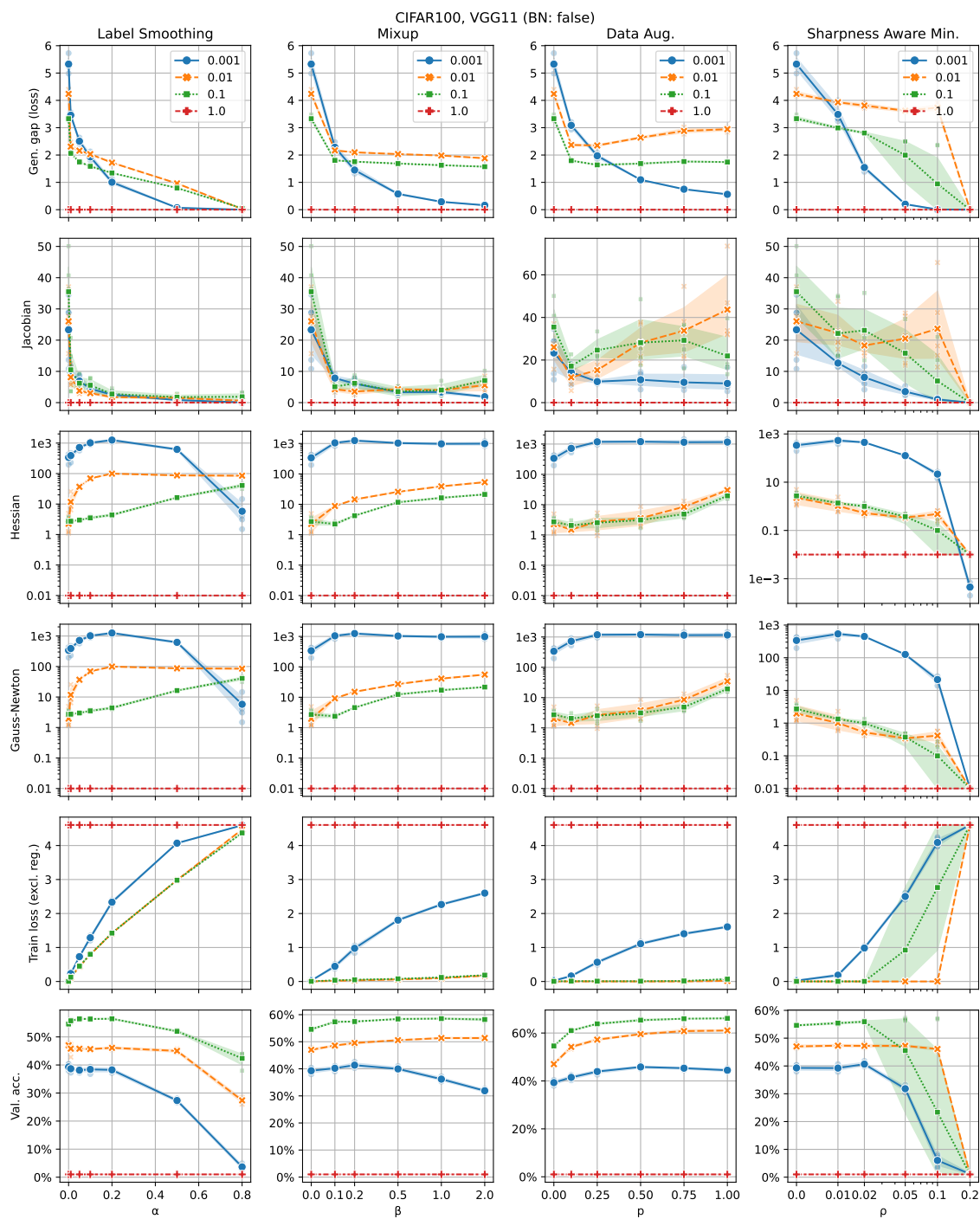


Figure 24: VGG11 without batch-norm on CIFAR100.

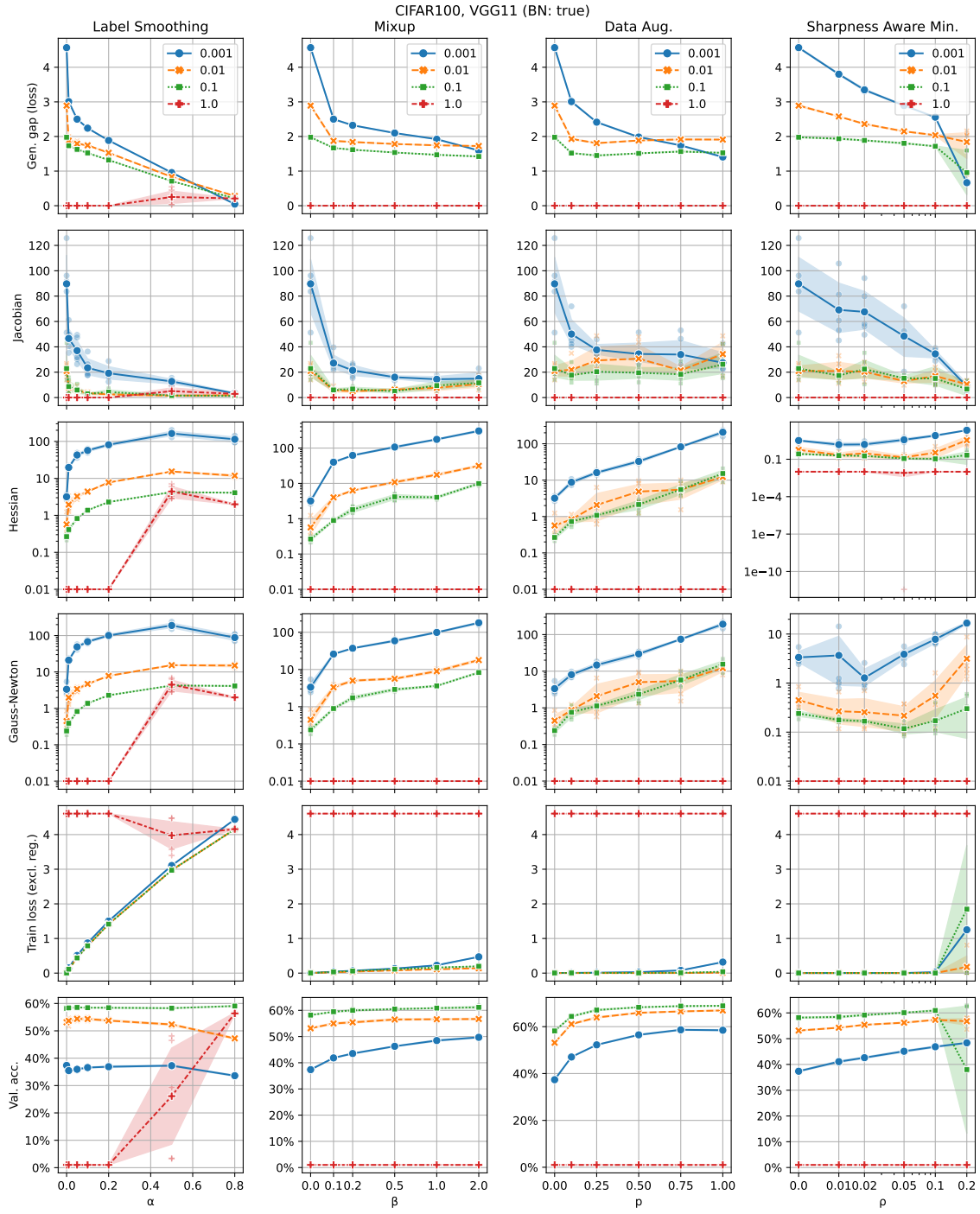


Figure 25: VGG11 with batch-norm on CIFAR100.



Figure 26: ResNet18 without batch-norm on CIFAR100.

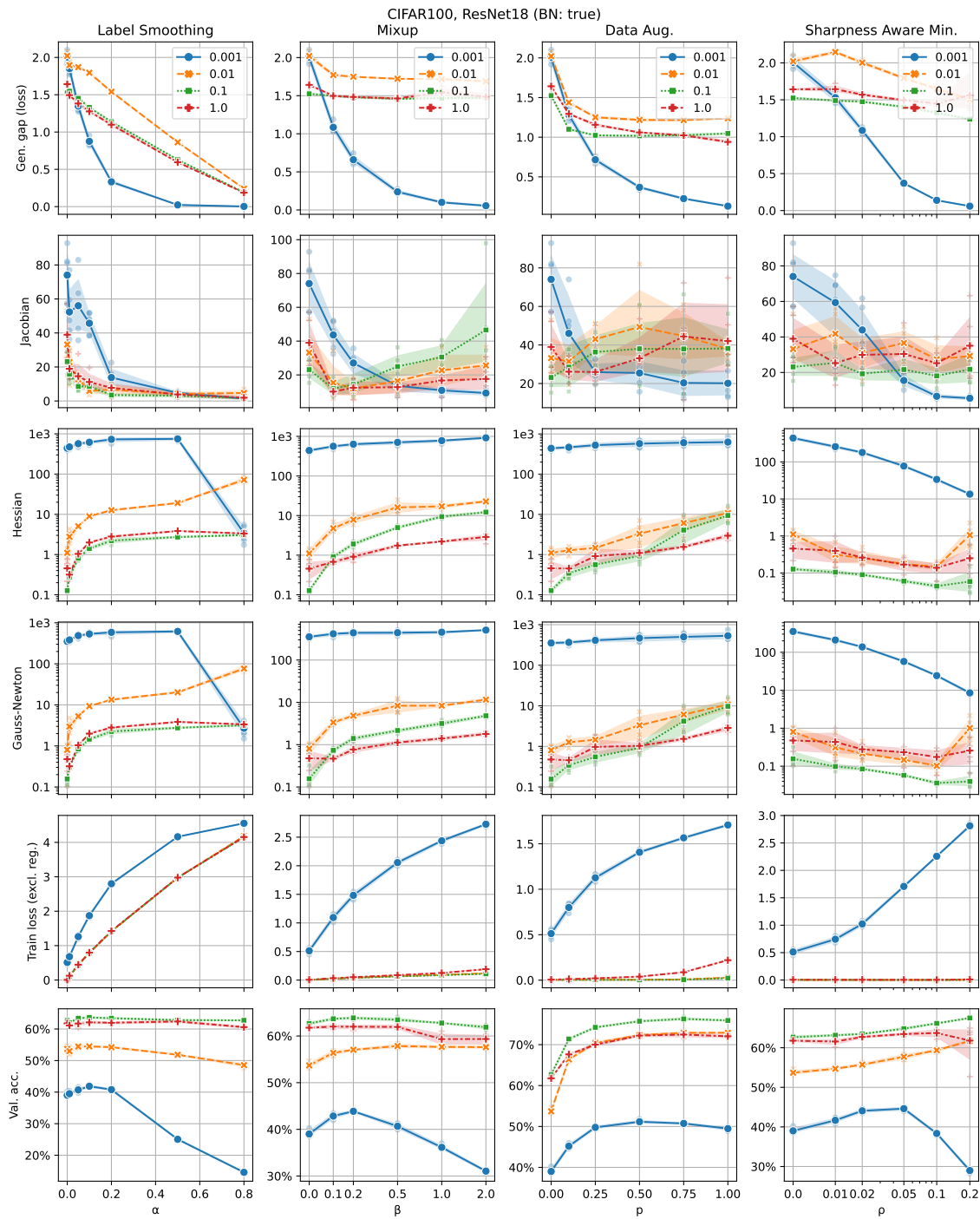


Figure 27: ResNet18 with batch-norm on CIFAR100.

C On mediating factors in the ansatz

In this section we provide classification and regression experiments to show that the relationship between loss curvature and input-output Jacobians is not always simple. Recall again Equation (4):

$$C D F_X D F_X^T C^T = C \left(\sum_{l=1}^L \left(J f_L \cdots J f_{l+1} D f_l D f_l^T J f_{l+1}^T \cdots J f_L^T \right) \right) C^T. \quad (33)$$

Due to the presence of the parameter derivatives $D f_l$ and the square root C of the Hessian of the cost function, as well as the absence of the first layer Jacobian in Equation (4), it is not always true that the magnitude of the Hessian and that of the model’s input-output Jacobian are correlated.

C.1 The cost function

We trained a VGG11, again using <https://github.com/chengyangfu/pytorch-vgg-cifar10/blob/master/vgg.py> with dropout layers removed and BN layers retained, on CIFAR10 using SGD with a batch size of 128, using differing degrees of label smoothing. Data was standardised so that each RGB channel has zero mean and unit standard deviation over all pixel coordinates and training samples. In Figure 28 we plot the sharpness and Jacobian norm every five iterations in the first epoch, observing exactly the same relationship between label smoothing, Jacobian norm and sharpness as predicted in Section 5. Refer to `vgg_sgd.py` in the supplementary material for the code to run these experiments.

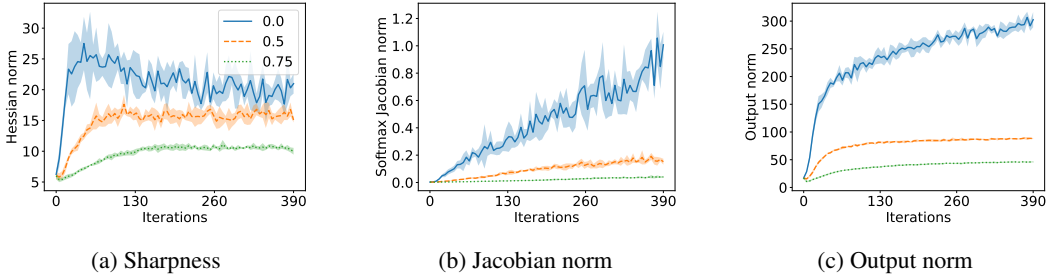


Figure 28: Effect of label smoothing on sharpness, Jacobian norm and output Frobenius norm during the first epoch of SGD. Line style indicates label smoothing. Exactly as in the full batch GD case, the smoother labels are associated with less severe increase of the Jacobian and less severe progressive sharpening.

Zooming out, however, to the full training run over 90 epochs, Figure 29 shows that the Hessian ultimately behaves markedly differently.

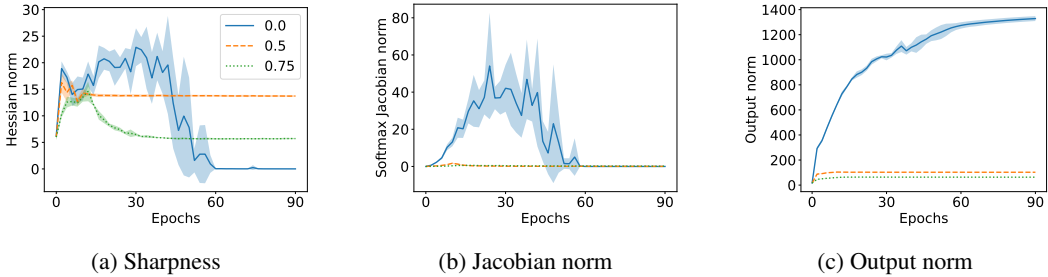


Figure 29: Effect of label smoothing on sharpness, Jacobian norm and output Frobenius norm during 90 epochs of SGD. Line style indicates label smoothing. When the output norm gets too large, the decay of the C terms in Equation (4) begins to overtake growth in Jacobian norm, ultimately causing the sharpness with unsmoothed labels to collapse to zero, whereas sharpness values for smoothed labels plateau at nonzero values. Note also the the Jacobian norm presented is the Jacobian norm of the softmaxed model: the growth in output norm therefore also drives the Jacobian norm to zero.

The reason for this is that the norm of the network output grows to infinity in the unsmoothed case (Figure 29), and this growth in output corresponds to a vanishing of the C term [12, Appendix C] corresponding to the cross-entropy cost in Equation (4). Note that this growth is also to blame for the eventual collapse of the norm of the Jacobian of the softmaxed model, since the derivative of softmax vanishes at infinity also. This vanishing of the Jacobian should not be interpreted as saying that the *Lipschitz* constant of the model has collapsed (since if this were the case the model would not be able to fit the data), but rather that the maximum Jacobian norm over training samples has ceased to be a good approximation of the Lipschitz constant of the model, due to the Jacobian itself having a large Lipschitz constant (Theorem 4.5).

C.2 The parameter derivatives

Recall the estimate

$$\|f\|_{Lip} \geq \frac{\|y_1 - y_2\| - 2\epsilon}{\|x_1 - x_2\|} \quad (34)$$

from Section 5. We have already shown that decreasing the distance between labels decreases the severity of Jacobian growth, and, in line with Ansatz 3.1, also reduces the severity of progressive sharpening. It seems that one could also manipulate this severity by increasing the distance between input data points, for instance by scaling all data by a constant.

We test this training simple, fully-connected, three layer networks, of width 200, with gradient descent on the first 5000 data points of CIFAR10. The data is standardised to have componentwise zero mean and unit standard deviation (measured across the whole dataset). The networks are initialised using the uniform distribution on the interval with endpoints $\pm 1/\sqrt{\text{in features}}$ (default PyTorch initialisation) and trained with a learning rate of 0.2 for 300 iterations using the cross-entropy cost. Both ReLU and tanh activations are considered. Refer to `small_network.py` in the supplementary material for the code we used to run these experiments.

We scale the inputs by factors of 0.5, 1.0 and 1.5, bringing the data closer together at 0.5 and further apart at 1.5. From Equation (34), we anticipate the Jacobian growth to go from most to least severe on these scalings respectively, with sharpness growth behaving similarly according to Ansatz 3.1. Surprisingly, we find that while Jacobian growth *is* more severe for the smaller scalings, the sharpness growth is *less*.

The reason for this is the effect of data scaling on the parameter derivatives Df_l in Equation (4). It is easily computed (cf. [34]) that the parameter derivative Df_l is simply determined by the matrix $f_{l-1}(X)$ of features from the previous layer. In the experimental settings we examined, these matrices are *larger* for the larger scaling values, which pushes the sharpness upwards despite Jacobian norm being smaller. Figures 30 and 31 show this phenomenon on ReLU and tanh activated networks respectively.

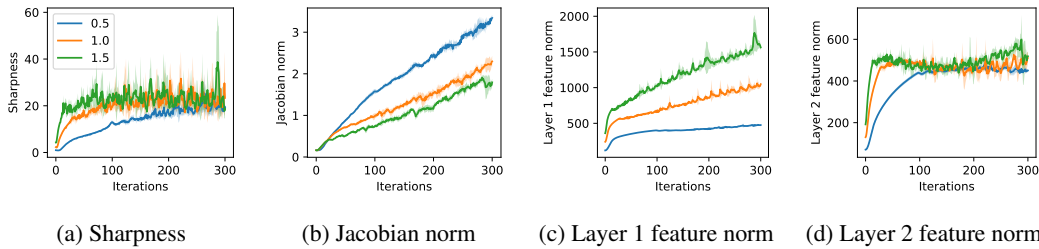


Figure 30: Effect of label smoothing on sharpness and Jacobian norm as inputs are scaled for a ReLU network (5 trials). Line style denotes input scaling factor. Scaling data closer together does cause more severe increase in Jacobian norm, but corresponds to *less* severe increase in sharpness. This is due to the data scaling increasing the spectral norm of the feature maps.

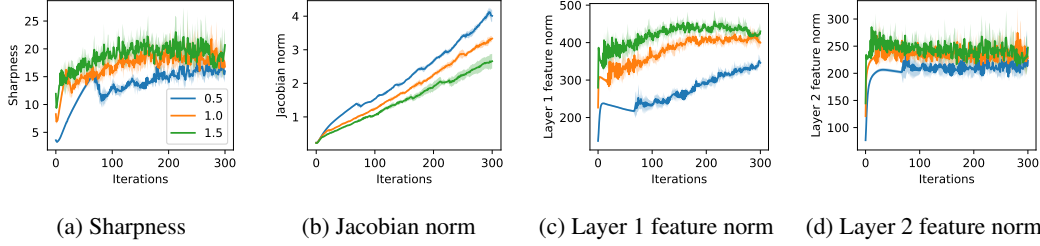


Figure 31: Effect of label smoothing on sharpness and Jacobian norm as inputs are scaled for a tanh network (5 trials). Line style denotes input scaling factor. Scaling data closer together does cause more severe increase in Jacobian norm, but corresponds to *less* severe increase in sharpness. This is due to the data scaling increasing the spectral norm of the feature maps.

C.3 The absence of the first layer Jacobian

We replicate the experiments of Section 8 in [42], where it is demonstrated that flatness of minimum is not correlated with model smoothness in a simple regression task, and point to a possible explanation of this within our theory.

In detail, a four layer MLP with either Gaussian or ReLU activations is trained to regress 8 points in \mathbb{R}^2 . Each network is trained from a *high frequency* initialisation (obtained from default PyTorch initialisation on the Gaussian network, and by pretraining on $\sin(6\pi x)$ for the ReLU network) to yield a non-smooth fit of the target data, and a *low frequency* initialisation (a wide Gaussian distribution for the Gaussian-activated network, and the default PyTorch initialisation for the ReLU network) to achieve a smooth fit of the data (see Figure 32). Refer to `coordinate_network.py` in the supplementary material for the code we used to run this experiment.

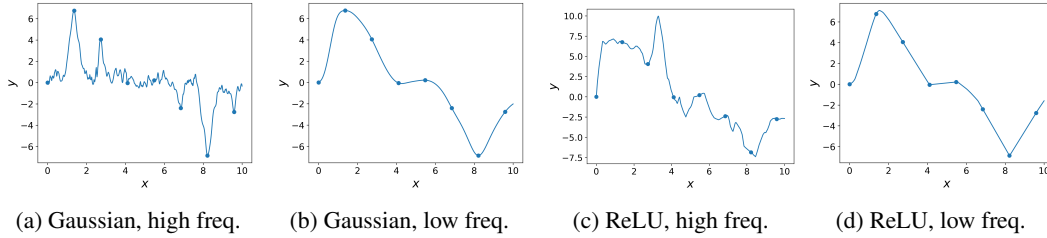


Figure 32: Interpolations of 8 points by networks started from high frequency and low frequency initialisations. The smoother interpolations have lower Jacobian norm over the training data.

The models were trained with gradient descent using a learning rate of $1e-4$ and momentum of 0.9, for 10000 epochs for the Gaussian networks and 100000 epochs for the ReLU networks. The pretraining for the high frequency ReLU initialisation was achieved using Adam with a learning rate of $1e-4$ and default PyTorch settings.

Table 1: Norm values for regression networks, replicating result of [42]

Norm	Gaussian activated network	
	High freq.	Low freq.
Jacobian norm	326.80 ± 182.51	84.07 ± 49.30
sharpness	13517.61 ± 3871.97	18353.54 ± 5751.98
First layer weight norm	9.09 ± 0.34	0.56 ± 0.02

Tables 1 and 2 record the means and standard deviations of loss Hessian and Jacobian norms of the Gaussian and ReLU networks over 10 trials. As in [42], we find that while the smoothly interpolating ReLU models on average land in flatter minima than the non-smoothly interpolating models, the opposite is true for the Gaussian networks.

Table 2: Norm values for regression networks, replicating result of [42]

Norm	ReLU activated network	
	High freq.	Low freq.
Jacobian norm	121.59 ± 80.79	42.96 ± 7.82
sharpness	14211.05 ± 4767.26	9922.45 ± 3330.37
First layer weight norm	9.61 ± 0.28	9.56 ± 0.16

Keeping in mind that the high variances in these numbers make it difficult to come to a conclusion about the trend, we propose a possible explanation for why the Gaussian activated networks do not appear to behave according to Ansatz 3.1. From Equation (4), the Hessian cannot be expected to relate to the Jacobian of the first layer. Note now the discrepancy between the first layer weight norms in the low frequency versus high frequency fits with the Gaussian networks, which does not occur for the ReLU networks. We believe this to be the cause of the larger Hessian for the low-frequency Gaussian models: with such a small weight matrix in the first layer, all higher layers must be relied upon to fit the data, making their Jacobians higher than they would otherwise need to be. These larger Jacobians would then push the Hessian magnitude higher by Equation (4).