

A SIMPLE CONNECTION FROM LOSS FLATNESS TO COMPRESSED REPRESENTATIONS IN NEURAL NETWORKS

Anonymous authors

Paper under double-blind review

ABSTRACT

1 Deep neural networks’ generalization capacity has been studied in a variety of
2 ways, including at least two distinct categories of approach: one based on the
3 shape of the loss landscape in parameter space, and the other based on the struc-
4 ture of the representation manifold in feature space (that is, in the space of unit
5 activities). These two approaches are related, but they are rarely studied together
6 and explicitly connected. Here, we present a simple analysis that makes such a
7 connection. We show that, in the last phase of learning of deep neural networks,
8 compression of the manifold of neural representations correlates with the flatness
9 of the loss around the minima explored by SGD. We show that this is predicted
10 by a relatively simple mathematical relationship: flatter loss gives a lower upper
11 bound on metrics of the compression of neural representations. Our results build
12 on the prior work of Ma and Ying, which shows how flatness (i.e., small eigenval-
13 ues of the loss Hessian) develops in late phases of learning and leads to robustness
14 to perturbations in network inputs. Moreover, we show there is no similarly di-
15 rect connection between local dimensionality and sharpness, suggesting that this
16 property may be controlled by different mechanisms than volume and hence may
17 play a complementary role in neural representations. Overall, we advance a dual
18 perspective on generalization in neural networks in both parameter and feature
19 space.

20 1 INTRODUCTION

21 Deep neural networks’ generalization capacity has been studied in many ways. Generalization is
22 a complex phenomenon influenced by myriad factors, including model architecture, dataset size
23 and diversity, and the specific task used to train a network. Researchers continue to develop new
24 techniques to enhance generalization (Elsayed et al., 2018; Galanti et al., 2023). From a theoretical
25 point of view, we can identify two distinct categories of approach. These are works that study
26 neural network generalization in the context of (a) properties of minima of the loss function that
27 learning algorithms find in parameter space (Dinh et al., 2017; Andriushchenko et al., 2023), and (b)
28 properties of the representations that optimized networks find in feature space – that is, in the space
29 of their neural activations (Ben-Shaul & Dekel, 2022; Ben-Shaul et al., 2023; Rangamani et al.,
30 2023; Pappan et al., 2020).

31 One of the most widely studied factors that influence generalization is the shape of the loss landscape
32 in parameter space. Empirical studies and theoretical analyses have shown that training deep neural
33 networks using stochastic gradient descent (SGD) with a small batch size and a large learning rate
34 often converges to flat and wide minima (Ma & Ying, 2021; Blanc et al., 2020; Geiger et al., 2021;
35 Li et al., 2022; Wu et al., 2018; Jastrzebski et al., 2018; Xie et al., 2021; Zhu et al., 2019). Flat
36 minima refer to regions in the loss landscape where the loss function has a relatively large basin:
37 put simply, the loss doesn’t change much in different directions around the minimum. Many works
38 conjecture that flat minima lead to a simpler model (shorter description length), and thus are less
39 likely to overfit and more likely to generalize well (Jastrzebski et al., 2018; Yang et al., 2023; Wu
40 et al., 2018). However, whether flatness positively correlates with the network’s generalization
41 capability remains unsettled (Dinh et al., 2017; Andriushchenko et al., 2023; Yang et al., 2021).
42 In particular, Dinh et al. (2017) argues that one can construct very sharp networks that generalize

43 well through reparametrization. However, more recent work (Andriushchenko et al., 2023) shows
44 that even reparametrization-invariant sharpness cannot capture the relationship between sharpness
45 and generalization.

46 In our work, we investigate how the sharpness of the loss function near learned solutions in param-
47 eter space influences local geometric features of neural representations. We demonstrate that as this
48 sharpness decreases and the minima become flatter, there is a set of mathematical bounds that imply
49 that the neural representation must undergo at least a specific, computable level of compression.
50 This process, which is related to previous results including the concept of neural collapse (Farrell
51 et al., 2022; Kothapalli et al., 2022; Zhu et al., 2021; Ansuini et al., 2019; Recanatesi et al., 2019;
52 Pappayan et al., 2020), refers to the emergence of a more compact and by some measures lower di-
53 mensional structure in the neural representation space. Compression in the feature space enables
54 networks to isolate the most crucial and discriminative features of input data. As a model becomes
55 less sensitive to small perturbations or noise in the input data, it gains increased robustness against
56 variations between training and test data. This simple and direct relationship between compression
57 and robustness creates a valuable lens into networks’ potential to generalize.

58 We find that bounds that apply to two different metrics of compression – volumetric ratio and maxi-
59 mum local sensitivity – include different terms, and therefore predict different levels of compression
60 for each. Moreover, we study the factors that contribute to the tightness of the bounds, or lack
61 thereof – and hence may allow representations to display trends that in practice appear to contradict
62 the theoretical predictions of the bounds. We also note that local dimensionality is a compression
63 metric of a distinct nature, and therefore does not necessarily correlate with sharpness. Taken to-
64 gether, this reveals that the impact of loss function sharpness on the neural representation is more
65 complex than a simple (and single) compression effect. These effects, despite their nuance, shed
66 light on the complex link between sharpness and generalization.

67 Throughout, we focus on the second, or final, stage of learning, which proceeds after SGD has al-
68 ready found parameters that give near-optimal performance (i.e., zero training *error*) on the training
69 data (Ma & Ying, 2021; Tishby & Zaslavsky, 2015; Ratzon et al., 2023). Here, additional learning
70 still occurs, which changes the properties of the solutions in both feature and parameter space in
71 very interesting ways.

72 Our work makes the following novel contributions:

- 73 • The paper identifies two representation space quantities that are bounded by sharpness
74 – volume compression and maximum local sensitivity (MLS) – and gives new explicit
75 formulas for these bounds that are reparametrization-invariant.
- 76 • The paper conducts empirical experiments with both VGG10 and MLP networks and finds
77 that volume compression and MLS are indeed strongly correlated with sharpness.
- 78 • The paper finds that sharpness, volume compression, and MLS are also correlated, if more
79 weakly, with test loss and hence generalization.

80 In these ways, we help reveal the interplay between key properties of trained neural networks in
81 parameter space and representation space. Specifically, we identify a sequence of equality condi-
82 tions for the bounds that link the volume and MLS of the neural representations to the sharpness in
83 parameter space. These conditions are helpful in explaining why there are the mixed results on the
84 relationship between sharpness and generalization in the literature, by looking through the additional
85 lens of the induced representations. Our findings altogether suggest that allied views into representa-
86 tion space offer a valuable dual perspective to that of parameter space landscapes for understanding
87 the effects of learning on generalization.

88 Our paper proceeds as follows. First, we review arguments of Ma & Ying (2021) that flatter minima
89 can constrain the gradient of the loss with respect to network inputs and extend the formulation to
90 the multidimensional input case (Sec. 2). Next, we prove that lower sharpness implies a lower upper
91 bound on two metrics of the compression of the representation manifold in feature space: the local
92 volume and the maximum local sensitivity (MLS) (Sec. 3.1, Sec. 3.2). We conclude our findings
93 with simulations that confirm our central theoretical results and show how they can be applied in
94 practice (Sec. 4).

95 **2 BACKGROUND AND SETUP**

96 Consider a feedforward neural network f with input data $\mathbf{x} \in \mathbb{R}^M$ and parameters $\boldsymbol{\theta}$. The output of
97 the network is:

$$\mathbf{y} = f(\mathbf{x}; \boldsymbol{\theta}), \quad (1)$$

98 with $\mathbf{y} \in \mathbb{R}^N$ ($N < M$). Consider a quadratic loss $L(\mathbf{y}, \mathbf{y}_{\text{true}}) = \frac{1}{2} \|\mathbf{y} - \mathbf{y}_{\text{true}}\|^2$ function of
99 the outputs and ground truth \mathbf{y}_{true} . In the following, we'll simply write $L(\mathbf{y})$, $L(f(\mathbf{x}, \boldsymbol{\theta}))$ or simply
100 $L(\boldsymbol{\theta})$ to highlight the dependence of the loss on the output, the network or its parameters.

101 During the last phase of learning, Ma and colleagues have recently argued that SGD appears to
102 regularize the sharpness of the loss (Li et al., 2022) (see also (Wu et al., 2018; Jastrzebski et al.,
103 2018; Xie et al., 2021; Zhu et al., 2019)). This is to say that the dynamics of SGD lead network
104 parameters to minima where the local loss landscape is flatter or wider. This is best captured by the
105 sharpness, measured by the sum of the eigenvalues of the Hessian:

$$S(\boldsymbol{\theta}) = \text{Tr}(H), \quad (2)$$

106 with $H = \nabla^2 L(\boldsymbol{\theta})$ being the Hessian. A solution with low sharpness is a flatter solution. Following
107 (Ma & Ying, 2021; Ratzon et al., 2023), we define $\boldsymbol{\theta}^*$ to be an “exact interpolation solution” on the
108 zero training loss manifold in the parameter space (the zero loss manifold in what follows), where
109 $f(\mathbf{x}_i, \boldsymbol{\theta}^*) = \mathbf{y}_i$ for all i 's (with $i \in \{1..n\}$ indexing the training set) and $L(\boldsymbol{\theta}^*) = 0$. On the zero
110 loss manifold, in particular, we have

$$S(\boldsymbol{\theta}^*) = \frac{1}{n} \sum_{i=1}^n \|\nabla_{\boldsymbol{\theta}} f(\mathbf{x}_i, \boldsymbol{\theta}^*)\|_F^2 \quad (3)$$

111 where $\|\cdot\|_F$ is the Frobenius norm. We give the proof of this equality in Appendix A. In practice,
112 the parameter $\boldsymbol{\theta}$ will never reach an exact interpolation solution due to the gradient noise of SGD,
113 however, Eq. (3) is a good enough approximation of the sharpness as long as we find an approximate
114 interpolation solution (Lemma. A.1).

115 In order to see why minimizing the sharpness of the solution leads to more compressed representa-
116 tions, we need to move from parameter space to input space. To do so we review the argument of Ma
117 & Ying (2021) that relates variations in input data \mathbf{x} and input weights. Let \mathbf{W} be the input weights
118 (the parameters of the first linear layer) of the network, and $\bar{\boldsymbol{\theta}}$ the rest of the parameters. Following
119 (Ma & Ying, 2021), as the weights \mathbf{W} multiply the inputs \mathbf{x} we have the following identities:

$$\begin{aligned} \|\nabla_{\mathbf{W}} f(\mathbf{W}\mathbf{x}; \bar{\boldsymbol{\theta}})\|_F &= \sqrt{\sum_{i,j,k} J_{jk}^2 x_i^2} = \|J\|_F \|\mathbf{x}\|_2 \geq \|J\|_2 \|\mathbf{x}\|_2 \\ \nabla_{\mathbf{x}} f(\mathbf{W}\mathbf{x}; \bar{\boldsymbol{\theta}}) &= \mathbf{W}^T J, \end{aligned} \quad (4)$$

120 where $J = \frac{\partial f(\mathbf{W}\mathbf{x}; \bar{\boldsymbol{\theta}})}{\partial(\mathbf{W}\mathbf{x})}$ is a complex expression as computed in, e.g., backpropagation. From Eq. (4)
121 and the sub-multiplicative property of the Frobenius norm and the matrix 2-norm, we have:

$$\begin{aligned} \|\nabla_{\mathbf{x}} f(\mathbf{W}\mathbf{x}; \bar{\boldsymbol{\theta}})\|_F &\leq \frac{\|\mathbf{W}\|_F}{\|\mathbf{x}\|_2} \|\nabla_{\mathbf{W}} f(\mathbf{W}\mathbf{x}; \bar{\boldsymbol{\theta}})\|_F, \\ \|\nabla_{\mathbf{x}} f(\mathbf{W}\mathbf{x}; \bar{\boldsymbol{\theta}})\|_2 &\leq \frac{\|\mathbf{W}\|_2}{\|\mathbf{x}\|_2} \|\nabla_{\mathbf{W}} f(\mathbf{W}\mathbf{x}; \bar{\boldsymbol{\theta}})\|_F. \end{aligned} \quad (5)$$

122 If the norms $\|\mathbf{W}\|_F$ or $\|\mathbf{W}\|_2$ and $\|\mathbf{x}\|_2$ are not excessively large or small respectively, these bounds
123 control the gradient with respect to inputs via the gradient with respect to weights. This in turn
124 reveals the impact of flatness in the loss function:

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \|\nabla_{\mathbf{x}} f(\mathbf{x}_i, \boldsymbol{\theta}^*)\|_F^k &\leq \frac{\|\mathbf{W}\|_F^k}{\min_i \|\mathbf{x}_i\|_2^k} \frac{1}{n} \sum_{i=1}^n \|\nabla_{\mathbf{W}} f(\mathbf{x}_i, \boldsymbol{\theta}^*)\|_F^k \\ &\leq \frac{\|\mathbf{W}\|_F^k}{\min_i \|\mathbf{x}_i\|_2^k} \frac{1}{n} \sum_{i=1}^n \|\nabla_{\boldsymbol{\theta}} f(\mathbf{x}_i, \boldsymbol{\theta}^*)\|_F^k. \end{aligned} \quad (6)$$

125 We define $G := \frac{1}{n} \sum_{i=1}^n \|\nabla_{\mathbf{x}} f(\mathbf{x}_i, \boldsymbol{\theta}^*)\|_F^2$ when $k = 2$. Similarly,

$$\frac{1}{n} \sum_{i=1}^n \|\nabla_{\mathbf{x}} f(\mathbf{x}_i, \boldsymbol{\theta}^*)\|_2^k \leq \frac{\|\mathbf{W}\|_2^k}{\min_i \|\mathbf{x}_i\|_2^k} \frac{1}{n} \sum_{i=1}^n \|\nabla_{\boldsymbol{\theta}} f(\mathbf{x}_i, \boldsymbol{\theta}^*)\|_F^k. \quad (7)$$

126 Thus, in (Ma & Ying, 2021), the effect of input perturbations is constrained by the sharpness of the
127 loss function. The flatter the minimum of the loss, the lower the effect of input space perturbations
128 on the network function $f(\mathbf{x}, \boldsymbol{\theta}^*)$ as determined by gradients.

129 3 FROM ROBUSTNESS TO INPUTS TO COMPRESSION OF REPRESENTATIONS

130 We now further analyze variations in the input and how they propagate through the network to shape
131 representations of sets of inputs. Although we only study the representations of the output of the
132 network here, our results apply to representations of any middle layer, through defining f to be
133 the transformation from input to the middle layer of interest. Overall, we focus on 3 key metrics
134 of network representations: local dimensionality, volumetric ratio, and maximum local sensitivity.
135 These quantities enable us to establish and evaluate the influence of input variations and, in turn,
136 sharpness on neural representation properties.

137 3.1 WHY SHARPNESS BOUNDS LOCAL VOLUMETRIC TRANSFORMATION IN 138 REPRESENTATION SPACE

139 Consider an input data point $\bar{\mathbf{x}}$ drawn from the training set: $\bar{\mathbf{x}} = \mathbf{x}_i$ for a specific $i \in \{1..n\}$. Let
140 the set of all possible perturbations around $\bar{\mathbf{x}}$ in input space be the ball $\mathcal{B}(\bar{\mathbf{x}})_\alpha \sim \mathcal{N}(\bar{\mathbf{x}}, \alpha \mathcal{I})$, where
141 α depends on the perturbation’s covariance, given as $C_{\mathcal{B}(\bar{\mathbf{x}})} = \alpha \mathcal{I}$, with \mathcal{I} as the identity matrix.
142 We’ll explore the network’s representation of inputs by measuring the expansion or contraction of
143 the ball $\mathcal{B}(\bar{\mathbf{x}})_\alpha$ as it propagates through the network. We first propagate the ball through the network
144 transforming each point \mathbf{x} into its image $f(\mathbf{x})$. Following a Taylor expansion for points within
145 $\mathcal{B}(\bar{\mathbf{x}})_\alpha$ as $\alpha \rightarrow 0$ we have:

$$f(\mathbf{x}) = f(\bar{\mathbf{x}}) + \nabla_{\mathbf{x}}(f(\bar{\mathbf{x}}, \boldsymbol{\theta}^*))(\mathbf{x} - \bar{\mathbf{x}}). \quad (8)$$

146 We can express the limit of the covariance matrix $C_{f(\mathcal{B}(\bar{\mathbf{x}}))}$ of the output $f(\mathbf{x})$ as

$$C_f^{\text{lim}} := \lim_{\alpha \rightarrow 0} C_{f(\mathcal{B}(\bar{\mathbf{x}})_\alpha)} = \alpha \nabla_{\mathbf{x}} f(\bar{\mathbf{x}}, \boldsymbol{\theta}^*) \nabla_{\bar{\mathbf{x}}}^T f(\bar{\mathbf{x}}, \boldsymbol{\theta}^*), \quad (9)$$

147 Our covariance expressions capture the distribution of points in $\mathcal{B}(\bar{\mathbf{x}})_\alpha$ as they go through the net-
148 work $f(\bar{\mathbf{x}}, \boldsymbol{\theta}^*)$.

149 Now we quantify how a network compresses its input volumes via the local volumetric ratio, be-
150 tween an hypercube of side length h at \mathbf{x} and its image under transformation f :

$$\begin{aligned} d \text{Vol}^{\text{ratio}}|_{f(\mathbf{x}, \boldsymbol{\theta}^*)} &= \lim_{h \rightarrow 0} \frac{\text{Vol}(f(\mathbf{x}, \boldsymbol{\theta}^*))}{\text{Vol}(\mathbf{x})} \\ &= \sqrt{\det(\nabla_{\mathbf{x}} f^T \nabla_{\mathbf{x}} f)} \end{aligned} \quad (10)$$

151 which is equal to the square root of the product of all positive eigenvalues of C_f^{lim} . Exploiting the
152 bound on the gradients derived earlier in Eq. (5), we derive a similar bound for the volumetric ratio:

$$\begin{aligned} d \text{Vol}^{\text{ratio}}|_{f(\mathbf{x}, \boldsymbol{\theta}^*)} &\leq \left(\frac{\text{Tr} \nabla_{\mathbf{x}} f^T \nabla_{\mathbf{x}} f}{N} \right)^{N/2} \\ &= N^{-N/2} \|\nabla_{\mathbf{x}} f(\mathbf{x}, \boldsymbol{\theta}^*)\|_F^N \end{aligned} \quad (11)$$

153 where the first line uses the inequality between arithmetic and geometric means and the second the
154 definition of the Frobenius norm. Introducing the averaged volumetric ratio across all input points
155 $dV^{\text{ratio}}(\boldsymbol{\theta}^*) = \frac{1}{n} \sum_{i=1}^n d \text{Vol}^{\text{ratio}}|_{f(\mathbf{x}_i, \boldsymbol{\theta}^*)}$, we obtain:

$$dV^{\text{ratio}}(\boldsymbol{\theta}^*) \leq \frac{N^{-N/2}}{n} \sum_{i=1}^n \|\nabla_{\mathbf{x}} f(\mathbf{x}_i, \boldsymbol{\theta}^*)\|_F^N \leq \frac{n^{\max(N/2-1, 0)} \|\mathbf{W}\|_F^N}{\min_i \|\mathbf{x}_i\|_2^N} \left(\frac{S(\boldsymbol{\theta}^*)}{N} \right)^{N/2}. \quad (12)$$

for all $N \geq 1$. A detailed derivation of the above inequality is given in Appendix B. Eq. (12) implies that flattened minima of the loss function in parameter space contribute to the compression of the data’s representation manifold. Our analysis demonstrates that these two phenomena are linked by the robustness properties of the network to input perturbations.

3.2 MAXIMUM LOCAL SENSITIVITY AS AN ALLIED METRIC TO TRACK NEURAL REPRESENTATION GEOMETRY

We observe that the equality condition in the first line of Eq. (11) rarely holds in practice, since to achieve equality, we need all singular values of the Jacobian matrix $\nabla_{\mathbf{x}} f$ to be identical. Our experiments in Sec. 4 show that the local dimensionality decreases rapidly with training onset, indicating that $\nabla_{\mathbf{x}} f^T \nabla_{\mathbf{x}} f$ has a non-uniform eigenspectrum. Moreover, the volume will decrease rapidly as the smallest eigenvalue vanishes. Thus, although sharpness upper bounds the volumetric ratio, it does not correlate well with it, nor does volumetric ratio give an accurate estimate of sharpness. Fortunately, considering only the maximum eigenvalue instead of the product alleviates this problem (recall that $\det(\nabla_{\mathbf{x}} f^T \nabla_{\mathbf{x}} f)$ in the definition Eq. (10) or volumetric ratio is the product of all eigenvalues): we define the maximum local sensitivity (MLS) to be the largest singular value of $\nabla_{\mathbf{x}} f$. The MLS is equivalently the matrix 2-norm of $\nabla_{\mathbf{x}} f$. Intuitively, it is the largest possible local change of $f(\mathbf{x})$ when the norm of the perturbation to \mathbf{x} is regularized. We denote the sample mean of MLS as $\overline{\text{MLS}}$. Given this definition, we obtain a bound of MLS using the Frobenius norm of the first linear layer, the quadratic mean of the input norm, and the sharpness.

$$\overline{\text{MLS}} = \frac{1}{n} \sum_{i=1}^n \|\nabla_{\mathbf{x}} f(\mathbf{x}_i, \theta^*)\|_2 \leq \|\mathbf{W}\|_2 \sqrt{\frac{1}{n} \sum_{i=1}^n \frac{1}{\|\mathbf{x}_i\|_2^2} S(\theta^*)^{1/2}}. \quad (13)$$

The derivation of the above bound is included in Appendix C, where we use Cauchy-Swartz inequality to tighten the bound in Eq. (7). As an alternative measure of compressed representations, we empirically show in Appendix D.2 that MLS has higher correlation with sharpness and test loss than the other two measures we consider in the feature space. We include more analysis of the tightness of this bound in Appendix D and discuss its connection to other works therein.

3.3 LOCAL DIMENSIONALITY IS TIED TO, BUT NOT BOUNDED BY, SHARPNESS

Now we introduce a local measure of dimensionality based on this covariance, the local Participation Ratio, given by:

$$D_{\text{PR}}(f(\bar{\mathbf{x}})) = \lim_{\alpha \rightarrow 0} \frac{\text{Tr}[C_f(\mathcal{B}(\bar{\mathbf{x}}))]^2}{\text{Tr}[(C_f(\mathcal{B}(\bar{\mathbf{x}}))]^2)} = \frac{\text{Tr}[C_f^{\text{lim}}]^2}{\text{Tr}[(C_f^{\text{lim}})^2]} \quad (14)$$

(cf. (Gao et al., 2017; Litwin-Kumar et al., 2017; Recanatesi et al., 2022)). This quantity can be averaged across a set of samples: $D_{\text{PR}}(\theta^*) = \frac{1}{n} \sum_{i=1}^n D_{\text{PR}}(f(\mathbf{x}_i))$. This quantity in some sense represents the sparseness of the eigenvalues of C_f^{lim} : if we let λ be all the eigenvalues of C_f^{lim} , then the local dimensionality can be written as $D_{\text{PR}} = (\|\lambda\|_1 / \|\lambda\|_2)^2$, which attains its maximum value when all eigenvalues are equal to each other, and its minimum when all but one eigenvalue is non-zero. Note that the quantity retains the same value when λ is arbitrarily scaled, therefore it is hard to find a relationship between local dimensionality and $\|\nabla_{\mathbf{x}} f(\mathbf{x}, \theta^*)\|_F^2$, which is basically $\|\lambda\|_1$.

4 EXPERIMENTS

4.1 SHARPNESS AND COMPRESSION: VERIFYING THE THEORY

The theoretical results derived above show that during the later phase of training – the interpolation phase – measures of compression of the network’s representation is upper bounded by a function of the sharpness of the loss function in parameter space. This links sharpness and compression of representation: the flatter is the loss landscape, the lower is the upper bound on the representation’s compression metrics.

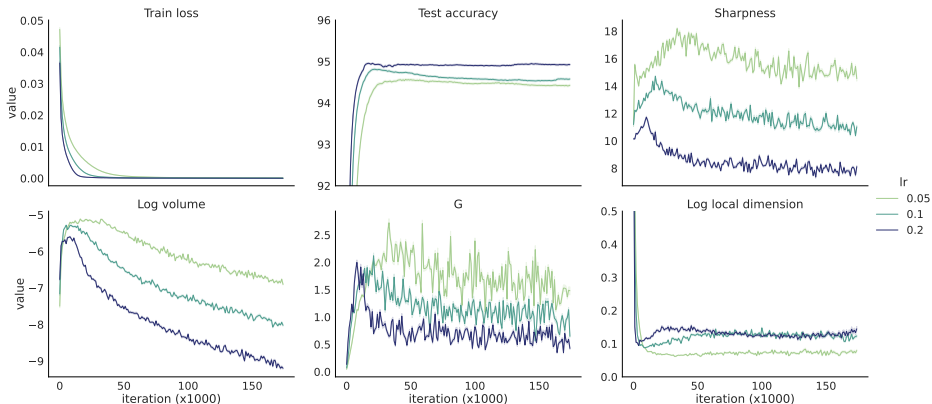


Figure 1: Trends in key variables across SGD training of the VGG10 network with fixed batch size (equal to 20) and varying learning rates (0.05, 0.1 and 0.2). After the loss is minimized (so that an approximate interpolation solution is found) sharpness and volumes decrease together. Moreover, higher learning rates lead to lower sharpness and hence stronger compression. From left to right: **train loss**, test accuracy, sharpness (square root of Eq. (3)), log volumetric ratio (Eq. (10)), left-hand side of Eq. (6) with $k = 2$ (axes titled G), and local dimensionality of the network output (Eq. (14)).

198 It remains to test in practice, however, whether these bounds are sufficiently tight so that a clear
 199 relationship between sharpness and representation collapse appears. As one such test, we ran the
 200 following experiment. We trained a network (Simonyan & Zisserman, 2015) to classify images from
 201 the CIFAR-10 dataset, and calculated the sharpness (Eq. (2)), the log volumetric ratio (Eq. (10)) and
 202 the left-hand side of Eq. (6) (the gradient with respect to the inputs, a quantity we term G in the
 203 figures below) during the training phase (Fig 1 and 2). We trained the network (VGG10) using SGD
 204 on images from 2 classes (out of 10) so that convergence to the interpolation regime, i.e. zero error,
 205 was faster. We explored the influence of two specific parameters that have a substantial effect on
 206 the network’s training: learning rate and batch size. For each pair of learning rate and batch size
 207 parameters, we computed all quantities at hand across 100 input samples and five different random
 208 initializations for network weights.

209 In the first set of experiments, we studied the link between a decrease in sharpness during the latter
 210 phases of training and volume compression (Fig. 1). We noticed that when the network reaches the
 211 interpolation regime, and the sharpness decreases, so does the volume. The quantity G similarly
 212 decreases. All these results were consistent across multiple learning rates for a fixed batch size (of
 213 20): specifically, for learning rates that gave lower values of sharpness, volume was lower as well.

214 We then repeated the experiments while keeping the learning rate fixed ($lr=0.1$) and varying the
 215 batch size. The same broadly consistent trends emerged linking a decrease in the sharpness to a
 216 compression in the representation volume (Fig. 2). However, we also find that while sharpness stops
 217 decreasing after about iteration $50 \cdot 10^3$ for batch size 32, the volume keeps decreasing as learning
 218 proceeds. This suggests that there may be other mechanisms at play, beyond sharpness, in driving
 219 the compression of volumes.

220 We repeat the experiments with an MLP trained on the FashionMNIST dataset (Fig. E.8 and
 221 Fig. E.7). Although the sharpness does not noticeably decrease at the end of the training, the sharp-
 222 ness has the same trend as G , which is consistent with our bound. The volume keeps decreasing after
 223 the sharpness plateaus, but it is also decreasing at a much slower rate, again matching our theory
 224 while suggesting that an additional factor is also involved in its decrease.

225 4.2 SHARPNESS AND COMPRESSION ON TEST SET DATA

226 Even though Eq. (3) is exact for interpolation solutions only (i.e., those with zero loss), we found
 227 that the test loss is small enough (Fig. 3) so that it should be a good approximation for test data as
 228 well. Therefore we analyzed our simulations to study trends in sharpness and volume for these held-
 229 out test data as well (Fig. 3). We discovered that this sharpness increased rather than diminished as

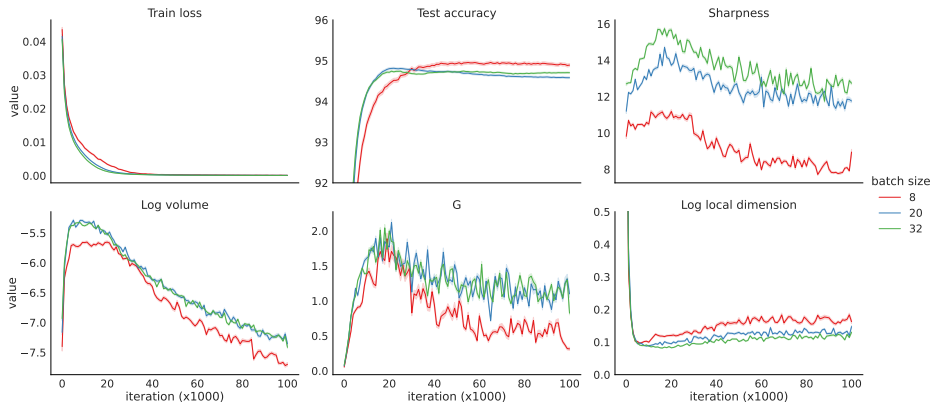


Figure 2: Trends in key variables across SGD training of the VGG10 network with fixed learning rate size (equal to 0.1) and varying batch size (8, 20, and 32). After the loss is minimized (so that an interpolation solution is found) sharpness and volumes decrease together. Moreover, lower batch sizes lead to lower sharpness and hence stronger compression. From left to right in row-wise order: train loss, test accuracy, sharpness (square root of Eq. (3)), log volumetric ratio (Eq. (10)), left-hand side of Eq. (6) with $k = 2$ (axes titled G), and local dimensionality of the network output (Eq. (14)).

230 a result of training. We hypothesized that sharpness could correlate with the difficulty of classifying
 231 testing points. This was supported by the fact that the sharpness of misclassified test data was even
 232 greater than that of all test data. Again we see that G has the same trend as the sharpness. Despite
 233 this increase in sharpness, the volume followed the same pattern as the training set. This suggests
 234 that compression in representation space is a robust phenomenon that can be driven by additional
 235 phenomena beyond sharpness. Nevertheless, the compression still is weaker for misclassified test
 236 samples that have higher sharpness than other test samples. Overall, these results emphasize an
 237 interesting distinction between how sharpness evolves for training vs. test data.

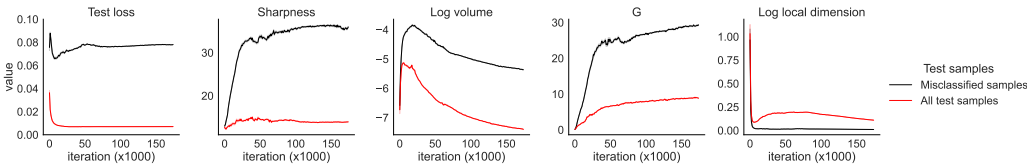


Figure 3: Trends in key variables across SGD training of the VGG10 network with fixed learning rate (equal to 0.1) and batch size (equal to 20) for samples of the test set. After the loss is minimized, we compute sharpness and volume on the test set. Moreover, the same quantities are computed separately over the entire test set or only on samples that are misclassified. In order from left to right in row-wise order: train loss, test loss, sharpness (Eq. (2)), log volumetric ratio (Eq. (10)), left-hand side of Eq. (6) with $k = 2$ (axes titled G), and local dimensionality of the network output (Eq. (14)).

238 4.3 SHARPNESS AND LOCAL DIMENSIONALITY

239 Lastly, we analyze the representation’s local dimensionality in a manner analogous to the analysis
 240 of volume and MLS. A priori, it is ambiguous whether the dimensionality of the data representation
 241 should increase or decrease as the volume is compressed. For instance, the volume could decrease
 242 while maintaining its overall form and symmetry, thus preserving its dimensionality. Alternatively,
 243 one or more of the directions in the relevant tangent space could be selectively compressed, leading
 244 to an overall reduction in dimensionality.

245 Figures 1 and 2 show our experiments computing the local dimensionality over the course of learn-
 246 ing. Here, we find that the local dimensionality of the representation decreases as the loss decreases
 247 to near 0, which is consistent with the viewpoint that the network compresses representations in

248 feature space as much as possible, retaining only the directions that code for task-relevant features
 249 (Berner et al., 2020; Cohen et al., 2020). However, the local dimensionality exhibits unpredictable
 250 behavior that cannot be explained by the sharpness once the network is near the zero-loss manifold
 251 and training continues. This discrepancy is consistent with the bounds established by our theory,
 252 which only bound the numerator of Eq. (14). It is also consistent with the property of local dimen-
 253 sionality that we described in Sec. 3.3 overall: it encodes the sparseness of the eigenvalues but it
 254 does not encode the magnitude of them. This shows how local dimensionality is a distinct quality
 255 of network representations compared with volume, and is driven by mechanisms that differ from
 256 sharpness alone. We emphasize that the dimensionality we study here is a local measure, on the
 257 finest scale around a point on the “global” manifold of unit activities; dimension on larger scales
 258 (i.e., across categories or large sets of task inputs (Farrell et al., 2022; Gao et al., 2017)) may show
 259 different trends.

260 5 CONCLUSION

261 This work presents a dual perspective, uniting views in both parameter and in feature space, of
 262 several key properties of trained neural networks that have been linked to their ability to generalize.
 263 We identify two representation space quantities that are bounded by sharpness – volume compression
 264 and maximum local sensitivity – and give new explicit formulas for these bounds. We conduct
 265 experiments with both VGG10 and MLP networks and find that the predictions of these bounds are
 266 born out for these networks, illustrating how MLS in particular is strongly correlated with sharpness.
 267 We also establish that sharpness, volume compression, and MLS are correlated, if more weakly, with
 268 test loss and hence generalization. Overall, we establish explicit links between sharpness properties
 269 in parameter spaces and compression and robustness properties in representation space.

270 By demonstrating both how these links can be tight, and how and when they may also become loose,
 271 we show that taking this dual perspective can bring more clarity to the often confusing question of
 272 what quantifies how well a network will generalize in practice. Indeed, many works, as reviewed in
 273 the introduction, have demonstrated how sharpness in parameter space can lead to generalization,
 274 but recent studies have established contradictory results. We show how looking at quantities not
 275 only in the parameter space (sharpness), but also in the feature space (compression, maximum local
 276 sensitivity, etc.) may help explain the wide range of results.

277 This said, we view our study as a starting point to open doors between two often-distinct perspectives
 278 on generalization in neural networks. Additional theoretical and experimental research is warranted
 279 to systematically investigate the implications of our findings, with a key area being further learning
 280 problems, such as predictive learning, beyond the classification tasks studied here. Nevertheless, we
 281 are confident that highly interesting and clarifying findings lie ahead at the interface between the
 282 parameter and representation space quantities explored here.

283 REFERENCES

- 284 Maksym Andriushchenko, Francesco Croce, Maximilian Müller, Matthias Hein, and Nicolas Flam-
 285 marion. A modern look at the relationship between sharpness and generalization. *arXiv preprint*
 286 *arXiv:2302.07011*, 2023.
- 287 Alessio Ansuini, Alessandro Laio, Jakob H. Macke, and Davide Zoccolan. Intrinsic dimension
 288 of data representations in deep neural networks. *Advances in Neural Information Processing*
 289 *Systems*, 32, 2019.
- 290 Ido Ben-Shaul and Shai Dekel. Nearest class-center simplification through intermediate layers. In
 291 *Topological, Algebraic and Geometric Learning Workshops 2022*, pp. 37–47. PMLR, 2022.
- 292 Ido Ben-Shaul, Ravid Shwartz-Ziv, Tomer Galanti, Shai Dekel, and Yann LeCun. Reverse engineer-
 293 ing self-supervised learning. *arXiv preprint arXiv:2305.15614*, 2023.
- 294 Julius Berner, Philipp Grohs, and Arnulf Jentzen. Analysis of the generalization error: Empirical
 295 risk minimization over deep artificial neural networks overcomes the curse of dimensionality in
 296 the numerical approximation of black–scholes partial differential equations. *SIAM Journal on*
 297 *Mathematics of Data Science*, 2(3):631–657, 2020. Publisher: SIAM.

- 298 Guy Blanc, Neha Gupta, Gregory Valiant, and Paul Valiant. Implicit regularization for deep neural
299 networks driven by an ornstein-uhlenbeck like process, 2020. URL [http://arxiv.org/
300 abs/1904.09080](http://arxiv.org/abs/1904.09080).
- 301 Uri Cohen, SueYeon Chung, Daniel D. Lee, and Haim Sompolinsky. Separability and geometry of
302 object manifolds in deep neural networks. *Nature communications*, 11(1):746, 2020. Publisher:
303 Nature Publishing Group UK London.
- 304 Lijun Ding, Dmitriy Drusvyatskiy, Maryam Fazel, and Zaid Harchaoui. Flat minima generalize for
305 low-rank matrix recovery, 2023.
- 306 Laurent Dinh, Razvan Pascanu, Samy Bengio, and Yoshua Bengio. Sharp minima can generalize
307 for deep nets. In *International Conference on Machine Learning*, pp. 1019–1028. PMLR, 2017.
- 308 Gamaleldin Elsayed, Dilip Krishnan, Hossein Mobahi, Kevin Regan, and Samy Bengio. Large
309 margin deep networks for classification. *Advances in neural information processing systems*, 31,
310 2018.
- 311 Matthew Farrell, Stefano Recanatesi, Timothy Moore, Guillaume Lajoie, and Eric Shea-Brown.
312 Gradient-based learning drives robust representations in recurrent neural networks by balancing
313 compression and expansion. *Nature Machine Intelligence*, 4(6):564–573, 2022. Publisher: Nature
314 Publishing Group UK London.
- 315 Tomer Galanti, Liane Galanti, and Ido Ben-Shaul. Comparative generalization bounds for deep
316 neural networks. *Transactions on Machine Learning Research*, 2023.
- 317 Peiran Gao, Eric Trautmann, Byron Yu, Gopal Santhanam, Stephen Ryu, Krishna Shenoy, and Surya
318 Ganguli. A theory of multineuronal dimensionality, dynamics and measurement. *BioRxiv*, pp.
319 214262, 2017.
- 320 Khashayar Gatmiry, Zhiyuan Li, Ching-Yao Chuang, Sashank Reddi, Tengyu Ma, and Stefanie
321 Jegelka. The inductive bias of flatness regularization for deep matrix factorization, 2023.
- 322 Mario Geiger, Leonardo Petrini, and Matthieu Wyart. Landscape and training regimes in deep
323 learning. *Physics Reports*, 924:1–18, 2021. ISSN 0370-1573. doi: 10.1016/j.physrep.
324 2021.04.001. URL [https://www.sciencedirect.com/science/article/pii/
325 S0370157321001290](https://www.sciencedirect.com/science/article/pii/S0370157321001290).
- 326 Stanisław Jastrzebski, Zachary Kenton, Devansh Arpit, Nicolas Ballas, Asja Fischer, Yoshua Ben-
327 gio, and Amos Storkey. Three Factors Influencing Minima in SGD, September 2018. URL
328 <http://arxiv.org/abs/1711.04623>. arXiv:1711.04623 [cs, stat].
- 329 Vignesh Kothapalli, Ebrahim Rasromani, and Vasudev Awatramani. Neural collapse: A review on
330 modelling principles and generalization. *arXiv preprint arXiv:2206.04041*, 2022.
- 331 Zhiyuan Li, Tianhao Wang, and Sanjeev Arora. What happens after SGD reaches zero loss? –a
332 mathematical framework, 2022. URL <http://arxiv.org/abs/2110.06914>.
- 333 Ashok Litwin-Kumar, Kameron Decker Harris, Richard Axel, Haim Sompolinsky, and LF Abbott.
334 Optimal degrees of synaptic connectivity. *Neuron*, 93(5):1153–1164, 2017.
- 335 Chao Ma and Lexing Ying. On linear stability of SGD and input-smoothness of neural networks,
336 2021. URL <http://arxiv.org/abs/2105.13462>.
- 337 Mor Shpigel Nacson, Kavya Ravichandran, Nathan Srebro, and Daniel Soudry. Implicit bias of the
338 step size in linear diagonal neural networks. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song,
339 Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International
340 Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*,
341 pp. 16270–16295. PMLR, 17–23 Jul 2022. URL [https://proceedings.mlr.press/
342 v162/nacson22a.html](https://proceedings.mlr.press/v162/nacson22a.html).
- 343 Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal
344 phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40):
345 24652–24663, 2020.

- 346 Akshay Rangamani, Marius Lindegaard, Tomer Galanti, and Tomaso A Poggio. Feature learning
347 in deep classifiers through intermediate neural collapse. In *International Conference on Machine*
348 *Learning*, pp. 28729–28745. PMLR, 2023.
- 349 Aviv Ratzon, Dori Derdikman, and Omri Barak. Representational drift as a result of implicit reg-
350 ularization, 2023. URL [https://www.biorxiv.org/content/10.1101/2023.05.](https://www.biorxiv.org/content/10.1101/2023.05.04.539512v3)
351 [04.539512v3](https://www.biorxiv.org/content/10.1101/2023.05.04.539512v3). Pages: 2023.05.04.539512 Section: New Results.
- 352 Stefano Recanatesi, Matthew Farrell, Madhu Advani, Timothy Moore, Guillaume Lajoie, and Eric
353 Shea-Brown. Dimensionality compression and expansion in deep neural networks. *arXiv preprint*
354 *arXiv:1906.00443*, 2019.
- 355 Stefano Recanatesi, Serena Bradde, Vijay Balasubramanian, Nicholas A Steinmetz, and Eric Shea-
356 Brown. A scale-dependent measure of system dimensionality. *Patterns*, 3(8), 2022.
- 357 Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image
358 recognition, 2015.
- 359 Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle, 2015.
360 URL <http://arxiv.org/abs/1503.02406>.
- 361 Kaiyue Wen, Zhiyuan Li, and Tengyu Ma. Sharpness minimization algorithms do not only minimize
362 sharpness to achieve better generalization, 2023.
- 363 Lei Wu, Chao Ma, and Weinan E. How SGD Selects the Global Minima in Over-parameterized
364 Learning: A Dynamical Stability Perspective. In *Advances in Neural Information Processing Sys-*
365 *tems*, volume 31. Curran Associates, Inc., 2018. URL [https://papers.nips.cc/paper_](https://papers.nips.cc/paper_files/paper/2018/hash/6651526b6fb8f29a00507de6a49ce30f-Abstract.html)
366 [files/paper/2018/hash/6651526b6fb8f29a00507de6a49ce30f-Abstract.](https://papers.nips.cc/paper_files/paper/2018/hash/6651526b6fb8f29a00507de6a49ce30f-Abstract.html)
367 [html](https://papers.nips.cc/paper_files/paper/2018/hash/6651526b6fb8f29a00507de6a49ce30f-Abstract.html).
- 368 Zeke Xie, Issei Sato, and Masashi Sugiyama. A Diffusion Theory For Deep Learning Dynamics:
369 Stochastic Gradient Descent Exponentially Favors Flat Minima, January 2021. URL <http://arxiv.org/abs/2002.03495>. arXiv:2002.03495 [cs, stat].
370
- 371 Ning Yang, Chao Tang, and Yuhai Tu. Stochastic gradient descent introduces an effective landscape-
372 dependent regularization favoring flat solutions. *Physical Review Letters*, 130(23):237101, 2023.
373 doi: 10.1103/PhysRevLett.130.237101. URL [https://link.aps.org/doi/10.1103/](https://link.aps.org/doi/10.1103/PhysRevLett.130.237101)
374 [PhysRevLett.130.237101](https://link.aps.org/doi/10.1103/PhysRevLett.130.237101). Publisher: American Physical Society.
- 375 Yaoqing Yang, Liam Hodgkinson, Ryan Theisen, Joe Zou, Joseph E Gonzalez, Kannan Ramchan-
376 dran, and Michael W Mahoney. Taxonomizing local versus global structure in neural network
377 loss landscapes. *Advances in Neural Information Processing Systems*, 34:18722–18733, 2021.
- 378 Zhanxing Zhu, Jingfeng Wu, Bing Yu, Lei Wu, and Jinwen Ma. The Anisotropic Noise in Stochastic
379 Gradient Descent: Its Behavior of Escaping from Sharp Minima and Regularization Effects, June
380 2019. URL <http://arxiv.org/abs/1803.00195>. arXiv:1803.00195 [cs, stat].
- 381 Zhihui Zhu, Tianyu Ding, Jinxin Zhou, Xiao Li, Chong You, Jeremias Sulam, and Qing Qu. A ge-
382 ometric analysis of neural collapse with unconstrained features. *Advances in Neural Information*
383 *Processing Systems*, 34:29820–29834, 2021.

384 **A PROOF OF EQ. (3)**

385 **Lemma A.1.** *If θ is an approximate interpolation solution, i.e. $\|f(\mathbf{x}_i, \theta) - y_i\| < \varepsilon$ for $i \in$
 386 $\{1, 2, \dots, n\}$, and second derivatives of the network function $\|\nabla_{\theta_j^2} f(\mathbf{x}_i, \theta)\| < M$ is bounded,
 387 then*

$$S(\theta^*) = \frac{1}{n} \sum_{i=1}^n \|\nabla_{\theta} f(\mathbf{x}_i, \theta^*)\|_F^2 + O(\varepsilon) \quad (15)$$

388 *Proof.* Using basic calculus we get

$$\begin{aligned} S(\theta) &= \text{Tr}(\nabla^2 L(\theta)) \\ &= \frac{1}{2n} \sum_{i=1}^n \text{Tr}(\nabla_{\theta}^2 \|f(\mathbf{x}_i, \theta) - y_i\|^2) \\ &= \frac{1}{2n} \sum_{i=1}^n \text{Tr} \nabla_{\theta} (2(f(\mathbf{x}_i, \theta) - y_i)^T \nabla_{\theta} f(\mathbf{x}_i, \theta)) \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m \frac{\partial}{\partial \theta_j} ((f(\mathbf{x}_i, \theta) - y_i)^T \nabla_{\theta} f(\mathbf{x}_i, \theta))_j \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m \frac{\partial}{\partial \theta_j} (f(\mathbf{x}_i, \theta) - y_i)^T \nabla_{\theta_j} f(\mathbf{x}_i, \theta) \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m \|\nabla_{\theta_j} f(\mathbf{x}_i, \theta)\|_2^2 + (f(\mathbf{x}_i, \theta) - y_i)^T \nabla_{\theta_j}^2 f(\mathbf{x}_i, \theta) \\ &= \frac{1}{n} \sum_{i=1}^n \|\nabla_{\theta} f(\mathbf{x}_i, \theta)\|_F^2 + \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i, \theta) - y_i)^T \nabla_{\theta_j}^2 f(\mathbf{x}_i, \theta). \end{aligned}$$

389 Therefore

$$\left| S(\theta) - \frac{1}{n} \sum_{i=1}^n \|\nabla_{\theta} f(\mathbf{x}_i, \theta)\|_F^2 \right| < \frac{1}{n} \sum_{i=1}^n |(f(\mathbf{x}_i, \theta) - y_i)^T \nabla_{\theta_j}^2 f(\mathbf{x}_i, \theta)| < M\varepsilon = O(\varepsilon). \quad (16)$$

390 \square

391 In other words, when the network reaches zero training error and enters the interpolation phase (i.e.
 392 it classifies all training data correctly), Eq. (3) will be a good enough approximation of the sharpness
 393 because the quadratic training loss is sufficiently small.

394 **B PROOF OF EQ. (12)**

395 We first show that Eq. (6) is correct. Because of Eq. (5), we have the first inequality of Eq. (6),

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \|\nabla_{\mathbf{x}} f(\mathbf{x}_i, \theta^*)\|_F^k &\leq \|\mathbf{W}\|_F^k \frac{1}{n} \sum_{i=1}^n \frac{\|\nabla_{\mathbf{w}} f(\mathbf{x}_i, \theta^*)\|_F^k}{\|\mathbf{x}_i\|_2^k} \\ &\leq \frac{\|\mathbf{W}\|_F^k}{\min_i \|\mathbf{x}_i\|_2^k} \frac{1}{n} \sum_{i=1}^n \|\nabla_{\mathbf{w}} f(\mathbf{x}_i, \theta^*)\|_F^k. \end{aligned} \quad (17)$$

396 Since the input weights \mathbf{W} is just a part of all the weights (θ) of the network, we have
 397 $\|\nabla_{\mathbf{w}} f(\mathbf{x}_i, \theta^*)\|_F^k \leq \|\nabla_{\theta} f(\mathbf{x}_i, \theta^*)\|_F^k$. Therefore

$$\frac{\|\mathbf{W}\|_F^k}{\min_i \|\mathbf{x}_i\|_2^k} \frac{1}{n} \sum_{i=1}^n \|\nabla_{\mathbf{w}} f(\mathbf{x}_i, \theta^*)\|_F^k \leq \frac{\|\mathbf{W}\|_F^k}{\min_i \|\mathbf{x}_i\|_2^k} \frac{1}{n} \sum_{i=1}^n \|\nabla_{\theta} f(\mathbf{x}_i, \theta^*)\|_F^k. \quad (18)$$

398 To show the correctness of Eq. (12), we discuss two cases.

399 **Case 1:** $k \geq 2$

400 **Lemma B.1.** For vector \mathbf{x} , $\|\mathbf{x}\|_p \geq \|\mathbf{x}\|_q$ for $1 \leq p \leq q \leq \infty$.

401 *Proof.* First we show that for $0 < k < 1$, we have $(|a| + |b|)^k \leq |a|^k + |b|^k$. It's trivial when either
 402 a or b is 0. So W.L.O.G, we can assume that $|a| < |b|$, and divide both sides by $|b|^k$. Therefore it
 403 suffices to show that for $0 < t < 1$, $(1+t)^k < t^k + 1$. Let $f(t) = (1+t)^k - t^k - 1$, then $f(0) = 0$, and
 404 $f'(t) = k(1+t)^{k-1} - kt^{k-1}$. Because $k-1 < 0$, $1+t > 1$ and $t < 1$, $t^{k-1} > (1+t)^{k-1}$. Therefore
 405 $f'(t) < 0$ and $f(t) < 0$ for $0 < t < 1$. Combining all cases, we have $(|a| + |b|)^k \leq |a|^k + |b|^k$ for
 406 $0 < k < 1$. By induction, we have $(\sum_n |a_n|)^k \leq \sum_n |a_n|^k$.

407 Now we can prove the lemma using the conclusion above,

$$\left(\sum_n |x_n|^q\right)^{1/q} = \left(\sum_n |x_n|^q\right)^{p/q \cdot 1/p} \leq \left(\sum_n (|x_n|^q)^{p/q}\right)^{1/p} = \left(\sum_n |x_n|^p\right)^{1/p} \quad (19)$$

408 □

409 Now take the x_i in above lemma to be $\|\nabla_{\theta} f(\mathbf{x}_i, \theta^*)\|_F^2$ and let $p = 1, q = k/2$, then we get

$$\left(\sum_{i=1}^n (\|\nabla_{\theta} f(\mathbf{x}_i, \theta^*)\|_F^2)^{k/2}\right)^{2/k} \leq \sum_{i=1}^n \|\nabla_{\theta} f(\mathbf{x}_i, \theta^*)\|_F^2. \quad (20)$$

410 Therefore,

$$\begin{aligned} \frac{\|\mathbf{W}\|_F^k}{\min_i \|\mathbf{x}_i\|_2^k} \frac{1}{n} \sum_{i=1}^n \|\nabla_{\theta} f(\mathbf{x}_i, \theta^*)\|_F^k &\leq \frac{n^{k/2-1} \|\mathbf{W}\|_F^k}{\min_i \|\mathbf{x}_i\|_2^k} \left(\frac{1}{n} \sum_{i=1}^n \|\nabla_{\theta} f(\mathbf{x}_i, \theta^*)\|_F^2\right)^{k/2} \\ &= \frac{n^{k/2-1} \|\mathbf{W}\|_F^k}{\min_i \|\mathbf{x}_i\|_2^k} S(\theta^*)^{k/2} \end{aligned} \quad (21)$$

411 **Case 2:** $1 \leq k < 2$

412 **Lemma B.2.** For vector $\mathbf{x} \in \mathbb{R}^n$, $\|\mathbf{x}\|_p \leq n^{1/p-1/q} \|\mathbf{x}\|_q$ for $1 \leq p \leq q \leq \infty$.

413 *Proof.* By Hölder's inequality, we have,

$$\sum_i |x_i|^p = \sum_i |x_i|^p \cdot 1 \leq \left(\sum_i |x_i|^q\right)^{p/q} \left(\sum_i 1\right)^{1-p/q} = n^{1-p/q} \|\mathbf{x}\|_q^p \quad (22)$$

414 Taking the p -th root on both sides gives us the desired inequality. □

415 Now take the x_i in above lemma to be $\|\nabla_{\theta} f(\mathbf{x}_i, \theta^*)\|_F$ and let $p = k, q = 2$, then we get

$$\left(\sum_{i=1}^n (\|\nabla_{\theta} f_i\|_F)^k\right)^{1/k} \leq n^{1/k-1/2} \left(\sum_{i=1}^n \|\nabla_{\theta} f_i\|_F^2\right)^{1/2}. \quad (23)$$

416 Therefore,

$$\begin{aligned} \frac{\|\mathbf{W}\|_F^k}{\min_i \|\mathbf{x}_i\|_2^k} \frac{1}{n} \sum_{i=1}^n \|\nabla_{\theta} f_i\|_F^k &\leq \frac{\|\mathbf{W}\|_F^k}{\min_i \|\mathbf{x}_i\|_2^k} \frac{n^{1-k/2}}{n} \left(\sum_{i=1}^n \|\nabla_{\theta} f_i\|_F^2\right)^{k/2} \\ &= \frac{\|\mathbf{W}\|_F^k}{\min_i \|\mathbf{x}_i\|_2^k} S(\theta^*)^{k/2}. \end{aligned} \quad (24)$$

417 Combining Eq. (21) and Eq. (24), we get Eq. (12).

418 C PROOF OF EQ. (13)

419 From Eq. (5), we get

$$\overline{\text{MLS}} = \frac{1}{n} \sum_{i=1}^n \|\nabla_{\mathbf{x}} f_i\|_2 \leq \|\mathbf{W}\|_2 \frac{1}{n} \sum_{i=1}^n \frac{\|\nabla_{\mathbf{w}} f(\mathbf{x}_i, \boldsymbol{\theta}^*)\|_F}{\|\mathbf{x}_i\|_2}. \quad (25)$$

420 Now Cauchy Swartz inequality tells us that

$$\left(\sum_{i=1}^n \frac{\|\nabla_{\mathbf{w}} f_i\|_F}{\|\mathbf{x}_i\|_2} \right)^2 \leq \left(\sum_{i=1}^n \frac{1}{\|\mathbf{x}_i\|_2^2} \right) \cdot \left(\sum_{i=1}^n \|\nabla_{\mathbf{w}} f_i\|_F^2 \right). \quad (26)$$

421 Therefore

$$\begin{aligned} \overline{\text{MLS}} &\leq \|\mathbf{W}\|_2 \sqrt{\frac{1}{n} \sum_{i=1}^n \frac{1}{\|\mathbf{x}_i\|_2^2}} \cdot \sqrt{\frac{1}{n} \sum_{i=1}^n \|\nabla_{\mathbf{w}} f_i\|_F^2} \\ &\leq \|\mathbf{W}\|_2 \sqrt{\frac{1}{n} \sum_{i=1}^n \frac{1}{\|\mathbf{x}_i\|_2^2}} \cdot S(\boldsymbol{\theta}^*)^{1/2}. \end{aligned} \quad (27)$$

422 D EMPIRICAL ANALYSIS OF THE BOUND

423 D.1 TIGHTNESS OF THE BOUND

424 In this section, we mainly explore the tightness of the bound in Eq. (13) for reasons discussed in
425 Sec. 3.2. First we rewrite Eq. (13) as

$$\begin{aligned} \overline{\text{MLS}} &= \frac{1}{n} \sum_{i=1}^n \|\nabla_{\mathbf{x}} f(\mathbf{x}_i, \boldsymbol{\theta}^*)\|_2 && := A \\ &\leq \frac{\|\mathbf{W}\|_2}{n} \sum_{i=1}^n \frac{\|\nabla_{\mathbf{w}} f(\mathbf{x}_i, \boldsymbol{\theta}^*)\|_F}{\|\mathbf{x}_i\|_2} && := B \\ &\leq \|\mathbf{W}\|_2 \sqrt{\frac{1}{n} \sum_{i=1}^n \frac{1}{\|\mathbf{x}_i\|_2^2}} \sqrt{\frac{1}{n} \sum_{i=1}^n \|\nabla_{\mathbf{w}} f(\mathbf{x}_i, \boldsymbol{\theta}^*)\|_F^2} && := C \\ &\leq \|\mathbf{W}\|_2 \sqrt{\frac{1}{n} \sum_{i=1}^n \frac{1}{\|\mathbf{x}_i\|_2^2}} S(\boldsymbol{\theta}^*)^{1/2} && := D \end{aligned} \quad (28)$$

426 Thus Eq. (13) consists of 3 different steps of relaxations. We analyze them one by one:

427 1. ($A \leq B$) The equality holds when $\|W^T J\|_2 = \|W\|_2 \|J\|_2$ and $\|J\|_F = \|J\|_2$, where
428 $J = \frac{\partial f(\mathbf{W}\mathbf{x}; \boldsymbol{\theta})}{\partial(\mathbf{W}\mathbf{x})}$. The former equality requires that W and J have the same left singular
429 vectors. The latter requires J to have zero singular values except for the largest singular
430 value. Since J depends on the specific neural network architecture and training process,
431 we test the tightness of this bound empirically (Fig. D.4).

432 2. ($B \leq C$) The equality requires $\frac{\|\nabla_{\mathbf{w}} f(\mathbf{x}_i, \boldsymbol{\theta}^*)\|_F}{\|\mathbf{x}_i\|_2}$ to be the same for all i . In other words, the
433 bound is tight when $\frac{\|\nabla_{\mathbf{w}} f(\mathbf{x}_i, \boldsymbol{\theta}^*)\|_F}{\|\mathbf{x}_i\|_2}$ does not vary too much from sample to sample.

434 3. ($C \leq D$) The equality holds if the model is linear, i.e. $\boldsymbol{\theta} = \mathbf{W}$.

435 We empirically verify the tightness of the above bounds in Fig. D.4

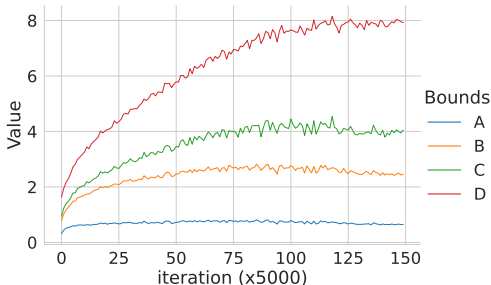


Figure D.4: **Empirical tightness of the bounds.** We empirically verify that the inequalities in Eq. (28) hold and test their tightness. The results are shown for a fully connected feedforward network trained on the FashionMNIST dataset. The quantities A, B, C, and D are defined in Eq. (28). We see that the gap between C and D is large compared to the gap between A and B or B and C. This indicates that partial sharpness $\|\nabla_{\mathbf{W}} f(\mathbf{x}_i, \boldsymbol{\theta}^*)\|_F$ (sensitivity of the loss w.r.t. only the input weights) is more indicative of the change in the maximum local sensitivity (A). Indeed, correlation analysis shows that bound C is positively correlated with MLS while bound D, perhaps surprisingly, is negatively correlated with MLS (Fig. D.6).

436 D.2 CORRELATION ANALYSIS

437 We empirically show how different metrics correlate with each other, and how these correlations
 438 can be predicted from our bounds. We train 20 VGG10 networks with different batch sizes, learning
 439 rates, and random initialization to classify images from the CIFAR-10 dataset, and plot pairwise
 440 scatter plots between 5 quantities at the end of the training: test loss, MLS, G (see Eq. (6)), log
 441 volume, sharpness and local dimensionality (Fig. D.5).

442 We find that

- 443 1. G and MLS are highly correlated and can be almost seen as the same quantity, scaled.
- 444 2. Although the bound in Eq. (12) is loose, log volume correlates well with sharpness and
 445 MLS.
- 446 3. Sharpness is positively correlated with the test loss, indicating that little reparametrization
 447 effect (Dinh et al., 2017) is happening during training, i.e. the network weights do not
 448 change too much during training. This is consistent with observations in Ma & Ying (2021).
- 449 4. MLS improves the correlation with the test loss over log volume and local dimensionality.
 450 This is consistent with the bound Eq. (13).

451 We repeat the analysis on an MLP trained on the FashionMNIST dataset, and observe the same
 452 phenomena (Fig. D.6).

453 D.3 CONNECTION TO OTHER WORKS

454 Our bound and its analysis are connected to many theoretical and experimental results. First of
 455 all, the right-hand side of Eq. (13) is related not only to the sharpness but also to the norm of the
 456 input weights. Therefore our bound takes into the effect of reparametrization, and is invariant under
 457 scaling of the input weights. This is consistent with the theoretical results in Dinh et al. (2017)
 458 which show that sharpness can be arbitrarily increased by reparametrization while the network can
 459 still generalize. Moreover, many works studied simplified linear models (Li et al., 2022; Ding et al.,
 460 2023; Nacson et al., 2022; Gatmiry et al., 2023), and showed that the flattest minima generalize well.
 461 Correspondingly, Eq. (28) shows that when the neural network is linear, the inequality between C
 462 and D becomes equality, and the flattest minima give the tightest bound on MLS. On the other hand,
 463 this also explains why sharpness does not always correlate with generalization when the network
 464 becomes more complicated (Wen et al., 2023; Andriushchenko et al., 2023): having weights that
 465 are other than the input weights makes the bound looser and more unpredictable. Experiments on

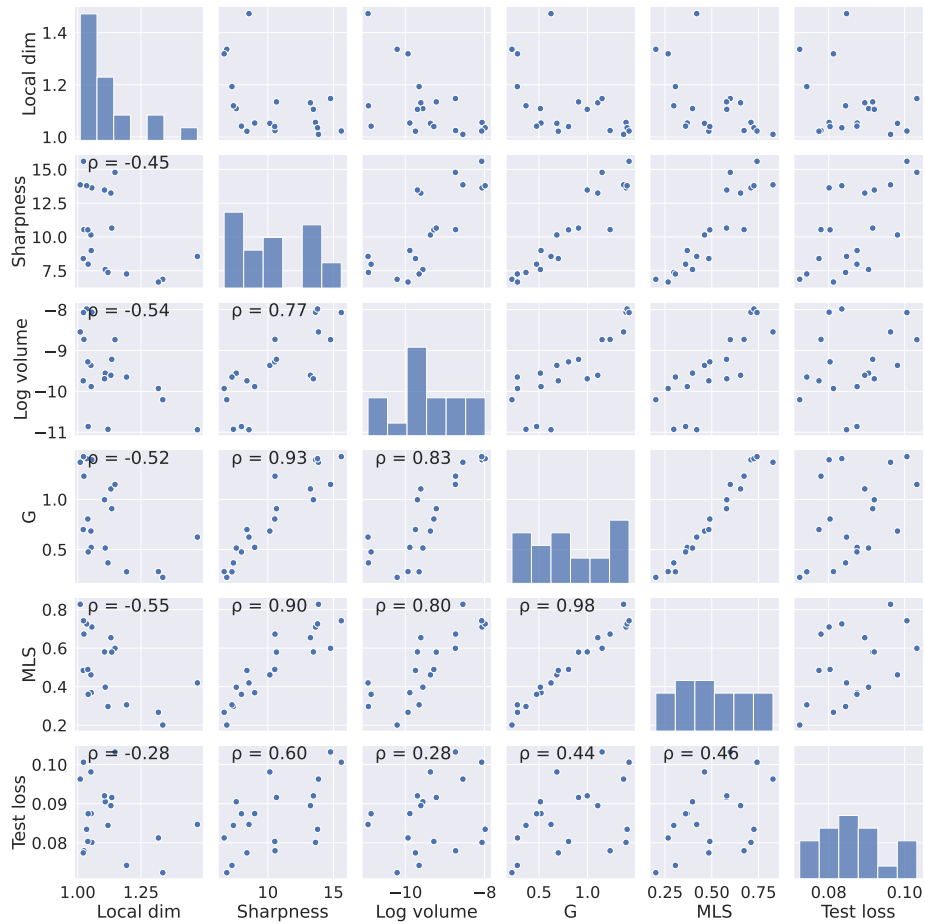


Figure D.5: Pairwise correlation among different metrics. We trained 20 different VGG10 networks using vanilla SGD with different learning rates, batch sizes, and random initializations and plot pairwise scatter plots between different quantities: local dimensionality, sharpness (square root of Eq. (3)), log volume (Eq. (10)), G (Eq. (6)), MLS (Eq. (13)) and test loss. The Pearson correlation coefficient ρ is shown in the top-left corner for each pair of quantities. See Appendix D.2 for a summary of the findings in this figure.

466 MLP show that the bound D in Eq. (28) can even be negatively correlated with MLS and test loss
 467 (Fig. D.6).

468 E ADDITIONAL EXPERIMENTS

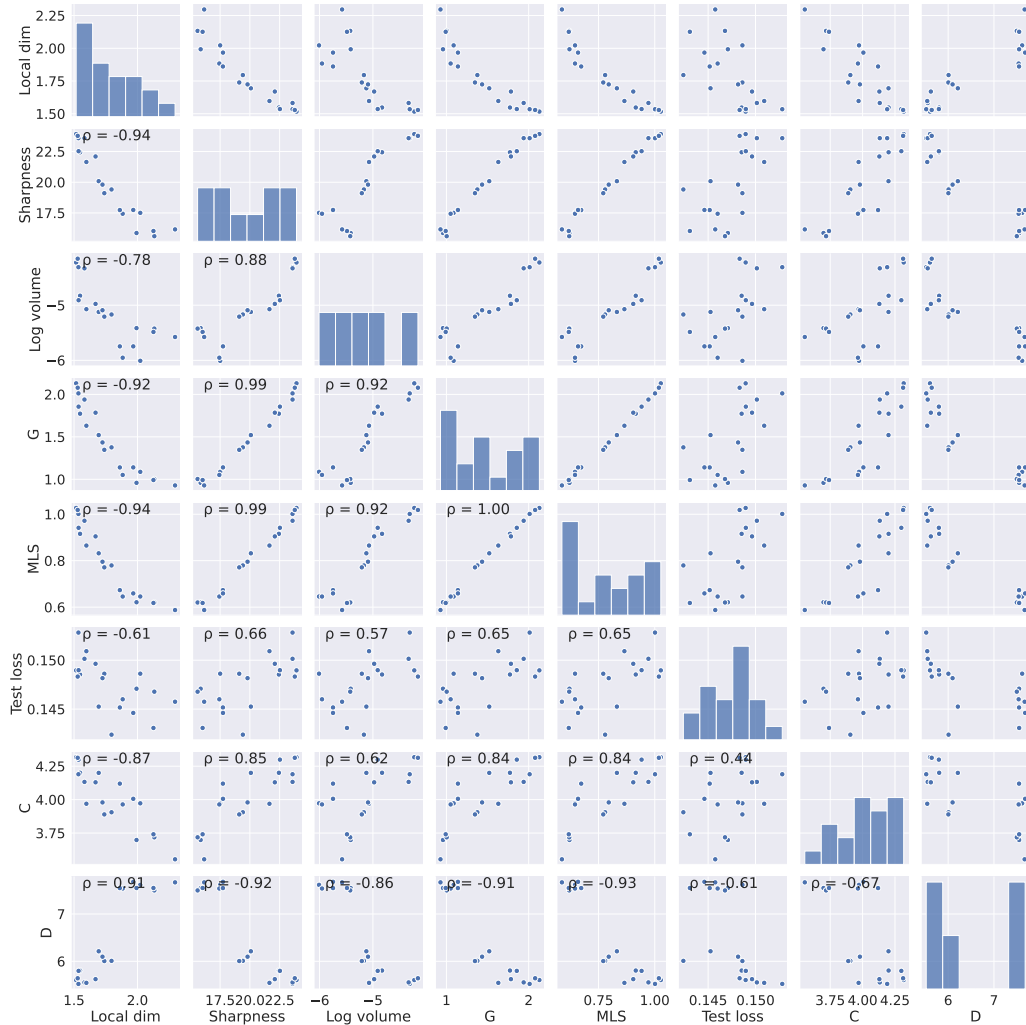


Figure D.6: Pairwise correlation among different metrics. We trained 20 different 4-layer MLPs using vanilla SGD with different learning rates, batch size, and random initializations and plot pairwise scatter plots between different quantities: local dimensionality, sharpness (square root of Eq. (3)), log volume (Eq. (10)), G (Eq. (6)), MLS (Eq. (13)), test loss and additionally bound C and D as defined in Eq. (28). The Pearson correlation coefficient ρ is shown in the top-left corner for each pair of quantities. See Appendix D.2 for a summary of the findings in this figure.

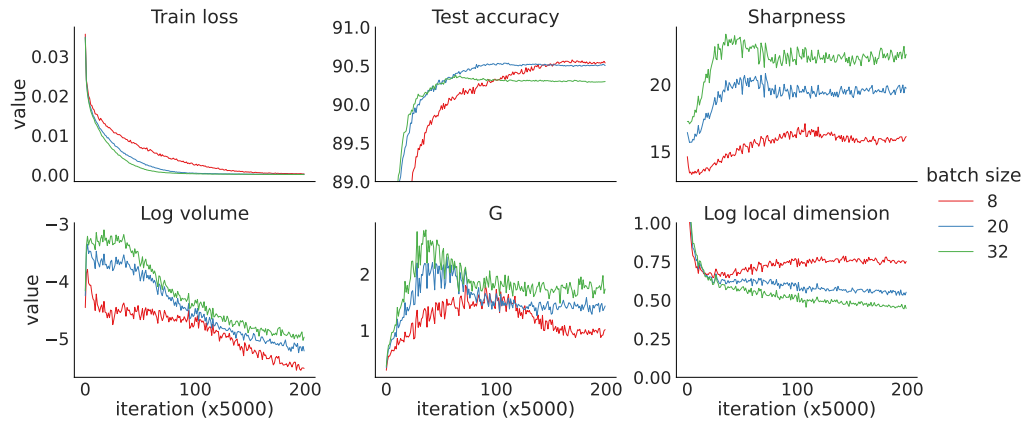


Figure E.7: Trends in key variables across SGD training of a 4-layer MLP with fixed learning rate (equal to 0.1) and varying batch size (8, 20, and 32). After minimizing the loss, lower batch sizes lead to lower sharpness and stronger compression. Moreover, G closely follows the trend of sharpness during the training. From left to right: train loss, test accuracy, sharpness (square root of Eq. (3)), log volumetric ratio (Eq. (10)), G (Eq. (6)), and local dimensionality of the network output (Eq. (14)).

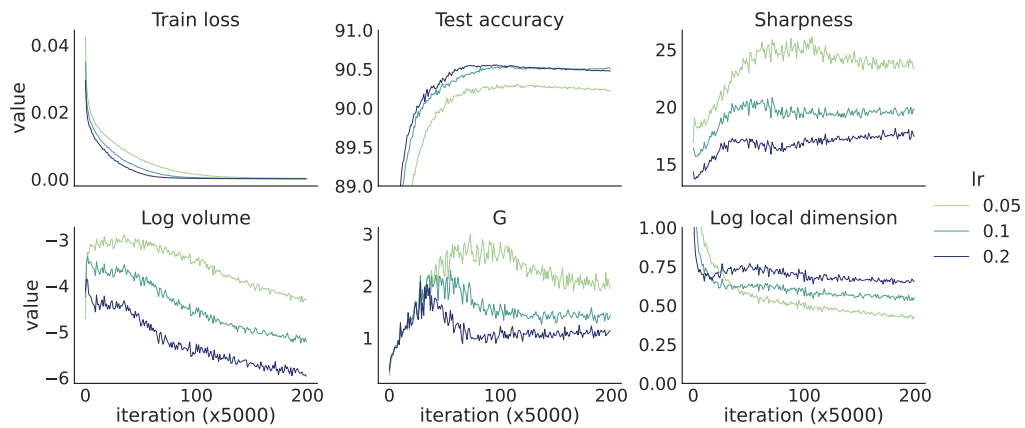


Figure E.8: Trends in key variables across SGD training of a 4-layer MLP with fixed batch size (equal to 20) and varying learning rates (0.05, 0.1 and 0.2). After the loss is minimized, higher learning rates lead to lower sharpness and hence stronger compression. Moreover, G closely follows the trend of sharpness during the training. From left to right: train loss, test accuracy, sharpness (square root of Eq. (3)), log volumetric ratio (Eq. (10)), G (Eq. (6)), and local dimensionality of the network output (Eq. (14)).