
Implicit Jacobian Regularization Weighted with Impurity of Probability Output

Sungyoon Lee¹ Jinseong Park² Jaewook Lee²

Abstract

The success of deep learning is greatly attributed to stochastic gradient descent (SGD), yet it remains unclear how SGD finds well-generalized models. We demonstrate that SGD has an implicit regularization effect on the logit-weight Jacobian norm of neural networks. This regularization effect is weighted with the *impurity* of the probability output, and thus it is active in a certain phase of training. Moreover, based on these findings, we propose a novel optimization method that explicitly regularizes the Jacobian norm, which leads to similar performance as other state-of-the-art sharpness-aware optimization methods.

1. Introduction

Deep learning has shown to great promise for many learning tasks in various areas. Numerous studies have aimed to understand how learning algorithms lead to the successful training of deep neural networks. In particular, it is crucial to understand the geometric properties of the loss landscape of neural networks and their interaction with gradient-based optimization methods, including stochastic gradient descent (SGD), along the training trajectory. These properties have been studied from the perspectives of both optimization (Gur-Ari et al., 2018; Jastrzębski et al., 2019; Ghorbani et al., 2019; Liu et al., 2020; Lewkowycz et al., 2020; Cohen et al., 2021) and generalization (Hochreiter & Schmidhuber, 1997; Keskar et al., 2017; Dinh et al., 2017; Jastrzębski et al., 2017; Wang et al., 2018; Chaudhari et al., 2019; Jiang et al., 2020; Barrett & Dherin, 2021; Smith et al., 2021).

We investigate the Hessian of the training loss (with respect to the model parameters) and its top eigenvalue (also called *sharpness*). The sharpness characterizes the dynamics of neural network training along the optimization trajectory

and appears to be correlated with the generalization capability. To be specific, the sharpness initially increases, and after reaching a certain value, the training dynamics become unstable and the iterate oscillates along the top eigenvector of the Hessian (Jastrzębski et al., 2019; 2020; Cohen et al., 2021). Moreover, the rapid increase in the sharpness during the early phase significantly impacts the final generalization performance (Achille et al., 2019; Golatkar et al., 2019; Jastrzębski et al., 2020; 2021; Lewkowycz et al., 2020; Fort et al., 2020). However, the Hessian of a deep neural network is high-dimensional which makes it difficult to analyze its eigensystem and the sharpness.

In this paper, we investigate the Hessian through two low-dimensional matrices: the *logit Hessian* and the *logit-weight Jacobian* (defined in Definition 3.1). Especially, we analyze the top eigenvalue of the logit Hessian. By doing so, we can provide a simple and intuitive explanation of the relation between the sharpness, the probability output of the classification model, and the Jacobian matrix. This enables us to understand how the sharpness of the loss landscape influences the learning dynamics and the generalization performance. We summarize the main contributions as follows:

- We provide connections between the top eigenvalue of the logit Hessian and the *impurity* of the probability output (Theorem 4.1 in Section 4.1 and 4.2).
- We derive a relation between the sharpness, the top eigenvalue of the logit Hessian, and the Jacobian norm (Theorem 5.1 in Section 5.1), which we call *sharpness-impurity-Jacobian relation*.
- We describe how the sharpness of the loss landscape influences the learning dynamics and the generalization performance (Section 5.2). In particular, we demonstrate that SGD has implicit effects on penalizing the Jacobian norm (*Implicit Jacobian Regularization*) during a certain phase of training (*Active Regularization Period*).
- We also evaluate the *Explicit Jacobian Regularization*, which outperforms state-of-the-art sharpness-aware optimization methods, SAM (Foret et al., 2021) and ASAM (Kwon et al., 2021) (Table 1 in Section 5.3).

¹Department of Computer Science, Hanyang University
²Department of Industrial Engineering, Seoul National University. Correspondence to: Sungyoon Lee <sungyoon-lee@hanyang.ac.kr>.

2. Related Work

We summarize some previous studies on the Hessian, learning dynamics, and generalization of neural networks.

Decomposition of the Hessian LeCun et al. (2012); Dauphin et al. (2014); Sagun et al. (2016; 2017) empirically found that the eigenvalue spectrum of the Hessian during training is composed of two parts: the bulk (concentrated around zero) and the outliers (scattered positively away from zero). Further, Pappas (2019; 2020) proposed a three-level hierarchical decomposition of the Hessian according to each label, logit coordinate, and example, using a well-known Gauss-Newton decomposition. Although we also use the Gauss-Newton decomposition, we focus on the top eigenvalue of the logit Hessian for which the predicted class is more important than the true label (detailed in Section 4.2).

Issues with the SDE modeling of SGD In many studies, SGD has been understood as a stochastic differential equation (SDE) in the limit of vanishing learning rate (Mandt et al., 2016; 2017; Hu et al., 2019; Li et al., 2017; 2019a; Smith & Le, 2018; Chaudhari & Soatto, 2018; Jastrzębski et al., 2017; Zhu et al., 2019; Park et al., 2019). However, Yaida (2019, Section 2.3.3) raised some theoretical concerns for such approximations, and Li et al. (2021) derived a sufficient condition for the SDE modeling to fail. Moreover, Barrett & Dherin (2021) argued that the SDE analysis within the limit of vanishing learning rate cannot explain the generalization benefits of finite learning rates, and proposed a modified gradient flow for finite learning rates. However, they still consider a continuous flow, which has a fundamental limitation in explaining the (discrete) oscillatory behavior with iterative catapult in a practical learning rate regime (Smith et al., 2021), as shown in the following paragraph.

Oscillatory catapult and the plateau of sharpness Xing et al. (2018) investigated the roles of the learning rate and batch size in the SGD dynamics. The authors observed that SGD explores the parameter space, bouncing between walls of valley-like regions. A large learning rate maintains a high valley height, and a small batch size induces gradient stochasticity. Both of them help exploring through the parameter space with different roles in the training dynamics. Jastrzębski et al. (2019) empirically investigated the evolution of the sharpness along the whole training trajectory of SGD. They observed the initial growth of the sharpness (dubbed *progressive sharpening* in Cohen et al. (2021)) as the loss decreases, reaching a maximum sharpness determined by learning rate and batch size, which then decreases towards the end of the training. Due to the progressive sharpening, the SGD step becomes excessively large compared to the shape of the loss landscape. This

is consistent with the valley-like structure shown in Xing et al. (2018). Lewkowycz et al. (2020) investigated simple theoretical models with solvable training dynamics. They showed that, in their setup with a large learning rate, the loss initially increases whereas the sharpness decreases, converging at a flat minimum. This mechanism is called *catapult mechanism*. Recently, Cohen et al. (2021) found that full-batch GD typically operates in a regime called the *Edge of Stability* where the sharpness can no longer increase and stays near a certain value, and the training loss behaves non-monotonically but decreases globally. This optimization behavior at the Edge of Stability can be seen as repeated catapult mechanisms. The authors explicitly marked the limit of the sharpness with $2/\eta$ where η is a learning rate. To describe the aforementioned evolution of the sharpness, Fort & Ganguli (2019) developed a theoretical model based on a random matrix modelling. To build a simple random model, they introduced assumptions regarding the gradients and Hessians, which are i.i.d. isotropic Gaussian with zero mean and varying variance during training. Whereas they focus on building a random model based on observation, we aim to explain the underlying mechanisms.

Implicit bias of SGD There have been many studies on the implicit bias of SGD (Neysshabur, 2017; Zhang et al., 2021; Soudry et al., 2018; Jastrzębski et al., 2020). Jastrzębski et al. (2021) empirically showed that SGD implicitly penalizes the trace of the Fisher Information Matrix (FIM). They also showed that the trace of FIM explodes during the early phase of training when using a small learning rate, which they call catastrophic Fisher explosion. Barrett & Dherin (2021) and Smith et al. (2021) demonstrated that SGD implicitly penalizes the norm of the total gradient and the non-uniformity of the mini-batch gradients. We demonstrate that the (logit-weight) Jacobian plays an important role in the generalization in each case.

3. Background

In this section, we provide some notations, basic equations, and definitions for the following sections. In this paper, we use the denominator layout notation for the vector derivatives, i.e., $\nabla_{\mathbf{v}} \mathbf{u} \equiv \left(\frac{\partial u_j}{\partial v_i} \right)_{ij} \in \mathbb{R}^{v \times u}$, where $\mathbf{u} : \mathbb{R}^v \rightarrow \mathbb{R}^u$ and $\mathbf{v} \in \mathbb{R}^v$.

We consider a problem of learning a C -class classifier that maps an input $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$ to a target label $y \in [C] = \{1, 2, \dots, C\}$. To this end, we build a parameterized model $f_{\theta} : \mathcal{X} \rightarrow \mathcal{Z} \subset \mathbb{R}^C$ with a model parameter $\theta \in \Theta \subset \mathbb{R}^m$ that outputs a logit vector $\mathbf{z} \equiv f_{\theta}(\mathbf{x}) \in \mathcal{Z} \subset \mathbb{R}^C$ (we often omit the dependence on \mathbf{x} and θ). The logit vector \mathbf{z} is then given as an input to the softmax function to yield a probability output $\mathbf{p} \equiv \text{softmax}(\mathbf{z}) \in \Delta^{C-1}$, where $\Delta^{C-1} \equiv \{\mathbf{p} \in [0, 1]^C : \mathbf{1}^T \mathbf{p} = 1, \mathbf{p} \geq \mathbf{0}\}$. We want

the model to match the most probable class c_1 to the true label y , where $c(x) \equiv \arg \text{sort}(\mathbf{p})$ in descending order. We exchangeably denote the probability value corresponding to the true label y as $p \equiv \mathbf{p}_y \in [0, 1]$. The cross-entropy loss, $l = l(\mathbf{z}, y) \in \mathbb{R}$, is equivalent to the negative log-likelihood $l \equiv -\log p$. We use the notations $\|\mathbf{v}\|$ for the Euclidean norm of a vector \mathbf{v} and $\|\mathbf{A}\| \equiv \max_{\|\mathbf{v}\|=1} \|\mathbf{A}\mathbf{v}\|$ for the Euclidean operator norm, $\|\mathbf{A}\|_F \equiv \sqrt{\sum_{i,j} |\mathbf{A}_{ij}|^2}$ for the Frobenius norm of a matrix \mathbf{A} , and the spectral norm $\|\mathbf{S}\|_\sigma$ and the trace $\text{tr}(\mathbf{S}) \equiv \sum_i \mathbf{S}_{ii}$ of a square matrix \mathbf{S} .

Starting with a simple computation of the derivatives of the softmax function in (1) (see Appendix C.1), we can easily derive the following equations in order:

$$\nabla_{\mathbf{z}} \mathbf{p} = \text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^T \in \mathbb{R}^{C \times C} \quad (1)$$

$$\nabla_{\mathbf{z}} p = [\nabla_{\mathbf{z}} \mathbf{p}]_{:,y} = p(\mathbf{e}^y - \mathbf{p}) \in \mathbb{R}^C \quad (2)$$

$$\nabla_{\mathbf{z}} l = \nabla_{\mathbf{z}} p \frac{\partial l}{\partial p} = p(\mathbf{e}^y - \mathbf{p}) \cdot -\frac{1}{p} = \mathbf{p} - \mathbf{e}^y \in \mathbb{R}^C \quad (3)$$

$$\nabla_{\mathbf{z}}^2 l = \nabla_{\mathbf{z}} (\nabla_{\mathbf{z}} l) = \nabla_{\mathbf{z}} (\mathbf{p} - \mathbf{e}^y) = \text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^T \quad (4)$$

where $\text{diag}(\mathbf{p}) \equiv (\delta_{ij} \mathbf{p}_i)_{ij} \in \mathbb{R}^{C \times C}$ is a diagonal matrix with \mathbf{p} as its diagonal entries, and $\mathbf{e}^i \equiv (\delta_{ij})_j \in \mathbb{R}^C$ is a one-hot vector with the i -th element as 1.

Next, the Hessian of the loss function l for the given example \mathbf{x} with respect to the model parameter can be expressed as follows:

$$\nabla_{\boldsymbol{\theta}}^2 l = \nabla_{\boldsymbol{\theta}} \mathbf{z} \nabla_{\mathbf{z}}^2 l \nabla_{\boldsymbol{\theta}} \mathbf{z}^T + \sum_{c=1}^C \nabla_{\boldsymbol{\theta}}^2 z_c \nabla_{z_c} l \quad (5)$$

$$\approx \nabla_{\boldsymbol{\theta}} \mathbf{z} \nabla_{\mathbf{z}}^2 l \nabla_{\boldsymbol{\theta}} \mathbf{z}^T \in \mathbb{R}^{m \times m} \quad (6)$$

using the well-known Gauss-Newton approximation (see Appendix C.2).

Now, we are ready to consider the training loss for the training set \mathcal{D} . We compute the total training loss over \mathcal{D} as $L \equiv \langle l \rangle$ which yields $\nabla L = \langle \nabla l \rangle$ and $\nabla^2 L = \langle \nabla^2 l \rangle$, where $\langle \cdot \rangle$ is the expectation over the empirical measure of the training set \mathcal{D} , i.e., $\langle \cdot \rangle \equiv \hat{\mathbb{E}}_{\mathcal{D}}[\cdot]$. We use the notation $\langle \cdot \rangle_{\mathcal{B}}$ when averaging over a mini-batch \mathcal{B} , i.e., $\langle \cdot \rangle_{\mathcal{B}} \equiv \hat{\mathbb{E}}_{\mathcal{B}}[\cdot]$. Following from (4) and (6), we define the Hessian \mathbf{H} for the total loss and its Gauss-Newton approximation matrix \mathbf{G} with matrices \mathbf{M} and \mathbf{J} as follows:

Definition 3.1. We call \mathbf{M} the logit Hessian, \mathbf{J} the Jacobian (of the logit function with respect to the model parameter), \mathbf{H} the Hessian, and \mathbf{G} the Gauss-Newton approximation, which are defined as follows:

$$\mathbf{M} \equiv \nabla_{\mathbf{z}}^2 l = \text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^T \in \mathbb{R}^{C \times C} \quad (7)$$

$$\mathbf{J} \equiv \nabla_{\boldsymbol{\theta}} \mathbf{z} \in \mathbb{R}^{m \times C} \quad (8)$$

$$\mathbf{H} \equiv \langle \nabla_{\boldsymbol{\theta}}^2 l \rangle \approx \langle \mathbf{J}\mathbf{M}\mathbf{J}^T \rangle \equiv \mathbf{G} \in \mathbb{R}^{m \times m} \quad (9)$$

It is interesting to note that while l is dependent on the true label y , the logit Hessian $\mathbf{M} = \nabla_{\mathbf{z}}^2 l$ is independent of y , as are \mathbf{J} , $\mathbf{J}\mathbf{M}\mathbf{J}^T$, and \mathbf{G} . From (9), although we often use the approximation $\|\mathbf{H}\|_\sigma \approx \|\mathbf{G}\|_\sigma$ as justified by Sagun et al. (2017) and Fort & Ganguli (2019), this approximation occasionally fails during the later phase of training when the top eigenvalues of the Gauss-Newton matrix are insufficiently isolated from the bulk near 0 (Papayan, 2018). Thus, we mainly focus on the early phase of training.

4. Logit Hessian and Impurity

In the previous section, we introduced the Gauss-Newton approximation \mathbf{G} of the Hessian \mathbf{H} , and decomposition of \mathbf{G} with the Jacobian \mathbf{J} and the logit Hessian \mathbf{M} , i.e., $\mathbf{G} = \langle \mathbf{J}\mathbf{M}\mathbf{J}^T \rangle$. Now, we focus on the logit Hessian \mathbf{M} and its eigendecomposition, estimate the top eigenvalue of \mathbf{M} with upper/lower bounds (Section 4.1), and explore the evolution of the top eigenvalue during training (Section 4.2).

4.1. Bounds on the Eigenvalues of the Logit Hessian

The lower-dimensional logit Hessian matrix $\mathbf{M} \in \mathbb{R}^{C \times C}$ is simple and fully characterized by only the probability vector \mathbf{p} as $\mathbf{M} = \text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^T$ in (7), however it turns out to be important for understanding the much higher-dimensional matrix $\mathbf{G} \in \mathbb{R}^{m \times m}$ ($C \ll m$). Because $\mathbf{M} = \text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^T$ is a rank-one modification of a simple diagonal matrix $\text{diag}(\mathbf{p})$, we can obtain its eigenvalues $\{\lambda^{(i)}\}_{i=1}^C$ and eigenvectors $\{\mathbf{q}^{(i)}\}_{i=1}^C$ from the theory of the rank-one modification of the eigenproblem (see, for example, Bunch et al. (1978); Golub (1973); Golub & Van Loan (2013)) where $\lambda^{(i)} \in \mathbb{R}$ is the i -th largest eigenvalue of \mathbf{M} and $\mathbf{q}^{(i)} \in \mathbb{R}^C$ is its corresponding eigenvector. We also use the same ordered index of $(i) \in [C]$ with parentheses for the probability output $\mathbf{p} \in \Delta^{C-1}$, i.e., $\mathbf{c}_i = (i)$ and $\mathbf{p}_{(1)} \geq \mathbf{p}_{(2)} \geq \dots \geq \mathbf{p}_{(C)}$, because this ordering is related to the eigenvalues $\{\lambda^{(i)}\}_{i=1}^C$ as shown in the following theorem.

Theorem 4.1 (Eigensystem of the logit Hessian \mathbf{M}). *The eigenvalues $\lambda^{(i)}$ ($\lambda^{(1)} \geq \lambda^{(2)} \geq \dots \geq \lambda^{(C)}$) and the corresponding eigenvectors $\mathbf{q}^{(i)}$ of the logit Hessian $\mathbf{M} = \nabla_{\mathbf{z}}^2 l = \text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^T$ satisfy the following properties:*

(a) *The eigenvalue $\lambda^{(i)}$ is the i -th largest solution of the following equation:*

$$v(\lambda) = 1 - \sum_{i=1}^C \mathbf{p}_i^2 (\mathbf{p}_i - \lambda)^{-1} = 0 \quad (10)$$

(b) *The eigenvector $\mathbf{q}^{(i)}$ is aligned with the direction of $(\text{diag}(\mathbf{p}) - \lambda^{(i)} \mathbf{I})^{-1} \mathbf{p}$*

(c) *$\mathbf{p}_{(i+1)} \leq \lambda^{(i)} \leq \mathbf{p}_{(i)}$ for $1 \leq i \leq C-1$, and $\lambda^{(C)} = 0$*

(d) $\frac{1}{2}\text{Gini}(\mathbf{p}_{(1)}) \leq \lambda^{(1)} \leq \text{Gini}(\mathbf{p}_{(1)})$, where $\text{Gini}(q) = 2q(1-q)$ is the Gini impurity for the binary $(q, 1-q)$.

We defer the proof to Appendix C.3. In the main text, we mainly focus on investigating the top eigenvalue $\lambda^{(1)}$ of \mathbf{M} by utilizing Theorem 4.1 (c) and (d), which provide the upper/lower bounds on $\lambda^{(1)}$ (Theorem 4.1 (a) and (b) are applied in Appendix). To be specific, the top eigenvalue $\lambda^{(1)}$ is bounded by $\frac{1}{2}\text{Gini}(\mathbf{p}_{(1)}) \leq \lambda^{(1)} \leq \text{Gini}(\mathbf{p}_{(1)})$, and thus we call $\lambda^{(1)}$ the *impurity* (of the probability output).

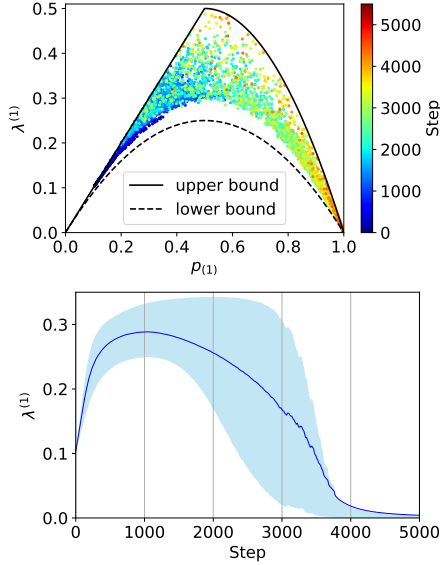


Figure 1: **Impurity $\lambda^{(1)}$ increases and then decreases as $p_{(1)}$ increases during training.** Top: The impurity is plotted against $p_{(1)}$ for a fixed example during training (from blue to red). We plot the upper bound $\min\{p_{(1)}, \text{Gini}(p_{(1)})\}$ (black solid line) and lower bound $\frac{1}{2}\text{Gini}(p_{(1)})$ (black dashed line) from Theorem 4.1 (c) and (d). Bottom: The impurity is plotted against the training step. The blue curve indicates its mean value $\langle \lambda^{(1)} \rangle$ and sky-blue area shows the 25%-75% quantile range for the training data. See Appendix E for detailed settings.

4.2. Evolution of Impurity

We explore the top eigenvalue $\lambda^{(1)}$ of \mathbf{M} (also referred to as impurity) during training. Figure 1 demonstrates the n-shaped evolution of the impurity, which increases in the beginning and then decreases in the later phase of training. We trained a model to zero training loss, and thus, for most of the examples, the probability p_y for the true class y eventually becomes the highest probability $p_{(1)}$. As the top probability $p_{(1)}$ increases from $1/C$ to 1 during training, the impurity starts from $\lambda^{(1)} \approx \frac{1}{C} \in [\frac{1}{2}\text{Gini}(\frac{1}{C}), \frac{1}{C}] = [\frac{C-1}{C^2}, \frac{1}{C}]$ (Theorem 4.1 (c) and (d)), and increases at the initial phase of the training, being lower bounded by $\frac{1}{2}\text{Gini}(p_{(1)}) =$

$p_{(1)}(1-p_{(1)})$, which increases for $p_{(1)} \in [0, 0.5]$. Then, $\lambda^{(1)}$ decreases as $p_{(1)}$ becomes larger than 0.5, which leads $\lambda^{(1)}$ to reach nearly 0 at the later phase because it is upper bounded by $\text{Gini}(p_{(1)}) = 2p_{(1)}(1-p_{(1)})$, which decreases for $p_{(1)} \in [0.5, 1]$. Note that Cohen et al. (2021) tried to estimate a similar value, but they use p_y , not $p_{(1)}$. Similarly, Pappayan (2019) decomposed \mathbf{G} into components using the label class information y . Thus, at the beginning of the training, the cluster members are not well-separated according to the true label y . We again emphasize that \mathbf{M} is independent of the label y , and the bounds on impurity $\lambda^{(1)}$ are well-described by the probability $p_{(1)}$ of the predicted class.

5. Implicit/Explicit Jacobian Regularization

In this section, based on the results of the previous sections, we aim to derive a relation between the sharpness, the impurity, and the Jacobian (Section 5.1), and answer how the sharpness of the loss landscape influences the learning dynamics and the generalization performance (Section 5.2 and 5.3). Detailed experimental settings for each figure and table are described in Appendix E.

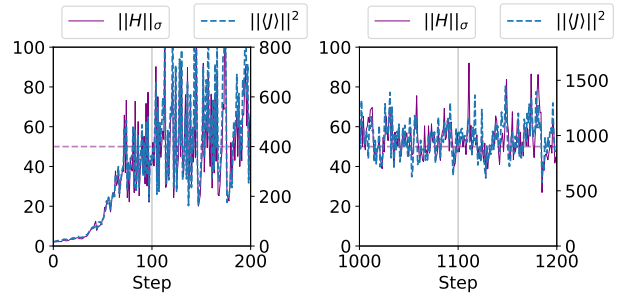


Figure 2: **Sharpness $\|H\|_\sigma$ and Jacobian norm $\|\langle J \rangle\|^2$.** The graphs show similar oscillating behaviors up to a factor $\hat{\lambda}^*$ which is locally constant and slowly changes during training (CIFAR-10, $\eta = 0.04$, $|\mathcal{B}| = 128$; Left: 0-200, Right: 1000-1200 steps). See Appendix I for more details. We highlighted $\|H\|_\sigma = 2/\eta$ with a dashed horizontal line.

5.1. Sharpness-Impurity-Jacobian Relation

We first take a closer look at the sharpness of the loss landscape during training and build a relation between the sharpness $\|H\|_\sigma$, the impurity $\lambda^{(1)}$, and the Jacobian \mathbf{J} . Because the Gauss-Newton matrix \mathbf{G} is known to approximate the true Hessian \mathbf{H} well, especially for the top eigenspace (Sagun et al., 2017; Fort & Ganguli, 2019; Pappayan, 2019), we can write the sharpness $\|H\|_\sigma$ as $\|H\|_\sigma \approx \|G\|_\sigma = \|\langle \mathbf{J} \mathbf{M} \mathbf{J}^T \rangle\|_\sigma$. This implies that the impurity $\|M\|_\sigma$ and squared Jacobian norm $\|J\|^2$ are highly correlated with the sharpness $\|H\|_\sigma$, as demonstrated in the following theorem. We defer the proof to Appendix C.4.

Theorem 5.1 (Sharpness-Impurity-Jacobian Relation). *For some lower bound $0 \leq \lambda^* \leq \lambda^{(1)}$ of the impurity for each $\mathbf{x} \in \mathcal{D}$, we have $\|\mathbf{G}\|_\sigma = \langle \lambda^* \|\mathbf{J}\|^2 \rangle$.*

In the next section, we will demonstrate that the sharpness is implicitly upper bounded (Proposition 5.2), and so is the per-example $\lambda^* \|\mathbf{J}\|^2$. Therefore, with a large λ^* , the Jacobian norm $\|\mathbf{J}\|^2$ is strongly regularized to a small value, i.e., λ^* in $\lambda^* \|\mathbf{J}\|^2$ acts as an adaptive regularization weight. Moreover, as the impurity $\lambda^{(1)}$ decreases, so does the lower bound λ^* , and the regularization effect diminishes. Now, because it is computationally inefficient to track $\|\mathbf{J}\|^2$ for every $\mathbf{x} \in \mathcal{D}$, we instead investigate $\|\langle \mathbf{J} \rangle\|^2$. We expect $\|\mathbf{H}\|_\sigma = \hat{\lambda}^* \|\langle \mathbf{J} \rangle\|^2$ for some $\hat{\lambda}^*$. Figure 2 shows that $\|\mathbf{H}\|_\sigma$ and $\|\langle \mathbf{J} \rangle\|^2$ have almost identical oscillating behaviors up to a factor $\hat{\lambda}^*$ which is locally constant and slowly changes during training (see Appendix I for further details).

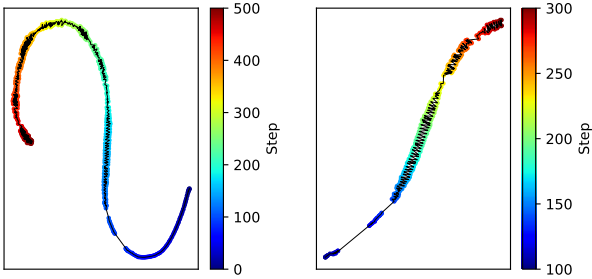


Figure 3: **Oscillatory catapult in the optimization trajectory $\{\theta^{(t)}\}$ (from blue to red) of full-batch GD.** Left: UMAP (McInnes et al., 2018) of the model parameters trained on CIFAR-10 for the first 500 steps. Right: Zoom-in of the oscillatory steps [100, 300]. After a few steps (~ 100), the sharpness reaches a threshold (see Figure 2) and the iterate shows an oscillatory behavior with an iterative catapult.

5.2. Implicit Jacobian Regularization

Now, we are ready to answer how the sharpness of the loss landscape influences the learning dynamics and generalization performance of neural networks.

Growing Jacobian and progressive sharpening during the early phase of training The weight norm $\|\theta\|$ increases, increasing the logit norm $\|z\|$ and minimizing the cross-entropy loss during training (Soudry et al., 2018) (see Appendix F for details). We hypothesize that this is one of the factors leading to an increase in the layerwise weight norms and the Jacobian norm. Thus, the progressive sharpening can be mainly attributed to the increase of the Jacobian norm (Theorem 5.1). For the MSE loss, Wang et al. (2022) proved the progressive sharpening along with the increase in the (squared) Jacobian norm for a two-layer neural network. We emphasize that the rate of increase in the Jacobian norm, however, varies throughout the training.

Oscillatory catapult and the plateau of sharpness As the sharpness increases in the beginning, the width of the valley of the loss landscape becomes narrower than the discrete step size of the SGD. After the sharpness reaches this threshold, the iterate starts to bounce off from one side of the valley to the other, which is then repeated (Xing et al., 2018; Jastrzębski et al., 2019). Figure 3 shows this oscillatory behavior with an iterative catapult after the sharpness reaches the threshold, using UMAP (McInnes et al., 2018). Owing to the catapult, the iterate cannot remain in a sharper area and returns to stability after the sharpness goes below the threshold (Lewkowycz et al., 2020; Damian et al., 2022). This oscillatory catapult and the plateau of the sharpness are attributed to the *discrete* dynamics of the gradient-based optimization with a finite learning rate and cannot be described through a *continuous* gradient flow. Figure 2 shows fine-grained patterns that the sharpness oscillates up and down around the threshold by the two conflicting effects: the Jacobian norm tends to increase the sharpness, and the self-stabilization reduces it again when the sharpness is over the threshold. Therefore, we can observe the plateau of the sharpness in a coarser scale (see Figure 4). We derive the threshold of the sharpness according to learning rate and batch size in the following proposition.

Proposition 5.2. *For SGD with a quadratic loss, the expected loss decreases when $\|\mathbf{H}\|_\sigma \leq \frac{2\rho_B}{\eta}$, where η is learning rate, B is batch size and $\rho_B \equiv \frac{\|\langle \nabla_{\theta^l} \rangle\|^2}{\mathbb{E}_{|B|=B}[\|\langle \nabla_{\theta^l} \rangle_B\|^2]} \leq 1$.*

The proof is deferred to Appendix C.5 with further discussion below. We suggest Proposition 5.2 as a modified version of the Edge of Stability proposed in Cohen et al. (2021) ($\rho_B = 1$ for full-batch), which generalizes to mini-batch SGD. As shown in Figure 2, this also approximately holds for the cross-entropy loss. Note that the threshold value $2\rho_B/\eta$ depends on learning rate η and batch size B , which is consistent with the results shown in Figure 5. To be specific, batch gradients $\langle \nabla_{\theta^l} \rangle_B$ are more scattered if B is smaller, which leads to a smaller ρ_B . Therefore, a smaller B and a larger η (dotted line) tend to lead to a lower threshold $2\rho_B/\eta$ (see the three purple lines of $\|\mathbf{H}\|_\sigma$ for different η and B in Figure 5).

Implicit Jacobian Regularization (IJR) Due to the catapult effect, the rate of increase in the Jacobian norm decreases. In other words, SGD implicitly penalizes the Jacobian norm since $\langle \lambda^* \|\mathbf{J}\|^2 \rangle \approx \|\mathbf{H}\|_\sigma \leq \frac{2\rho_B}{\eta}$. This implicit Jacobian regularization (IJR) effect begins after the sharpness reaches the threshold. In addition, because the lower bound $\lambda^* \leq \lambda^{(1)}$ acts as a regularization coefficient, the effect diminishes as the impurity $\lambda^{(1)}$ decreases with increasing $p_{(1)} \geq 0.5$ during the later phase (see Figures 1, 4, and 5). This explains why the behavior of the sharpness in the early phase of the training (where $\lambda^{(1)}$ is not small)

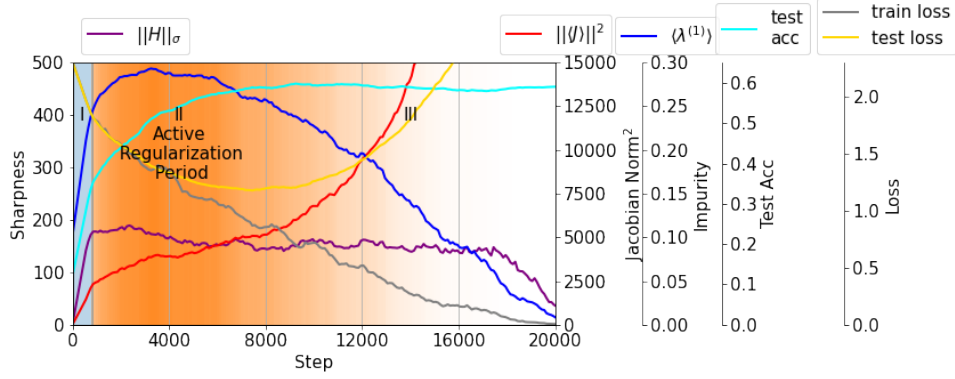


Figure 4: **Three phases of Implicit Jacobian Regularization (IJR)**. The Jacobian norm (red) (I) initially rapidly increases before the sharpness (purple) reaches near the threshold, (II) is actively regularized with a gentle slope, and (III) again increases quickly as the regularization effect diminishes (as the regularization weight $\lambda^* \leq \lambda^{(1)}$ decreases) with the slope being gradually steeper. We call phase II the Active Regularization Period (ARP). SGD gradually progresses from phase II to phase III (strong regularization with a high $\langle \lambda^{(1)} \rangle$ indicated with a dark orange color). Thus, we do not explicitly separate these two phases, but we may arbitrarily mark the end of the ARP when $\langle \lambda^{(1)} \rangle$ decreases below a certain value (e.g., 0.25).

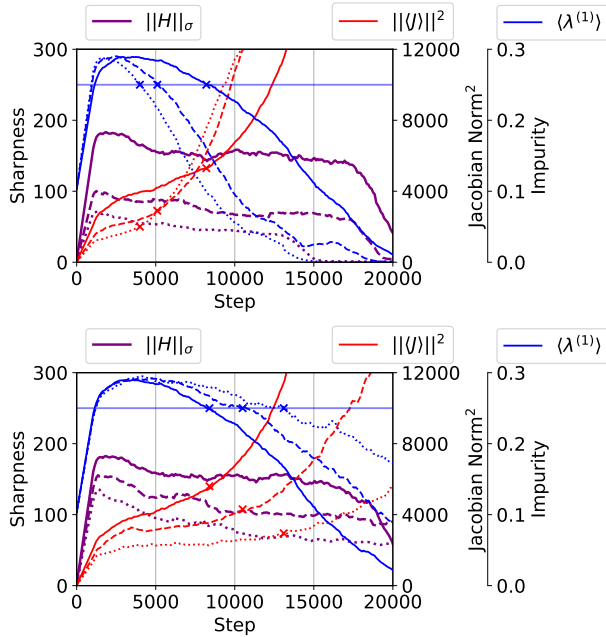


Figure 5: **The effect of IJR varies depending on (η, B)** . We used SGD (solid/dashed/dotted lines) with fixed $B = 128$ and different $\eta = 0.01/0.02/0.03$ (Top), and with fixed $\eta = 0.01$ and different $B = 128/64/32$ (Bottom) on CIFAR-10. We arbitrarily mark the end of the ARP (where $\langle \lambda^{(1)} \rangle = 0.25$) and the corresponding Jacobian norm with "x". Training with a large learning rate and a small batch size (dotted lines) penalizes the Jacobian norm (red) more strongly with lower limits of the sharpness (purple) in the ARP. The curves are smoothed for visual clarity. See Figure 2 for a finer structure.

seems to highly impact the final generalization (Achille et al., 2019; Golatkar et al., 2019; Jastrzębski et al., 2020; 2021; Lewkowycz et al., 2020; Fort et al., 2020).

Figure 4 shows that SGD has implicit regularization effects on the Jacobian norm in a certain period. There are three phases in which the Jacobian norm (I) initially rapidly increases before the sharpness reaches near the threshold, (II) is actively regularized with a gentle slope, and (III) again increases fast as the regularization effect diminishes (as the implicit regularization weight $\lambda^* \leq \lambda^{(1)}$ decreases) with the slope becoming gradually steeper. We call the second phase *Active Regularization Period (ARP)*.

The evolution of the sharpness is highly affected by learning rate and batch size (Jastrzębski et al., 2019; 2020; Lewkowycz et al., 2020; Cohen et al., 2021). As Figure 5 shows, when comparing the three red lines (solid/dashed/dotted) with different learning rates (0.01/0.02/0.03) and batch sizes (128/64/32), training with a larger learning rate and a smaller batch size (dotted lines) encourages a stronger implicit regularization on the Jacobian norm (see the red "x" marks) with a lower threshold $2\rho_B/\eta$ of the sharpness in the ARP. This could explain why training with a large learning rate and a small batch size yields a better generalization (Keskar et al., 2017; Hoffer et al., 2017; Li et al., 2019b; Jastrzębski et al., 2020).

Jacobian norm and Generalization There has been a line of work on the norm-based capacity control and generalization (Neyshabur et al., 2014; 2015a;b; 2017; 2018; Bartlett et al., 2017; Dziugaite & Roy, 2017; Nagarajan & Kolter, 2019). Similarly, IJR in SGD controls the Jacobian norm and the capacity of the model. Moreover, as the Jacobian norm is directly related to the sharpness (Theorem 5.1),

Table 1: **Explicit Jacobian Regularization (EJR) enhances the test accuracy under various settings.** We report the improvement ($\Delta\text{Acc.}$) and Error Reduction Rate (ERR) when trained with EJR, in comparison to the standard SGD. First, we train a simple 6-layer CNN (SimpleCNN) on CIFAR-10 without data augmentation until convergence. Second, we train WRN-28-2 and WRN-28-10 with an efficient variant of EJR and compare it with the state-of-the-art methods, SAM (Foret et al., 2021) and ASAM (Kwon et al., 2021). We use different settings for SimpleCNN and WRN (see Appendix E for details). We report the mean \pm std of three independent runs. ASAM does not work under the full-batch setting (test accuracy around 20%). More results are found in Appendix Q.

| Dataset | Model | Hyperparameters | | | Test Accuracy | | | | $\Delta\text{Acc.}$ (%p) | ERR (%) |
|-----------|-----------|-----------------|------------------|----------------------------|------------------|------------------|----------------------------------|----------------------------------|-----------------------------|----------------------------------|
| | | B | η | λ_{reg}/ρ_{reg} | SGD | SAM | ASAM | EJR | | |
| CIFAR-10 | SimpleCNN | 128 | 0.003 | $\lambda_{reg} = 0.01$ | 68.92 \pm 0.13 | 74.15 \pm 0.01 | 69.41 \pm 0.03 | 75.61\pm0.03 | +6.69 | 21.53 |
| | | | 0.01 | | 71.08 \pm 0.21 | 74.18 \pm 0.08 | 70.84 \pm 0.29 | 75.64\pm0.02 | +4.56 | 15.77 |
| | | | 0.03 | | 72.16 \pm 0.36 | 74.10 \pm 0.13 | 71.51 \pm 0.22 | 75.71\pm0.10 | +3.55 | 12.75 |
| | | | 0.1 | | 72.59 \pm 0.29 | 74.94 \pm 0.29 | 72.65 \pm 0.18 | 75.49\pm0.05 | +2.90 | 10.58 |
| | | | 0.3 | | 69.17 \pm 0.75 | 71.09 \pm 0.06 | 70.91 \pm 0.43 | 74.98\pm0.64 | +5.81 | 18.85 |
| | | 50000 | 0.01 | 69.10 \pm 0.05 | 64.80 \pm 0.06 | - | 73.99\pm0.01 | +4.89 | 15.83 | |
| | | | 0.03 | 69.72 \pm 0.07 | 64.80 \pm 0.07 | - | 74.19\pm0.03 | +4.47 | 14.76 | |
| | | | 0.1 | 69.34 \pm 0.39 | 65.28 \pm 0.04 | - | 73.88\pm0.11 | +4.54 | 14.81 | |
| | | | 0.3 | 64.54 \pm 0.39 | 66.26 \pm 0.11 | - | 71.86\pm0.20 | +7.32 | 20.64 | |
| | | | WRN-28-2 | 256 | 0.1 | $\rho_{reg} = 2$ | 95.46 \pm 0.03 | 96.17 \pm 0.06 | 96.16 \pm 0.14 | 96.25\pm0.04 |
| WRN-28-10 | 128 | 0.1 | $\rho_{reg} = 2$ | 96.21 \pm 0.03 | 96.98 \pm 0.06 | 97.13 \pm 0.06 | 97.21\pm0.04 | +1.00 | 26.39 | |
| CIFAR-100 | WRN-28-2 | 256 | 0.1 | $\rho_{reg} = 4$ | 75.32 \pm 0.17 | 78.12 \pm 0.31 | 78.93 \pm 0.16 | 79.15\pm0.39 | +3.83 | 15.52 |
| | WRN-28-10 | 128 | 0.1 | $\rho_{reg} = 4$ | 80.72 \pm 0.28 | 83.23 \pm 0.17 | 83.65 \pm 0.09 | 83.94\pm0.10 | +3.22 | 16.70 |

it can provide connections between the sharpness and the norm-based capacity control. We argue that IJR is one of the main reasons why SGD finds well-generalized minima.

5.3. Explicit Jacobian Regularization (EJR)

To further investigate and boost the effectiveness of IJR, we *explicitly* regularize the Jacobian norm. We expect improvements in the generalization when introducing EJR. This supports the effectiveness of IJR in that it efficiently controls the capacity of the model and helps find better-generalized minima. However, it is computationally hard to back-propagate through the computation graph of the operator norm $\|\langle \mathbf{J} \rangle\|^2$ for a practical neural network even with a simple iterative method (see Algorithm 2 in Appendix E). Thus, we instead penalize an upper bound, that is, the Frobenius norm $\|\langle \mathbf{J} \rangle\|_F^2 (\geq \|\langle \mathbf{J} \rangle\|^2)$, with the regularization coefficient λ_{reg}/C , i.e., we minimize $L + \lambda_{reg}\|\langle \mathbf{J} \rangle\|_F^2/C$. The Frobenius regularization term can be efficiently computed with an unbiased estimator $\|\langle \mathbf{J} \rangle\|_F^2 = C\mathbb{E}_{\mathbf{u} \sim U(\mathbb{S}^{C-1})}[\|\langle \mathbf{J} \rangle \mathbf{u}\|^2] = C\mathbb{E}_{\mathbf{u} \sim U(\mathbb{S}^{C-1})}[\|\nabla_{\theta} \langle \mathbf{u}^T \mathbf{z} \rangle\|^2]$, where \mathbf{u} is randomly drawn from the unit hypersphere \mathbb{S}^{C-1} . Because the batch-size is large enough, we efficiently use a single sample $\mathbf{u} \sim U(\mathbb{S}^{C-1})$ for each batch, as suggested in Hoffman et al. (2019). Table 1 (SimpleCNN) shows clear improvements in the test accuracy of a simple 6-layer CNN (SimpleCNN) when introducing EJR. We follow a similar setting to Jaszczyński et al. (2021).

Furthermore, we also propose an efficient variant of EJR for

larger networks with another regularization coefficient ρ_{reg} , which update the model parameter $\theta^{(t)}$ using two gradient steps like SAM (Foret et al., 2021) as follows:

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_{\theta} L(\hat{\theta}^{(t)}), \quad (11)$$

where

$$\hat{\theta}^{(t)} \equiv \theta^{(t)} + \rho_{reg} \delta^{(t)} / \|\delta^{(t)}\|, \quad (12)$$

$$\delta^{(t)} \equiv \langle \mathbf{J} \rangle \mathbf{u}^{(t)} = \nabla_{\theta} \langle \mathbf{z}^T \mathbf{u}^{(t)} \rangle. \quad (13)$$

Here, \mathbf{J} and \mathbf{z} are evaluated at $\theta^{(t)}$, and $\mathbf{u}^{(t)} \sim U(\mathbb{S}^{C-1})$ is randomly sampled for each step. If we use the gradient ascent step $\delta^{(t)} = \nabla_{\theta} L(\theta^{(t)}) = \langle \mathbf{J}(\mathbf{p} - \mathbf{e}^y) \rangle = \nabla_{\theta} \langle \mathbf{z}^T (\mathbf{p} - \mathbf{e}^y) \rangle$ instead of (13), we can obtain SAM. We will shortly show that EJR is comparable to SAM, which implies that a random direction $\mathbf{u}^{(t)}$ is as good as the specific direction of $\mathbf{p} - \mathbf{e}^y$, if not better. Note that our perturbation $\delta^{(t)}$ in (13) has a 1/2 chance of increasing the loss (\mathbf{u} and $-\mathbf{u}$ have the same probability to be sampled). Therefore, we may conclude that SAM is successful not because it solves a minimax problem, but its ascent step is in the column space of the Jacobian, thereby regularizing the Jacobian norm. See Appendix Q for the details.

Figure 6 shows that the efficient version of EJR (dotted line) successfully mitigates overfitting of the model, especially after each learning rate decay. Table 1 (WRN) shows that it outperforms the state-of-the-art sharpness-aware optimization methods such as SAM (Foret et al., 2021) and ASAM (Kwon et al., 2021) on WideResNet (WRN) (Zagoruyko

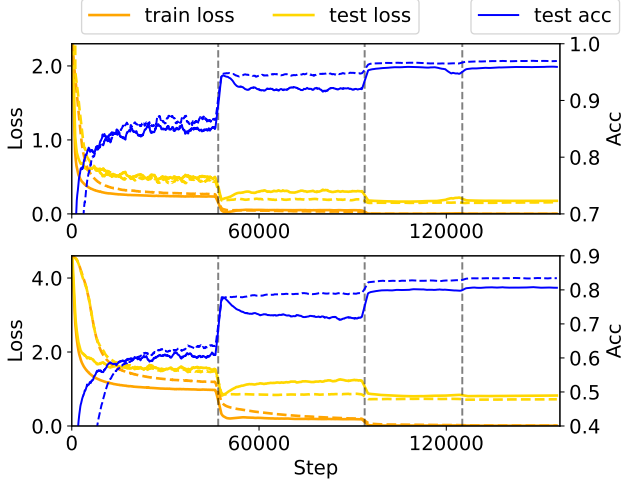


Figure 6: **Efficient EJR (dashed lines) in comparison to SGD (solid lines) for WRN-28-10 on CIFAR-10/CIFAR-100 (Top/Bottom).** EJR effectively mitigates the overfitting, especially after each learning rate decay (undesirable decrease/increase of test accuracy/loss) at steps = $\lceil 50k/128 \rceil \times 400 \times [0.3, 0.6, 0.8]$.

& Komodakis, 2016) with the same computational complexity (we run SGD for twice longer epochs than the others). We follow a similar setting to Kwon et al. (2021). We also evaluate the results of Smith et al. (2021) (regularization on $\|\langle \nabla_{\theta} l \rangle_{\mathcal{B}}\|^2$) and Hoffman et al. (2019) (regularization on $\|\nabla_{\mathbf{x}} z\|_F^2$) on WRN, but their performance is about 95.4-95.6/96.2-96.3 (WRN-28-2/WRN-28-10) on CIFAR-10, which is similar to SGD (95.46/96.21), but is not competitive with SAM (96.17/96.98), ASAM (96.16/97.13), and **EJR (96.25/97.21)**. Note that the logit-input Jacobian (Novak et al., 2018; Hoffman et al., 2019) has nothing to do with the implicit bias of SGD in terms of the sharpness and the Edge of Stability, unlike the logit-weight Jacobian.

Connections between the Jacobian and Fisher/Gradient Penalty Our explanation of the implicit bias of SGD may extend to the catastrophic Fisher explosion (Jastrzebski et al., 2021) with \mathbf{G} instead of the Fisher Information Matrix (FIM). The trace of \mathbf{G} can be written as follows:

$$\begin{aligned} \text{tr}(\mathbf{G}) &= \langle \text{tr}(\mathbf{J}\mathbf{M}\mathbf{J}^T) \rangle = \langle \sum_{i=1}^{C-1} \lambda^{(i)} \|\mathbf{J}\mathbf{q}^{(i)}\|^2 \rangle \\ &\approx \langle \text{tr}(\mathbf{M}) \|\mathbf{J}\|_F^2 \rangle / C \end{aligned} \quad (14)$$

where we assume $\|\mathbf{J}\mathbf{q}^{(i)}\|^2 \approx \|\mathbf{J}\|_F^2 / C$ since $\sum_{i=1}^C \|\mathbf{J}\mathbf{q}^{(i)}\|^2 = \|\mathbf{J}\|_F^2$ (see Figure 11 (left) in Appendix D for empirical evidence). Here, the trace of the logit Hessian \mathbf{M} can be equivalently written as a C -class Gini impurity, i.e., $\text{tr}(\mathbf{M}) = \sum_{i=1}^C \mathbf{p}_i(1 - \mathbf{p}_i) = 1 - \sum_{i=1}^C \mathbf{p}_i^2 \equiv \text{Gini}^C(\mathbf{p})$, which is $\frac{C-1}{C}$ for the initial uniform distribution and zero for a one-hot probabil-

ity. Thus, penalizing $\text{tr}(\mathbf{G})$ induces the effect of penalizing $\|\mathbf{J}\|_F$, especially in the early phase of training with large C -class impurity $\text{Gini}^C(\mathbf{p})$. Thus, as Jastrzebski et al. (2021) argued, Fisher Penalty on the trace of the FIM improves the generalization performance by limiting the memorization, and thus the Jacobian regularization may have similar effects. Moreover, because $\nabla_{\theta} l(\mathbf{z}, \hat{y}) = \nabla_{\theta} z \nabla_z l(\mathbf{z}, \hat{y}) = \mathbf{J}(\mathbf{p} - \mathbf{e}^{\hat{y}})$, the trace of the FIM they approximated is simply $\|\mathbb{E}_{\mathbf{x} \sim \mathcal{B}} \mathbb{E}_{\hat{y} \sim \mathcal{P}} [\nabla_{\theta} l(\mathbf{z}, \hat{y})]\|^2 = \|\mathbb{E}_{\mathbf{x} \sim \mathcal{B}} \mathbb{E}_{\hat{y} \sim \mathcal{P}} [\mathbf{J}(\mathbf{p} - \mathbf{e}^{\hat{y}})]\|^2$ with a single sample \hat{y} , the gradient norm penalty (Barrett & Dherin, 2021; Smith et al., 2021) is $\|\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{B}} [\mathbf{J}(\mathbf{p} - \mathbf{e}^y)]\|^2$, and the EJR regularizer is $\|\mathbf{J}\|_F^2 = C \mathbb{E}_{\mathbf{u} \sim U(\mathbb{S}^{C-1})} \|\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{B}} [\mathbf{J}\mathbf{u}]\|^2$. In each case, the Jacobian \mathbf{J} plays an important role in the generalization.

Adversarially robust generalization We also apply EJR to train a model that is robust to adversarial attacks (Szegedy et al., 2013). Adversarially robust training suffers from robust overfitting (Rice et al., 2020) and we expect EJR to help improve robust generalization. We test the effectiveness of EJR with AT (Madry et al., 2018) and TRADES (Zhang et al., 2019). Table 2 shows that the proposed method can obtain better robust generalization against the adversarial attacks. We report the best results with the regularization parameter tuning. We follow a similar setting to Pang et al. (2020). It outperforms the baselines (Madry et al., 2018; Zhang et al., 2019) in terms of the robust accuracy against AA (Croce & Hein, 2020). See Appendix R for details.

Table 2: **Adversarial training with efficient EJR on CIFAR-10 compared with AT (Madry et al., 2018) and TRADES (Zhang et al., 2019).** We report standard accuracy (Std) and robust accuracy ($\epsilon = 8/255$) against PGD-20 and AutoAttack (AA) (Croce & Hein, 2020).

| Method | PreAct ResNet-18 | | | WRN-34-10 | | |
|------------|------------------|--------|--------------|-----------|--------|--------------|
| | Std | PGD-20 | AA | Std | PGD-20 | AA |
| AT | 82.45 | 52.85 | 48.76 | 86.99 | 52.20 | 49.83 |
| AT-EJR | 82.32 | 53.58 | 49.14 | 86.85 | 57.82 | 53.73 |
| TRADES | 82.34 | 52.83 | 49.06 | 83.62 | 57.08 | 53.29 |
| TRADES-EJR | 81.56 | 53.05 | 49.56 | 83.97 | 57.48 | 53.85 |

6. Conclusion

We investigated the Hessian using the Jacobian and the top eigenvalue of the logit Hessian. By doing so, we provided a simple and intuitive explanation on the relation between the sharpness of the loss landscape, the learning dynamics of the gradient-based optimization methods, and the generalization performance of neural networks. We hope this research can help answer other intriguing questions regarding SGD.

Acknowledgements

S. Lee was supported by the research fund of Hanyang University (HY-202300000000552). This work was also partly supported by the Institute of Information & communications Technology Planning & Evaluation (IITP) (No.RS-2023-002206284, Artificial intelligence for prediction of structure-based protein interaction reflecting physicochemical principles and No.2022-0-00984, Development of Artificial Intelligence Technology for Personalized Plug-and-Play Explanation and Verification of Explanation) and by the National Research Foundation of Korea (NRF) (No. RS-2023-00244896 and No. 2019R1A2C2002358) grant funded by the Korean government (MSIT).

References

- Achille, A., Rovere, M., and Soatto, S. Critical learning periods in deep networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=BkeStsCckQ>.
- Barrett, D. and Dherin, B. Implicit gradient regularization. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=3q5IqUrkcF>.
- Bartlett, P. L., Foster, D. J., and Telgarsky, M. J. Spectrally-normalized margin bounds for neural networks. In *Advances in Neural Information Processing Systems*, pp. 6240–6249, 2017.
- Bunch, J. R., Nielsen, C. P., and Sorensen, D. C. Rank-one modification of the symmetric eigenproblem. *Numerische Mathematik*, 31(1):31–48, 1978.
- Chaudhari, P. and Soatto, S. Stochastic gradient descent performs variational inference, converges to limit cycles for deep networks. In *International Conference on Learning Representations*, 2018.
- Chaudhari, P., Choromanska, A., Soatto, S., LeCun, Y., Baldassi, C., Borgs, C., Chayes, J., Sagun, L., and Zecchina, R. Entropy-sgd: Biasing gradient descent into wide valleys. *Journal of Statistical Mechanics: Theory and Experiment*, 2019(12):124018, 2019.
- Cohen, J., Kaur, S., Li, Y., Kolter, J. Z., and Talwalkar, A. Gradient descent on neural networks typically occurs at the edge of stability. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=jh-rTtvkGeM>.
- Croce, F. and Hein, M. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *International conference on machine learning*, pp. 2206–2216. PMLR, 2020.
- Damian, A., Nichani, E., and Lee, J. D. Self-stabilization: The implicit bias of gradient descent at the edge of stability. *arXiv preprint arXiv:2209.15594*, 2022.
- Dauphin, Y. N., Pascanu, R., Gulcehre, C., Cho, K., Ganguli, S., and Bengio, Y. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. *Advances in neural information processing systems*, 27, 2014.
- Dinh, L., Pascanu, R., Bengio, S., and Bengio, Y. Sharp minima can generalize for deep nets. In *International Conference on Machine Learning*, pp. 1019–1028. PMLR, 2017.
- Dziugaite, G. K. and Roy, D. M. Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. *arXiv preprint arXiv:1703.11008*, 2017.
- Foret, P., Kleiner, A., Mobahi, H., and Neyshabur, B. Sharpness-aware minimization for efficiently improving generalization. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=6TmlmposlrM>.
- Fort, S. and Ganguli, S. Emergent properties of the local geometry of neural loss landscapes. *arXiv preprint arXiv:1910.05929*, 2019.
- Fort, S., Dziugaite, G. K., Paul, M., Kharaghani, S., Roy, D. M., and Ganguli, S. Deep learning versus kernel learning: an empirical study of loss landscape geometry and the time evolution of the neural tangent kernel. *Advances in Neural Information Processing Systems*, 33: 5850–5861, 2020.
- Fulton, W. Eigenvalues, invariant factors, highest weights, and schubert calculus. *Bulletin of the American Mathematical Society*, 37(3):209–249, 2000.
- Gershgorin, S. A. Über die abgrenzung der eigenwerte einer matrix. *Известия Российской академии наук. Серия математическая*, 7(6): 749–754, 1931.
- Ghorbani, B., Krishnan, S., and Xiao, Y. An investigation into neural net optimization via hessian eigenvalue density. In *International Conference on Machine Learning*, pp. 2232–2241. PMLR, 2019.
- Golatkar, A. S., Achille, A., and Soatto, S. Time matters in regularizing deep networks: Weight decay and data augmentation affect early learning dynamics, matter little near convergence. *Advances in Neural Information Processing Systems*, 32, 2019.

- Golub, G. H. Some modified matrix eigenvalue problems. *Siam Review*, 15(2):318–334, 1973.
- Golub, G. H. and Van Loan, C. F. *Matrix computations*, volume 3. JHU press, 2013.
- Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., and He, K. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- Gur-Ari, G., Roberts, D. A., and Dyer, E. Gradient descent happens in a tiny subspace. *arXiv preprint arXiv:1812.04754*, 2018.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Hochreiter, S. and Schmidhuber, J. Flat minima. *Neural computation*, 9(1):1–42, 1997.
- Hoffer, E., Hubara, I., and Soudry, D. Train longer, generalize better: closing the generalization gap in large batch training of neural networks. *Advances in neural information processing systems*, 30, 2017.
- Hoffman, J., Roberts, D. A., and Yaida, S. Robust learning with jacobian regularization. *arXiv preprint arXiv:1908.02729*, 2019.
- Hu, W., Li, C. J., Li, L., and Liu, J.-G. On the diffusion approximation of nonconvex stochastic gradient descent. *Annals of Mathematical Sciences and Applications*, 4(1), 2019.
- Jastrzębski, S., Kenton, Z., Arpit, D., Ballas, N., Fischer, A., Bengio, Y., and Storkey, A. Three factors influencing minima in sgd. *arXiv preprint arXiv:1711.04623*, 2017.
- Jastrzębski, S., Kenton, Z., Ballas, N., Fischer, A., Bengio, Y., and Storkey, A. On the relation between the sharpest directions of DNN loss and the SGD step length. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=SkgeAj05t7>.
- Jastrzębski, S., Szymczak, M., Fort, S., Arpit, D., Tabor, J., Cho*, K., and Geras*, K. The break-even point on optimization trajectories of deep neural networks. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=r1g87C4KwB>.
- Jastrzębski, S., Arpit, D., Astrand, O., Kerg, G. B., Wang, H., Xiong, C., Socher, R., Cho, K., and Geras, K. J. Catastrophic fisher explosion: Early phase fisher matrix impacts generalization. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 4772–4784. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/jastrzebski21a.html>.
- Jiang, Y., Neyshabur, B., Mobahi, H., Krishnan, D., and Bengio, S. Fantastic generalization measures and where to find them. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=SJgIPJBFvH>.
- Keskar, N. S., Mudigere, D., Nocedal, J., Smelyanskiy, M., and Tang, P. T. P. On large-batch training for deep learning: Generalization gap and sharp minima. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL <https://openreview.net/forum?id=HloyRlYgg>.
- Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. *Master’s thesis, Department of Computer Science, University of Toronto*, 2009.
- Kwon, J., Kim, J., Park, H., and Choi, I. K. Asam: Adaptive sharpness-aware minimization for scale-invariant learning of deep neural networks, 2021.
- LeCun, Y. A., Bottou, L., Orr, G. B., and Müller, K.-R. Efficient backprop. In *Neural networks: Tricks of the trade*, pp. 9–48. Springer, 2012.
- Lee, S. and Jang, C. A new characterization of the edge of stability based on a sharpness measure aware of batch gradient distribution. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=bH-kCY6LdKg>.
- Lewkowycz, A., Bahri, Y., Dyer, E., Sohl-Dickstein, J., and Gur-Ari, G. The large learning rate phase of deep learning: the catapult mechanism. *arXiv preprint arXiv:2003.02218*, 2020.
- Li, C. J., Li, L., Qian, J., and Liu, J.-G. Batch size matters: A diffusion approximation framework on nonconvex stochastic gradient descent. *stat*, 1050:22, 2017.
- Li, Q., Tai, C., and Weinan, E. Stochastic modified equations and dynamics of stochastic gradient algorithms i: Mathematical foundations. *The Journal of Machine Learning Research*, 20(1):1474–1520, 2019a.
- Li, Y., Wei, C., and Ma, T. Towards explaining the regularization effect of initial large learning rate in training neural networks. *Advances in Neural Information Processing Systems*, 32, 2019b.

- Li, Z., Malladi, S., and Arora, S. On the validity of modeling sgd with stochastic differential equations (sdes). *Advances in Neural Information Processing Systems*, 34, 2021.
- Liu, C., Zhu, L., and Belkin, M. Toward a theory of optimization for over-parameterized systems of non-linear equations: the lessons of deep learning. *arXiv preprint arXiv:2003.00307*, 2020.
- Liu, H., Simonyan, K., and Yang, Y. DARTS: Differentiable architecture search. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=S1eYHoC5FX>.
- Loshchilov, I. and Hutter, F. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.
- Mandt, S., Hoffman, M., and Blei, D. A variational analysis of stochastic gradient algorithms. In *International conference on machine learning*, pp. 354–363. PMLR, 2016.
- Mandt, S., Hoffman, M. D., and Blei, D. M. Stochastic gradient descent as approximate bayesian inference. *Journal of Machine Learning Research*, 18:1–35, 2017.
- McInnes, L., Healy, J., Saul, N., and Großberger, L. Umap: Uniform manifold approximation and projection. *Journal of Open Source Software*, 3(29):861, 2018. doi: 10.21105/joss.00861. URL <https://doi.org/10.21105/joss.00861>.
- Müller, R., Kornblith, S., and Hinton, G. E. When does label smoothing help? *Advances in neural information processing systems*, 32, 2019.
- Nagarajan, V. and Kolter, J. Z. Generalization in deep networks: The role of distance from initialization. *arXiv preprint arXiv:1901.01672*, 2019.
- Neyshabur, B. Implicit regularization in deep learning. *arXiv preprint arXiv:1709.01953*, 2017.
- Neyshabur, B., Tomioka, R., and Srebro, N. In search of the real inductive bias: On the role of implicit regularization in deep learning. *arXiv preprint arXiv:1412.6614*, 2014.
- Neyshabur, B., Salakhutdinov, R., and Srebro, N. Path-sgd: Path-normalized optimization in deep neural networks. *arXiv preprint arXiv:1506.02617*, 2015a.
- Neyshabur, B., Tomioka, R., and Srebro, N. Norm-based capacity control in neural networks. In *Conference on Learning Theory*, pp. 1376–1401. PMLR, 2015b.
- Neyshabur, B., Bhojanapalli, S., and Srebro, N. A pac-bayesian approach to spectrally-normalized margin bounds for neural networks. *arXiv preprint arXiv:1707.09564*, 2017.
- Neyshabur, B., Li, Z., Bhojanapalli, S., LeCun, Y., and Srebro, N. Towards understanding the role of over-parametrization in generalization of neural networks. *arXiv preprint arXiv:1805.12076*, 2018.
- Novak, R., Bahri, Y., Abolafia, D. A., Pennington, J., and Sohl-Dickstein, J. Sensitivity and generalization in neural networks: an empirical study. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=HJC2SszZCW>.
- Pang, T., Yang, X., Dong, Y., Su, H., and Zhu, J. Bag of tricks for adversarial training. In *International Conference on Learning Representations*, 2020.
- Papayan, V. The full spectrum of deepnet Hessians at scale: Dynamics with sgd training and sample size. *arXiv preprint arXiv:1811.07062*, 2018.
- Papayan, V. Measurements of three-level hierarchical structure in the outliers in the spectrum of deepnet Hessians. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 5012–5021. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/papayan19a.html>.
- Papayan, V. Traces of class/cross-class structure pervade deep learning spectra. *Journal of Machine Learning Research*, 21(252):1–64, 2020.
- Park, D., Sohl-Dickstein, J., Le, Q., and Smith, S. The effect of network width on stochastic gradient descent and generalization: an empirical study. In *International Conference on Machine Learning*, pp. 5042–5051. PMLR, 2019.
- Rice, L., Wong, E., and Kolter, Z. Overfitting in adversarially robust deep learning. In *International Conference on Machine Learning*, pp. 8093–8104. PMLR, 2020.
- Robbins, H. and Monro, S. A stochastic approximation method. *The annals of mathematical statistics*, pp. 400–407, 1951.
- Sagun, L., Bottou, L., and LeCun, Y. Eigenvalues of the Hessian in deep learning: Singularity and beyond. *arXiv preprint arXiv:1611.07476*, 2016.

- Sagun, L., Evci, U., Guney, V. U., Dauphin, Y., and Bottou, L. Empirical analysis of the hessian of over-parametrized neural networks. *arXiv preprint arXiv:1706.04454*, 2017.
- Schraudolph, N. N. Fast curvature matrix-vector products for second-order gradient descent. *Neural computation*, 14(7):1723–1738, 2002.
- Simonyan, K. and Zisserman, A. Very deep convolutional networks for large-scale image recognition. In Bengio, Y. and LeCun, Y. (eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1409.1556>.
- Smith, S. L. and Le, Q. V. A bayesian perspective on generalization and stochastic gradient descent. In *International Conference on Learning Representations*, 2018.
- Smith, S. L., Dherin, B., Barrett, D., and De, S. On the origin of implicit regularization in stochastic gradient descent. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=rq_Qr0clHyo.
- Soudry, D., Hoffer, E., Nacson, M. S., Gunasekar, S., and Srebro, N. The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878, 2018.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- Wang, H., Keskar, N. S., Xiong, C., and Socher, R. Identifying generalization properties in neural networks. *arXiv preprint arXiv:1809.07402*, 2018.
- Wang, Z., Li, Z., and Li, J. Analyzing sharpness along GD trajectory: Progressive sharpening and edge of stability. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=thgItcQrJ4y>.
- Weyl, H. Das asymptotische verteilungsgesetz der eigenwerte linearer partieller differentialgleichungen (mit einer anwendung auf die theorie der hohlraumstrahlung). *Mathematische Annalen*, 71(4):441–479, 1912.
- Xing, C., Arpit, D., Tsirigotis, C., and Bengio, Y. A walk with sgd. *arXiv preprint arXiv:1802.08770*, 2018.
- Yaida, S. Fluctuation-dissipation relations for stochastic gradient descent. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=SkNksoRctQ>.
- Yao, Z., Gholami, A., Keutzer, K., and Mahoney, M. W. Py-hessian: Neural networks through the lens of the hessian. In *2020 IEEE International Conference on Big Data (Big Data)*, pp. 581–590. IEEE, 2020.
- Zagoruyko, S. and Komodakis, N. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.
- Zhang, H., Yu, Y., Jiao, J., Xing, E., El Ghaoui, L., and Jordan, M. Theoretically principled trade-off between robustness and accuracy. In *International Conference on Machine Learning*, pp. 7472–7482, 2019.
- Zhu, Z., Wu, J., Yu, B., Wu, L., and Ma, J. The anisotropic noise in stochastic gradient descent: Its behavior of escaping from sharp minima and regularization effects. In *International Conference on Machine Learning*, pp. 7654–7663. PMLR, 2019.

A. Notations

We summarize the notations for a quick reference.

| | |
|--|---|
| $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$ | input |
| $y \in [C] = \{1, 2, \dots, C\}$ | the corresponding label |
| $f_{\boldsymbol{\theta}} : \mathcal{X} \rightarrow \mathcal{Z} \subset \mathbb{R}^C$ | parameterized model |
| $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^m$ | model parameter |
| $\mathbf{z} = f_{\boldsymbol{\theta}}(\mathbf{x}) \in \mathcal{Z} \subset \mathbb{R}^C$ | logit vector |
| $\mathbf{p} = \text{softmax}(\mathbf{z}) \in \Delta^{C-1}$ | probability output |
| $\Delta^{C-1} = \{\mathbf{p} \in [0, 1]^C : \mathbf{1}^T \mathbf{p} = 1, \mathbf{p} \geq \mathbf{0}\}$ | probability simplex |
| $p = \mathbf{p}_y$ | probability value |
| $\mathbf{c} = \arg \text{sort}(\mathbf{p})$ | class output |
| $\ell = \ell(\mathbf{z}, y) \in \mathbb{R}$ | cross-entropy loss |
| $\ \mathbf{v}\ = \sqrt{\sum_i \mathbf{v}_i^2}$ | Euclidean norm of a vector \mathbf{v} |
| $\ \mathbf{A}\ = \max_{\ \mathbf{v}\ =1} \ \mathbf{A}\mathbf{v}\ $ | Euclidean operator norm of a matrix \mathbf{A} |
| $\ \mathbf{A}\ _F = \sqrt{\sum_{i,j} \mathbf{A}_{ij} ^2}$ | Frobenius norm of a matrix \mathbf{A} |
| $\ \mathbf{S}\ _{\sigma}$ | spectral norm of a square matrix \mathbf{S} |
| $\text{tr}(\mathbf{S}) = \sum_i \mathbf{S}_i i$ | trace of a square matrix \mathbf{S} |
| $\text{diag}(\mathbf{v}) = (\delta_{ij} \mathbf{v}_i)_{ij}$ | diagonal matrix with \mathbf{v} as its diagonal entries |
| $\mathbf{e}^i = (\delta_{ij})_j$ | one-hot vector with the i -th element as 1 |
| \mathcal{D} | training set |
| $L = \langle \ell \rangle$ | total training loss |
| $\langle \cdot \rangle = \langle \cdot \rangle_{\mathcal{D}} = \hat{\mathbb{E}}_{\mathcal{D}}[\cdot]$ | expectation over \mathcal{D} |
| $\langle \cdot \rangle_{\mathcal{B}} = \hat{\mathbb{E}}_{\mathcal{B}}[\cdot]$ | expectation over \mathcal{B} |
| \mathcal{B} | mini-batch |
| $\mathbf{H} = \langle \nabla_{\boldsymbol{\theta}}^2 \ell \rangle \in \mathbb{R}^{m \times m}$ | Hessian |
| $\mathbf{M} = \nabla_{\mathbf{z}}^2 \ell \in \mathbb{R}^{C \times C}$ | logit Hessian |
| $\mathbf{J} = \nabla_{\boldsymbol{\theta}} \mathbf{z} \in \mathbb{R}^{m \times C}$ | (logit-weight) Jacobian |
| $\mathbf{G} = \langle \mathbf{J} \mathbf{M} \mathbf{J}^T \rangle \in \mathbb{R}^{m \times m}$ | Gauss-Newton approximation |
| $\lambda^{(i)}$ | the i -th largest eigenvalues of \mathbf{M} |
| $\mathbf{q}^{(i)}$ | the corresponding eigenvector |
| $(i) = \mathbf{c}_i \in [C]$ | the ordered index |
| $\text{Gini}(q) = 2q(1 - q)$ | (binary) Gini impurity |
| $\text{Gini}^C(\mathbf{p}) = 1 - \sum_{i=1}^C \mathbf{p}_i^2$ | (C -ary) Gini impurity |
| $\eta > 0$ | learning rate |
| $B = \mathcal{B} $ | batch size |

B. Summary

The concepts in the paper, weight Hessian \mathbf{H} , logit Hessian \mathbf{M} , impurity $\lambda^{(1)} \equiv \|\mathbf{M}\|_\sigma$, sharpness $\|\mathbf{H}\|_\sigma$, logit Jacobian \mathbf{J} and Gauss-Newton matrix \mathbf{G} , connect with each other as follows:

$$\|\mathbf{H}\|_\sigma \stackrel{(i)}{\approx} \|\mathbf{G}\|_\sigma \stackrel{(ii)}{=} \langle \lambda^* \|\mathbf{J}\|^2 \rangle \stackrel{(iii)}{\leq} \langle \lambda^{(1)} \|\mathbf{J}\|^2 \rangle \stackrel{(iv)}{=} \langle \|\mathbf{M}\|_\sigma \|\mathbf{J}\|^2 \rangle,$$

where (i), (ii), (iii), (iv) come from the Gauss-Newton approximation, Theorem 5.1, $\lambda^* \leq \lambda^{(1)}$, and $\lambda^{(1)} \equiv \|\mathbf{M}\|_\sigma$, respectively.

C. Proofs

C.1. Proof of (1)

$$\nabla_{\mathbf{z}} \mathbf{p} = \text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^T \in \mathbb{R}^{C \times C} \quad (1)$$

Proof. By definition of the softmax function,

$$\mathbf{p}_j = [\text{softmax}(\mathbf{z})]_j = \frac{\exp(\mathbf{z}_j)}{\sum_k \exp(\mathbf{z}_k)} = \exp(\mathbf{z}_j) s^{-1} \quad (15)$$

where $s = \sum_k \exp(\mathbf{z}_k)$, we have

$$\nabla_{\mathbf{z}_i} \mathbf{p}_j = \begin{cases} -\exp(\mathbf{z}_j) s^{-2} \exp(\mathbf{z}_i) = -\mathbf{p}_i \mathbf{p}_j, & \text{if } i \neq j \\ -\exp(\mathbf{z}_i) s^{-2} \exp(\mathbf{z}_i) + \exp(\mathbf{z}_i) s^{-1} = -\mathbf{p}_i^2 + \mathbf{p}_i, & \text{if } i = j \end{cases} \quad (16)$$

which leads to $\nabla_{\mathbf{z}} \mathbf{p} = (\nabla_{\mathbf{z}_i} \mathbf{p}_j)_{ij} = -\mathbf{p}\mathbf{p}^T + \text{diag}(\mathbf{p})$. \square

C.2. Proof of (6)

We provide a self-contained proof of the first part (6) (see, for example, Schraudolph (2002) for more).

$$\nabla_{\boldsymbol{\theta}}^2 l = \nabla_{\boldsymbol{\theta}} \mathbf{z} \nabla_{\mathbf{z}}^2 l \nabla_{\boldsymbol{\theta}} \mathbf{z}^T + \sum_{c=1}^C \nabla_{\boldsymbol{\theta}}^2 \mathbf{z}_c \nabla_{\mathbf{z}_c} l \quad (6)$$

Proof. We apply the chain rule to obtain the following equations:

$$(\nabla_{\boldsymbol{\theta}}^2 l)_{i,j} = \frac{\partial}{\partial \boldsymbol{\theta}_j} \frac{\partial l}{\partial \boldsymbol{\theta}_i} = \frac{\partial}{\partial \boldsymbol{\theta}_j} \left(\sum_{c=1}^C \frac{\partial l}{\partial \mathbf{z}_c} \frac{\partial \mathbf{z}_c}{\partial \boldsymbol{\theta}_i} \right) \quad (17)$$

$$= \sum_{c=1}^C \frac{\partial}{\partial \boldsymbol{\theta}_j} \left(\frac{\partial l}{\partial \mathbf{z}_c} \frac{\partial \mathbf{z}_c}{\partial \boldsymbol{\theta}_i} \right) \quad (18)$$

$$= \sum_{c=1}^C \frac{\partial}{\partial \boldsymbol{\theta}_j} \left(\frac{\partial l}{\partial \mathbf{z}_c} \right) \frac{\partial \mathbf{z}_c}{\partial \boldsymbol{\theta}_i} + \sum_{c=1}^C \frac{\partial l}{\partial \mathbf{z}_c} \frac{\partial^2 \mathbf{z}_c}{\partial \boldsymbol{\theta}_j \partial \boldsymbol{\theta}_i} \quad (19)$$

$$= \sum_{c=1}^C \left(\sum_{k=1}^C \frac{\partial \mathbf{z}_c}{\partial \boldsymbol{\theta}_j} \frac{\partial^2 l}{\partial \mathbf{z}_c \partial \mathbf{z}_k} \right) \frac{\partial \mathbf{z}_c}{\partial \boldsymbol{\theta}_i} + \sum_{c=1}^C \frac{\partial^2 \mathbf{z}_c}{\partial \boldsymbol{\theta}_j \partial \boldsymbol{\theta}_i} \frac{\partial l}{\partial \mathbf{z}_c} \quad (20)$$

which leads to the conclusion. \square

Remark C.1. Assuming $\nabla_{\mathbf{z}_c} l$ are uncorrelated with $\nabla_{\boldsymbol{\theta}} \mathbf{z}_c^2$, we can approximate $\nabla_{\boldsymbol{\theta}}^2 l \approx \nabla_{\boldsymbol{\theta}} \mathbf{z}^T \nabla_{\mathbf{z}}^2 l \nabla_{\boldsymbol{\theta}} \mathbf{z}$ (Sagun et al., 2017; Fort & Ganguli, 2019).

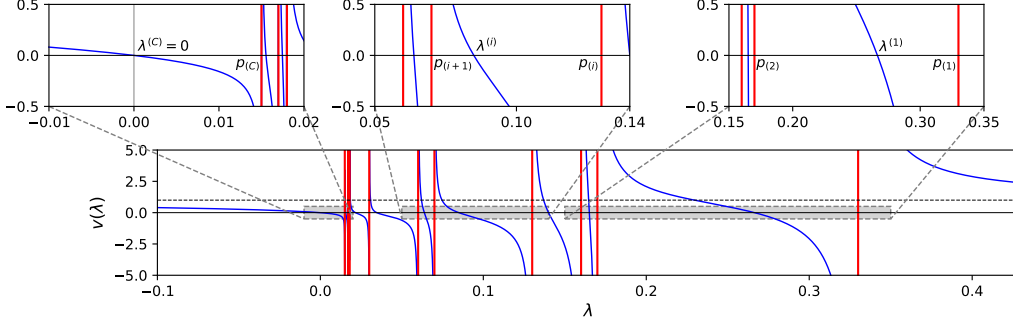


Figure 7: **Graph of secular function $v(\lambda)$ (blue curves) in (10) for some p having zeros at the eigenvalues $\{\lambda^{(i)}\}_{i=1}^C$ of $M = \nabla_z^2 l$.** We highlighted the singularities $\lambda = p^{(i)}$ with red vertical lines. The figure illustrates Theorem 4.1 (a) and (c).

C.3. Proof of Theorem 4.1

Theorem 4.1 (restated). The eigenvalues $\lambda^{(i)}$ ($\lambda^{(1)} \geq \lambda^{(2)} \geq \dots \geq \lambda^{(C)}$) and the corresponding eigenvectors $\mathbf{q}^{(i)}$ of the logit Hessian $M = \nabla_z^2 l = \text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^T$ satisfy the following properties:

(a) The eigenvalue $\lambda^{(i)}$ is the i -th largest solution of the following equation:

$$v(\lambda) = 1 - \sum_{i=1}^C \frac{\mathbf{p}_i^2}{\mathbf{p}_i - \lambda} = 0 \quad (21)$$

(b) The eigenvector $\mathbf{q}^{(i)}$ is aligned with the direction of $(\text{diag}(\mathbf{p}) - \lambda^{(i)}\mathbf{I})^{-1}\mathbf{p}$

(c) $\mathbf{p}_{(i+1)} \leq \lambda^{(i)} \leq \mathbf{p}_{(i)}$ for $1 \leq i \leq C-1$, and $\lambda^{(C)} = 0$

(d) $\frac{1}{2}\text{Gini}(\mathbf{p}_{(1)}) \leq \lambda^{(1)} \leq \text{Gini}(\mathbf{p}_{(1)})$ where $\text{Gini}(q) = 2q(1-q)$ is the Gini impurity for the binary case $(q, 1-q)$.

Proof. The eigenvalues $\lambda^{(i)}$ of $M = \text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^T$ are the zeros of the following characteristic polynomial:

$$\phi_M(\lambda) = \det(\text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^T - \lambda\mathbf{I}) \quad (22)$$

$$= \det(\text{diag}(\mathbf{p}) - \lambda\mathbf{I}) \det(\mathbf{I} - (\text{diag}(\mathbf{p}) - \lambda\mathbf{I})^{-1}\mathbf{p}\mathbf{p}^T) \quad (23)$$

$$= \prod_{i=1}^C (\mathbf{p}_i - \lambda) \left(1 - \sum_{j=1}^C \frac{\mathbf{p}_j^2}{\mathbf{p}_j - \lambda} \right) \quad (24)$$

where the second equality follows from $\mathbf{A} - \mathbf{p}\mathbf{p}^T = \mathbf{A}(\mathbf{I} - \mathbf{A}^{-1}\mathbf{p}\mathbf{p}^T)$ with the matrix $\mathbf{A} = \text{diag}(\mathbf{p}) - \lambda\mathbf{I}$, and the third inequality holds because $\det(\mathbf{I} + \mathbf{u}\mathbf{v}^T) = 1 + \mathbf{u}^T\mathbf{v}$ for vectors \mathbf{u} and \mathbf{v} . Then it is equivalent to solving the following equation:

$$v(\lambda) = 1 - \sum_{i=1}^C \frac{\mathbf{p}_i^2}{\mathbf{p}_i - \lambda} = 0 \quad (25)$$

which implies (a). Note that this result also implies (c).

Next, to prove (b), put $\mathbf{A} \equiv \text{diag}(\mathbf{p}) - \lambda\mathbf{I}$ and $\mathbf{q} \equiv \mathbf{A}^{-1}\mathbf{p}$. Then it is required to show that $(M - \lambda\mathbf{I})\mathbf{q} = (\mathbf{A} - \mathbf{p}\mathbf{p}^T)\mathbf{q} = \mathbf{0}$ for the eigenvalues $\lambda = \lambda^{(i)}$. We have

$$(\mathbf{A} - \mathbf{p}\mathbf{p}^T)\mathbf{q} = (\mathbf{A} - \mathbf{p}\mathbf{p}^T)\mathbf{A}^{-1}\mathbf{p} = \mathbf{p} - \mathbf{p}\mathbf{p}^T\mathbf{A}^{-1}\mathbf{p} \quad (26)$$

Here, $(\mathbf{p}\mathbf{p}^T\mathbf{A}^{-1}\mathbf{p})_i = \sum_{j,k} \mathbf{p}_i\mathbf{p}_j\mathbf{A}_{jk}^{-1}\mathbf{p}_k = \sum_{j,k} \mathbf{p}_i\mathbf{p}_j\delta_{jk}(\mathbf{p}_j - \lambda)^{-1}\mathbf{p}_k = \sum_k \mathbf{p}_i\mathbf{p}_k(\mathbf{p}_k - \lambda)^{-1}\mathbf{p}_k = \mathbf{p}_i \sum_k \mathbf{p}_k^2 / (\mathbf{p}_k - \lambda) = \mathbf{p}_i$. The last equality holds for the eigenvalues $\lambda = \lambda^{(i)}$ which follows from (a).

Now, we want to prove the statement (c). Since

$$\lambda^{(i)}(\mathbf{C}) \leq \lambda^{(j)}(\mathbf{A}) + \lambda^{(k)}(\mathbf{B}) \text{ if } k + j - i = 1 \quad (27)$$

$$\lambda^{(i)}(\mathbf{C}) \geq \lambda^{(j)}(\mathbf{A}) + \lambda^{(k)}(\mathbf{B}) \text{ if } k + j - i = C \quad (28)$$

for $\mathbf{C} = \mathbf{A} + \mathbf{B} \in \mathbb{R}^{C \times C}$ where $\lambda^{(i)}(\mathbf{D})$ is the i -th largest eigenvalue of a matrix \mathbf{D} (Weyl, 1912; Fulton, 2000), we can get $\lambda^{(i)}(\mathbf{C}) \leq \lambda^{(i)}(\mathbf{A}) + \lambda^{(1)}(\mathbf{B})$ and $\lambda^{(i)}(\mathbf{C}) \geq \lambda^{(i+1)}(\mathbf{A}) + \lambda^{(C-1)}(\mathbf{B})$. Thus, for $\mathbf{A} = \text{diag}(\mathbf{p})$ and $\mathbf{B} = -\mathbf{p}\mathbf{p}^T$, we can get

$$\mathbf{p}_{(i+1)} \leq \lambda^{(i)}(\mathbf{M}) \leq \mathbf{p}_{(i)} \text{ for } 1 \leq i \leq C - 1 \quad (29)$$

since $\lambda^{(i)}(\mathbf{A}) = \mathbf{p}_{(i)}$, $\lambda^{(i+1)}(\mathbf{A}) = \mathbf{p}_{(i+1)}$ and $\lambda^{(1)}(\mathbf{B}) = \lambda^{(C-1)}(\mathbf{B}) = 0$. Moreover, since $\mathbf{M}\mathbf{1} = \mathbf{p} - \mathbf{p}\mathbf{p}^T\mathbf{1} = \mathbf{p} - \mathbf{p}\sum_i \mathbf{p}_i = \mathbf{0}$, the smallest eigenvalue is $\lambda^{(C)} = 0$.

Lastly, we prove the statement (d). From the Gershgorin circle theorem (Gershgorin, 1931), we have

$$\lambda^{(1)} \in \bigcup_i \mathbb{B}(\mathbf{M}_{ii}, \sum_{j \neq i} |\mathbf{M}_{ij}|) = \mathbb{B}(\mathbf{p}_{(1)}(1 - \mathbf{p}_{(1)}), \mathbf{p}_{(1)}(1 - \mathbf{p}_{(1)})) = [0, 2\mathbf{p}_{(1)}(1 - \mathbf{p}_{(1)})] \quad (30)$$

which implies $\lambda^{(1)} \leq 2\mathbf{p}_{(1)}(1 - \mathbf{p}_{(1)})$. Note that $\mathbf{p}_{(1)}(1 - \mathbf{p}_{(1)}) \geq \mathbf{p}_{(i)}(1 - \mathbf{p}_{(i)})$ since $g(t) = t(1 - t)$ is increasing for $0 \leq t \leq 0.5$. In detail, if $\mathbf{p}_{(1)} \geq 0.5$, since $\mathbf{p}_{(i)} \leq 1 - \mathbf{p}_{(1)} \leq 0.5$, we have $g(\mathbf{p}_{(i)}) \leq g(1 - \mathbf{p}_{(1)}) = g(\mathbf{p}_{(1)})$. Otherwise ($\mathbf{p}_{(1)} < 0.5$), since $\mathbf{p}_{(i)} \leq \mathbf{p}_{(1)}$, it leads to the same inequality $g(\mathbf{p}_{(i)}) \leq g(\mathbf{p}_{(1)})$. With the Rayleigh principle, we can express the largest eigenvalue as $\lambda^{(1)} = \max_{\|\mathbf{u}\|_2=1} \mathbf{u}^T \mathbf{M} \mathbf{u}$, and thus $e^{(1)T} \mathbf{M} e^{(1)} = \mathbf{M}_{(1)(1)} = \mathbf{p}_{(1)}(1 - \mathbf{p}_{(1)}) \leq \lambda^{(1)}$. \square

C.4. Proof of Theorem 5.1

Theorem 5.1 (restated). For some lower bound $0 \leq \lambda^* \leq \lambda^{(1)}$ of the impurity for each $\mathbf{x} \in \mathcal{D}$, we have

$$\|\mathbf{G}\|_\sigma = \langle \lambda^* \|\mathbf{J}\|^2 \rangle. \quad (31)$$

Proof. We start with the Rayleigh principle:

$$\|\langle \mathbf{J} \mathbf{M} \mathbf{J}^T \rangle\|_\sigma = \max_{\|\mathbf{q}\|=1} \mathbf{q}^T \langle \mathbf{J} \mathbf{M} \mathbf{J}^T \rangle \mathbf{q} = \max_{\|\mathbf{q}\|=1} \langle \mathbf{q}^T \mathbf{J} \mathbf{M} \mathbf{J}^T \mathbf{q} \rangle \quad (32)$$

Since $\mathbf{M} = \sum_i \lambda^{(i)} \mathbf{q}^{(i)} \mathbf{q}^{(i)T}$, we can continue by putting $\mathbf{v} = \mathbf{J}^T \mathbf{q}$, and then

$$(32) = \max_{\|\mathbf{q}\|=1} \langle \mathbf{v}^T \mathbf{M} \mathbf{v} \rangle = \max_{\|\mathbf{q}\|=1} \left\langle \sum_i \lambda^{(i)} (\mathbf{q}^{(i)T} \mathbf{v})^2 \right\rangle \quad (33)$$

Then by putting $\tilde{\lambda} = \sum_i \gamma^{(i)} \lambda^{(i)}$ with $\gamma^{(i)} = (\mathbf{q}^{(i)T} \mathbf{v})^2 / \sum_i (\mathbf{q}^{(i)T} \mathbf{v})^2 \geq 0$ ($\sum_i \gamma_i = 1$),

$$(33) = \max_{\|\mathbf{q}\|=1} \langle \tilde{\lambda} \sum_i (\mathbf{q}^{(i)T} \mathbf{v})^2 \rangle = \max_{\|\mathbf{q}\|=1} \langle \tilde{\lambda} \sum_i (\mathbf{q}^{(i)T} \mathbf{J}^T \mathbf{q})^2 \rangle \quad (34)$$

$$= \max_{\|\mathbf{q}\|=1} \langle \tilde{\lambda} \sum_i (\mathbf{q}^{(i)T} (\mathbf{J}^T \mathbf{q}))^2 \rangle \quad (35)$$

Since $\{\mathbf{q}^{(i)}\}_{i=1}^C$ is an orthonormal basis of \mathbb{R}^C (eigenvectors of a symmetric matrix \mathbf{M}), we have the following by putting $\lambda^* = \frac{\|\mathbf{J}^T \mathbf{q}^*\|^2}{\|\mathbf{J}\|^2} \tilde{\lambda}$ with $\mathbf{q}^* = \arg \max_{\|\mathbf{q}\|=1} \langle \tilde{\lambda} \|\mathbf{J}^T \mathbf{q}\|^2 \rangle$,

$$(35) = \max_{\|\mathbf{q}\|=1} \langle \tilde{\lambda} \|\mathbf{J}^T \mathbf{q}\|^2 \rangle = \langle \tilde{\lambda} \|\mathbf{J}^T \mathbf{q}^*\|^2 \rangle = \langle \lambda^* \|\mathbf{J}\|^2 \rangle \quad (36)$$

Because of the definition of λ^* and $\tilde{\lambda}$ with $\|\mathbf{q}^*\| = 1$ and $\sum_i \gamma^{(i)} = 1$ ($\gamma^{(i)} \geq 0$), we have

$$\lambda^* \leq \tilde{\lambda} \leq \lambda^{(1)}. \quad (37)$$

\square

C.5. Proof of Proposition 5.2

Proposition 5.2 (restated). *For SGD with a quadratic loss, the expected loss decreases when $\|\mathbf{H}\|_\sigma \leq \frac{2\rho_B}{\eta}$, where η is learning rate, B is batch size and $\rho_B \equiv \frac{\|\langle \nabla_\theta l \rangle_B\|^2}{\mathbb{E}_{|B|=B}[\|\langle \nabla_\theta l \rangle_B\|^2]} \leq 1$.*

Proof. Put $\mathbf{g} = \nabla_\theta l(\boldsymbol{\theta})$. With SGD update $\boldsymbol{\theta}' = \boldsymbol{\theta} - \eta \langle \mathbf{g} \rangle_B$ and $L' = L(\boldsymbol{\theta}')$, we can obtain

$$L' = L - \eta \nabla L^T \langle \mathbf{g} \rangle_B + \frac{\eta^2}{2} \langle \mathbf{g} \rangle_B^T \mathbf{H} \langle \mathbf{g} \rangle_B \quad (38)$$

Therefore, the expected loss at the next step is

$$\mathbb{E}[L'] = L - \eta \nabla L^T \mathbb{E}[\langle \mathbf{g} \rangle_B] + \frac{\eta^2}{2} \mathbb{E}[\langle \mathbf{g} \rangle_B^T \mathbf{H} \langle \mathbf{g} \rangle_B] \quad (39)$$

$$\leq L - \eta \|\nabla L\|^2 + \frac{\eta^2}{2} \mathbb{E}[\|\langle \mathbf{g} \rangle_B\|^2] \|\mathbf{H}\|_\sigma \quad (40)$$

When $\|\mathbf{H}\|_\sigma \leq \frac{2}{\eta} \frac{\|\nabla L\|^2}{\mathbb{E}[\|\langle \mathbf{g} \rangle_B\|^2]} = \frac{2\rho_B}{\eta}$, i.e., $-\eta \|\nabla L\|^2 + \frac{\eta^2}{2} \mathbb{E}[\|\langle \mathbf{g} \rangle_B\|^2] \|\mathbf{H}\|_\sigma \leq 0$, the expected loss decreases, and the iterate stays within the quadratic basin and does not diverge. \square

Remark C.2. A necessary and sufficient condition for the expected loss to decrease is that $-\eta \|\nabla L\|^2 + \frac{\eta^2}{2} \mathbb{E}[\langle \mathbf{g} \rangle_B^T \mathbf{H} \langle \mathbf{g} \rangle_B] \leq 0$ from (39), i.e., $\frac{\mathbb{E}[\langle \mathbf{g} \rangle_B^T \mathbf{H} \langle \mathbf{g} \rangle_B]}{\mathbb{E}[\|\langle \mathbf{g} \rangle_B\|^2]} \leq \frac{2\rho_B}{\eta}$. Thus, when the batch gradients are aligned with the sharpest direction, the condition is equivalent to $\|\mathbf{H}\|_\sigma \leq \frac{2\rho_B}{\eta}$.

Remark C.3. We can qualitatively analyze the norm of the batch gradient $\langle \mathbf{g} \rangle_B$ and the ratio ρ_B with respect to the batch size B . First, if $B = N = |\mathcal{D}|$ (full-batch), then $\rho_B = 1$ by definition. Second, if B is large enough, then $\rho_B \approx 1$ since $\mathbb{E}_{|B|=B}[\|\langle \mathbf{g} \rangle_B\|^2] \approx \|\langle \mathbf{g} \rangle\|^2$ (cf. $B = 512$ in Figure 24). Now, we consider the third case of $B \ll N$. Put $\mathbf{g}_B = \langle \mathbf{g} \rangle_B$ and $\mathbf{g}_D = \langle \mathbf{g} \rangle$. Then, we have

$$\mathbb{E}[\|\mathbf{g}_B - \mathbf{g}_D\|^2] = \mathbb{E}[(\mathbf{g}_B - \mathbf{g}_D)^T (\mathbf{g}_B - \mathbf{g}_D)] = \mathbb{E}[\mathbf{g}_B^T \mathbf{g}_B - \mathbf{g}_B^T \mathbf{g}_D - \mathbf{g}_D^T \mathbf{g}_B + \mathbf{g}_D^T \mathbf{g}_D] = \mathbb{E}[\|\mathbf{g}_B\|^2] - \|\mathbf{g}_D\|^2 \quad (41)$$

and

$$\mathbb{E}[\|\mathbf{g}_B - \mathbf{g}_D\|^2] = \frac{1}{B} \frac{N - B}{N - 1} \mathbb{E}[\|\mathbf{g} - \mathbf{g}_D\|^2] \quad (42)$$

(see Appendix A in Smith et al. (2021) for the detailed proof). Therefore,

$$\frac{\mathbb{E}[\|\langle \mathbf{g} \rangle_B\|^2]}{\|\langle \mathbf{g} \rangle\|^2} = \frac{\mathbb{E}[\|\mathbf{g}_B\|^2]}{\|\mathbf{g}_D\|^2} = \frac{\|\mathbf{g}_D\|^2 + \mathbb{E}[\|\mathbf{g}_B - \mathbf{g}_D\|^2]}{\|\mathbf{g}_D\|^2} = 1 + \frac{\mathbb{E}[\|\mathbf{g}_B - \mathbf{g}_D\|^2]}{\|\mathbf{g}_D\|^2} \quad (43)$$

$$= 1 + \frac{1}{B} \frac{N - B}{N - 1} \frac{\mathbb{E}[\|\mathbf{g} - \mathbf{g}_D\|^2]}{\|\mathbf{g}_D\|^2} = 1 + \frac{1}{B} A = \frac{A + B}{B} \geq 1 \quad (44)$$

which becomes larger as B decreases ($B \ll N$) where $A = \frac{N - B}{N - 1} \frac{\mathbb{E}[\|\mathbf{g} - \mathbf{g}_D\|^2]}{\|\mathbf{g}_D\|^2} > 0$. In other words, the smaller the batch size B , the smaller the ratio $\rho_B = \frac{\|\langle \mathbf{g} \rangle_B\|^2}{\mathbb{E}[\|\langle \mathbf{g} \rangle_B\|^2]} = \frac{B}{A + B}$. Thus, if the batch gradients are diverse such that $\frac{\mathbb{E}[\|\mathbf{g} - \mathbf{g}_D\|^2]}{\|\mathbf{g}_D\|^2} \gg B$ and $A \gg B$, then we have $\rho_B \propto B$ which leads to the sharpness $\propto \frac{B}{\eta}$, a similar result with the linear scaling rule (Goyal et al., 2017). We refer the readers to Lee & Jang (2023) for further detailed analysis.

D. Gradient descent in the top Hessian subspace

Gur-Ari et al. (2018) showed that the gradient of the loss quickly converges to a tiny subspace spanned by a few top eigenvectors of the Hessian after a short training. Then, the top Hessian subspace does not evolve much, which implies gradient descent happens in a tiny subspace. However, the underlying mechanism has not been fully understood.

Direction of $\mathbf{q}^{(i)}$ and two salient elements We investigate the direction of the eigenvector $\mathbf{q}^{(i)}$ ($1 \leq i \leq C - 1$) of M . The eigenvector

$$\mathbf{q}^{(i)} = \alpha \left(\frac{\mathbf{p}_j}{\mathbf{p}_j - \lambda^{(i)}} \right)_j \quad (45)$$

can be obtained from Theorem 4.1 (b) for some $\alpha > 0$. Here, the magnitude of the denominator $|\mathbf{p}_j - \lambda^{(i)}|$ is small for the two indices $j = (i), (i + 1)$, and is large for the others. This is because the eigenvalue $\lambda^{(i)}$ lies between $\mathbf{p}_{(i+1)}$ and $\mathbf{p}_{(i)}$ (Theorem 4.1 (c)). Therefore, the eigenvector $\mathbf{q}^{(i)}$ has a relatively large positive value in $\mathbf{q}_{(i)}^{(i)}$ and a large negative value in $\mathbf{q}_{(i+1)}^{(i)}$ compared to the other components.

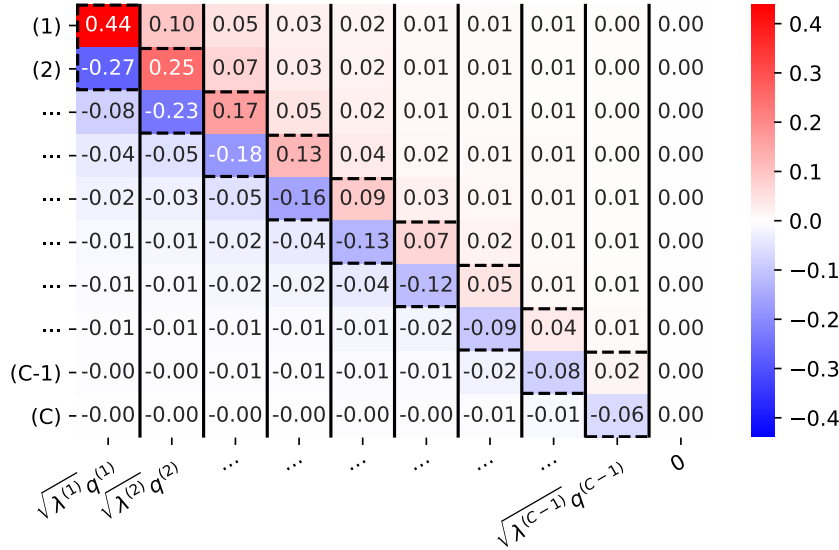


Figure 8: Heatmap of the matrix $\mathbf{Q}\Lambda^{1/2} = [\sqrt{\lambda^{(1)}}\mathbf{q}^{(1)}; \dots; \sqrt{\lambda^{(C-1)}}\mathbf{q}^{(C-1)}; \mathbf{0}]$ averaged over the training set \mathcal{D} where $M = \mathbf{Q}\Lambda\mathbf{Q}^T$. Each column of $\mathbf{Q}\Lambda^{1/2}$ visualizes the color-encoded direction of $\mathbf{q}^{(i)}$ multiplied by $\sqrt{\lambda^{(i)}}$. We highlighted the elements $\mathbf{q}_{(i)}^{(i)}$ and $\mathbf{q}_{(i+1)}^{(i)}$ with the dashed boxes for $0 \leq i \leq C - 1$ (see Theorem 4.1 (b)).

Figure 8 shows the directions of $\mathbf{q}^{(i)}$ with the heatmap of the matrix $\mathbf{Q}\Lambda^{1/2}$ where $\mathbf{Q}\Lambda^{1/2} = [\sqrt{\lambda^{(1)}}\mathbf{q}^{(1)}; \dots; \sqrt{\lambda^{(C-1)}}\mathbf{q}^{(C-1)}; \mathbf{0}] \in \mathbb{R}^{C \times C}$ for $M = \mathbf{Q}\Lambda\mathbf{Q}^T$. As expected, considering each column of $\mathbf{Q}\Lambda^{1/2}$, the eigenvector $\mathbf{q}^{(i)}$ is colored in red (+) at $\mathbf{q}_{(i)}^{(i)}$ and in blue (-) at $\mathbf{q}_{(i+1)}^{(i)}$ for $1 \leq i \leq C - 1$. The two salient elements are highlighted with the dashed boxes.

Direction of $\mathbf{J}\mathbf{q}^{(i)}$ and margin maximization In light of the previous discussion, the direction of

$$\mathbf{J}\mathbf{q}^{(i)} = (\mathbf{J}\mathbf{q}^{(i)})|_{\theta=\theta^{(t)}} = \nabla_{\theta} \left(\mathbf{q}^{(i)}(\theta^{(t)})^T \mathbf{z}(\theta) \right) |_{\theta=\theta^{(t)}} \in \mathbb{R}^m \quad (46)$$

is approximately a direction maximizing $\mathbf{q}_{(i)}^{(i)}\mathbf{z}_{(i)} + \mathbf{q}_{(i+1)}^{(i)}\mathbf{z}_{(i+1)}$ at the current parameter $\theta^{(t)} \in \Theta$ because the other terms are relatively small. In other words, it tends to maximize the margin $\mathbf{z}_{(i)} - \mathbf{z}_{(i+1)}$ in the logit space \mathcal{Z} between the two classes (i) and $(i + 1)$ (see Figure 8). In particular, $\mathbf{J}\mathbf{q}^{(1)}$ is approximately a direction that maximizes the margin between the most likely class and the second most likely class.

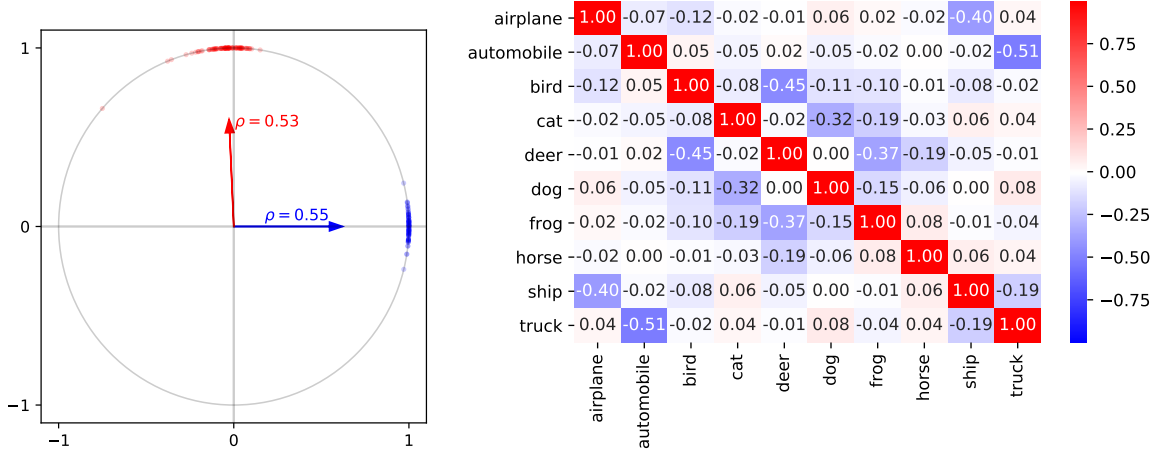


Figure 9: **There are C clusters of $\mathbf{m} = \sqrt{\lambda^{(1)}} \mathbf{J} \mathbf{q}^{(1)}$ according to the most likely class (not the true class).** **Left: (Within-class similarity)** Directional data of \mathbf{m} from \mathcal{D}_i and \mathcal{D}_j for the two classes, dog (blue) and automobile (red). They are projected onto the 2D-plane spanned by the two mean vectors indicated with the arrows. We highlight the MRL ρ for each class. The directional data of \mathbf{m} are concentrated within the class but separated from each other. **Right: (Between-class dissimilarity)** Cosine similarities between each pair of $\{\bar{\mathbf{m}}^i\}$. They are mostly orthogonal, but some pairs are even negatively aligned, for example, automobile and truck. This is because the examples predicted to be automobile mostly have the second most probable class as truck.

Clustering of $\mathbf{J} \mathbf{q}^{(1)}$ and the most probable class We first define following subsets of the training set according to the most probable class (and the second most one): $\mathcal{D}_i = \{\mathbf{x} \in \mathcal{D} : \mathbf{c}_1(\mathbf{x}) = i\} \subset \mathcal{D}$ and $\mathcal{D}_{ij} = \{\mathbf{x} \in \mathcal{D} : \mathbf{c}_1(\mathbf{x}) = i, \mathbf{c}_2(\mathbf{x}) = j\} \subset \mathcal{D}$ for $i \neq j \in [C]$. Note that $\mathcal{D}_i = \bigcup_{j \neq i} \mathcal{D}_{ij}$. Given two examples from \mathcal{D}_{ij} , their $\mathbf{J} \mathbf{q}^{(1)}$ are expected to be highly aligned to each other. This is because the direction of $\mathbf{J} \mathbf{q}^{(1)}$ is approximately a direction of maximizing the margin and of learning the features to discriminate the class i from the class j . Moreover, two examples from \mathcal{D}_i also have highly-aligned $\mathbf{J} \mathbf{q}^{(1)}$. Figure 9 (Left) shows the concentration of the directional data of $\mathbf{J} \mathbf{q}^{(1)}$ from \mathcal{D}_i . We also compute the mean resultant length (MRL) to measure the concentration. The MRL ρ of the directional variable $V \in \mathbb{S}^{m-1} \equiv \{\mathbf{v} \in \mathbb{R}^m : \|\mathbf{v}\| = 1\}$ defined as $\rho \equiv \|\mathbb{E}[V]\| \in [0, 1]$ indicates how V is distributed (the higher, the more concentrated).

Now, we focus on $\mathbf{m} \equiv \sqrt{\lambda^{(1)}} \mathbf{J} \mathbf{q}^{(1)}$ as the other $\lambda^{(i)}$ -terms are dominated by the $\lambda^{(1)}$ -term after a few epochs (see Appendix M for details). Then, we follow a similar approach from Papyan (2019) and provide the following equation:

$$\langle \mathbf{m} \mathbf{m}^T \rangle = \sum_{i=1}^C \gamma_i \langle \mathbf{m} \mathbf{m}^T \rangle_{\mathcal{D}_i} = \sum_{i=1}^C \gamma_i (\bar{\mathbf{m}}^i \bar{\mathbf{m}}^{iT} + \langle (\mathbf{m} - \bar{\mathbf{m}}^i)(\mathbf{m} - \bar{\mathbf{m}}^i)^T \rangle_{\mathcal{D}_i}) \quad (47)$$

where $\gamma_i = |\mathcal{D}_i|/|\mathcal{D}|$ and $\bar{\mathbf{m}}^i = \langle \mathbf{m} \rangle_{\mathcal{D}_i}$. Here, the covariance term $\langle (\mathbf{m} - \bar{\mathbf{m}}^i)(\mathbf{m} - \bar{\mathbf{m}}^i)^T \rangle_{\mathcal{D}_i}$ is weak as \mathbf{m} is concentrated within \mathcal{D}_i , and thus we can roughly approximate \mathbf{G} with $\sum_{i=1}^C \gamma_i \bar{\mathbf{m}}^i \bar{\mathbf{m}}^{iT}$. This implies that the top eigensubspace of the Hessian highly overlaps with the at most C -dimensional subspace spanned by $\{\bar{\mathbf{m}}^i\}_{i=1}^C$. Figure 9 (Right) demonstrates that the mean vectors $\{\bar{\mathbf{m}}^i\}_{i=1}^C$ are well separated from each other. This also implies the outliers in the Hessian spectrum (Sagun et al., 2016; 2017).

Why gradient descent happens mostly in the top Hessian subspace? Given input \mathbf{x} , after the model becomes to correctly predict the true label y , the gradient descent direction $-\mathbf{g} = \mathbf{J}(e^y - \mathbf{p})$ used in the training tends to be highly aligned with $\mathbf{J} \mathbf{q}^{(1)}$. This is because $e^y - \mathbf{p}$ and $\mathbf{q}^{(1)}$ both have similar direction. They have positive values $1 - p_y$ and $q_{(1)}^{(1)}$ in $y(=1)$ -th element, negative value $-p_i$ and $q_i^{(1)}$ in the others, and especially large negative value for the second most probable class $i = (2)$. Figure 10 (Middle) shows the cosine similarity between the gradient descent direction $-\mathbf{g}$ and $\mathbf{J} \mathbf{q}^{(1)}$. As expected, they are highly aligned with the cosine similarity near 1 as the two vectors $e^y - \mathbf{p}$ and $\mathbf{q}^{(1)}$ become more aligned to each other.

Next, we move on to the subspace $\mathcal{S} \equiv \text{span}(\{\bar{\mathbf{m}}^i\}_{i=1}^C)$ spanned by $\{\bar{\mathbf{m}}^i\}_{i=1}^C$. As each $\mathbf{m} = \sqrt{\lambda^{(1)}} \mathbf{J} \mathbf{q}^{(1)}$ is highly aligned with $-\mathbf{g}$, it is reasonably expected that the total gradient $\mathbf{g}^{\mathcal{D}}$ lies in the subspace \mathcal{S} . To measure how much the vector $\mathbf{v} \neq \mathbf{0}$

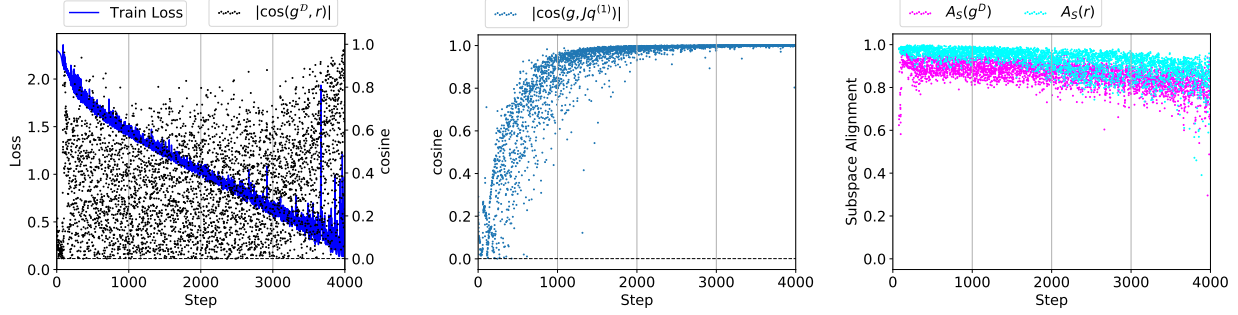


Figure 10: **Left: Total gradient $g^D = \langle g \rangle$ is aligned with the top eigenvector r of the Hessian H at each step during training (Jastrzębski et al., 2019; Gur-Ari et al., 2018).** They have large cosine similarities considering that they are very high-dimensional. We highlighted the cosine value for random m -dimensional vectors in Θ with the dashed horizontal line (about $1e-3$). **Middle: $Jq^{(1)}$ (or m) is highly aligned with the gradient g for given example at each step during training.** They have cosine similarities near 1 as the model becomes to correctly predict the true label. See Figure 1 (Right) together. **Right: Total gradient g^D and the top eigenvector r of the Hessian H mostly lie in the at most C -dimensional subspace \mathcal{S} spanned by $\{\bar{m}^i\}_{i=1}^C$.** The subspace alignment measure $A_{\mathcal{S}}$ is defined in (48).

is aligned with the subspace \mathcal{S} , we define the cosine similarity between the vector v and its projection $P_{\mathcal{S}}(v)$ onto the subspace \mathcal{S} as follows:

$$A_{\mathcal{S}}(v) \equiv \cos(v, P_{\mathcal{S}}(v)) = \|P_{\mathcal{S}}(v)\|/\|v\| \quad (48)$$

In particular, $A_{\mathcal{S}}(v) = 0$ means $v \in \mathcal{S}^{\perp}$ and $A_{\mathcal{S}}(v) = 1$ means $v \in \mathcal{S}$. Figure 10 (Right) shows the high alignment of the total gradient g^D in the subspace $\mathcal{S} \equiv \text{span}(\{\bar{m}^i\}_{i=1}^C)$ although the subspace \mathcal{S} is of dimension at most $C \ll m$. Moreover, since G can be roughly approximated by $\sum_{i=1}^C \gamma^i \bar{m}^i \bar{m}^{iT}$, the top eigenvector r of the Hessian H is also highly aligned with the subspace \mathcal{S} .

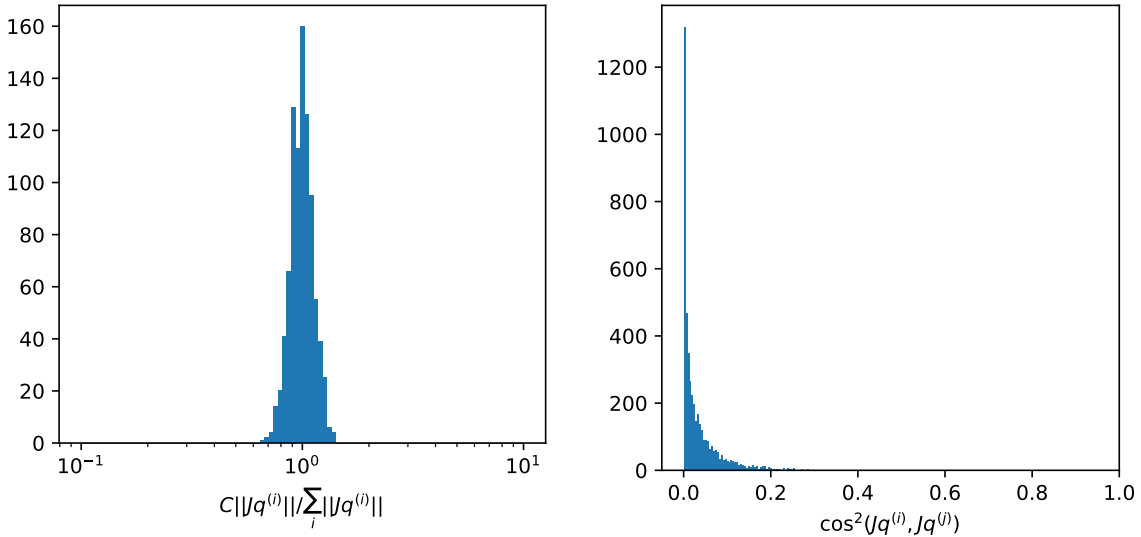


Figure 11: **Histograms of $\frac{\|Jq^{(i)}\|}{\sum_j \|Jq^{(j)}\|/C}$ and $\cos^2(Jq^{(i)}, Jq^{(j)})$.**

E. Experimental settings

We summarize the experimental settings for each Figure and Table in the main text in Table 3 and 4.

E.1. Data

We use the CIFAR-10 dataset ((Krizhevsky & Hinton, 2009), <https://www.cs.toronto.edu/~kriz/cifar.html>) and the MNIST dataset which have $C = 10$ number of classes. We also conduct some experiments on the CIFAR-100 dataset with the number of classes $C = 100$. We sometimes do not use the data augmentation for training (1) not to introduce the randomness in the training loss and (2) to allow the training loss to converge to a small value.

E.2. Network architectures

We use the following models: VGG-11 (VGG) (Simonyan & Zisserman, 2015) without batch-normalization, VGG for CIFAR-100 (VGG-CIFAR-100), ResNet-20 (ResNet) (He et al., 2016) without batch-normalization, a 6-layer CNN (6CNN), SimpleCNN used in Jastrzębski et al. (2021) (SimpleCNN) two 3-layer fully-connected networks (3FCN-CIFAR and 3FCN-MNIST), and two WRNs (WRN-28-2/WRN-28-10) for CIFAR-10 and CIFAR-100 with the number of model parameters, $m = 9750922, 9797092, 268346, 511926, 361706, 656810, 199210, 1467610/36479194, 1479220/36536884$, respectively.

We use a modified version of the implementation of VGG-11 from <https://github.com/chengyangfu/pytorch-vgg-cifar10/blob/master/vgg.py> without the dropout layers and ResNet-20 from https://github.com/locuslab/edge-of-stability/blob/github/src/resnet_cifar.py. We change the last linear layer for the CIFAR-100 dataset. The 6CNN model can be expressed in the Pytorch code as follows:

```
nn.Sequential(
    nn.Conv2d(3, 32, 3, stride=1, padding=1, bias=False)
    nn.ReLU(),
    nn.Conv2d(32, 32, 4, stride=2, padding=1, bias=False)
    nn.ReLU(),
    nn.Conv2d(32, 64, 3, stride=1, padding=1, bias=False)
    nn.ReLU(),
    nn.Conv2d(64, 64, 4, stride=2, padding=1, bias=False)
    nn.ReLU(),
    nn.Flatten(),
    nn.Linear(4096, 100, bias=True),
    nn.ReLU(),
    nn.Linear(100, 10, bias=True),
)
```

and the 3FCN architecture is as follows:

```
nn.Sequential(
    nn.Flatten(),
    nn.Linear(n, 200, bias=True),
    nn.ReLU(),
    nn.Linear(200, 200, bias=True),
    nn.ReLU(),
    nn.Linear(200, 10, bias=True),
)
```

where $n=784$ for 3FCN-MNIST, and $n=3072$ for 3FCN-CIFAR (the same one used in Cohen et al. (2021)).

E.3. Hyperparameters

SGD (Robbins & Monro, 1951) with the learning rate η can be expressed as follows:

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta \langle \mathbf{g}(\boldsymbol{\theta}^{(t)}) \rangle_{B^{(t)}} \quad (49)$$

where $\theta^{(t)}$ is the model parameter, $\mathcal{B}^{(t)} \subset \mathcal{D}$ is the training batch at t -th step, and $\mathbf{g}(\theta) = \nabla_{\theta} l$. For a simple CNN, we use the simplest form of SGD as described in (49) without momentum, weight decay, and learning rate decay. However, for a larger network, we use the momentum of 0.9, the weight decay of 0.0005, and two types of learning rate schedule, the multi-step scheduler with the learning rate decay of 0.2 in [30%, 60%, 80%] of the whole training epochs and the cosine annealing (Loshchilov & Hutter, 2016).

E.4. Hessian

When computing the top eigenvalue $\|\mathbf{H}\|_{\sigma}$ and the corresponding eigenvector \mathbf{r} of the Hessian \mathbf{H} , we use the tool developed in PyHessian ((Yao et al., 2020), <https://github.com/amirgholami/PyHessian>, MIT License) based on power iteration method using a small subset (5-25%) of training dataset \mathcal{D} .

E.5. Power iteration algorithm

Even though we have the secular function $v(\lambda)$ in (10) and the algorithms for computing the eigenvectors (Bunch et al., 1978), we use the power iteration in Algorithm 1 to get the top eigenvalue $\lambda^{(1)}$ of the logit Hessian $\mathbf{M} \in \mathbb{R}^{C \times C}$ since we can run the algorithm for a mini-batch in parallel. Then, we can compute the corresponding top eigenvector $\mathbf{q}^{(1)}$ from Theorem 4.1 (b). To compute the second largest eigenvalue, we apply the power iteration to $\mathbf{M}' \equiv \mathbf{M} - \lambda^{(1)} \mathbf{q}^{(1)} \mathbf{q}^{(1)T}$ instead of \mathbf{M} after computing $\lambda^{(1)}$ and $\mathbf{q}^{(1)}$.

Algorithm 1 Power iteration

input matrix \mathbf{M} , maximum iteration n_{\max} , tolerance bound ϵ

output the spectral norm $\lambda^{(1)} = \|\mathbf{M}\|_{\sigma}$ of the matrix \mathbf{M}

Initialize $\mathbf{u} \in \mathbb{R}^C$ with a random vector.

$i \leftarrow 0$

repeat

$\mathbf{v} \leftarrow \mathbf{M}\mathbf{u} / \|\mathbf{M}\mathbf{u}\|$

$\mathbf{u} \leftarrow \mathbf{M}^T \mathbf{v} / \|\mathbf{M}^T \mathbf{v}\|$

$i \leftarrow i + 1$

until it converges within the tolerance bound ϵ or $i \geq n_{\max}$

return $\mathbf{v}^T \mathbf{M} \mathbf{u}$

We also compute the operator norm of the Jacobian $\|\langle \mathbf{J} \rangle\|$ with the power iteration as in Algorithm 2. It requires $(C+1)$ -times scalar function differentiations with respect to θ for each iteration.

Algorithm 2 Power iteration for the Jacobian

input logit function \mathbf{z}_{θ} (not matrix $\langle \mathbf{J} \rangle$), maximum iteration n_{\max} , tolerance bound ϵ

output the operator norm $\|\langle \mathbf{J} \rangle\|$ of the Jacobian

Initialize $\mathbf{u} \in \mathbb{R}^C$ with a random vector.

$i \leftarrow 0$

repeat

Compute $\langle \mathbf{J} \rangle \mathbf{u} = \nabla_{\theta} \langle \mathbf{u}^T \mathbf{z} \rangle$ ($1 \times$ scalar function differentiation)

$\mathbf{v} \leftarrow \langle \mathbf{J} \rangle \mathbf{u} / \|\langle \mathbf{J} \rangle \mathbf{u}\|$

Compute $\langle \mathbf{J} \rangle^T \mathbf{v} = (\mathbf{v}^T \langle \mathbf{J} \rangle)^T = [\mathbf{v}^T \nabla_{\theta} \langle \mathbf{z}_1 \rangle, \dots, \mathbf{v}^T \nabla_{\theta} \langle \mathbf{z}_C \rangle]^T$ ($C \times$ scalar function differentiations)

$\mathbf{u} \leftarrow \langle \mathbf{J} \rangle^T \mathbf{v} / \|\langle \mathbf{J} \rangle^T \mathbf{v}\|$

$i \leftarrow i + 1$

until it converges within the tolerance bound ϵ or $i \geq n_{\max}$

return $\mathbf{v}^T (\langle \mathbf{J} \rangle \mathbf{u})$

Table 3: Experimental settings for each Figure and Table in the main text (see Table 4 together)

| Figure/Table | Dataset | Model | Batch Size | Initial lr | Others |
|----------------------|-----------|---|------------|------------|-----------|
| Figure 1 | CIFAR-10 | 6CNN | 50000 | 0.04 | Setting 1 |
| Figure 2 | CIFAR-10 | 6CNN | 128 | 0.04 | Setting 1 |
| Figure 3 | CIFAR-10 | 6CNN | 50000 | 0.04 | Setting 1 |
| Figure 4 | CIFAR-10 | 6CNN | 128 | 0.01 | Setting 1 |
| Figure 5 | CIFAR-10 | 6CNN | - | - | Setting 1 |
| Figure 6 (Left) | CIFAR-10 | WRN-28-10 | 128 | 0.1 | Setting 2 |
| Figure 6 (Right) | CIFAR-100 | WRN-28-10 | 128 | 0.1 | Setting 2 |
| Table 1 | CIFAR-10 | SimpleCNN | - | - | Setting 1 |
| Table 1 ^a | - | WRN-28-10/WRN-28-2 ^b | 128/256 | 0.1 | Setting 3 |
| Table 2 | CIFAR-10 | PreAct ResNet-18/WRN-34-10 ^c | 128 | 0.1 | Setting 4 |

^a For SAM, we use $\rho = 0.05/0.1$ for CIFAR-10/100.

For ASAM, we use $\rho = 0.5/1$ and $\eta = 0.01$ for CIFAR-10/100.

For EJR, we use $\rho = 2/4$ for CIFAR-10/100.

^b For WRN-28-2, we use a different setting from Kwon et al. (2021) that $B = 256$ and epochs of 1600 (not $B = 128$ and epochs of 200).

^b For PreAct ResNet-18, we train the model for 110 epochs.

For WRN-34-10, we train the model for 120 epochs.

Table 4: Additional experimental settings for each Figure and Table in the main text (see Table 3 together)

| Settings | Data Aug. | Label Smoothing (Müller et al., 2019) | Epochs | lr scheduler | Momentum | Weight Decay | Description |
|-----------|--------------------|---------------------------------------|----------------------------|-------------------------------|----------|--------------|---|
| Setting 1 | None | None | until convergence (< 100k) | constant | None | None | Simple SGD |
| Setting 2 | Basic ^a | 0.1 | 400 | multi-step decay ^b | 0.9 | 0.0005 | used in SAM implementation ^c |
| Setting 3 | Basic | 0.1 | 200/1600 | cosine annealing | 0.9 | 0.0005 | used in ASAM (Kwon et al., 2021) |
| Setting 4 | Basic | None | 110/120 | multi-step decay ^d | 0.9 | 0.0005 | used in Pang et al. (2020) ^{e,f} |

^a random crop (with 4-pixel padding) and random horizontal flip

^b lr decay of 0.2 at $[0.3, 0.6, 0.8] \times \text{epochs}$

^c <https://github.com/davda54/sam>

^d lr decay of 0.1 at $[100, 105]$ for 110 epochs

lr decay of 0.1 at $[100, 110]$ for 120 epochs

^e BN mode when crafting adversarial examples for training: train for AT and eval for TRADES

^f ReLU activation

F. The tendency of the Jacobian norm to increase

The weight norm $\|\theta^{(t)}\|$ increases (Figure 12 (the third one)) in order to increase the logit norm $\|z(\theta^{(t)})\|$ (Figure 12 (the leftmost figure)) and to minimize the cross-entropy loss during training (Soudry et al., 2018). This also leads to the increase in the layerwise weight norms (Figure 14) and the Jacobian norm (Figure 5). We ran the experiments on the CIFAR-10 dataset and 6CNN with learning rate $\eta = 0.04$ and batch size 50000 (full-batch).

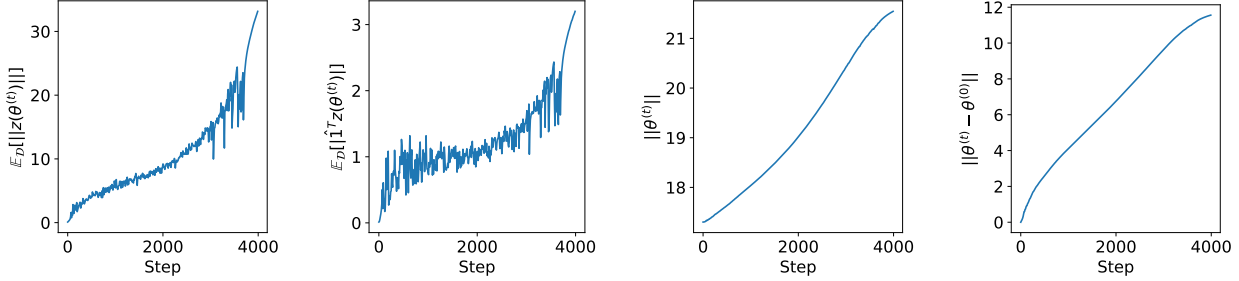


Figure 12: (Left to Right): the logit norm $\langle \|z(\theta^{(t)})\| \rangle$, the absolute value of the logit sum $\langle |\hat{1}^T z(\theta^{(t)})| \rangle$, the weight norm $\|\theta^{(t)}\|$, the distance from the initial weight $\|\theta^{(t)} - \theta^{(0)}\|$ during training (every 10 steps). See together with Figure 5 (Left, solid red line) and Figure 4.

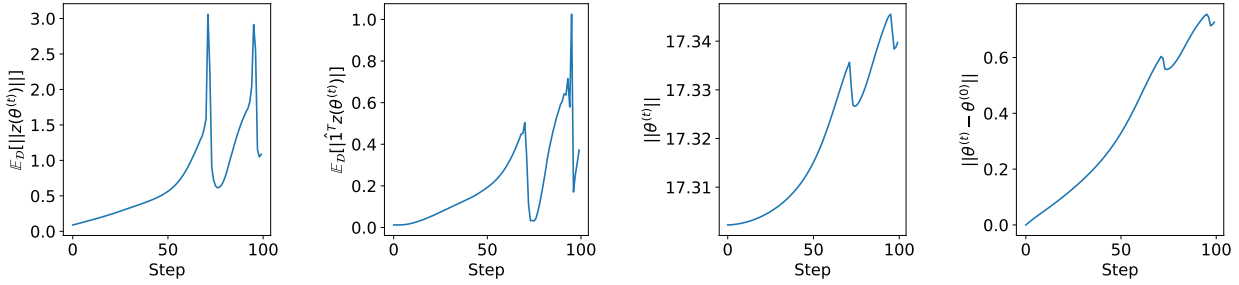


Figure 13: (Left to Right): the logit norm $\langle \|z(\theta^{(t)})\| \rangle$, the absolute value of the logit sum $\langle |\hat{1}^T z(\theta^{(t)})| \rangle$, the weight norm $\|\theta^{(t)}\|$, the distance from the initial weight $\|\theta^{(t)} - \theta^{(0)}\|$ in the early phase of training (every step).

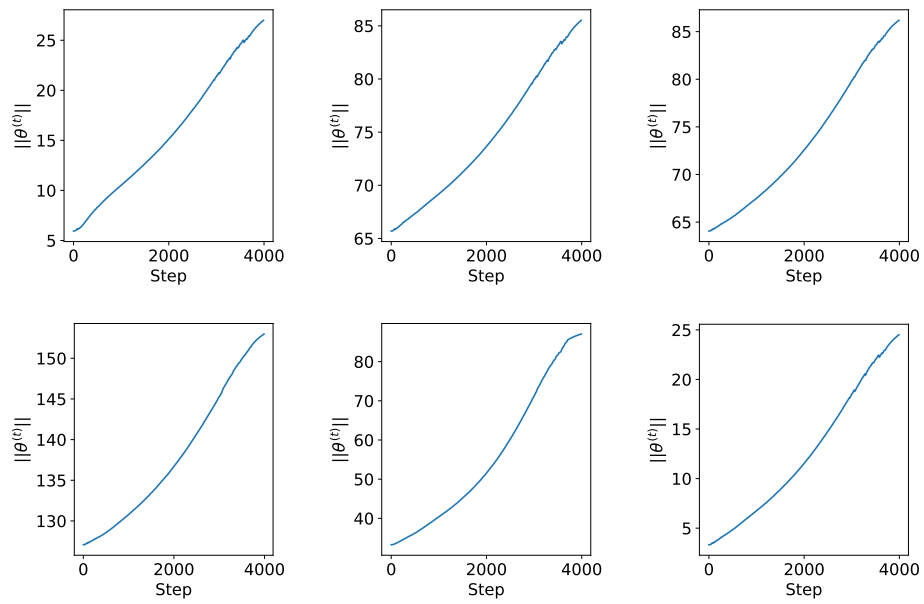


Figure 14: The layerwise weight norms (6 layers from left to right and from top to bottom) of the 6CNN model in the early phase of training (every 10 step).

$$\mathbf{G.} \quad \|\mathbf{H}\|_{\sigma} \propto \langle \lambda^{(1)} \rangle \|\langle \mathbf{J} \rangle \hat{\mathbf{1}}\|^2$$

Surprisingly, we empirically observed that $\|\mathbf{H}\|_{\sigma} \propto \langle \lambda^{(1)} \rangle \|\langle \mathbf{J} \rangle \hat{\mathbf{1}}\|^2$ during training where $\hat{\mathbf{1}} = \mathbf{1}/\sqrt{C} \in \mathbb{R}^C$. Since it requires further theoretical grounding, we left it as future work. We observe this relation for a variety of learning rates (Figure 15), network architectures (Figure 16 and 17), batch sizes (Figure 18 and 19), and datasets (Figure 20 and 21). At least, $\|\mathbf{H}\|_{\sigma}$ and $\langle \lambda^{(1)} \rangle \|\langle \mathbf{J} \rangle \hat{\mathbf{1}}\|^2$, they increase and decrease together. We emphasize that Figure 21 shows the case when the sharpness did not reach the limit $2/\eta$. This is because the learning rate is relatively low and it is easy to train a model for the MNIST dataset, and thus $\langle \lambda^{(1)} \rangle$ decreases before the sharpness reaches the limit.

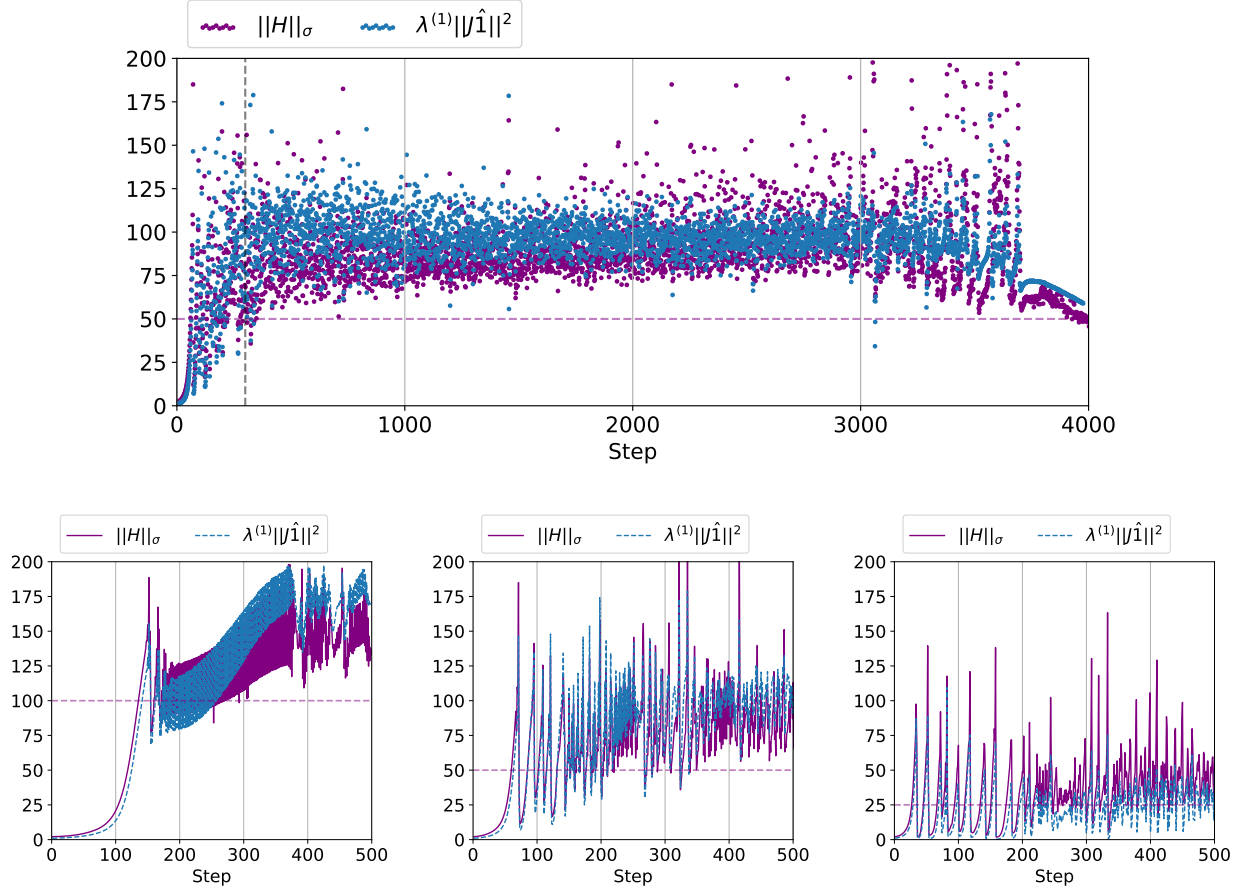


Figure 15: **The relation $\|\mathbf{H}\|_{\sigma} \propto \langle \lambda^{(1)} \rangle \|\langle \mathbf{J} \rangle \hat{\mathbf{1}}\|^2$ on the CIFAR-10 dataset and the 6CNN model trained using full-batch GD with different learning rates $\eta = 0.02/0.04/0.08$ (from left to right). Bottom figures are plotted for the early phase of training. Top figure is plotted for $\eta = 0.02$. Curves are plotted for every step.**

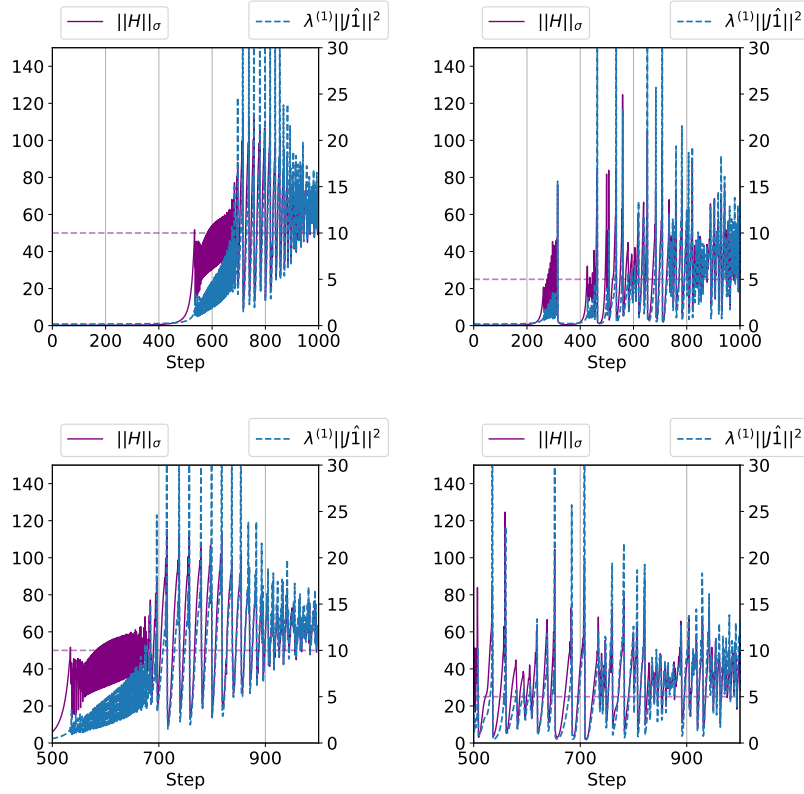


Figure 16: **The relation $\|H\|_\sigma \propto \langle \lambda^{(1)} \rangle \|\langle J \rangle \hat{\mathbf{1}}\|^2$ on the CIFAR-10 dataset and the VGG model trained using full-batch GD with different learning rates $\eta = 0.04/0.08$ (left/right) in the early phase of training.** Curves are plotted for every step (Top: 0-1000 stpes and Bottom: 500-1000 stpes).

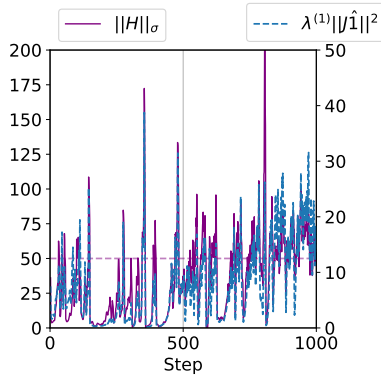


Figure 17: **The relation $\|H\|_\sigma \propto \langle \lambda^{(1)} \rangle \|\langle J \rangle \hat{\mathbf{1}}\|^2$ on the CIFAR-10 dataset and the ResNet model trained using SGD with learning rate $\eta = 0.04$ and batch sizes $|\mathcal{B}^{(t)}| = 128$ during training (0-1000 steps).** Curves are plotted for every step.

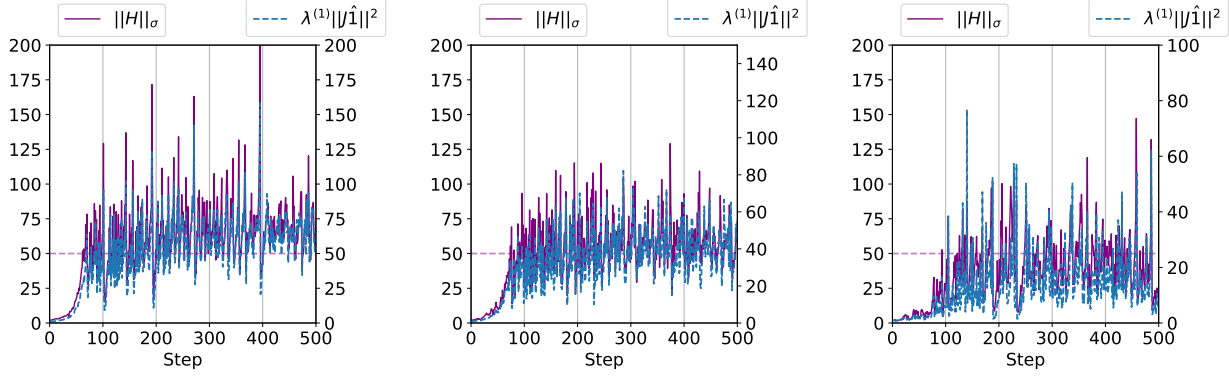


Figure 18: **The relation $\|H\|_\sigma \propto \langle \lambda^{(1)} \rangle \|\langle J \rangle \hat{\mathbf{1}}\|^2$ on the CIFAR-10 dataset and the 6CNN model trained using SGD with fixed learning rate $\eta = 0.04$ and different batch sizes $|\mathcal{B}^{(t)}| = 512/128/32$ (from left to right) in the initial phase (0-500 steps).** Note that the proportionality constant may change according to the batch size (the smaller the batch size, the larger the proportionality constant). Curves are plotted for every step.

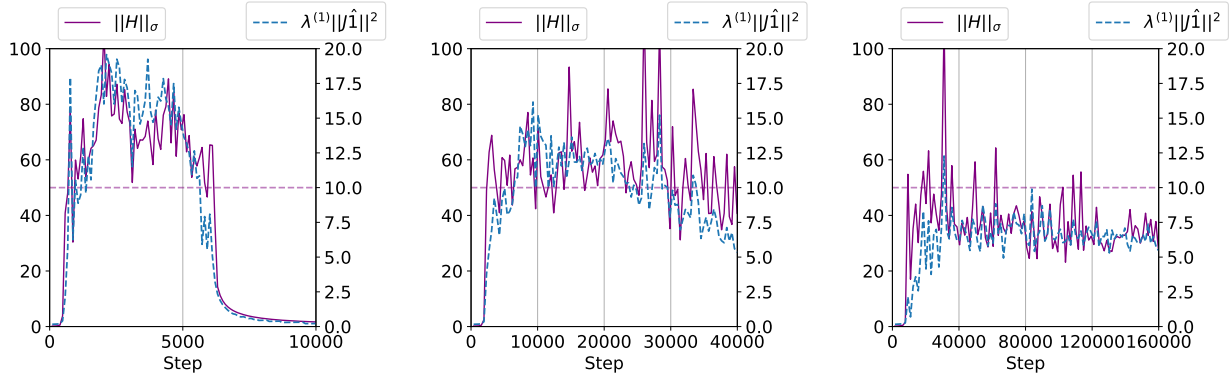


Figure 19: **The relation $\|H\|_\sigma \propto \langle \lambda^{(1)} \rangle \|\langle J \rangle \hat{\mathbf{1}}\|^2$ on the CIFAR-10 dataset and the VGG model trained using SGD with fixed learning rate $\eta = 0.04$ and different batch sizes $|\mathcal{B}^{(t)}| = 512/128/32$ (from left to right) during training (0-100 epochs).** Curves are plotted for every n steps ($n = 97/388/1552$) where $97 = \lfloor 50000/512 \rfloor$.

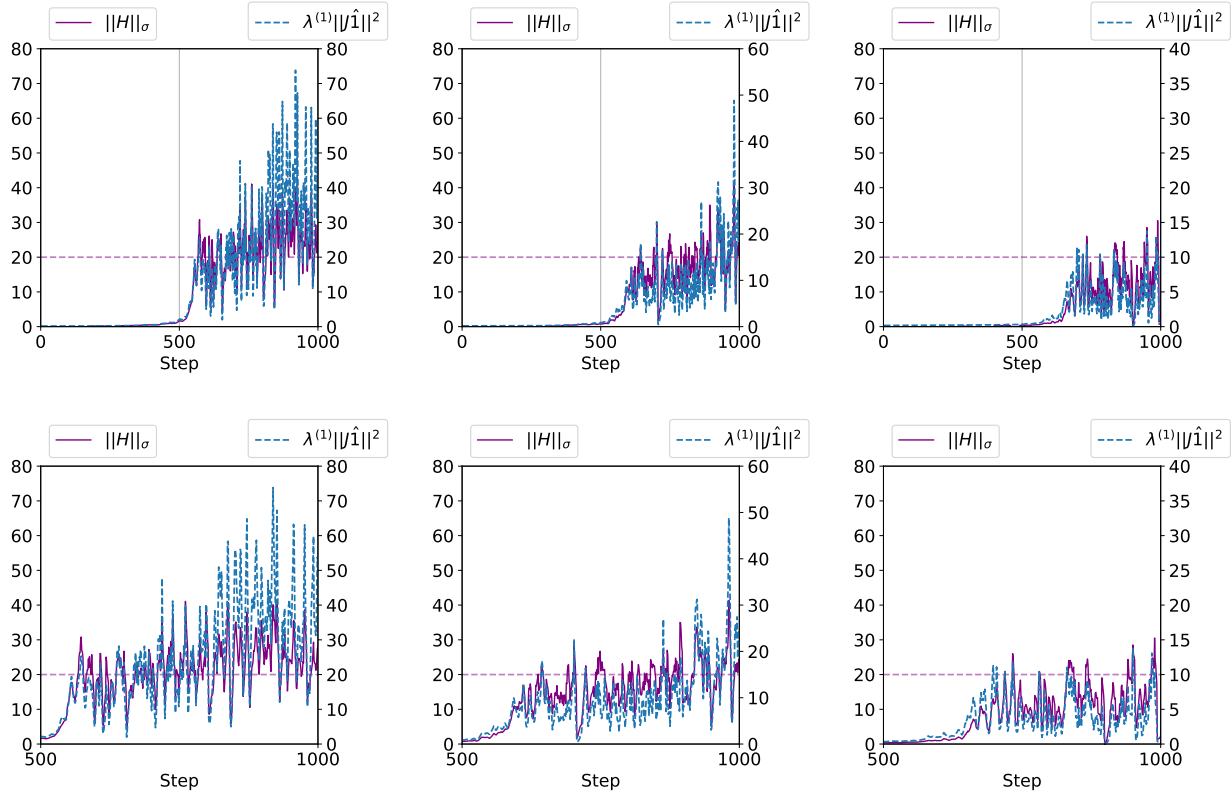


Figure 20: **The relation $\|H\|_\sigma \propto \langle \lambda^{(1)} \rangle \|\langle J \rangle \hat{\mathbf{1}}\|^2$ on the CIFAR-100 dataset and the VGG model trained using SGD with fixed learning rate $\eta = 0.1$ and different batch sizes $|\mathcal{B}^{(t)}| = 128/64/32$ (from left to right) in the early phase of training.** Note that the proportionality constant may change according to the batch size (the smaller the batch size, the larger the proportionality constant). Curves are plotted for every step (Top: 0-1000 stpes and Bottom: 500-1000 stpes).

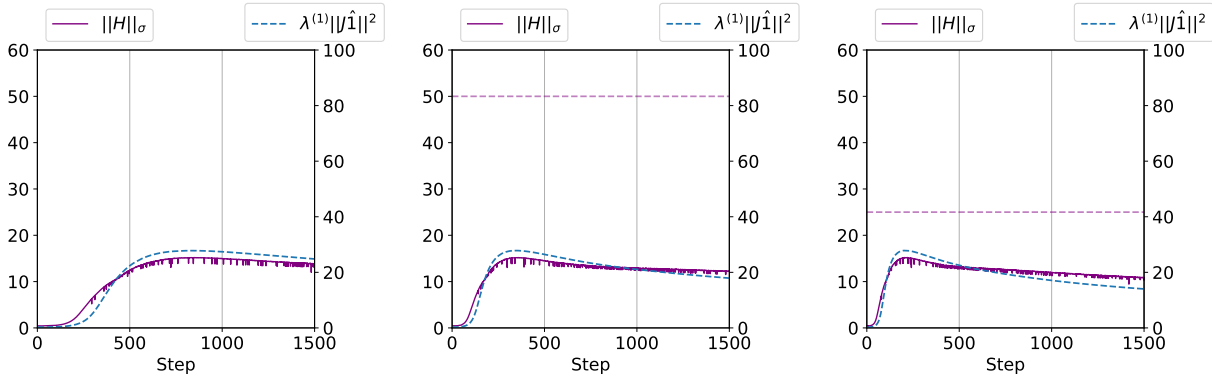


Figure 21: **The relation $\|H\|_\sigma \propto \langle \lambda^{(1)} \rangle \|\langle J \rangle \hat{\mathbf{1}}\|^2$ on the MNIST dataset and the 3FCN-MNIST model trained using full-batch GD with different learning rates $\eta = 0.02/0.04/0.08$ (from left to right) in the early phase of training.** Note that the sharpness did not reach the limit $2/\eta$ (the dashed horizontal line). Curves are plotted for every step.

H. Implicit Regularization on $\|\langle \mathbf{J} \rangle \hat{\mathbf{1}}\|^2$

We further investigate the relation $\|\mathbf{H}\|_\sigma \propto \langle \lambda^{(1)} \rangle \|\langle \mathbf{J} \rangle \hat{\mathbf{1}}\|^2$ in the previous section and its effect on $\|\langle \mathbf{J} \rangle \hat{\mathbf{1}}\|^2$. For the early phase of training, it is hard to see the initial rapid growth of the sharpness in this smoothed curves and when exactly the regularization begins to activate. We refer the readers to the previous section, Appendix G, for the fine-grained analysis of the early phase of training. We provide plots with different settings. Again, we observe that training with a larger learning rate and a smaller batch size limits $\|\langle \mathbf{J} \rangle \hat{\mathbf{1}}\|^2$ with a smaller value (dotted red lines) in the Active Regularization Period. Curves are smoothed for visual clarity.

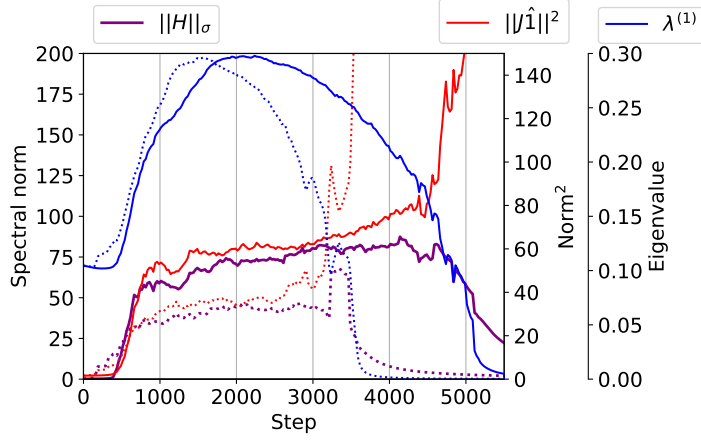


Figure 22: The evolution of $\|\mathbf{H}\|_\sigma$, $\|\langle \mathbf{J} \rangle \hat{\mathbf{1}}\|^2$, and $\langle \lambda^{(1)} \rangle$ on the CIFAR-10 dataset and the VGG model trained using full-batch GD with the learning rates $\eta = 0.04/0.08$ (solid/dotted lines). See the Figure 5 caption together.

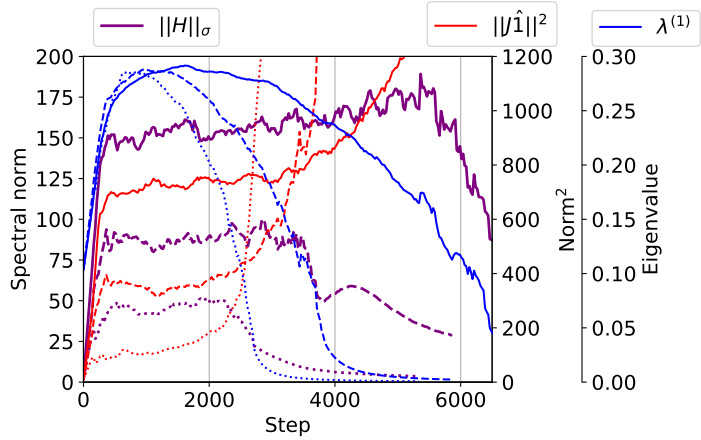


Figure 23: The evolution of $\|\mathbf{H}\|_\sigma$, $\|\langle \mathbf{J} \rangle \hat{\mathbf{1}}\|^2$, and $\langle \lambda^{(1)} \rangle$ on the CIFAR-10 dataset and the 6CNN model trained using full-batch GD with the learning rates $\eta = 0.02/0.04/0.08$ (solid/dashed/dotted lines). See the Figure 5 caption together.

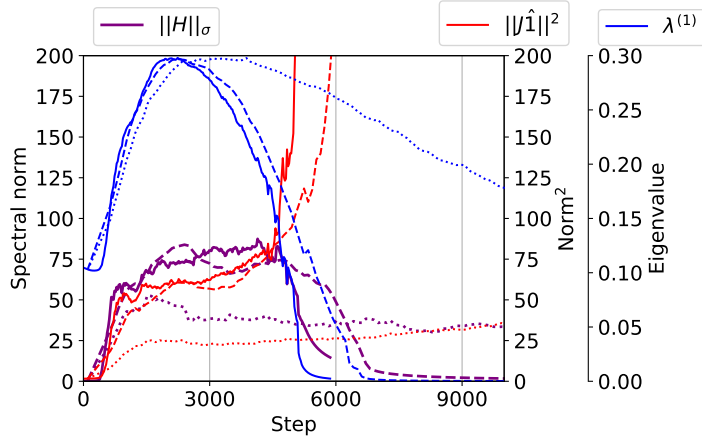


Figure 24: The evolution of $\|H\|_\sigma$, $\|\langle J \rangle \hat{\mathbf{1}}\|^2$, and $\langle \lambda^{(1)} \rangle$ on the CIFAR-10 dataset and the VGG model trained using SGD with the fixed learning rate $\eta = 0.04$ and the batch sizes $|\mathcal{B}^{(t)}| = 50000(\text{GD})/512/32$ (solid/dashed/dotted lines). Training with a batch size 512 shows similar evolutions to the GD training. See the Figure 5 caption together.

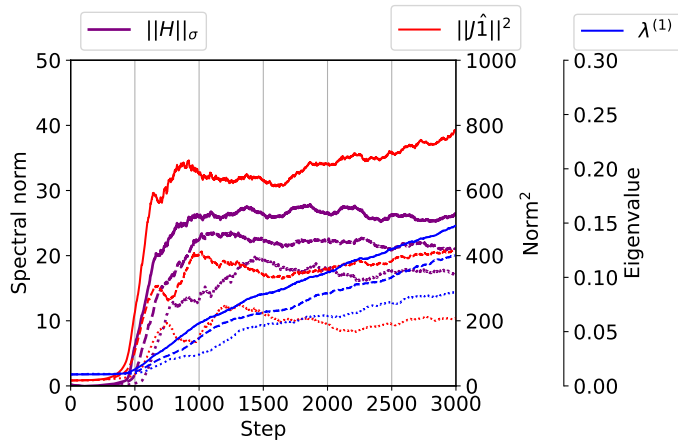


Figure 25: The evolution of $\|H\|_\sigma$, $\|\langle J \rangle \hat{\mathbf{1}}\|^2$, and $\langle \lambda^{(1)} \rangle$ on the CIFAR-100 dataset and the VGG model trained using SGD with the fixed learning rate $\eta = 0.1$ and the batch sizes $|\mathcal{B}^{(t)}| = 128/64/32$ (solid/dashed/dotted lines). See the Figure 5 caption together.

I. Figure 2 (Sharpness and Jacobian norm)

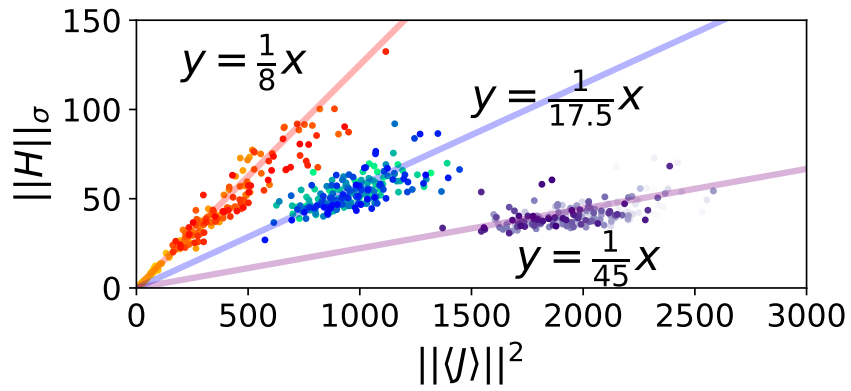


Figure 26: The sharpness $\|H\|_\sigma$ and the Jacobian norm $\|\langle J \rangle\|^2$ show similar behavior up to a factor $\hat{\lambda}^*$ which is locally constant and slowly changes during training (CIFAR-10, $\eta = 0.04$, $B = 128$). $\hat{\lambda}^* \approx 1/8, 1/17.5, 1/45$ (red/blue/purple) for 0-200/1000-1200/4000-4200 steps

J. Figure 3 (Visualization of the optimization trajectory)

We use UMAP (McInnes et al., 2018) to visualize the optimization trajectory $\{\theta^{(t)}\}_{t \in [T]} \subset \Theta \subset \mathbb{R}^m$ in a 2D space. As GD enters into the Edge of Stability (Cohen et al., 2021), it oscillates in a direction nearly orthogonal to its global descent direction (Xing et al., 2018). Note that GD may not enter the Edge of Stability as shown in Figure 27 (Bottom).

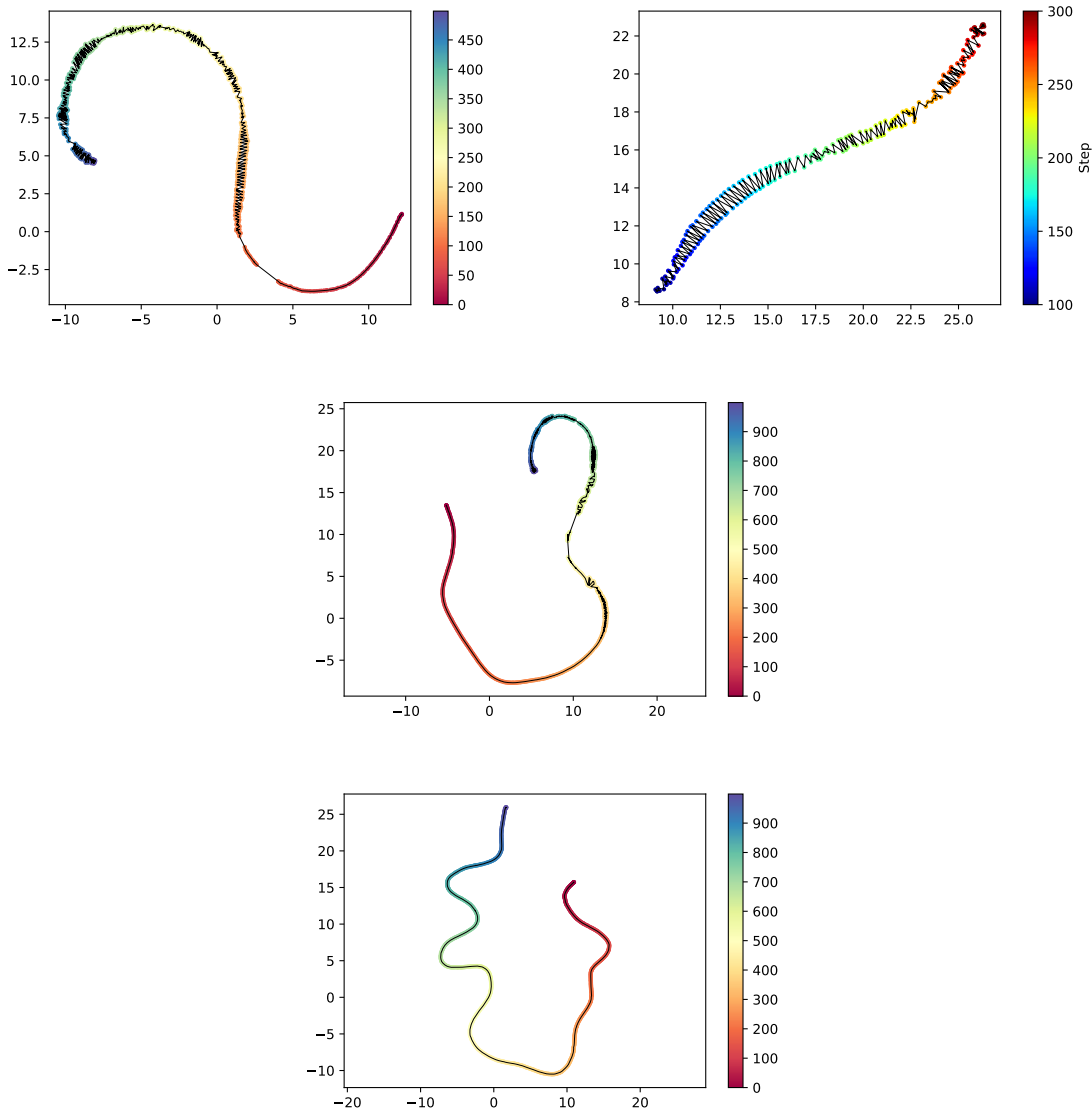


Figure 27: **Visualization of the optimization trajectory using UMAP.** UMAP on the CIFAR-10 dataset trained with 6CNN for the first 500 steps (Top Left) and for 100-300 steps (Top Right), on the CIFAR-10 dataset trained with VGG for the first 1000 steps (Middle), and on the MNIST dataset trained with 3FCN-MNIST for the first 1000 steps (Bottom) from red to blue.

K. Figure 4 (Active Regularization Period)

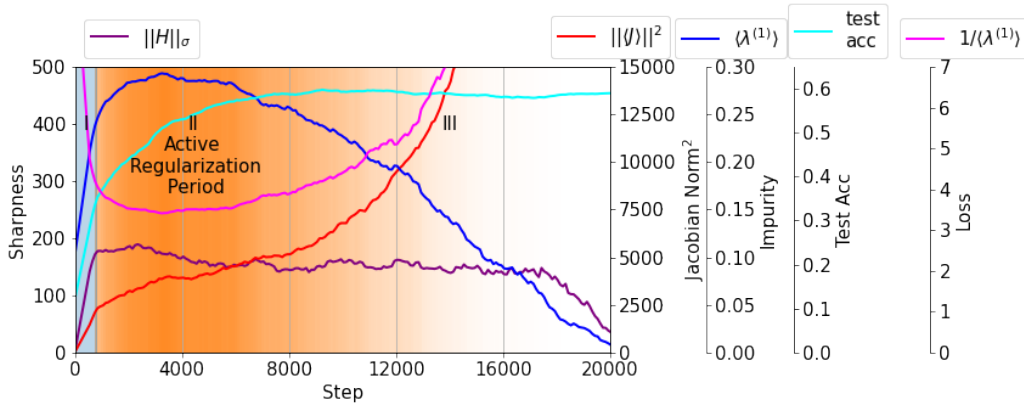


Figure 28: **Three phases of Implicit Jacobian Regularization (IJR).** Figure 4 with $1/\langle\lambda^{(1)}\rangle$. As $1/\langle\lambda^{(1)}\rangle$ increases, the Jacobian norm increases faster in phase III.

L. Figure 6 (Efficient EJR)

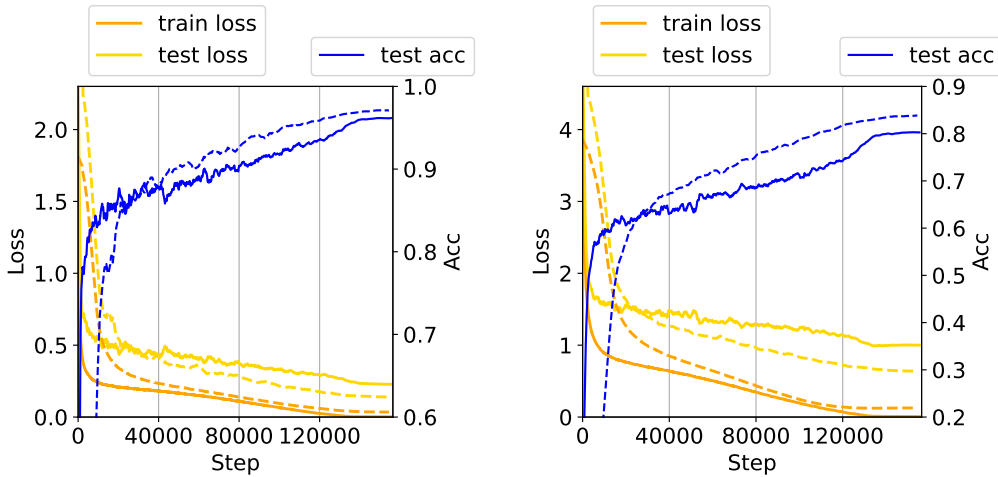


Figure 29: **The effectiveness of efficient EJR (dashed lines) compared to SGD (solid lines) trained with Setting 3 in Table 4 (cosine annealing) for WRN-28-10 on CIFAR-10 (Left) and CIFAR-100 (Right).** Unlike Figure 6, we train the model with EJR for 200 epochs and with SGD for 400 epochs and thus we plot with different x-axis to match this difference ($\times 2$ for EJR).

M. Figure 8 ($Q\Lambda^{1/2}$)

Figure 30 shows the matrix $Q\Lambda^{1/2}$ (Figure 8) for different training steps. Figure 8 and Figure 30 (Top Right) are plotted at the equivalent step ($t = 1000$). Figure 30 demonstrates that $\lambda^{(1)}$ becomes more dominant than the others as training progresses. The argument that there are two salient elements in each $q^{(i)}$ in its (i) - and $(i + 1)$ -th elements is empirically shown to be valid throughout the training.

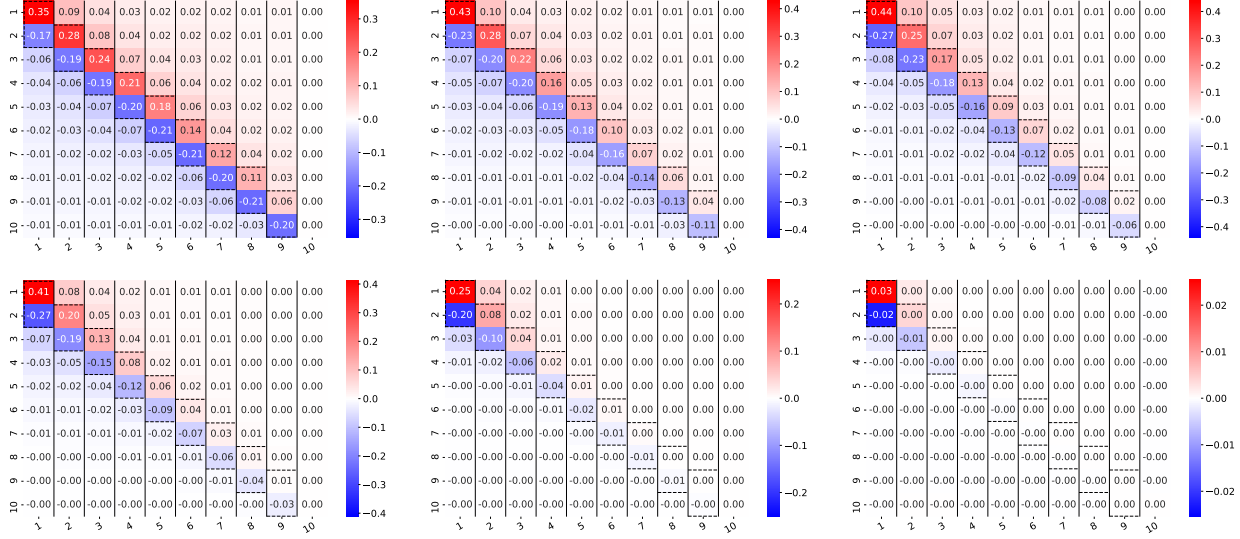


Figure 30: **Evolution of $Q\Lambda^{1/2}$ for some training steps during the training.** They are visualized for 100/500/1000/2000/4000/6000 steps from left to right and from top to bottom.

N. Figure 9 (Clusters of m around each \bar{m}^i)

Figure 31 shows Figure 9 (Left) for some different class pairs (i, j) with negative cosine similarity, i.e., $\cos(\bar{m}^i, \bar{m}^j) < 0$. We use the model trained for $t = 1000$ steps. Figure 32 shows Figure 9 (Right) for different training steps. Figure 33 shows the evolution of the Mean Resultant Length (MRL) ρ in Figure 9 (Left) during training.

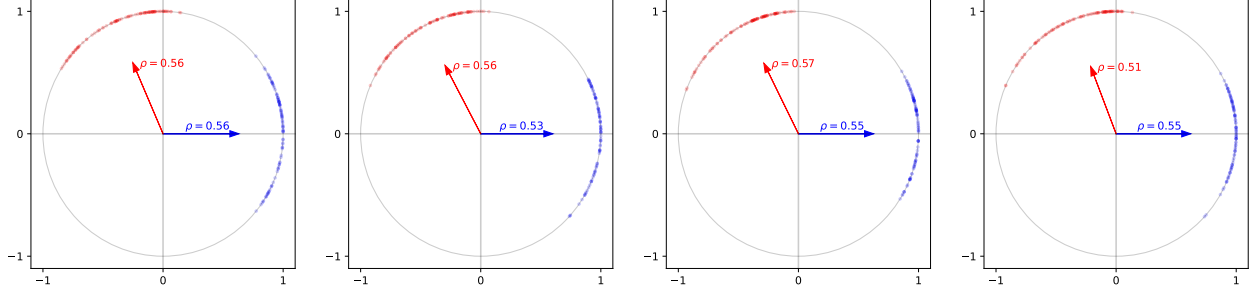


Figure 31: Directional data of m from \mathcal{D}_i and \mathcal{D}_j . They are visualized for $(i, j) = (\text{airplane, ship})/(\text{automobile, truck})/(\text{dog, cat})/(\text{deer, bird})$ from left to right. See the Figure 9 caption.

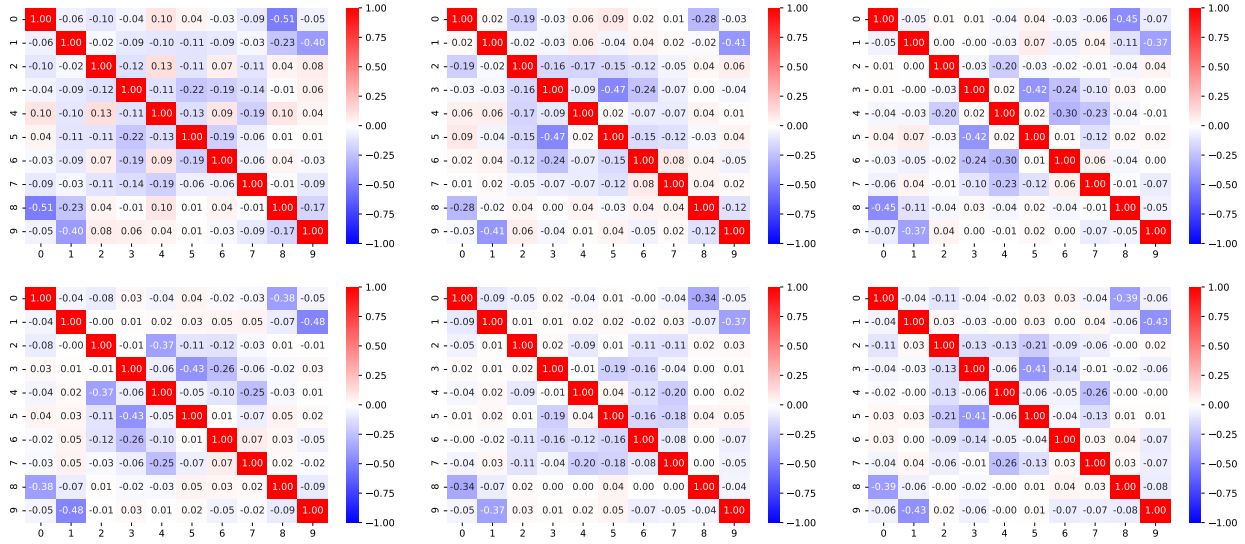


Figure 32: Cosine similarities of $\{\bar{m}^i\}_{i=1}^C$ during training. They are visualized for step=500/1000/1500/2000/3000/4000 from left to right and from top to bottom. See the Figure 9 caption.

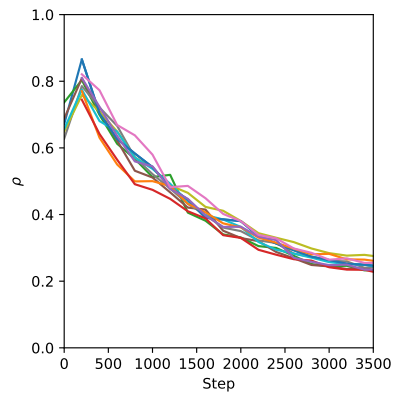


Figure 33: **The Mean Resultant Lengths (MRL) ρ for each \mathcal{D}_i .** Note that ρ is not defined for first few steps because some \mathcal{D}_i are empty.

O. Figure 10 (gradient descent happens mostly in the top Hessian subspace)

Figure 34 shows similar results with different settings with VGG and learning rate $\eta = 0.08$.

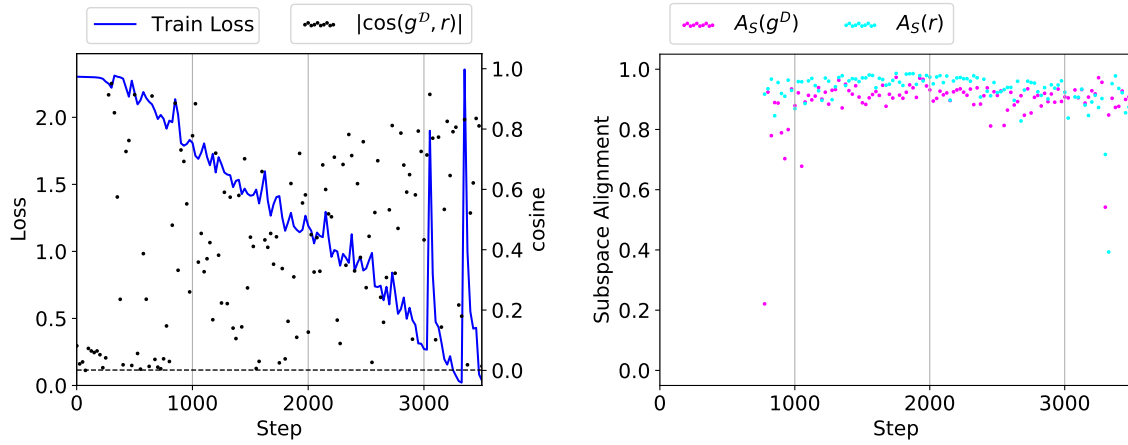


Figure 34: **(Left) Alignment between the two vectors $g^{\mathcal{D}}$ and r , and (Right) alignment of $g^{\mathcal{D}}$, r in the subspace $\mathcal{S} = \text{span}(\{\bar{m}^i\}_{i=1}^C)$ using VGG with $\eta = 0.08$.** See the Figure 10 caption. They are plotted for every 25 steps. Note that $A_{\mathcal{S}}$ is not defined for first few steps (about 0-800 steps) because some \mathcal{D}_i are empty.

P. Analysis of the MSE loss

In the main text, we focus on the cross-entropy loss. Here, we briefly analyze the MSE loss, $l = \frac{1}{2}\|z - e^y\|^2$. Then, we have $M = \nabla_z^2 l = I$, $\lambda^{(1)} = \|M\|_\sigma = 1$ and $G = \langle JJ^T \rangle$. It leads to the same conclusion as in Theorem 5.1:

$$\|G\|_\sigma = \|\langle JJ^T \rangle\|_\sigma \leq \langle \|JJ^T\|_\sigma \rangle = \langle \|J\|^2 \rangle \quad (50)$$

We empirically observed that $\|H\|_\sigma \propto \|\langle J \rangle\|^2$ as shown in Figure 35.

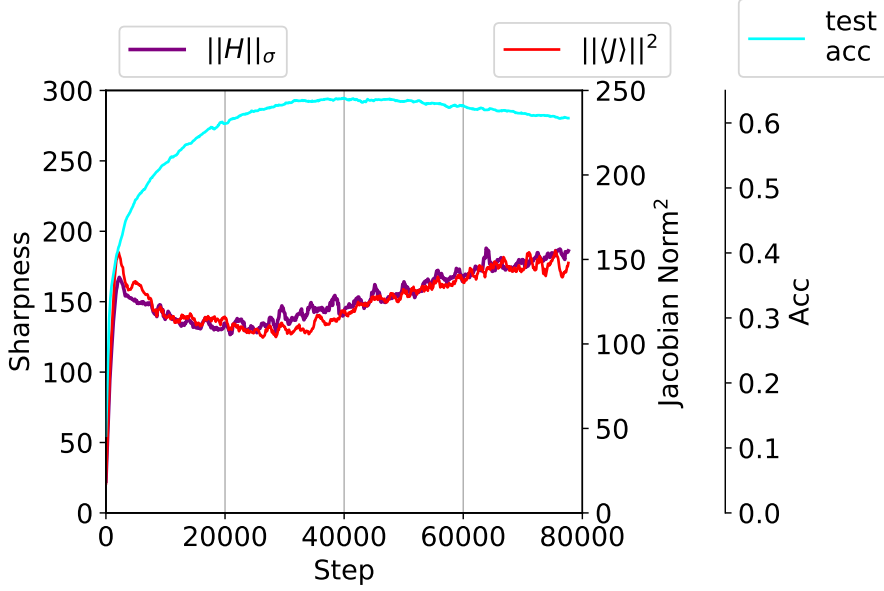


Figure 35: The sharpness $\|H\|_\sigma$ and the Jacobian norm $\|\langle J \rangle\|^2$ during training with the MSE loss

Q. Details of Explicit Jacobian Regularization (EJR)

We first propose a simple form of EJR with the regularized loss as follows:

$$\tilde{L}(\boldsymbol{\theta}) = L(\boldsymbol{\theta}) + \lambda_{reg} \|\langle \mathbf{J} \rangle\|_F^2 / C \quad (51)$$

and update the model parameter as

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta (\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^{(t)}) + \lambda_{reg} \nabla_{\boldsymbol{\theta}} \|\langle \mathbf{J} \rangle\|_F^2 / C) \quad (52)$$

Efficient EJR However, this requires to build a computational graph for $\|\langle \mathbf{J} \rangle\|_F^2$ which is inefficient for a large network (e.g. WRN-28-10). To this end, we propose an efficient variant of EJR. SVD of the Jacobian is $\langle \mathbf{J} \rangle = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T$ where $\mathbf{U} \in \mathbb{R}^{m \times m}$, $\boldsymbol{\Sigma} \in \mathbb{R}^{m \times C}$, $\mathbf{V} \in \mathbb{R}^{C \times C}$ where

$$\Sigma_{i,j} = \begin{cases} \sigma_i & \text{if } i = j \in [C] \\ 0 & \text{if } i \neq j \end{cases}$$

and $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_C \geq 0$. Then, we have $\|\langle \mathbf{J} \rangle\| = \sigma_1$ and $\|\langle \mathbf{J} \rangle\|_F^2 = \sum_i \sigma_i^2$. We propose to update the model parameter as follows:

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta \nabla_{\boldsymbol{\theta}} L(\hat{\boldsymbol{\theta}}^{(t)}) \quad (53)$$

where

$$\hat{\boldsymbol{\theta}} = \boldsymbol{\theta} + \tilde{\rho}_{reg} \langle \mathbf{J} \rangle \mathbf{u} \quad (54)$$

and $\tilde{\rho}_{reg} = \rho_{reg} / \|\langle \mathbf{J} \rangle \mathbf{u}\|$ with $\mathbf{u} \sim \mathbb{S}^{C-1}$ as in Foret et al. (2021); Liu et al. (2019). This perturbation in $\hat{\boldsymbol{\theta}}$ is highly aligned with the principal left-singular vectors of the Jacobian and the principal eigenspace of the Hessian. Therefore, the update of (53) leads to minimize the largest singular values σ_i of the Jacobian and the sharpness as well. As shown in Table 1, this method outperforms the state-of-the-art sharpness-aware optimization methods like SAM (Foret et al., 2021) and ASAM (Kwon et al., 2021) which can be attributed to the principal eigensubspace awareness.

Efficient EJR.v2 We propose another variant using \hat{L} instead of L in (53) as follows:

$$\hat{L} = L + \mu_{reg} \mathbf{u}^T \langle \mathbf{z} \rangle \quad (55)$$

We can approximate \hat{L} as follows:

$$\hat{L}(\hat{\boldsymbol{\theta}}) \approx \hat{L}(\boldsymbol{\theta}) + (\nabla_{\boldsymbol{\theta}} \hat{L}(\boldsymbol{\theta}))^T (\tilde{\rho}_{reg} \langle \mathbf{J} \rangle \mathbf{u}) \quad (56)$$

$$= \hat{L}(\boldsymbol{\theta}) + (\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}))^T (\tilde{\rho}_{reg} \langle \mathbf{J} \rangle \mathbf{u}) + (\nabla_{\boldsymbol{\theta}} \mu_{reg} \mathbf{u}^T \langle \mathbf{z} \rangle)^T (\tilde{\rho}_{reg} \langle \mathbf{J} \rangle \mathbf{u}) \quad (57)$$

$$= \hat{L}(\boldsymbol{\theta}) + (\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}))^T (\tilde{\rho}_{reg} \langle \mathbf{J} \rangle \mathbf{u}) + (\mu_{reg} \langle \mathbf{J} \rangle \mathbf{u})^T (\tilde{\rho}_{reg} \langle \mathbf{J} \rangle \mathbf{u}) \quad (58)$$

$$= \hat{L}(\boldsymbol{\theta}) + (\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}))^T (\tilde{\rho}_{reg} \langle \mathbf{J} \rangle \mathbf{u}) + \mu_{reg} \rho_{reg} \|\langle \mathbf{J} \rangle \mathbf{u}\| \quad (59)$$

$$\approx \hat{L}(\boldsymbol{\theta}) + (\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}))^T (\tilde{\rho}_{reg} \langle \mathbf{J} \rangle \mathbf{u}) + \mu_{reg} \rho_{reg} \|\langle \mathbf{J} \rangle\|_F / \sqrt{C} \quad (60)$$

We used the first-order Taylor expansion of $\hat{L}(\boldsymbol{\theta} + \hat{\rho}_{reg} \langle \mathbf{J} \rangle \mathbf{u})$ in (56). We expect additional effect of minimizing $\mu_{reg} \mathbf{u}^T \langle \mathbf{z} \rangle + \mu_{reg} \rho_{reg} \|\langle \mathbf{J} \rangle\|_F / \sqrt{C} \approx \mu_{reg} \rho_{reg} \|\langle \mathbf{J} \rangle\|_F / \sqrt{C}$ compared to the first variant.

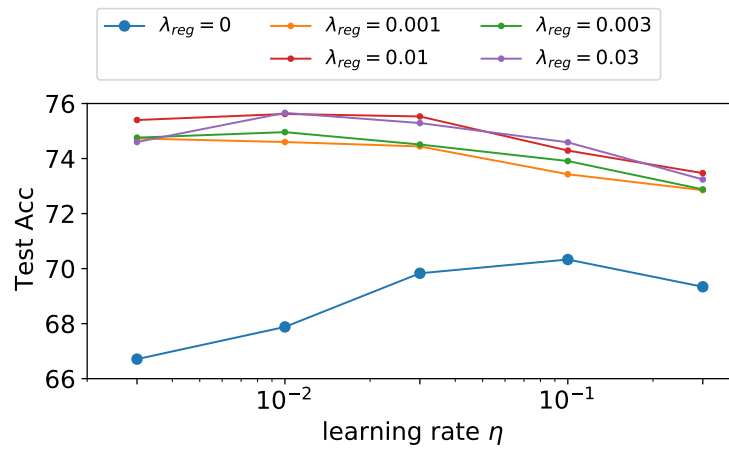


Figure 36: We plot the test accuracy for different learning rates η and regularization coefficients λ_{reg} . The models are trained with batch size of $|\mathcal{B}| = 128$ on CIFAR-10.

Table 5: **Effectiveness of Efficient EJR**. We report improvement ($\Delta\text{Acc.}$) and Error Reduction Rate (**ERR**) on CIFAR-10 when trained with EJR (+EJR), compared to the standard training (**Baseline**). We used Setting 1 and Setting 2 in Table 4 for SimpleCNN and WRN-28-10, respectively.

| Dataset | Network Architecture | Batch Size | lr | Reg. param. | Test Accuracy | | $\Delta\text{Acc.}$ (%p) | ERR (%) |
|---------------------------|---------------------------|-----------------------|--------------------|-------------------------|--|------------------|-----------------------------|------------|
| | | | | | Baseline | +EJR | | |
| CIFAR-10 | SimpleCNN | 128 | 0.003 | $\lambda_{reg} = 0.01$ | 66.71 | 75.40 | +8.69 | 26.10 |
| | | | 0.01 | $\lambda_{reg} = 0.03$ | 67.88 | 75.66 | +7.78 | 24.22 |
| | | | 0.03 | $\lambda_{reg} = 0.01$ | 69.83 | 75.53 | +5.70 | 18.89 |
| | | | 0.1 | $\lambda_{reg} = 0.03$ | 70.33 | 74.59 | +4.26 | 14.36 |
| | | | 0.3 | $\lambda_{reg} = 0.01$ | 69.34 | 73.47 | +4.13 | 13.47 |
| | | 50000 (full-batch) | 0.01 | $\lambda_{reg} = 0.001$ | 66.81 | 74.43 | +7.62 | 22.96 |
| | | | 0.03 | | 67.72 | 74.31 | +6.59 | 20.42 |
| | | | 0.1 | | 67.53 | 73.69 | +6.16 | 18.97 |
| | | | 0.3 | | 61.08 | 72.15 | +11.07 | 28.44 |
| | | | 0.1 | | $\rho_{reg} = 0.5$ $\rho_{reg} = 1$ $\rho_{reg} = 2$ $\rho_{reg} = 3$ | 95.93 \pm 0.15 | 96.44 | +0.51 |
| | 0.1 | 96.57 | +0.64 | 15.72 | | | | |
| | 0.1 | 96.62 | +0.69 | 16.95 | | | | |
| | 0.1 | 96.30 | +0.37 | 9.09 | | | | |
| | WRN-28-10 (200 epochs) | 128 | 0.1 | $\rho_{reg} = 0.5$ | 96.10 \pm 0.05 | 96.65 | +0.55 | 14.10 |
| $\rho_{reg} = 1$ | | | | 96.78 | | +0.68 | 17.44 | |
| $\rho_{reg} = 2$ | | | | 97.07 | | +0.97 | 24.87 | |
| $\rho_{reg} = 3$ | | | | 96.79 | | +0.69 | 17.69 | |
| 128 | | 0.1 | $\rho_{reg} = 0.5$ | 80.29 \pm 0.25 | 80.42 | +0.13 | 0.66 | |
| | | | $\rho_{reg} = 1$ | | 81.11 | +0.82 | 4.16 | |
| | | | $\rho_{reg} = 2$ | | 81.50 | +1.21 | 6.14 | |
| | | | $\rho_{reg} = 3$ | | 82.51 | +2.22 | 11.26 | |
| WRN-28-10 (400 epochs) | 128 | 0.1 | $\rho_{reg} = 4$ | 80.69 \pm 0.21 | 82.65 | +2.36 | 11.97 | |
| | | | $\rho_{reg} = 5$ | | 82.31 | +2.02 | 10.25 | |
| | | | $\rho_{reg} = 6$ | | 82.03 | +1.74 | 8.83 | |
| | | | $\rho_{reg} = 1$ | | 82.55 | +1.86 | 9.63 | |
| | | | $\rho_{reg} = 2$ | | 82.51 | +1.82 | 9.42 | |
| | | | $\rho_{reg} = 3$ | | 82.84 | +2.15 | 11.13 | |
| | 128 | 0.1 | $\rho_{reg} = 4$ | 83.35 | +2.66 | 13.78 | | |
| | | | $\rho_{reg} = 5$ | 83.73 | +3.04 | 15.74 | | |
| | | | $\rho_{reg} = 6$ | 83.16 | +2.47 | 12.79 | | |
| | | | $\rho_{reg} = 7$ | 83.12 | +2.43 | 12.58 | | |
| WRN-28-10 (400 epochs) | 128 | 0.1 | $\rho_{reg} = 8$ | 81.16 | +0.47 | 2.43 | | |

Table 6: **Effectiveness of Efficient EJR.v2.** We report improvement ($\Delta\text{Acc.}$) and Error Reduction Rate (**ERR**) on CIFAR-10 when trained with the second variant of EJR (**+EJR.v2**), compared to the standard training (**Baseline**). We used Setting 2 in Table 4

| Dataset | Network Architecture | Batch Size | lr | Reg. param. | Test Accuracy | | $\Delta\text{Acc.}$ (%p) | ERR (%) |
|----------|---------------------------|------------|-----------------------------------|-------------------------------------|------------------|--------------|-----------------------------|------------|
| | | | | | Baseline | +EJR.v2 | | |
| CIFAR-10 | WRN-28-10 (400 epochs) | 128 | 0.1 | $\rho_{reg} = 2, \mu_{reg} = 0.001$ | 96.10 \pm 0.05 | 97.28 | +1.18 | 30.26 |
| | | | | $\rho_{reg} = 2, \mu_{reg} = 0.003$ | | 97.32 | +1.23 | 31.28 |
| | | | | $\rho_{reg} = 2, \mu_{reg} = 0.01$ | | 97.33 | +1.23 | 31.54 |
| | | | | $\rho_{reg} = 2, \mu_{reg} = 0.03$ | | 97.38 | +1.28 | 32.82 |
| | | | | $\rho_{reg} = 2, \mu_{reg} = 0.1$ | | 97.17 | +1.07 | 27.44 |
| | | | | $\rho_{reg} = 1, \mu_{reg} = 0.01$ | | 97.07 | +0.97 | 24.87 |
| | | | | $\rho_{reg} = 1, \mu_{reg} = 0.02$ | | 97.34 | +1.24 | 31.79 |
| | | | | $\rho_{reg} = 1, \mu_{reg} = 0.06$ | | 97.33 | +1.23 | 31.54 |
| | | | $\rho_{reg} = 1, \mu_{reg} = 0.1$ | | 97.26 | +1.16 | 29.74 | |

R. Adversarial Robustness

We also apply EJR to train a model that is robust to adversarial attacks (Szegedy et al., 2013). Adversarially robust training suffers from robust overfitting (Rice et al., 2020) and we expect EJR to help to improve robust generalization. We test the effectiveness of EJR with AT (Madry et al., 2018) and TRADES (Zhang et al., 2019). Table 2 shows that the proposed method can obtain better robust generalization against the adversarial attacks. We report the best result with the regularization parameter tuning. We follow a similar setting to Pang et al. (2020).

We tune the regularization parameter ρ_{reg} as in Table 7 to improve the results and to report the best results in Table 2. However, we did not perform an exhaustive investigation due to limited computational resources, so there is still room for further improvement.

Table 7: Adversarial training with efficient EJR on CIFAR-10 compared with AT (Madry et al., 2018) and TRADES (Zhang et al., 2019). We report Standard Accuracy and robust accuracy ($\epsilon = 8/255$) against PGD-20 and AutoAttack (AA) (Croce & Hein, 2020).

| Model | Method | Epoch | ρ_{reg} | β_{TRADES} | Standard | PGD-20 | AA |
|---------------------|-----------------|-------------|--------------|-------------------------|----------|--------|--------------|
| PreAct ResNet-18 | AT Baseline | 110 | - | - | 82.45 | 52.85 | 48.76 |
| | AT-EJR | 110 | 0.05 | - | 82.32 | 53.58 | 49.14 |
| | | 110 | 0.1 | - | 81.45 | 53.61 | 48.77 |
| | TRADES Baseline | 110 | - | 6 | 82.34 | 52.83 | 49.06 |
| | | 110 | - | 10 | 79.20 | 53.19 | 49.07 |
| | TRADES-EJR | 110 | 0.1 | 6 | 81.56 | 53.05 | 49.56 |
| | | 110 | 0.1 | 10 | 80.43 | 53.44 | 49.00 |
| | | 110 | 0.2 | 6 | 81.87 | 53.19 | 49.37 |
| | | 110 | 0.2 | 10 | 80.68 | 53.37 | 49.21 |
| | WRN-34-10 | AT Baseline | 120 | - | - | 86.99 | 52.20 |
| AT-EJR | | 120 | 0.05 | - | 87.16 | 55.89 | 52.70 |
| | | 120 | 0.1 | - | 86.85 | 57.82 | 53.73 |
| TRADES Baseline | | 120 | - | 6 | 85.23 | 55.68 | 52.46 |
| | | 120 | - | 10 | 83.62 | 57.08 | 53.29 |
| TRADES-EJR | | 120 | 0.02 | 6 | 85.77 | 56.97 | 53.62 |
| | | 120 | 0.02 | 10 | 83.61 | 57.08 | 53.73 |
| | | 120 | 0.1 | 6 | 85.91 | 56.34 | 53.18 |
| | | 120 | 0.1 | 10 | 83.97 | 57.48 | 53.85 |