

Gradient Norm Aware Minimization Seeks First-Order Flatness and Improves Generalization

Xingxuan Zhang[†], Renzhe Xu[†], Han Yu, Hao Zou, Peng Cui*

Department of Computer Science, Tsinghua University

xingxuanzhang@hotmail.com, xrz199721@gmail.com

yuh21@mails.tsinghua.edu.cn, zouh18@mails.tsinghua.edu.cn, cuip@tsinghua.edu.cn

Abstract

Recently, flat minima are proven to be effective for improving generalization and sharpness-aware minimization (SAM) achieves state-of-the-art performance. Yet the current definition of flatness discussed in SAM and its follow-ups are limited to the zeroth-order flatness (i.e., the worst-case loss within a perturbation radius). We show that the zeroth-order flatness can be insufficient to discriminate minima with low generalization error from those with high generalization error both when there is a single minimum or multiple minima within the given perturbation radius. Thus we present first-order flatness, a stronger measure of flatness focusing on the maximal gradient norm within a perturbation radius which bounds both the maximal eigenvalue of Hessian at local minima and the regularization function of SAM. We also present a novel training procedure named Gradient norm Aware Minimization (GAM) to seek minima with uniformly small curvature across all directions. Experimental results show that GAM improves the generalization of models trained with current optimizers such as SGD and AdamW on various datasets and networks. Furthermore, we show that GAM can help SAM find flatter minima and achieve better generalization. The code is available at <https://github.com/xxgege/GAM>.

1. Introduction

Current neural networks have achieved promising results in a wide range of fields [39, 57, 59, 73, 79–81, 84], yet they are typically heavily over-parameterized [2, 4]. Such heavy overparameterization leads to severe overfitting and poor generalization to unseen data when the model is learned simply with common loss functions (e.g., cross-entropy) [29]. Thus effective training algorithms are required to limit the negative effects of overfitting training data and find gen-

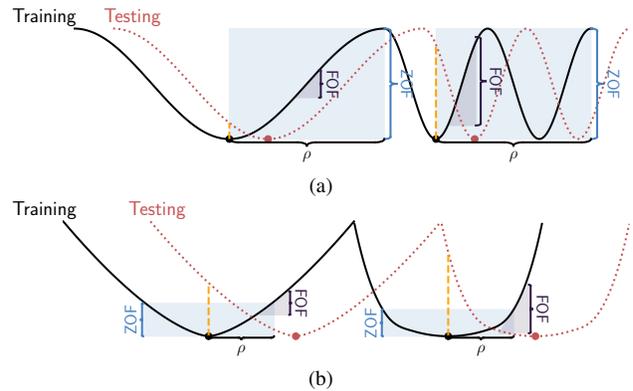


Figure 1. The comparison of the zeroth-order flatness (ZOF) and first-order flatness (FOF). Given a perturbation radius ρ , ZOF can fail to indicate generalization error both when there are multiple minima (1a) and a single minimum (1b) in the radius while FOF remains discriminative. The height of blue rectangles in curly brackets is the value of ZOF and the height of gray triangles (which indicates the slope) is the value of FOF. In Figure 1a, when ρ is large and enough to cover multiple minima, ZOF can not measure the fluctuation frequency while FOF prefers the flatter valley which has a smaller gradient norm. When ρ is small and covers only a single minimum, the maximum loss in ρ can be misleading as it can be misaligned with the uptrend of loss. As shown in Figure 1b, ZOF prefers the valley on the right, which has a larger generalization error (the orange dotted line), while FOF prefers the left one.

eralizable solutions.

Many studies try to improve model generalization by modifying the training procedure, such as batch normalization [28], dropout [25], and data augmentation [14, 74, 78]. Especially, some works discuss the connection between the geometry of the loss landscape and generalization [20, 23, 29]. A branch of effective approaches, sharpness-Aware Minimization (SAM) [20] and its variants [17, 18, 37, 47, 52, 83], minimizes the worst-case loss within a perturbation radius, which we call zeroth-order flatness. It is proven that optimizing the zeroth-order flatness leads to lower generalization error and achieves state-of-the-art per-

[†]Equal contribution, *Corresponding author

formance on various image classification tasks [20, 42, 86].

Optimizing the worst case, however, relies on a reasonable choice of perturbation radius ρ . As a prefixed hyperparameter in SAM or a hyperparameter under parameter re-scaling in its variants, such as ASAM [42], ρ can not always be a perfect choice in the whole training process. We show that the zeroth-order flatness may fail to indicate the generalization error with a given ρ . As in Figure 1a, when ρ covers multiple minima, the zeroth-order flatness (SAM) can not measure the fluctuation frequency. When there is a single minimum within ρ , as in Figure 1b the observation radius is limited and the maximum loss in ρ can be misaligned with the uptrend of loss. So zeroth-order flatness can be misleading and the knowledge of loss gradient is required for generalization error minimization.

To address this problem, we introduce first-order flatness, which controls the maximum gradient norm in the neighborhood of minima. We show that the first-order flatness is stronger than the zeroth-order flatness as the loss intensity of the loss fluctuation can be bounded by the maximum gradient. When the perturbation radius covers multiple minima, which we show is quite common in practice, the first-order flatness discriminates more drastic jitters from real flat valleys, as in Figure 1a. When the perturbation radius is small and covers only one minimum, the first-order flatness demonstrates the trend of loss gradient and can help indicate generalization error. We further show that the first-order flatness directly controls the maximal eigenvalue of Hessian of the training loss, which is a proper sharpness/flatness measure indicating the loss uptrend under an adversarial perturbation to the weights [34–36].

To optimize the first-order flatness in deep model training, we propose Gradient norm Aware Minimization (GAM), which approximates the maximum gradient norm with stochastic gradient ascent and Hessian-vector products to avoid the materialization of the Hessian matrix.

We summarize our contributions as follows.

- We present first-order flatness, which measures the largest gradient norm in the neighborhood of minima. We show that the first-order flatness is stronger than current zeroth-order flatness and it controls the maximum eigenvalue of Hessian.
- We propose a novel training procedure, GAM, to simultaneously optimize prediction loss and first-order flatness. We analyze the generalization error and the convergence of GAM.
- We empirically show that GAM considerably improves model generalization when combined with current optimizers such as SGD and AdamW across a wide range of datasets and networks. We show that GAM further improves the generalization of models trained with SAM.
- We empirically validate that GAM indeed finds flatter optima with lower Hessian spectra.

2. Related Works

Optimizer Some studies [20, 68] have demonstrated that current optimization approaches, such as SGD [53], Adam [38], AdamW [49] and others [19, 46] affect generalization. Some previous literature finds that Adam is more vulnerable to sharp minima than SGD [64], which results in worse generalization ability [22, 26, 67]. Some following works [10, 50, 68, 76] propose generalizable optimizers to address this problem. However, it can be a trade-off between generalization ability and convergence speed [19, 36, 46, 68, 76]. Different tasks and network architectures may agree with different optimizers (e.g., SGD is often chosen for ResNet [24] while AdamW [49] for ViTs [16]). Thus selecting a proper optimizer is critical while the understanding of its relationship to model generalization remains nascent [20].

Flat Minima and Generalization Many recent works show that flatter minima lead to better generalization [32, 36, 36, 55, 86]. Recently, [35] thoroughly reviews the literature related to generalization and sharpness of minima. It highlights the role of maximum Hessian eigenvalue in deciding the sharpness of minima [36, 63]. And there also have been several simple strategies to achieve a smaller maximum Hessian eigenvalue, such as choosing a large learning rate [12, 31, 44] and smaller batch size [30, 44, 61]. Sharpness-Aware Minimization (SAM) [20] and its variants [17, 18, 37, 42, 47, 52, 83, 86] are representative training algorithm to seek flat minima for better generalization. However, their definition of flatness is limited to zeroth-order flatness. In this paper, we present first-order flatness, a stronger flatness measure to learn better generalization. It is shown that discrete steps of gradient descent regularize deep models implicitly by penalizing the gradient descent trajectories with large loss gradients and this implicit regularization helps to find flat minima [7]. [82] proposes to directly control the gradient norm. They focus on the gradient norm at each training step, while we propose to penalize the maximum gradient norm in the neighborhood of minima and show the connection between our regularizer and the largest eigenvalue of Hessian and generalization error.

3. Preliminaries

Notations Let \mathcal{X} and \mathcal{Y} be the sample space and label space, respectively. Let \mathcal{D} denote the training distribution on $\mathcal{X} \times \mathcal{Y}$ and $S = \{(x_i, y_i)\}_{i=1}^n$ denote the training dataset with n data-points drawn independently from \mathcal{D} . Let $\theta \in \Theta \subseteq \mathbb{R}^d$ denote the parameters of the model. In addition, we use $B(\theta, \rho)$ to denote the open ball of radius $\rho > 0$ centered at the point θ in the Euclidean space, i.e., $B(\theta, \rho) = \{\theta' : \|\theta - \theta'\| < \rho\}$ ¹.

Let $\ell : \Theta \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ be the per-data-point loss

¹We use $\|\cdot\|$ to denote the L2 norm throughout the paper.

function. Let $\hat{L}(\boldsymbol{\theta}) = \sum_{i=1}^n \ell(\boldsymbol{\theta}, x_i, y_i)$ and $L(\boldsymbol{\theta}) = \mathbb{E}_{(x,y) \sim \mathcal{D}}[\ell(\boldsymbol{\theta}, x, y)]$ denote the empirical loss function and population-level loss function, respectively. We assume $\hat{L}(\boldsymbol{\theta})$ and $L(\boldsymbol{\theta})$ are twice differentiable throughout the paper. $\nabla L(\boldsymbol{\theta})$ and $\nabla^2 L(\boldsymbol{\theta})$ ($\nabla \hat{L}(\boldsymbol{\theta})$ and $\nabla^2 \hat{L}(\boldsymbol{\theta})$) are the derivative and Hessian matrix of the function $L(\cdot)$ ($\hat{L}(\cdot)$) at point $\boldsymbol{\theta}$, respectively. Besides, for any $\boldsymbol{\theta} \in \Theta$, we use $\nabla \|\nabla \hat{L}(\boldsymbol{\theta})\|$ to represent the gradient of function $\|\nabla \hat{L}(\cdot)\|$ at point $\boldsymbol{\theta}$. In addition, we use $L^{\text{oracle}}(\boldsymbol{\theta})$ to denote an oracle loss function and it can be chosen as empirical loss function $\hat{L}(\boldsymbol{\theta})$, $\hat{L}(\boldsymbol{\theta})$ with the weight decay regularization, and other common loss functions.

3.1. Zeroth-Order Flatness

The most popular mathematical definitions of flatness considers the maximal loss value within a radius [20, 36], which we call the zeroth-order flatness. We follow the loss function proposed in SAM:

$$L^{\text{sam}}(\boldsymbol{\theta}) = \hat{L}(\boldsymbol{\theta}) + \max_{\boldsymbol{\theta}' \in B(\boldsymbol{\theta}, \rho)} \left(\hat{L}(\boldsymbol{\theta}') - \hat{L}(\boldsymbol{\theta}) \right). \quad (1)$$

The second term in the right-hand side of Equation (1) can be considered as a measure of the zeroth-order flatness.

Definition 3.1 (ρ -zeroth-order flatness). For any $\rho > 0$, the ρ -zeroth-order flatness $R_\rho^{(0)}(\boldsymbol{\theta})$ of function $\hat{L}(\boldsymbol{\theta})$ at a point $\boldsymbol{\theta}$ is defined as

$$R_\rho^{(0)}(\boldsymbol{\theta}) \triangleq \max_{\boldsymbol{\theta}' \in B(\boldsymbol{\theta}, \rho)} \left(\hat{L}(\boldsymbol{\theta}') - \hat{L}(\boldsymbol{\theta}) \right), \quad \forall \boldsymbol{\theta} \in \Theta. \quad (2)$$

Here ρ is the perturbation radius that controls the magnitude of the neighborhood.

Intuitively, we name the term zeroth-order flatness because it measures the gap between the maximum loss value and the current point. As a measure of accumulation of gradients, zeroth-order flatness can be insufficient to indicate the generalization loss as shown in Section 4.2. In this paper, we propose a novel first-order flatness measure and compare these two flatness notions in Section 4.2.

4. First-order Flatness and Optimization

In this section, we introduce the first-order flatness and the corresponding minimizer for optimization. In Section 4.1, we formulate the first-order flatness and show its connection with the maximal eigenvalue of the Hessian. Afterward, we discuss the relationship between the zeroth-order and first-order flatness in Section 4.2. In Section 4.3, we present the optimization framework based on the first-order flatness as shown in Algorithm 1. We further provide a generalization bound with respect to the empirical loss, the first-order flatness, and high order terms, indicating that optimizing the first-order flatness improves generalization abilities. We then prove the convergence of the algorithm.

4.1. First-order Flatness

We first introduce the formulation of the first-order flatness, which measures the maximal gradient norm in the neighbourhood of a point $\boldsymbol{\theta} \in \Theta$.

Definition 4.1 (ρ -first-order flatness). For any $\rho > 0$, the ρ -first-order flatness $R_\rho^{(1)}(\boldsymbol{\theta})$ of function $\hat{L}(\boldsymbol{\theta})$ at a point $\boldsymbol{\theta}$ is defined as

$$R_\rho^{(1)}(\boldsymbol{\theta}) \triangleq \rho \cdot \max_{\boldsymbol{\theta}' \in B(\boldsymbol{\theta}, \rho)} \left\| \nabla \hat{L}(\boldsymbol{\theta}') \right\|, \quad \forall \boldsymbol{\theta} \in \Theta. \quad (3)$$

Here ρ is the perturbation radius that controls the magnitude of the neighbourhood.

Intuitively, the first-order flatness entails that the loss function $\hat{L}(\boldsymbol{\theta})$ should not change drastically in the neighbourhood of $\boldsymbol{\theta}$ so that the largest gradient norm of loss is constrained.

We then discuss the relationship between the first-order flatness and the maximal eigenvalue of the Hessian matrix $\nabla^2 \hat{L}(\boldsymbol{\theta}^*)$ (denoted as $\lambda_{\max}(\nabla^2 \hat{L}(\boldsymbol{\theta}^*))$). λ_{\max} is proven to be a proper measure of the curvature of minima [35, 36] and is closely related to generalization abilities [11, 30, 63]. As another definition of flatness in related works [9, 44], λ_{\max} is widely accepted yet hard to calculate. We show in the following lemma that given a radius ρ , the first-order flatness controls λ_{\max} , which reinforces the validity of the first-order flatness.

Lemma 4.1. *Let $\boldsymbol{\theta}^*$ be a local minimum of \hat{L} . Suppose \hat{L} can be second-order Taylor approximated in the neighbourhood $B(\boldsymbol{\theta}^*, \rho)^2$, i.e., $\forall \boldsymbol{\theta} \in B(\boldsymbol{\theta}^*, \rho)$, $\hat{L}(\boldsymbol{\theta}) = \hat{L}(\boldsymbol{\theta}^*) + (\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \nabla^2 \hat{L}(\boldsymbol{\theta}^*)(\boldsymbol{\theta} - \boldsymbol{\theta}^*)/2$. Then*

$$\lambda_{\max} \left(\nabla^2 \hat{L}(\boldsymbol{\theta}^*) \right) = \frac{R_\rho^{(1)}(\boldsymbol{\theta}^*)}{\rho^2}. \quad (4)$$

Since the maximal eigenvalue of Hessian matrices is usually difficult to approximate and optimize directly [71, 72], the first-order flatness becomes a proper surrogate of λ_{\max} .

4.2. Comparison with Zeroth-order Flatness

We compare the first-order flatness with the zeroth-order flatness. We first show that $R_\rho^{(0)}(\boldsymbol{\theta})$ in Equation (2) is bounded by $R_\rho^{(1)}(\boldsymbol{\theta})$ in Equation (3).

Proposition 4.2. *For any $\boldsymbol{\theta} \in \Theta$, $R_\rho^{(0)}(\boldsymbol{\theta})$ is bounded by $R_\rho^{(1)}(\boldsymbol{\theta})$, i.e., $R_\rho^{(1)}(\boldsymbol{\theta}) \geq R_\rho^{(0)}(\boldsymbol{\theta})$.*

Thus a smaller $R_\rho^{(1)}$ also leads to a smaller $R_\rho^{(0)}$, indicating that $R_\rho^{(1)}$ is a stronger flatness measure than $R_\rho^{(0)}$.

²The second order Taylor approximation assumption is commonly adopted in optimization-related literature [51, 66, 68, 77] to analyze the properties near critical points.

Proposition 4.2 gives an explanation that the first-order flatness covers wider scenarios compared with the zeroth-order flatness.

We present scenarios where the zeroth-order flatness fails to indicate generalization error while the first-order flatness remains discriminative in Figure 1. The gap between a local minimum and the largest loss in ρ can be considered as an accumulation of gradients across the trajectory while the largest gradient norm measures the maximum ascent rate, which may indicate the trends of loss outside of ρ .

When ρ is large, there probably exist several other local minima in the neighborhood $B(\theta^*, \rho)$ as shown in Figure 1a. This case is common in practice as shown in Section 5.1. In addition, when the number of local minimum in $B(\theta^*, \rho)$ becomes larger, θ^* is expected to become sharper since the valley of θ^* becomes narrower. However, the zeroth-order flatness $R_\rho^{(0)}$ only measures the maximal gap of the loss function in $B(\theta^*, \rho)$ and fails to distinguish the cases when the number of local minimums varies. By contrast, the maximal gradient norm in $B(\theta^*, \rho)$ increases when the number of local minima is larger, indicating that the first-order flatness can successfully characterize the sharpness in this case.

When ρ only covers a single minimum, as shown in Figure 1b, the zeroth-order flatness in ρ can be misleading since the observation radius is insufficient to measure the loss trend with the maximum loss. The first-order flatness can help to learn more about the loss trend.

From the perspective of flatness, the zeroth-order flatness focuses on the average gradient within a radius while the first-order flatness measures the maximum gradient. Intuitively, the combination of the zeroth-order and first-order captures a more comprehensive picture of the loss landscape. Furthermore, as discussed in the following Section 4.3, minimizers for both flatness measures adopt the first-order approximation to calculate the maxima within a radius. This may be the reason that the combination of the two flatness measures achieves the best performance as shown in Section 5.

4.3. Gradient Norm Aware Minimization

In this subsection, we propose a novel Gradient norm Aware Minimization (GAM) framework to incorporate the first-order flatness $R_\rho^{(1)}(\theta)$ into optimization procedures.

Specifically, suppose we could obtain an oracle loss function $L^{\text{oracle}}(\theta)$ and calculate its gradient $\nabla L^{\text{oracle}}(\theta)$. $L^{\text{oracle}}(\theta)$ can be chosen as the empirical loss function $\hat{L}(\theta)$ and the empirical loss function with other regularizations (such as the weight decay and the zeroth-order flatness as shown in Definition 3.1).

Generalization analysis We first derive a generalization bound *w.r.t.* the first-order flatness in Proposition 4.3.

Proposition 4.3. *Suppose the per-data-point loss function ℓ is differentiable and bounded by M . Fix $\rho > 0$ and $\theta \in \Theta$. Then with probability at least $1 - \delta$ over training set S generated from the distribution \mathcal{D} ,*

$$\begin{aligned} & \mathbb{E}_{\epsilon_i \sim N(0, \rho^2 / (\sqrt{d} + \sqrt{\log n})^2)} [L(\theta + \epsilon)] \\ & \leq \hat{L}(\theta) + R_\rho^{(1)}(\theta) + \frac{M}{\sqrt{n}} \\ & + \sqrt{\frac{\frac{1}{4}d \log \left(1 + \frac{\|\theta\|^2 (\sqrt{d} + \sqrt{\log n})^2}{d\rho^2} \right) + \frac{1}{4} + \log \frac{n}{\delta} + 2 \log(6n + 3d)}{n - 1}}. \end{aligned} \quad (5)$$

Remark. The left-hand side of Equation (5) is close to the population-level loss function $L(\theta)$ since the numbers of samples n and parameters d are often large. As a result, ignoring high-order terms, the population-level loss $L(\theta)$ is bounded by the empirical loss $\hat{L}(\theta)$ and the first-order flatness $R_\rho^{(1)}(\theta)$, which motivates us to use $R_\rho^{(1)}(\theta)$ as a regularizer to help improve the generalization abilities of models.

Inspired by Lemma 4.1 and Proposition 4.3, the overall loss function is given by

$$L^{\text{overall}}(\theta) = L^{\text{oracle}}(\theta) + \alpha R_\rho^{(1)}(\theta), \quad (6)$$

where α is a hyperparameter that determines the strength of regularization. The gradient of the loss function $L^{\text{overall}}(\theta)$ is given by $\nabla L^{\text{overall}}(\theta) = \nabla L^{\text{oracle}}(\theta) + \alpha \nabla R_\rho^{(1)}(\theta)$. Using similar techniques in [20], GAM approximates $\nabla R_\rho^{(1)}(\theta)$ by

$$\begin{aligned} \nabla R_\rho^{(1)}(\theta) & \approx \rho \cdot \nabla \left\| \nabla \hat{L}(\theta^{\text{adv}}) \right\|, \quad \theta^{\text{adv}} = \theta + \rho \cdot \frac{\mathbf{f}}{\|\mathbf{f}\|}, \\ \mathbf{f} & = \nabla \left\| \nabla \hat{L}(\theta) \right\|. \end{aligned} \quad (7)$$

Details of the derivation of $\nabla R_\rho^{(1)}(\theta)$ can be found in Appendix A. Notice that

$$\forall \theta \in \Theta, \quad \nabla \left\| \nabla \hat{L}(\theta) \right\| = \frac{\nabla^2 \hat{L}(\theta) \cdot \nabla \hat{L}(\theta)}{\|\nabla \hat{L}(\theta)\|}. \quad (8)$$

As a result, Equation (7) can be calculated efficiently by the Hessian vector product. The pseudocode of the whole optimization procedure is shown in Algorithm 1.

Convergence analysis We further analyze the convergence properties of GAM. Firstly, we introduce the Lipschitz smoothness, which is common adopted in optimization-related literature [1, 70, 86].

Definition 4.2. A function $J : \Theta \rightarrow \mathbb{R}$ is γ -Lipschitz smooth if

$$\forall \theta_1, \theta_2 \in \Theta, \quad \left\| \nabla J(\theta_1) - \nabla J(\theta_2) \right\| \leq \gamma \|\theta_1 - \theta_2\|. \quad (9)$$

With Definition 4.2, we could prove the convergence property of GAM as shown in Theorem 4.4.

Theorem 4.4. Suppose $L^{\text{oracle}}(\theta)$ is γ_1 -Lipschitz smooth and $\hat{L}(\theta)$ is γ_2 -Lipschitz smooth. Suppose $|L^{\text{oracle}}(\theta)|$ is bounded by M . For any timestamp $t \in \{0, 1, \dots, T\}$ and any $\theta \in \Theta$, suppose we can obtain noisy and bounded observations $g_t^{\text{loss}}(\theta)$, $g_t^{\text{norm}}(\theta)$, and $\tilde{g}_t^{\text{loss}}(\theta)$ of $\nabla \hat{L}(\theta)$, $\nabla \|\nabla \hat{L}(\theta)\|$, and $\nabla L^{\text{oracle}}(\theta)$ such that

$$\begin{aligned} \mathbb{E}[g_t^{\text{loss}}(\theta)] &= \nabla \hat{L}(\theta), \|g_t^{\text{loss}}(\theta)\| \leq G^{\text{loss}}, \|g_t^{\text{norm}}(\theta)\| \leq G^{\text{norm}}, \\ \mathbb{E}[\tilde{g}_t^{\text{loss}}(\theta)] &= \nabla L^{\text{oracle}}(\theta), \|\tilde{g}_t^{\text{loss}}(\theta)\| \leq \tilde{G}^{\text{loss}}. \end{aligned} \quad (10)$$

Then with learning rate $\eta_t = \eta_0/\sqrt{t}$ and perturbation radius $\rho_t = \rho_0/\sqrt{t}$, GAM could obtain

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\|\nabla L^{\text{overall}}(\theta_t)\|^2 \right] \leq \frac{C_1 + C_2 \log T}{\sqrt{T}}, \quad (11)$$

for some constants C_1 and C_2 that only depend on $\gamma, G^{\text{loss}}, G^{\text{norm}}, \tilde{G}^{\text{loss}}, M, \eta_0, \rho_0$, and α . Here $\nabla L^{\text{overall}}(\theta_t) = \nabla L^{\text{oracle}}(\theta_t) + \alpha \nabla R_{\rho_t}^{(1)}(\theta_t)$ and $\nabla R_{\rho_t}^{(1)}(\theta_t)$ is approximated in Equation (7).

Remark. The assumptions in Theorem 4.4 are common and standard when analyzing convergence of non-convex functions via SGD-based methods [38, 56, 86]. In addition, the requirements on $L^{\text{oracle}}(\theta)$ (i.e., $L^{\text{oracle}}(\theta)$ is Lipschitz smooth and we can obtain unbiased and bounded observations of $\nabla L^{\text{oracle}}(\theta)$) are mild and common. For example, when the empirical loss function $\hat{L}(\theta)$ satisfies the constraints, it is easy to check that $\hat{L}(\theta)$ with the weight decay regularization also meets the requirements.

5. Experiments

We empirically show that the case discussed in Section 4.2 is common in practice. Then we evaluate GAM with random initialization on various state-of-the-art models and the transfer learning setting on various datasets. We show the Hessian spectra of GAM at convergence and discuss the computation overhead of GAM with the considerable improvement of model generalization.

5.1. The Density of Local Minima

To investigate the number of local minima within the perturbation radius, we train 3 ResNet-18 models with SAM on CIFAR-100 with proper hyperparameters for 200 epochs. The perturbation radius is set to 0.1 as suggested by [20]. We load the checkpoints at convergence for evaluation. We randomly generate 100 perturbation directions with the same size as the model weights for each model. For each direction, we repeatedly add a perturbation with the norm of 0.01 along the selected direction 10 times. We calculate the

Algorithm 1 Gradient norm Aware Minimization (GAM)

- 1: **Input:** Batch size b , Learning rate η_t , Perturbation radius ρ_t , Trade-off coefficient α , Small constant ξ
 - 2: $t \leftarrow 0, \theta_0 \leftarrow$ initial parameters
 - 3: **while** θ_t not converged **do**
 - 4: Sample W_t from the training data with b instances
 - 5: $\mathbf{h}_t^{\text{loss}} \leftarrow \nabla L^{\text{oracle}}(\theta_t)$ ▷ Calculate the oracle loss gradient $\nabla L^{\text{oracle}}(\theta_t)$
 - 6: $\mathbf{f}_t \leftarrow \nabla^2 \hat{L}_{W_t}(\theta_t) \cdot \frac{\nabla \hat{L}_{W_t}(\theta_t)}{\|\nabla \hat{L}_{W_t}(\theta_t)\| + \xi}$
 - 7: $\theta_t^{\text{adv}} \leftarrow \theta_t + \rho_t \cdot \frac{\mathbf{f}_t}{\|\mathbf{f}_t\| + \xi}$
 - 8: $\mathbf{h}_t^{\text{norm}} \leftarrow \rho_t \cdot \nabla^2 \hat{L}_{W_t}(\theta_t^{\text{adv}}) \cdot \frac{\nabla \hat{L}_{W_t}(\theta_t^{\text{adv}})}{\|\nabla \hat{L}_{W_t}(\theta_t^{\text{adv}})\| + \xi}$ ▷
 - 9: Calculate the norm gradient $\nabla R_{\rho_t}^{(1)}(\theta_t)$
 - 10: $\theta_{t+1} \leftarrow \theta_t - \eta_t(\mathbf{h}_t^{\text{loss}} + \alpha \mathbf{h}_t^{\text{norm}})$
 - 11: $t \leftarrow t + 1$
 - 12: **end while**
 - 13: **return** θ_t
-

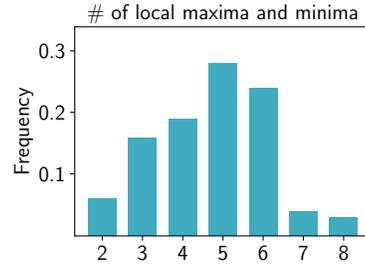


Figure 2. The distribution of numbers of local minima and maxima within the perturbation radius ρ after convergence.

training loss after each addition and report the distribution of the number of local maxima and minima along each perturbation direction within the perturbation radius ρ of 0.1. As shown in Figure 2, we find more than 1 local minima within ρ for most of the directions, indicating that the case is common in practice. As discussed in Section 4.2, zeroth-order flatness fails to tell the sharpness caused by multiple minima while the first-order flatness measure increases as the sharpness grow.

5.2. Training from Scratch

5.2.1 CIFAR-10 and CIFAR-100

We conduct experiments on CIFAR-10 and CIFAR-100 [41] with ResNets [24], WideResNet [75], ResNeXt [65], PyramidNet [21] and Vision Transformers (ViTs) [16]. All the models are trained for 200 epochs from scratch. We evaluate GAM both with basic data augmentations (i.e., horizontal flip, padding by four pixels, and random crop) and advanced data augmentation including cutout regularization [15], RandAugment [13] and AutoAugment [14].

GAM has two hyperparameters, ρ and α . We conduct a

Table 1. Results of GAM with state-of-the-art models on CIFAR-10 and CIFAR-100. The best results are highlighted in bold font.

Model	Aug	CIFAR-10				CIFAR-100			
		SGD	SGD + GAM	SAM	SAM + GAM	SGD	SGD + GAM	SAM	SAM + GAM
ResNet18	Basic	95.32 \pm 0.13	96.17 \pm 0.21	96.10 \pm 0.20	96.75 \pm 0.18	78.32 \pm 0.32	79.53 \pm 0.30	79.27 \pm 0.16	80.45 \pm 0.25
ResNet18	Cutout	95.99 \pm 0.13	96.46 \pm 0.20	96.64 \pm 0.13	96.99 \pm 0.23	78.73 \pm 0.13	79.89 \pm 0.31	79.43 \pm 0.15	80.80 \pm 0.14
ResNet18	RA	96.07 \pm 0.07	96.52 \pm 0.09	96.64 \pm 0.17	97.06 \pm 0.13	78.62 \pm 0.32	79.82 \pm 0.24	79.71 \pm 0.15	80.97 \pm 0.29
ResNet18	AA	96.13 \pm 0.05	96.71 \pm 0.07	96.75 \pm 0.08	97.17 \pm 0.08	78.88 \pm 0.15	80.56 \pm 0.21	80.58 \pm 0.25	81.59 \pm 0.24
ResNet101	Basic	96.35 \pm 0.08	96.98 \pm 0.11	96.82 \pm 0.16	97.20 \pm 0.15	80.47 \pm 0.13	82.21 \pm 0.40	82.03 \pm 0.12	83.13 \pm 0.07
ResNet101	Cutout	96.56 \pm 0.18	97.22 \pm 0.05	97.07 \pm 0.08	97.36 \pm 0.24	80.53 \pm 0.30	82.36 \pm 0.24	81.60 \pm 0.35	83.40 \pm 0.13
ResNet101	RA	96.68 \pm 0.25	97.33 \pm 0.30	97.12 \pm 0.18	97.40 \pm 0.23	80.60 \pm 0.28	82.40 \pm 0.31	82.19 \pm 0.34	83.28 \pm 0.20
ResNet101	AA	96.78 \pm 0.14	97.39 \pm 0.18	97.18 \pm 0.11	97.42 \pm 0.1	81.83 \pm 0.37	83.19 \pm 0.15	82.44 \pm 0.47	83.94 \pm 0.23
WRN28_2	Basic	94.82 \pm 0.07	95.69 \pm 0.13	95.47 \pm 0.08	95.85 \pm 0.08	75.45 \pm 0.25	77.21 \pm 0.31	77.04 \pm 0.18	77.69 \pm 0.20
WRN28_2	Cutout	95.70 \pm 0.20	96.41 \pm 0.18	96.22 \pm 0.13	96.39 \pm 0.22	76.80 \pm 0.45	78.58 \pm 0.24	78.04 \pm 0.43	79.33 \pm 0.12
WRN28_2	RA	95.75 \pm 0.16	96.35 \pm 0.13	96.22 \pm 0.08	96.49 \pm 0.20	76.73 \pm 0.27	78.66 \pm 0.03	77.88 \pm 0.29	78.96 \pm 0.13
WRN28_2	AA	95.44 \pm 0.06	95.98 \pm 0.09	96.07 \pm 0.08	96.44 \pm 0.09	77.35 \pm 0.02	79.05 \pm 0.10	78.64 \pm 0.23	79.50 \pm 0.21
WRN28_10	Basic	95.73 \pm 0.10	96.61 \pm 0.15	96.78 \pm 0.80	97.29 \pm 0.11	81.40 \pm 0.13	83.45 \pm 0.09	83.41 \pm 0.04	84.31 \pm 0.06
WRN28_10	Cutout	96.74 \pm 0.03	96.97 \pm 0.05	97.35 \pm 0.16	97.56 \pm 0.12	81.53 \pm 0.40	83.69 \pm 0.08	82.38 \pm 0.15	84.43 \pm 0.13
WRN28_10	RA	97.14 \pm 0.04	96.83 \pm 0.03	97.58 \pm 0.07	97.49 \pm 0.03	81.65 \pm 0.18	83.84 \pm 0.09	82.79 \pm 0.06	84.68 \pm 0.13
WRN28_10	AA	96.93 \pm 0.12	97.05 \pm 0.04	97.48 \pm 0.06	97.67 \pm 0.08	81.99 \pm 0.11	84.02 \pm 0.18	83.84 \pm 0.30	84.81 \pm 0.21
PyramidNet110	Basic	96.19 \pm 0.11	97.11 \pm 0.14	97.26 \pm 0.05	97.51 \pm 0.09	82.74 \pm 0.12	84.91 \pm 0.09	85.01 \pm 0.09	85.25 \pm 0.06
PyramidNet110	Cutout	96.82 \pm 0.09	97.32 \pm 0.21	97.49 \pm 0.06	97.91 \pm 0.14	83.31 \pm 0.21	85.20 \pm 0.19	84.90 \pm 0.03	85.46 \pm 0.10
PyramidNet110	RA	97.15 \pm 0.21	97.80 \pm 0.22	97.60 \pm 0.09	98.01 \pm 0.10	84.04 \pm 0.19	86.47 \pm 0.14	85.33 \pm 0.27	85.64 \pm 0.20
PyramidNet110	AA	97.11 \pm 0.01	97.85 \pm 0.02	97.61 \pm 0.14	97.95 \pm 0.10	84.48 \pm 0.03	85.92 \pm 0.03	85.69 \pm 0.17	86.35 \pm 0.18

grid search over $\{0.05, 0.1, 0.2, 0.5, 1.0, 2.0\}$ to tune ρ and $\{0.1, 0.2, 0.5, 1.0, 2.0, 3.0, \dots, 10.0\}$ for α using 10% of the training data as a validation set. The selection of hyperparameters is in Appendix C.5.

As a gradient regularizer, GAM can be integrated with current optimizers such as SGD and Adam [36]. We also show that GAM can be combined with sharpness-aware training procedures such as SAM. As shown in Section 4.2, the GAM term bounds the regularization term in SAM. Yet the practical implementations of GAM and SAM rely on first-order Taylor expansion of different objective functions (GAM approximates the maximum gradient norm while SAM approximates the maximum loss). We empirically show that the combination of GAM and SAM outperforms both of them, indicating that they may strengthen each other with omitted items.

As shown in Table 1, GAM improves generalization for all models on CIFAR-10 and CIFAR-100. When combined with SGD, GAM achieves considerably higher test accuracy compared with SGD. Moreover, GAM further improves generalization when combined with SAM. For example, GAM improves SAM performance by 1.18% and 1.10% on CIFAR-100 with ResNet-18 and ResNet-101, respectively, which are noticeable margins. Other experimental results are in Appendix C.1.

5.2.2 ImageNet

We use ResNet50, ResNet101 [24], ViT-S/32 and ViT-B/32 [16] for evaluations on ImageNet [58] to evaluate GAM on

large scale data. For ResNet, we use SGD with momentum=0.9 as the base optimizer for both GAM and SAM. For ViT, we use the AdamW optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$. We train ResNets for 90 epochs and ViTs for 300 epochs following [16]. We set the batch size to 256, learning rate to 0.1, and weight decay to 0.0001. The learning rate is decayed using a cosine schedule.

As shown in Table 2, GAM consistently improves SGD performance on ImageNet for both ResNets and ViTs. GAM also further improves the model generalization compared with SAM. The combination of GAM and SAM outperforms both SGD and SAM by a noticeable margin.

5.3. Transfer Learning

Transfer learning shows the generalization of models when trained on sufficient labeled data and finetuned on a small dataset [85]. We show that GAM improves generalization on all datasets in this setting.

We consider Stanford Cars [40], CIFAR-10, CIFAR-100 [41], Oxford_IIIT_Pets [54] and Food101 [8] for this setting. We apply SGD, SAM, and GAM to finetuning EfficientNet-b0 [62] and Swin-Transformer-t [48] on these datasets. Both EfficientNet-b0 and Swin-Transformer-t are pretrained on ImageNet.

We use ImageNet pretrained weights of EfficientNet-b0 and Swin-t except for the last linear layer for classification. Following previous works, we train for 40k steps since our batch size is 128. The initial learning rate is set to 2e-3 with cosine learning rate decay. Weight decay is set to 1e-5. We do not use any data augmentations for Stanford Cars, Ox-

Table 2. Results of GAM with ResNet50 on ImageNet.

Model	Dataset	Base Opt	Base + GAM	SAM	SAM + GAM
ResNet50	Top-1	76.01 \pm 0.19	76.59 \pm 0.15	76.47 \pm 0.11	76.86 \pm 0.15
ResNet50	Top-5	92.75 \pm 0.08	93.10 \pm 0.08	93.07 \pm 0.05	93.22 \pm 0.06
ResNet101	Top-1	77.69 \pm 0.08	78.45 \pm 0.10	78.35 \pm 0.12	78.70 \pm 0.12
ResNet101	Top-5	93.76 \pm 0.09	94.09 \pm 0.12	94.02 \pm 0.06	94.15 \pm 0.12
ViT-S/32	Top-1	68.26 \pm 0.22	69.95 \pm 0.16	69.73 \pm 0.05	70.15 \pm 0.18
ViT-S/32	Top-5	87.39 \pm 0.19	88.11 \pm 0.26	87.91 \pm 0.30	88.23 \pm 0.18
ViT-B/32	Top-1	71.15 \pm 0.14	73.58 \pm 0.06	73.10 \pm 0.18	73.70 \pm 0.10
ViT-B/32	Top-5	90.12 \pm 0.07	91.15 \pm 0.19	91.03 \pm 0.06	91.50 \pm 0.16

Table 3. Results of GAM for finetuning EfficientNet-b0 and Swin Transformers on various datasets.

Dataset	EfficientNet-b0				Swin-t			
	SGD	SGD + GAM	SAM	SAM + GAM	AdamW	AdamW + GAM	SAM	SAM + GAM
Stanford Cars	82.14	83.50	83.21	83.98	83.50	84.90	83.55	85.29
CIFAR-10	86.26	87.37	86.95	87.97	91.32	92.06	91.77	92.55
CIFAR-100	63.75	64.85	64.29	65.03	72.88	73.78	73.99	74.30
Oxford_IIIT_Pets	91.03	91.80	91.65	91.96	93.49	93.87	93.59	94.03
Food101	82.54	82.69	82.57	83.01	86.38	86.89	86.64	87.03

ford_IIIT_Pets and Food101. For CIFAR datasets, we employ the same data augmentations as previous experiments.

As seen in Table 3, GAM once again brings generalization improvement for SGD, AdamW, and SAM on both EfficientNet-b0 and Swin-t. For example, GAM improves AdamW by 1.2% on Stanford Cars with Swin-t and 1.11% on CIFAR-10 with EfficientNet-b0.

Moreover, we leave the experiments of robustness to label noise in Appendix C.2.

5.4. Top Eigenvalues of Hessian and Hessian Trace

Lemma 4.1 shows that the GAM term can be an equivalent measure of the maximum eigenvalue of the Hessian, which is a well-known measure of flatness/sharpness. Thus optimizing the GAM term decreases the maximum eigenvalue of the Hessian and leads to flatter minima. To empirically validate that GAM finds optima with low curvature, we present the Hessian spectra of SGD, SAM, and GAM. We consider the maximum eigenvalue of Hessian and the Hessian trace, which measures the expected loss increase under random perturbations to the weights [35] as the measures of flatness. We empirically show that GAM significantly decreases both the maximum eigenvalue and the trace of Hessian during training compared with SGD and SAM, and thus finds flatter minima.

We compute the Hessian spectra of ResNet-18 trained on CIFAR-100 for 200 epochs with SGD, SAM, SGD + GAM, and SAM + GAM. We use power iteration [72] to compute the top eigenvalues of Hessian and Hutchinson’s method [5, 6, 71] to compute the Hessian trace. We report the histogram of the distribution of the top-50 Hessian

eigenvalues for each method.

As shown in Figure 3, the model trained with SGD has a higher maximum Hessian eigenvalue and Hessian trace at convergence compared to the middle of training, indicating that optimizing directly with cross-entropy loss does not contribute to the lower Hessian spectra. In contrast, GAM leads to lower Hessian spectra and thus flatter minima. Moreover, GAM helps to reduce both top eigenvalues and the Hessian trace when combined with SAM, where Hessian spectra at convergence are lower than other methods. We show visualizations of landscapes of SGD, SAM, and GAM in Section 5.6.

5.5. Computation Overhead

As discussed in Section 4.3, the GAM term can be easily calculated via the Hessian vector product, which is an efficient approach to calculating the dot product between the Hessian and a vector without the need to calculate the entire Hessian. However, it can still introduce extra computation when calculated in each iteration. To accelerate the training with GAM, we investigate applying GAM to only a few iterations in each epoch. Surprisingly, we show that only several iterations of learning with GAM (with higher α compared with applying GAM to all iterations) improve model generalization considerably. As shown in Figure 4, with approximately 1/20 of iterations, GAM considerably improves test accuracy for both SGD and SAM on CIFAR-10 and CIFAR-100. When applying GAM to 1/10 iterations of training, it shows similar effectiveness to applying GAM to all the iterations, while the extra computational cost for GAM is less than 25% of the original

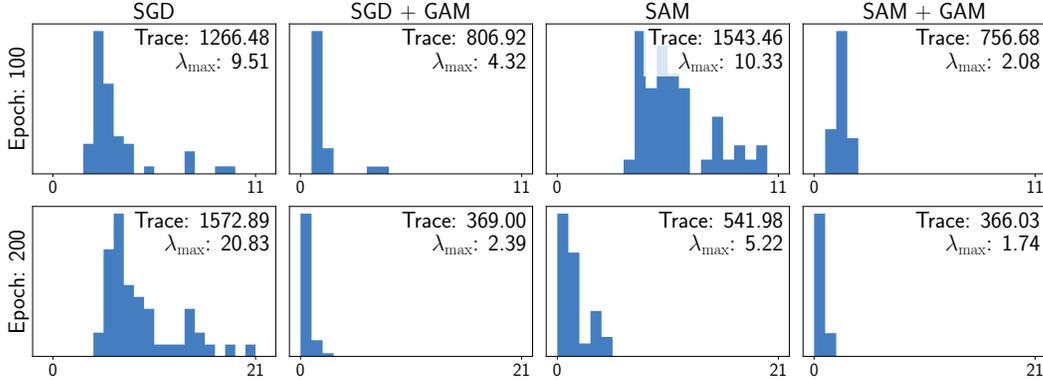


Figure 3. The distribution of top eigenvalues and the trace of Hessian at epoch 100 and 200 on CIFAR-100 with SGD, SGD + GAM, SAM, or SAM + GAM.

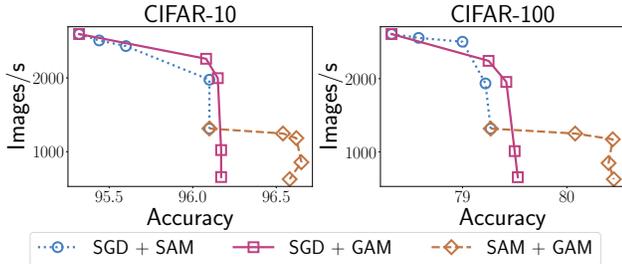


Figure 4. Accuracy and training speed of training with different ratios $\{0, 0.05, 0.1, 0.5, 1\}$ from upper left to lower right, see details in Appendix C.3) of iterations using GAM. Numbers in parentheses indicate the ratio of the training speed compared with the vanilla base optimizer SGD/SAM.

cost. GAM outperforms SAM with lower computation overhead and achieves significant improvement when combined with SGD (the red line in the figure). When combined with SAM, GAM also improves generalization with low computation cost. Thus the computation overhead of GAM can be easily controlled. The optimization of first-order flatness can be further accelerated by approximation of second-order gradient with first-order gradient and the details are in Appendix D.

5.6. Visualization of Landscapes

We visualize the loss landscapes of models trained with SGD, SGD + GAM, SAM, SAM + GAM of the ResNet-18 model on CIFAR-100 following [45]. All the models are trained with the same hyperparameters for 200 epochs as described in Section 5.2.1. As shown in Figure 5, GAM consistently helps SGD and SAM find flatter minima.

6. Discussions

We show that the most popular definitions of flatness, which we call the zeroth-order flatness, can be insufficient to indicate generalization error. Thus we propose first-order flatness, a stronger flatness measure that bounds both the

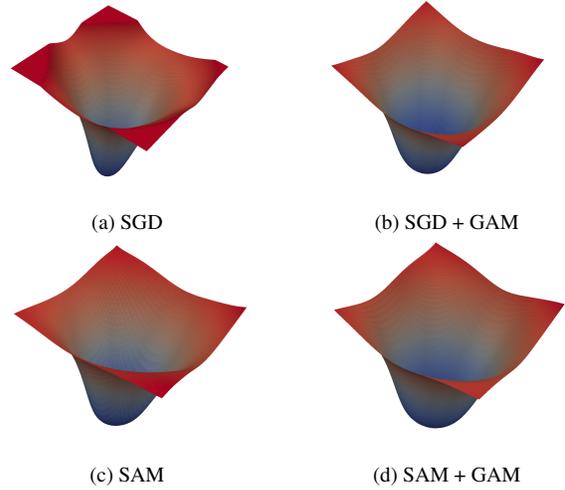


Figure 5. Visualization of loss landscape for SGD, SGD + GAM, SAM, SAM + GAM.

maximum eigenvalue of Hessian and the zeroth-order flatness. We also propose a novel Gradient norm Aware Minimization (GAM) to optimize the first-order flatness. We empirically show that GAM considerably improves generalization for SGD, AdamW, and SAM.

Despite the empirical effectiveness of GAM, adopting the first-order flatness for generalization has the following limitations which could lead to potential future work. First, a theoretical explanation of whether a stronger flatness measure is better for generalization is vital for selecting flatness measures in practice. Second, the contribution to generalization of combining the zeroth-order and first-order flatness requires a thorough theoretical analysis.

Acknowledgement

This work was supported in part by National Key R&D Program of China (No. 2018AAA0102004, No. 2020AAA0106300), National Natural Science Foundation of China (No. U1936219, 62141607), Beijing Academy of Artificial Intelligence (BAAI).

References

- [1] Zeyuan Allen-Zhu and Yuanzhi Li. Neon2: Finding local minima via first-order oracles. *Advances in Neural Information Processing Systems*, 31, 2018. 4
- [2] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In *International Conference on Machine Learning*, pages 242–252. PMLR, 2019. 1
- [3] Eric Arazo, Diego Ortego, Paul Albert, Noel O’Connor, and Kevin McGuinness. Unsupervised label noise modeling and loss correction. In *International conference on machine learning*, pages 312–321. PMLR, 2019. 16
- [4] Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit acceleration by over-parameterization. In *International Conference on Machine Learning*, pages 244–253. PMLR, 2018. 1
- [5] Haim Avron and Sivan Toledo. Randomized algorithms for estimating the trace of an implicit symmetric positive semi-definite matrix. *Journal of the ACM (JACM)*, 58(2):1–34, 2011. 7
- [6] Zhaojun Bai, Gark Fahey, and Gene Golub. Some large-scale matrix computation problems. *Journal of Computational and Applied Mathematics*, 74(1-2):71–89, 1996. 7
- [7] David GT Barrett and Benoit Dherin. Implicit gradient regularization. *arXiv preprint arXiv:2009.11162*, 2020. 2
- [8] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative components with random forests. In *European conference on computer vision*, pages 446–461. Springer, 2014. 6, 19
- [9] Pratik Chaudhari, Anna Choromanska, Stefano Soatto, Yann LeCun, Carlo Baldassi, Christian Borgs, Jennifer Chayes, Levent Sagun, and Riccardo Zecchina. Entropy-sgd: Biasing gradient descent into wide valleys. *Journal of Statistical Mechanics: Theory and Experiment*, 2019(12):124018, 2019. 3
- [10] Jinghui Chen, Dongruo Zhou, Yiqi Tang, Ziyang Yang, Yuan Cao, and Quanquan Gu. Closing the generalization gap of adaptive gradient methods in training deep neural networks. *arXiv preprint arXiv:1806.06763*, 2018. 2
- [11] Xiangning Chen, Cho-Jui Hsieh, and Boqing Gong. When vision transformers outperform resnets without pre-training or strong data augmentations. *arXiv preprint arXiv:2106.01548*, 2021. 3, 19
- [12] Jeremy M. Cohen, Simran Kaur, Yuanzhi Li, J. Zico Kolter, and Ameet Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. 2
- [13] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation policies from data. *arXiv preprint arXiv:1805.09501*, 2018. 5, 16
- [14] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 702–703, 2020. 1, 5, 16
- [15] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017. 5, 16
- [16] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020. 2, 5, 6, 16
- [17] Jiawei Du, Hanshu Yan, Jiashi Feng, Joey Tianyi Zhou, Liangli Zhen, Rick Siow Mong Goh, and Vincent YF Tan. Efficient sharpness-aware minimization for improved training of neural networks. *arXiv preprint arXiv:2110.03141*, 2021. 1, 2
- [18] Jiawei Du, Daquan Zhou, Jiashi Feng, Vincent YF Tan, and Joey Tianyi Zhou. Sharpness-aware training for free. *arXiv preprint arXiv:2205.14083*, 2022. 1, 2
- [19] John Duchi, Elad Hazan, and Yoram Singer. Adaptive sub-gradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011. 2
- [20] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. In *International Conference on Learning Representations*, 2021. 1, 2, 3, 4, 5, 13, 14, 16
- [21] Dongyoon Han, Jiwhan Kim, and Junmo Kim. Deep pyramidal residual networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5927–5935, 2017. 5
- [22] Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International conference on machine learning*, pages 1225–1234. PMLR, 2016. 2
- [23] Haowei He, Gao Huang, and Yang Yuan. Asymmetric valleys: Beyond sharp and flat local minima. *Advances in neural information processing systems*, 32, 2019. 1
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 2, 5, 6
- [25] Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012. 1
- [26] Sepp Hochreiter and Jürgen Schmidhuber. Simplifying neural nets by discovering flat minima. *Advances in neural information processing systems*, 7, 1994. 2
- [27] Gao Huang, Zhuang Liu, Geoff Pleiss, Laurens Van Der Maaten, and Kilian Weinberger. Convolutional networks with dense connectivity. *IEEE transactions on pattern analysis and machine intelligence*, 2019. 16
- [28] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015. 1
- [29] Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights

- leads to wider optima and better generalization. *arXiv preprint arXiv:1803.05407*, 2018. **1**
- [30] Stanisław Jastrzebski, Zachary Kenton, Devansh Arpit, Nicolas Ballas, Asja Fischer, Yoshua Bengio, and Amos Storkey. Three factors influencing minima in SGD. *arXiv preprint arXiv:1711.04623*, 2017. **2, 3**
- [31] Stanisław Jastrzebski, Maciej Szymczak, Stanislav Fort, Devansh Arpit, Jacek Tabor, Kyunghyun Cho, and Krzysztof J. Geras. The break-even point on optimization trajectories of deep neural networks. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. **2**
- [32] Zhiwei Jia and Hao Su. Information-theoretic local minima characterization and regularization. In *International Conference on Machine Learning*, pages 4773–4783. PMLR, 2020. **2**
- [33] Lu Jiang, Di Huang, Mason Liu, and Weilong Yang. Beyond synthetic noise: Deep learning on controlled noisy labels. In *International Conference on Machine Learning*, pages 4804–4815. PMLR, 2020. **16**
- [34] Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. *arXiv preprint arXiv:1912.02178*, 2019. **2**
- [35] Simran Kaur, Jeremy Cohen, and Zachary C Lipton. On the maximum hessian eigenvalue and generalization. *arXiv preprint arXiv:2206.10654*, 2022. **2, 3, 7**
- [36] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In *International Conference on Learning Representations*, 2017. **2, 3, 6**
- [37] Taero Kim, Sungjun Lim, and Kyungwoo Song. Sharpness-aware minimization for worst case optimization. *arXiv preprint arXiv:2210.13533*, 2022. **1, 2**
- [38] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *International Conference on Learning Representations*, 2015. **2, 5**
- [39] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016. **1**
- [40] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 554–561, 2013. **6, 19**
- [41] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *Citeseer*, 2009. **5, 6, 16, 19**
- [42] Jungmin Kwon, Jeongseop Kim, Hyunseo Park, and In Kwon Choi. Asam: Adaptive sharpness-aware minimization for scale-invariant learning of deep neural networks. In *International Conference on Machine Learning*, pages 5905–5914. PMLR, 2021. **2, 16**
- [43] Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, pages 1302–1338, 2000. **14**
- [44] Aitor Lewkowycz, Yasaman Bahri, Ethan Dyer, Jascha Sohl-Dickstein, and Guy Gur-Ari. The large learning rate phase of deep learning: the catapult mechanism. *arXiv preprint arXiv:2003.02218*, 2020. **2, 3**
- [45] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. *Advances in neural information processing systems*, 31, 2018. **8**
- [46] Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. On the variance of the adaptive learning rate and beyond. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. **2**
- [47] Yong Liu, Siqi Mai, Xiangning Chen, Cho-Jui Hsieh, and Yang You. Towards efficient and scalable sharpness-aware minimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12360–12370, 2022. **1, 2**
- [48] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. **6**
- [49] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. **2**
- [50] Liangchen Luo, Yuanhao Xiong, Yan Liu, and Xu Sun. Adaptive gradient methods with dynamic bound of learning rate. *arXiv preprint arXiv:1902.09843*, 2019. **2**
- [51] Stephan Mandt, Matthew D Hoffman, and David M Blei. Stochastic gradient descent as approximate bayesian inference. *Journal of Machine Learning Research*, 18:1–35, 2017. **3**
- [52] Peng Mi, Li Shen, Tianhe Ren, Yiyi Zhou, Xiaoshuai Sun, Rongrong Ji, and Dacheng Tao. Make sharpness-aware minimization stronger: A sparsified perturbation approach. *arXiv preprint arXiv:2210.05177*, 2022. **1, 2**
- [53] Yu E Nesterov. A method for solving the convex programming problem with convergence rate. In *Dokl. Akad. Nauk SSSR*, volume 269, pages 543–547, 1983. **2**
- [54] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3498–3505. IEEE, 2012. **6, 19**
- [55] Henning Petzka, Michael Kamp, Linara Adilova, Cristian Sminchisescu, and Mario Boley. Relative flatness and generalization. *Advances in Neural Information Processing Systems*, 34:18420–18432, 2021. **2**
- [56] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2018. **5**
- [57] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015. **1**
- [58] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy,

- Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015. 6
- [59] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 1
- [60] Sidak Pal Singh and Dan Alistarh. Woodfisher: Efficient second-order approximation for neural network compression. *Advances in Neural Information Processing Systems*, 33:18098–18109, 2020. 19
- [61] Samuel L. Smith and Quoc V. Le. A bayesian perspective on generalization and stochastic gradient descent. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. 2
- [62] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019. 6
- [63] Yeming Wen, Kevin Luk, Maxime Gazeau, Guodong Zhang, Harris Chan, and Jimmy Ba. An empirical study of large-batch stochastic gradient descent with structured covariance noise. *arXiv preprint arXiv:1902.08234*, 2019. 2, 3
- [64] Ashia C Wilson, Rebecca Roelofs, Mitchell Stern, Nati Srebro, and Benjamin Recht. The marginal value of adaptive gradient methods in machine learning. *Advances in neural information processing systems*, 30, 2017. 2
- [65] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017. 5, 16
- [66] Zeke Xie, Issei Sato, and Masashi Sugiyama. A diffusion theory for deep learning dynamics: Stochastic gradient descent exponentially favors flat minima. In *International Conference on Learning Representations*, 2021. 3
- [67] Zeke Xie, Qian-Yuan Tang, Yunfeng Cai, Mingming Sun, and Ping Li. On the power-law spectrum in deep learning: A bridge to protein science. *arXiv preprint arXiv:2201.13011*, 2022. 2
- [68] Zeke Xie, Xinrui Wang, Huishuai Zhang, Issei Sato, and Masashi Sugiyama. Adaptive inertia: Disentangling the effects of adaptive learning rate and momentum. In *International Conference on Machine Learning*, pages 24430–24459. PMLR, 2022. 2, 3
- [69] Haoyi Xiong, Ruosi Wan, Jian Zhao, Zeyu Chen, Xingjian Li, Zhanxing Zhu, and Jun Huan. Grod: Deep learning with gradients orthogonal decomposition for knowledge transfer, distillation, and adversarial training. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 16(6):1–25, 2022. 20
- [70] Yi Xu, Rong Jin, and Tianbao Yang. First-order stochastic algorithms for escaping from saddle points in almost linear time. *Advances in neural information processing systems*, 31, 2018. 4
- [71] Zhewei Yao, Amir Gholami, Kurt Keutzer, and Michael W Mahoney. Pyhessian: Neural networks through the lens of the hessian. In *2020 IEEE international conference on big data (Big data)*, pages 581–590. IEEE, 2020. 3, 7
- [72] Zhewei Yao, Amir Gholami, Qi Lei, Kurt Keutzer, and Michael W Mahoney. Hessian-based analysis of large batch training and robustness to adversaries. *Advances in Neural Information Processing Systems*, 31, 2018. 3, 7
- [73] Tom Young, Devamanyu Hazarika, Soujanya Poria, and Erik Cambria. Recent trends in deep learning based natural language processing. *IEEE Computational Intelligence Magazine*, 13(3):55–75, 2018. 1
- [74] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6023–6032, 2019. 1
- [75] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016. 5
- [76] Manzil Zaheer, Sashank Reddi, Devendra Sachan, Satyen Kale, and Sanjiv Kumar. Adaptive methods for nonconvex optimization. *Advances in neural information processing systems*, 31, 2018. 2
- [77] Guodong Zhang, Lala Li, Zachary Nado, James Martens, Sushant Sachdeva, George Dahl, Chris Shallue, and Roger B Grosse. Which algorithmic choices matter at which batch sizes? insights from a noisy quadratic model. *Advances in neural information processing systems*, 32, 2019. 3
- [78] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018. 1
- [79] Xingxuan Zhang, Feng Cheng, and Shilin Wang. Spatio-temporal fusion based convolutional sequence learning for lip reading. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 713–722, 2019. 1
- [80] Xingxuan Zhang, Peng Cui, Renzhe Xu, Linjun Zhou, Yue He, and Zheyang Shen. Deep stable learning for out-of-distribution generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5372–5382, 2021. 1
- [81] Xingxuan Zhang, Linjun Zhou, Renzhe Xu, Peng Cui, Zheyang Shen, and Haoxin Liu. Towards unsupervised domain generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4910–4920, 2022. 1
- [82] Yang Zhao, Hao Zhang, and Xiuyuan Hu. Penalizing gradient norm for efficiently improving generalization in deep learning. In *International Conference on Machine Learning*, pages 26982–26992. PMLR, 2022. 2, 16, 19
- [83] Qihuang Zhong, Liang Ding, Li Shen, Peng Mi, Juhua Liu, Bo Du, and Dacheng Tao. Improving sharpness-aware minimization with fisher mask for better generalization on language models. *arXiv preprint arXiv:2210.05497*, 2022. 1, 2
- [84] Daquan Zhou, Bingyi Kang, Xiaojie Jin, Linjie Yang, Xiaochen Lian, Zihang Jiang, Qibin Hou, and Jiashi Feng. Deepvit: Towards deeper vision transformer. *arXiv preprint arXiv:2103.11886*, 2021. 1

- [85] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1):43–76, 2020. [6](#)
- [86] Juntang Zhuang, Boqing Gong, Liangzhe Yuan, Yin Cui, Hartwig Adam, Nicha C Dvornek, James s Duncan, Ting Liu, et al. Surrogate gap minimization improves sharpness-aware training. In *International Conference on Learning Representations*, 2022. [2](#), [4](#), [5](#), [19](#), [20](#)

A. Omitted details in Section 4

A.1. Derivation of Equation (7)

We follow the steps in [20] to approximate

$$\nabla R^{(1)}(\boldsymbol{\theta}) = \rho \cdot \nabla_{\boldsymbol{\theta}} \max_{\boldsymbol{\epsilon} \in B(0, \rho)} \left\| \nabla \hat{L}(\boldsymbol{\theta} + \boldsymbol{\epsilon}) \right\|. \quad (12)$$

We first conduct the first-order Taylor expansion of $\|\nabla \hat{L}(\boldsymbol{\theta} + \boldsymbol{\epsilon})\|$ and get that

$$\begin{aligned} \boldsymbol{\epsilon}^*(\boldsymbol{\theta}) &= \arg \max_{\boldsymbol{\epsilon} \in B(0, \rho)} \left\| \nabla \hat{L}(\boldsymbol{\theta} + \boldsymbol{\epsilon}) \right\| \approx \arg \max_{\boldsymbol{\epsilon} \in B(0, \rho)} \left\| \nabla \hat{L}(\boldsymbol{\theta}) \right\| + \left(\nabla \left\| \nabla \hat{L}(\boldsymbol{\theta}) \right\| \right)^\top \boldsymbol{\epsilon} \\ &= \arg \max_{\boldsymbol{\epsilon} \in B(0, \rho)} \left(\nabla \left\| \nabla \hat{L}(\boldsymbol{\theta}) \right\| \right)^\top \boldsymbol{\epsilon} = \frac{\rho \cdot \mathbf{f}}{\|\mathbf{f}\|}, \end{aligned} \quad (13)$$

where $\mathbf{f} = \nabla \|\nabla \hat{L}(\boldsymbol{\theta})\|$. As a result, by letting $\boldsymbol{\theta}^{\text{adv}} = \boldsymbol{\theta} + \boldsymbol{\epsilon}^*(\boldsymbol{\theta})$,

$$\nabla R^{(1)}(\boldsymbol{\theta}) \approx \rho \cdot \nabla_{\boldsymbol{\theta}} \left\| \nabla \hat{L}(\boldsymbol{\theta} + \boldsymbol{\epsilon}^*(\boldsymbol{\theta})) \right\| = \rho \cdot \nabla \left\| \nabla \hat{L}(\boldsymbol{\theta}^{\text{adv}}) \right\| + \rho \cdot \frac{d\boldsymbol{\epsilon}^*(\boldsymbol{\theta})}{d\boldsymbol{\theta}} \cdot \nabla \left\| \nabla \hat{L}(\boldsymbol{\theta}^{\text{adv}}) \right\|. \quad (14)$$

In addition, similar to [20], we further drop the second-order term to accelerate the computation. Finally, the derivative $\nabla R^{(1)}(\boldsymbol{\theta})$ is given by

$$\nabla R^{(1)}(\boldsymbol{\theta}) \approx \rho \cdot \nabla \left\| \nabla \hat{L}(\boldsymbol{\theta}^{\text{adv}}) \right\|, \quad \boldsymbol{\theta}^{\text{adv}} = \boldsymbol{\theta} + \rho \cdot \frac{\mathbf{f}}{\|\mathbf{f}\|}, \quad \mathbf{f} = \nabla \left\| \nabla \hat{L}(\boldsymbol{\theta}) \right\|. \quad (15)$$

B. Proofs

B.1. Proof of Lemma 4.1

Proof. By assumption, we have that for all $\boldsymbol{\theta} \in B(\boldsymbol{\theta}^*, \rho)$,

$$\hat{L}(\boldsymbol{\theta}) = \hat{L}(\boldsymbol{\theta}^*) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^\top \left(\nabla^2 \hat{L}(\boldsymbol{\theta}^*) \right) (\boldsymbol{\theta} - \boldsymbol{\theta}^*). \quad (16)$$

In addition,

$$\nabla \hat{L}(\boldsymbol{\theta}) = \left(\nabla^2 \hat{L}(\boldsymbol{\theta}^*) \right) (\boldsymbol{\theta} - \boldsymbol{\theta}^*). \quad (17)$$

As a result,

$$\max_{\boldsymbol{\theta} \in B(\boldsymbol{\theta}^*, \rho)} \left\| \nabla \hat{L}(\boldsymbol{\theta}) \right\| = \max_{\boldsymbol{\theta} \in B(\boldsymbol{\theta}^*, \rho)} \left\| \left(\nabla^2 \hat{L}(\boldsymbol{\theta}^*) \right) (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \right\| = \rho \left\| \nabla^2 \hat{L}(\boldsymbol{\theta}^*) \right\| = \rho \lambda_{\max} \left(\nabla^2 \hat{L}(\boldsymbol{\theta}^*) \right). \quad (18)$$

Now the claim follows. \square

B.2. Proof of Proposition 4.2

Proof. Suppose $\boldsymbol{\epsilon}^* = \arg \max_{\boldsymbol{\epsilon} \in B(0, \rho)} \hat{L}(\boldsymbol{\theta} + \boldsymbol{\epsilon})$. Then $R^{(0)}(\boldsymbol{\theta}) = \hat{L}(\boldsymbol{\theta} + \boldsymbol{\epsilon}^*) - \hat{L}(\boldsymbol{\theta})$. According to the mean value theorem, there exists a constant $0 \leq c \leq 1$ such that

$$\hat{L}(\boldsymbol{\theta} + \boldsymbol{\epsilon}^*) - \hat{L}(\boldsymbol{\theta}) = \left(\nabla \hat{L}(\boldsymbol{\theta} + c \cdot \boldsymbol{\epsilon}^*) \right)^\top \boldsymbol{\epsilon}^*. \quad (19)$$

As a result, by the Cauchy–Schwarz inequality,

$$\begin{aligned} R^{(0)}(\boldsymbol{\theta}) &= \hat{L}(\boldsymbol{\theta} + \boldsymbol{\epsilon}^*) - \hat{L}(\boldsymbol{\theta}) = \left(\nabla \hat{L}(\boldsymbol{\theta} + c \cdot \boldsymbol{\epsilon}^*) \right)^\top \boldsymbol{\epsilon}^* \leq \left\| \nabla \hat{L}(\boldsymbol{\theta} + c \cdot \boldsymbol{\epsilon}^*) \right\| \|\boldsymbol{\epsilon}^*\| \\ &\leq \max_{\boldsymbol{\epsilon} \in B(0, \rho)} \left\| \nabla \hat{L}(\boldsymbol{\theta} + \boldsymbol{\epsilon}) \right\| \cdot \rho = R^{(1)}(\boldsymbol{\theta}). \end{aligned} \quad (20)$$

\square

B.3. Proof of Proposition 4.3

Proof. Define $h(\boldsymbol{\theta}) = \max_{\boldsymbol{\theta}' \in B(\boldsymbol{\theta}, \rho)} \left\| \nabla \hat{L}(\boldsymbol{\theta}') \right\|$. Fix $\sigma = \rho / (\sqrt{d} + \sqrt{\log n})$, following the proof of Theorem 1 in [20], we can obtain that with probability at least $1 - \delta$,

$$\mathbb{E}_{\epsilon_i \sim N(0, \sigma^2)} [L(\boldsymbol{\theta} + \epsilon)] \leq \mathbb{E}_{\epsilon_i \sim N(0, \sigma^2)} [\hat{L}(\boldsymbol{\theta} + \epsilon)] + \sqrt{\frac{\frac{1}{4}d \log \left(1 + \frac{\|\boldsymbol{\theta}\|_2^2}{d\sigma^2}\right) + \frac{1}{4} + \log \frac{n}{\delta} + 2 \log(6n + 3d)}{n-1}}. \quad (21)$$

Since $\epsilon_i \sim N(0, \sigma^2)$, $\|\epsilon\|^2/\sigma^2$ has a chi-square distribution. As a result, according to [43, Lemma 1], we have that for any $t > 0$,

$$\mathbb{P} \left(\|\epsilon\|^2/\sigma^2 - d \geq 2\sqrt{dt} + 2t \right) \leq \exp(-t). \quad (22)$$

By letting $t = \frac{1}{2} \log n$, we can get that with probability at least $1 - 1/\sqrt{n}$,

$$\|\epsilon\|^2 \leq \sigma^2 \left(d + \sqrt{2d \log n} + \log n \right) \leq \sigma^2 \left(\sqrt{d} + \sqrt{\log n} \right)^2 = \rho^2. \quad (23)$$

As a result,

$$\begin{aligned} & \mathbb{E}_{\epsilon_i \sim N(0, \sigma^2)} [\hat{L}(\boldsymbol{\theta} + \epsilon)] \\ & \leq \mathbb{E}_{\epsilon_i \sim N(0, \sigma^2)} [\hat{L}(\boldsymbol{\theta} + \epsilon) \mid \|\epsilon\| \leq \rho] \mathbb{P}(\|\epsilon\| \leq \rho) + \mathbb{E}_{\epsilon_i \sim N(0, \sigma^2)} [\hat{L}(\boldsymbol{\theta} + \epsilon) \mid \|\epsilon\| > \rho] \mathbb{P}(\|\epsilon\| > \rho) \\ & \leq \mathbb{E}_{\epsilon_i \sim N(0, \sigma^2)} [\hat{L}(\boldsymbol{\theta} + \epsilon) \mid \|\epsilon\| \leq \rho] + \frac{M}{\sqrt{n}}. \end{aligned} \quad (24)$$

According to the mean value theorem and Cauchy–Schwarz inequality, for any ϵ such that $\|\epsilon\| < \rho$, there exists a constant $0 \leq c \leq 1$, such that

$$\hat{L}(\boldsymbol{\theta} + \epsilon) = \hat{L}(\boldsymbol{\theta}) + \left(\nabla \hat{L}(\boldsymbol{\theta} + c\epsilon) \right)^\top \epsilon \leq \hat{L}(\boldsymbol{\theta}) + \left\| \nabla \hat{L}(\boldsymbol{\theta} + c\epsilon) \right\| \cdot \|\epsilon\| \leq \hat{L}(\boldsymbol{\theta}) + h(\boldsymbol{\theta})\rho = \hat{L}(\boldsymbol{\theta}) + R^{(1)}(\boldsymbol{\theta}). \quad (25)$$

Now the claim follows from Equations (21), (24), and (25). \square

B.4. Proof of Theorem 4.4

Proof. Observe that

$$\begin{aligned} \left\| \nabla L^{\text{overall}}(\boldsymbol{\theta}_t) \right\|^2 &= \left\| \nabla L^{\text{oracle}}(\boldsymbol{\theta}_t) + \alpha \rho_t \cdot \nabla \left\| \nabla \hat{L}(\boldsymbol{\theta}_t^{\text{adv}}) \right\| \right\|^2 \\ &\leq 2 \left(\left\| \nabla L^{\text{oracle}}(\boldsymbol{\theta}_t) \right\|^2 + \left\| \alpha \rho_t \cdot \nabla \left\| \nabla \hat{L}(\boldsymbol{\theta}_t^{\text{adv}}) \right\| \right\|^2 \right). \end{aligned} \quad (26)$$

The claim follows from Propositions B.1 and B.2. \square

Proposition B.1. *Assume the conditions in Theorem 4.4 hold (with parameters $\gamma_1, \gamma_2, G^{\text{loss}}, G^{\text{norm}}, \tilde{G}^{\text{loss}}, M, \eta_0, \rho_0, \alpha$). Then with learning rate $\eta_t = \eta_0/\sqrt{t}$ and perturbation radius $\rho_t = \rho_0/\sqrt{t}$, Algorithm 1 could obtain*

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \nabla L^{\text{oracle}}(\boldsymbol{\theta}_t) \right\|^2 \right] \leq \frac{C'_1 + C'_2 \log T}{\sqrt{T}} \quad (27)$$

for some constants C'_1 and C'_2 that only depend on $\gamma_1, \gamma_2, G^{\text{loss}}, G^{\text{norm}}, \tilde{G}^{\text{loss}}, M, \eta_0, \rho_0, \alpha$.

Proof. By definition, we have $\mathbf{h}_t^{\text{loss}} = \tilde{g}_t^{\text{loss}}(\boldsymbol{\theta}_t)$ and $\mathbf{h}_t^{\text{norm}} = g_t^{\text{norm}}(\boldsymbol{\theta}_t^{\text{adv}})$. By assumption,

$$\begin{aligned} L^{\text{oracle}}(\boldsymbol{\theta}_{t+1}) &\leq L^{\text{oracle}}(\boldsymbol{\theta}_t) + (\nabla L^{\text{oracle}}(\boldsymbol{\theta}_t))^\top (\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t) + \frac{\gamma_1}{2} \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|^2 \\ &= L^{\text{oracle}}(\boldsymbol{\theta}_t) - \eta_t (\nabla L^{\text{oracle}}(\boldsymbol{\theta}_t))^\top \left(\mathbf{h}_t^{\text{loss}} + \alpha \rho_t \mathbf{h}_t^{\text{norm}} \right) + \frac{\gamma_1 \eta_t^2}{2} \left\| \mathbf{h}_t^{\text{loss}} + \alpha \rho_t \mathbf{h}_t^{\text{norm}} \right\|^2. \end{aligned} \quad (28)$$

Take the expectation conditioned on the observations till timestamp t . By the assumption $\mathbb{E}[\mathbf{h}_t^{\text{loss}}] = \mathbb{E}[\tilde{g}_t^{\text{loss}}(\boldsymbol{\theta}_t)] = \nabla L^{\text{oracle}}(\boldsymbol{\theta}_t)$ and $\mathbb{E}[\mathbf{h}_t^{\text{norm}}] = \mathbb{E}[g_t^{\text{norm}}(\boldsymbol{\theta}_t^{\text{adv}})]$, we can obtain that

$$\begin{aligned} & \mathbb{E} [L^{\text{oracle}}(\boldsymbol{\theta}_{t+1})] - L^{\text{oracle}}(\boldsymbol{\theta}_t) \\ & \leq -\eta_t \|\nabla L^{\text{oracle}}(\boldsymbol{\theta}_t)\|^2 - \eta_t \rho_t \alpha (\nabla L^{\text{oracle}}(\boldsymbol{\theta}_t))^\top \mathbb{E} [g_t^{\text{norm}}(\boldsymbol{\theta}_t^{\text{adv}})] + \frac{\gamma_1 \eta_t^2}{2} \|\mathbf{h}_t^{\text{loss}} + \alpha \rho_t \mathbf{h}_t^{\text{norm}}\|^2 \end{aligned} \quad (29)$$

We have

$$-\eta_t \rho_t \alpha (\nabla L^{\text{oracle}}(\boldsymbol{\theta}_t))^\top \mathbb{E} [g_t^{\text{norm}}(\boldsymbol{\theta}_t^{\text{adv}})] \leq \eta_t \rho_t \alpha \|\nabla L^{\text{oracle}}(\boldsymbol{\theta}_t)\| \|\mathbb{E} [g_t^{\text{norm}}(\boldsymbol{\theta}_t^{\text{adv}})]\| \leq \eta_t \rho_t \alpha \tilde{G}^{\text{loss}} G^{\text{norm}}. \quad (30)$$

In addition,

$$\mathbb{E} \left[\|\mathbf{h}_t^{\text{loss}} + \alpha \mathbf{h}_t^{\text{norm}}\|^2 \right] \leq 2\mathbb{E} \left[\|\mathbf{h}_t^{\text{loss}}\|^2 \right] + 2\alpha^2 \mathbb{E} \left[\|\mathbf{h}_t^{\text{norm}}\|^2 \right] \leq 2 \left(\tilde{G}^{\text{loss}} \right)^2 + 2\alpha^2 \left(G^{\text{norm}} \right)^2. \quad (31)$$

Combining Equations (29), (30), and (31), we can get that

$$\eta_t \|\nabla L^{\text{oracle}}(\boldsymbol{\theta}_t)\|^2 \leq -\mathbb{E} [L^{\text{oracle}}(\boldsymbol{\theta}_{t+1})] + L^{\text{oracle}}(\boldsymbol{\theta}_t) + \eta_t \rho_t Z_1 + \eta_t^2 Z_2 \quad (32)$$

for some constants Z_1 and Z_2 that only depend on $\gamma_1, \gamma_2, G^{\text{loss}}, G^{\text{norm}}, \tilde{G}^{\text{loss}}, \alpha$. Now perform telescope sum and take the expectations at each step, we can obtain that

$$\sum_{t=1}^T \eta_t \|\nabla L^{\text{oracle}}(\boldsymbol{\theta}_t)\|^2 \leq -\mathbb{E} [L^{\text{oracle}}(\boldsymbol{\theta}_{T+1})] + L^{\text{oracle}}(\boldsymbol{\theta}_1) + Z_1 \sum_{t=1}^T \eta_t \rho_t + Z_2 \sum_{t=1}^T \eta_t^2. \quad (33)$$

By letting $\eta_t = \eta_0/\sqrt{t}$ and $\rho_t = \rho_0/\sqrt{t}$, we can get that

$$\begin{aligned} \frac{\eta_0}{\sqrt{T}} \sum_{t=1}^T \|\nabla L^{\text{oracle}}(\boldsymbol{\theta}_t)\|^2 & \leq \sum_{t=1}^T \eta_t \|\nabla L^{\text{oracle}}(\boldsymbol{\theta}_t)\|^2 \\ & \leq -\mathbb{E} [L^{\text{oracle}}(\boldsymbol{\theta}_{T+1})] + L^{\text{oracle}}(\boldsymbol{\theta}_1) + Z_1 \sum_{t=1}^T \eta_t \rho_t + Z_2 \sum_{t=1}^T \eta_t^2 \\ & \leq 2M + Z_1 \eta_0 \rho_0 \sum_{t=1}^T \frac{1}{t} + Z_2 \eta_0^2 \sum_{t=1}^T \frac{1}{t} \\ & \leq Z_4 + Z_5 \log T \end{aligned} \quad (34)$$

for some constants Z_4 and Z_5 that only depend on $\gamma_1, \gamma_2, G^{\text{loss}}, G^{\text{norm}}, \tilde{G}^{\text{loss}}, M, \eta_0, \rho_0, \alpha$. Divide the two sides of the equation by $\eta_0 \sqrt{T}$ and the claim follows. \square

Proposition B.2. Assume the conditions in Theorem 4.4 hold (with parameters $\gamma_1, \gamma_2, G^{\text{loss}}, G^{\text{norm}}, \tilde{G}^{\text{loss}}, M, \eta_0, \rho_0, \alpha$). Then with perturbation radius $\rho_t = \rho_0/\sqrt{t}$, Algorithm 1 could obtain

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\|\alpha \rho_t \cdot \nabla \|\nabla \hat{L}(\boldsymbol{\theta}^{\text{adv}})\|\|^2 \right] \leq \frac{C_1'' + C_2'' \log T}{\sqrt{T}} \quad (35)$$

for some constants C_1'' and C_2'' that only depend on γ_2, ρ_0, α .

Proof. For any $t \in \{1, 2, \dots, T\}$,

$$\begin{aligned} & \mathbb{E} \left[\|\alpha \rho_t \cdot \nabla \|\nabla \hat{L}(\boldsymbol{\theta}^{\text{adv}})\|\|^2 \right] = \alpha^2 \rho_t^2 \mathbb{E} \left[\|\nabla \|\nabla \hat{L}(\boldsymbol{\theta}^{\text{adv}})\|\|^2 \right] \\ & = \alpha^2 \rho_t^2 \mathbb{E} \left[\left\| \nabla^2 \hat{L}(\boldsymbol{\theta}_t^{\text{adv}}) \cdot \frac{\nabla \hat{L}(\boldsymbol{\theta}_t^{\text{adv}})}{\|\nabla \hat{L}(\boldsymbol{\theta}_t^{\text{adv}})\|} \right\|^2 \right] \leq \alpha^2 \rho_t^2 \mathbb{E} \left[\|\nabla^2 \hat{L}(\boldsymbol{\theta}_t^{\text{adv}})\| \left\| \frac{\nabla \hat{L}(\boldsymbol{\theta}_t^{\text{adv}})}{\|\nabla \hat{L}(\boldsymbol{\theta}_t^{\text{adv}})\|} \right\|^2 \right] \\ & \leq \alpha^2 \rho_t^2 \mathbb{E}[\gamma_2] = \alpha^2 \rho_t^2 \gamma_2. \end{aligned} \quad (36)$$

Table 4. Results of GAM with state-of-the-art models on CIFAR-10 and CIFAR-100. The best results are highlighted in bold font.

Model	Aug	CIFAR-10				CIFAR-100			
		SGD	SGD + GAM	SAM	SAM + GAM	SGD	SGD + GAM	SAM	SAM + GAM
DenseNet121	Basic	91.16 \pm 0.13	92.35 \pm 0.14	92.19 \pm 0.20	92.72 \pm 0.30	69.25 \pm 0.40	70.48 \pm 0.27	70.44 \pm 0.19	71.16 \pm 0.25
DenseNet121	Cutout	91.85 \pm 0.17	92.93 \pm 0.26	92.35 \pm 0.16	93.30 \pm 0.17	70.17 \pm 0.31	71.47 \pm 0.23	70.89 \pm 0.15	71.80 \pm 0.07
DenseNet121	RA	91.59 \pm 0.16	92.37 \pm 0.20	92.32 \pm 0.29	92.97 \pm 0.27	69.65 \pm 0.36	70.10 \pm 0.26	70.49 \pm 0.16	71.43 \pm 0.17
DenseNet121	AA	92.65 \pm 0.10	94.17 \pm 0.25	92.96 \pm 0.19	94.05 \pm 0.22	70.53 \pm 0.17	72.25 \pm 0.19	71.34 \pm 0.20	72.90 \pm 0.17
ResNeXt29-32x4d	Basic	95.75 \pm 0.31	96.46 \pm 0.25	96.32 \pm 0.36	96.90 \pm 0.24	79.45 \pm 0.29	81.67 \pm 0.26	81.35 \pm 0.12	82.93 \pm 0.25
ResNeXt29-32x4d	Cutout	96.20 \pm 0.37	97.82 \pm 0.24	96.44 \pm 0.21	97.85 \pm 0.27	80.56 \pm 0.20	82.62 \pm 0.33	82.49 \pm 0.25	83.58 \pm 0.09
ResNeXt29-32x4d	RA	95.86 \pm 0.28	97.17 \pm 0.26	96.75 \pm 0.35	97.79 \pm 0.19	79.88 \pm 0.12	81.75 \pm 0.23	82.26 \pm 0.15	83.02 \pm 0.22
ResNeXt29-32x4d	AA	96.58 \pm 0.18	97.46 \pm 0.15	97.38 \pm 0.25	97.58 \pm 0.16	80.47 \pm 0.13	82.02 \pm 0.19	81.52 \pm 0.26	83.35 \pm 0.09
ViT-S/16	Basic	95.27 \pm 0.23	97.21 \pm 0.14	96.85 \pm 0.25	97.58 \pm 0.20	79.52 \pm 0.36	83.35 \pm 0.28	82.77 \pm 0.29	84.30 \pm 0.25
ViT-S/16	Cutout	95.36 \pm 0.22	97.53 \pm 0.17	97.10 \pm 0.25	97.85 \pm 0.10	79.36 \pm 0.21	83.59 \pm 0.28	82.86 \pm 0.14	84.53 \pm 0.16
ViT-S/16	RA	95.59 \pm 0.19	97.44 \pm 0.31	97.18 \pm 0.12	97.59 \pm 0.11	79.96 \pm 0.22	83.80 \pm 0.20	83.36 \pm 0.13	84.66 \pm 0.23
ViT-S/16	AA	96.40 \pm 0.30	97.82 \pm 0.21	97.52 \pm 0.25	97.97 \pm 0.13	80.35 \pm 0.06	84.02 \pm 0.18	83.54 \pm 0.19	85.20 \pm 0.26

By letting $\rho_t = \rho_0/\sqrt{t}$,

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\left\| \alpha \rho_t \cdot \nabla \left\| \nabla \hat{L}(\theta^{\text{adv}}) \right\|^2 \right\| \right] \leq \frac{1}{T} \alpha^2 \gamma_2 \rho_0^2 \sum_{t=1}^T \frac{1}{t} \leq \frac{C_1'' + C_2'' \log T}{\sqrt{T}} \quad (37)$$

for some constants C_1'' and C_2'' that only depend on γ_2, ρ_0, α . □

C. More Experimental Results and Details

C.1. More Results on Training from Scratch

Due to space limitation, we omit the experimental results on CIFAR-10 and CIFAR-100 [41] with ResNeXt [65], DenseNet [27] and ViTs [16] in Section 5.2.1 in the main paper and report them in Section C.1.1. Then we report the results of robustness to label noise in Section C.2.

C.1.1 CIFAR-10 and CIFAR-100

We report the omitted results on CIFAR-10 and CIFAR-100 with ResNeXt, DenseNet and ViTs. As described in the main paper, all the models are trained for 200 epochs from scratch. We evaluate GAM both with basic data augmentations (i.e., horizontal flip, padding by four pixels, and random crop) and advanced data augmentation including cutout regularization [15], RandAugment [13] and AutoAugment [14]. The hyperparameters, ρ and α are searched with the same approach described in the main paper.

Results are shown in Table 4. GAM consistently improves generalization for all models. We observe the same results as in the main paper. When combined with SGD, GAM achieves considerably higher test accuracy compared with SGD. And GAM also achieves improvement when combined with SAM.

Comparison with GNP GNP can be considered as a special case of GAM where ρ is set to 0. We compare GAM with GNP [82] in on CIFAR-100 in Table 5. We follow hyperparameters searching and choice of GNP in its original paper. GAM consistently outperforms GNP by noticeable margins.

C.2. Robustness to Label Noise

It is observed that sharpness-aware minimization methods are robust to perturbations to label noise [20, 42]. Here we assess the degree of robustness that GAM provides to label noise.

Following [20, 42], we measure the effectiveness of GAM in the classical noisy-label setting for CIFAR-10. A fraction of the training data labels is randomly flipped [33] while the test data remains unmodified. We train a ResNet32 for 200 epochs following [33]. Hyperparameter settings for all the models are the same as that of previous CIFAR experiments. Following [3, 42], we report the best results during the training instead of the results at the end of the training.

Table 5. Comparison with GNP on CIFAR-100. *-BA* indicates the basic data augmentation, *-CU* indicates cutout regularization, *-RA* indicates RandAugment, and *-AA* indicates AutoAugment.

	Res18-BA	Res18-CU	Res18-RA	Res18-AA	Res101-BA	Res101-CU	Res101-RA	Res101-AA
SGD	78.32 \pm 0.32	78.73 \pm 0.13	78.62 \pm 0.32	78.88 \pm 0.15	80.47 \pm 0.13	80.53 \pm 0.30	80.60 \pm 0.28	81.83 \pm 0.37
GNP + SGD	78.80 \pm 0.40	79.29 \pm 0.15	79.21 \pm 0.27	80.29 \pm 0.05	81.17 \pm 0.29	81.10 \pm 0.14	81.31 \pm 0.88	82.53 \pm 0.25
GAM + SGD	79.53 \pm 0.30	79.89 \pm 0.31	79.82 \pm 0.24	80.56 \pm 0.21	82.21 \pm 0.40	82.36 \pm 0.24	82.40 \pm 0.31	83.19 \pm 0.15
	WRN10-BA	WRN10-CU	WRN10-RA	WRN10-AA	Pyr110-BA	Pyr110-CU	Pyr110-RA	Pyr110-AA
SGD	81.40 \pm 0.13	81.53 \pm 0.40	81.65 \pm 0.18	81.99 \pm 0.11	82.74 \pm 0.12	83.31 \pm 0.21	84.04 \pm 0.19	84.48 \pm 0.03
GNP + SGD	82.30 \pm 0.05	82.54 \pm 0.19	82.99 \pm 0.39	83.58 \pm 0.32	83.99 \pm 0.27	84.46 \pm 0.16	84.47 \pm 0.08	84.83 \pm 0.21
GAM + SGD	83.45 \pm 0.09	83.69 \pm 0.08	83.84 \pm 0.09	84.02 \pm 0.18	84.91 \pm 0.09	85.20 \pm 0.19	86.47 \pm 0.14	85.92 \pm 0.03

Table 6. Test accuracy of ResNet32 on CIFAR-10 with label noise.

Noise Rate (%)	Base Opt	Base + GAM	SAM	SAM + GAM
0%	95.71	96.55	96.25	96.88
20%	92.05	94.20	93.85	94.74
40%	88.89	92.17	90.85	92.57
60%	83.17	88.49	87.37	89.65
80%	63.16	73.82	70.65	76.33

Table 7. Accuracy and training speed of training with different ratios of iterations using GAM. Superscripts indicate the ratio of iterations in each epoch is trained with GAM (e.g., GAM^{0.05} indicates that 5% of iterations are trained with GAM, while the remaining iterations are trained with the basic optimizer). Numbers in parentheses indicate the ratio of the training speed compared with the vanilla base optimizer SGD/SAM. We mark runs whose training speed is lower than 50% of the basic optimizer in red and others in green. Please note that the speed of SAM is about 50% w.r.t SGD’s speed. Thus when combined with SGD, green markers indicate that the speed of GAM under the corresponding ratio is faster than SAM.

CIFAR-10		SGD	SGD + GAM ^{0.05}	SGD + GAM ^{0.1}	SGD + GAM ^{0.5}	SGD + GAM ¹
	Accuracy	95.32	96.08	96.15	96.17	96.17
	Images/s	2,593 (100%)	2,258 (87%)	1,996 (77%)	1,023 (39%)	658 (25%)
		SAM	SAM + GAM ^{0.05}	SAM + GAM ^{0.1}	SAM + GAM ^{0.5}	SAM + GAM ¹
	Accuracy	96.10	96.54	96.62	96.65	96.58
	Images/s	1,314 (100%)	1,247 (95%)	1,184 (90%)	858 (65%)	629 (48%)
CIFAR-100		SGD	SGD + GAM ^{0.05}	SGD + GAM ^{0.1}	SGD + GAM ^{0.5}	SGD + GAM ¹
	Accuracy	78.32	79.25	79.42	79.50	79.53
	Images/s	2,609 (100%)	2,243 (86%)	1,955 (75%)	1,011 (39%)	655 (25%)
		SAM	SAM + GAM ^{0.05}	SAM + GAM ^{0.1}	SAM + GAM ^{0.5}	SAM + GAM ¹
	Accuracy	79.27	80.08	80.44	80.40	80.45
	Images/s	1,318 (100%)	1,251 (95%)	1,172 (89%)	848 (64%)	628 (48%)

We report test accuracies for SGD, SAM, SGD+GAM, and SAM+GAM obtained from 3 independent runs for each label noise level in Table 6. As seen in Table 6, GAM shows a high degree of robustness to label noise. GAM consistently improves the robustness to label noise for both SGD and SAM.

C.3. Detailed Results and Discussions about Computation Overhead

Here we report the detailed results of the trade-off between computation overhead and test accuracy of GAM. As discussed in Section 4.3 and 5.5 in the main paper, the GAM term can be easily calculated via the Hessian vector product, which is an efficient approach to calculating the dot product between the Hessian and a vector without the need to calculate the entire Hessian. And we also notice that only several iterations of learning with GAM (with higher α compared with applying GAM

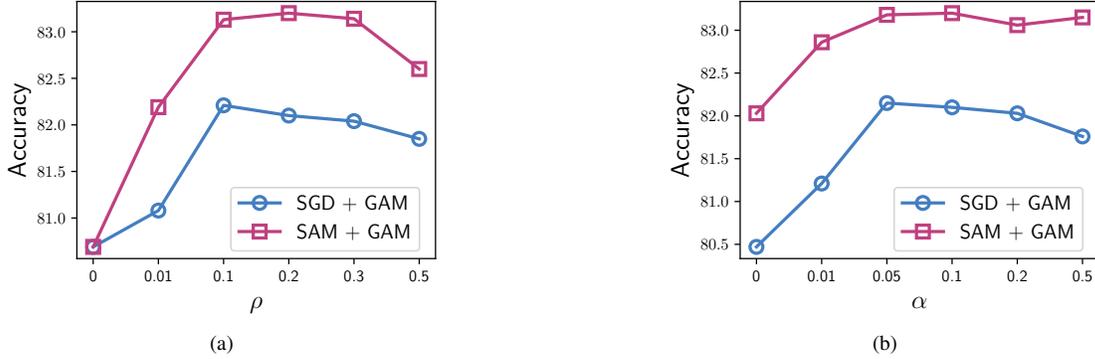


Figure 6. The influence of hyperparameters ρ and α on the performance of ResNet101 on CIFAR-100.

Table 8. Hyperparameters for Algorithm 1 on CIFAR-10 and CIFAR-100 datasets.

Model	Learning Rate	Weight Decay	Base Optimizer	Epochs	LR Schedule
ResNet18	0.1	0.0005	SGD	200	Cosine
ResNet101	0.1	0.0005	SGD	200	Cosine
WRN28_2	0.1	0.0005	SGD	200	Cosine
WRN28_10	0.1	0.0005	SGD	200	Cosine
PyramidNet110	0.05	0.0005	SGD	200	Cosine
DenseNet121	0.1	0.001	SGD	200	Cosine
ResNeXt29-32x4d	0.1	0.0005	SGD	200	Cosine

to all iterations) improve model generalization considerably. As seen in Table 7, applying GAM to 1/10 iterations of training shows similar generalization performance to applying GAM to all the iterations, while the extra computational cost for GAM is less than 25% of the original cost. Thus the computation overhead of GAM can be easily controlled.

C.4. Ablation Study

GAM has two hyperparameters, ρ , and α . We analyze the influence of the choice of them in the following subsection C.4.1 and C.4.2. The results are shown in Figure 6.

C.4.1 Influence of ρ

ρ controls the step length of gradient ascent in GAM. When ρ is set to 0, GAM degenerates into a naive regularizer constraining the gradient norm at each training step. We plot the performance of ResNet101 on CIFAR-100 with varying ρ . For experiments where GAM is combined with SAM, we apply the same ρ for both GAM and SAM. As shown in Figure 6a, GAM with ρ larger than 0 outperforms GAM without gradient ascent, showing that the gradient ascent is necessary for GAM. Moreover, GAM consistently outperforms SGD with different ρ . GAM also consistently improves SAM’s performance with various ρ , indicating that the first-order flatness can improve the generalization ability of zero-order flatness.

C.4.2 Influence of α

α controls the strength of GAM penalty. When α is set to 0 GAM degenerates into the basic optimizer (SGD or SAM). We show the performance of GAM with varying α in Figure 6b. Compared with SGD, GAM shows considerable improvement with varying α . The improvement of GAM under various α is also observed when combined with SAM.

C.5. Training Details and Selection of Hyperparameters

C.5.1 Training Details of Training from Scratch Experiments

We search hyperparameters, including learning rate and weight decay for all the models unless otherwise noted. For ResNets, we conduct a grid search of learning rate in {0.01, 0.1, 1.0} and weight decay in {0.0001, 0.0005, 0.001, 0.01, 0.1}. The batch size is set to 128 for all models. For Vits, we search the learning rate in {1e-3, 3e-3, 1e-2, 3e-3}, and weight decay in

Table 9. Hyperparameters for Algorithm 1 on ImageNet.

Model	ρ	α	Learning Rate	Weight Decay	Base Optimizer	Epochs	LR Schedule
ResNet50	0.2	0.1	0.1	0.0001	SGD	90	Cosine
ResNet101	0.2	0.1	0.1	0.0001	SGD	90	Cosine
ViT-S/32	0.3	0.5	0.0003	0.3	AdamW	300	Cosine
ViT-B/32	0.3	0.5	0.0003	0.3	AdamW	300	Cosine

{0.001, 0.01, 0.1}. We adopt SGD with momentum = 0.9 for ResNets and AdamW with $\beta_1 = 0.9$, $\beta_2 = 0.999$ for ViTs. We train ResNets for 90 epochs, and train ViTs for 300 epochs following [11, 86]. We first search for the optimal learning rate and weight decay for training with basic optimizers and keep them fixed for SAM and GAM. We search ρ in {0.05, 0.1, 0.2, 0.5, 1.0, 2.0} for both SAM and GAM and search α in {0.1, 0.2, 0.5, 1.0, 2.0, 3.0, ..., 10.0} for GAM. We set ρ to 0.04 for CIFAR-10 and 0.1 for CIFAR-100. We set α to 0.3 for ResNet-18 and 0.1 for other models. We report the best selection of hyperparameters for each individual model in Table 8 and Table 9.

C.5.2 Training Details of Transfer Learning Experiments

We finetune the models on downstream datasets including Stanford Cars [40], CIFAR-10, CIFAR-100 [41], Oxford_IIIT_Pets [54] and Food101 [8] from the weights pretrained on ImageNet. For EfficientNet-b0, we adopt SGD with momentum = 0.9. For Swin-t, we adopt AdamW with $\beta_1 = 0.9$, $\beta_2 = 0.999$. We train all the models for 40k steps and the batch size is 128. The initial learning rate is 2e-3 and the cosine learning rate decay is used. Weight decay is set to 1e-5.

D. Further Acceleration of GAM

Acceleration Optimizing the gradient of $R_\rho^{(1)}(\theta)$ according to Equations (7) and (8) requires the Hessian vector product operation, which can still introduce considerable extra computation when the model is large. Inspired by [60, 82], one can approximate $\nabla\|\nabla\hat{L}(\theta)\|$ with first-order gradient as follows.

$$\forall \theta \in \Theta, \quad \nabla\|\nabla\hat{L}(\theta)\| \approx \frac{\nabla\hat{L}\left(\theta + \rho' \cdot \frac{\nabla\hat{L}(\theta)}{\|\nabla\hat{L}(\theta)\|}\right) - \nabla\hat{L}(\theta)}{\rho'}, \quad (38)$$

where ρ' is a small constant. Thus we can further accelerate GAM by applying Equation (38) to the θ^{adv} term in Equation (7) as follows,

$$\theta^{\text{adv}} \approx \theta + \rho \cdot \frac{\nabla\hat{L}\left(\theta + \rho' \cdot \frac{\nabla\hat{L}(\theta)}{\|\nabla\hat{L}(\theta)\|}\right) - \nabla\hat{L}(\theta)}{\left\|\nabla\hat{L}\left(\theta + \rho' \cdot \frac{\nabla\hat{L}(\theta)}{\|\nabla\hat{L}(\theta)\|}\right) - \nabla\hat{L}(\theta)\right\|} = \theta + \rho \cdot \frac{\mathbf{g}_1 - \mathbf{g}_0}{\|\mathbf{g}_1 - \mathbf{g}_0\|}, \quad (39)$$

where

$$\mathbf{g}_0 = \nabla\hat{L}(\theta), \quad \mathbf{g}_1 = \nabla\hat{L}(\tilde{\theta}_1), \quad \text{and} \quad \tilde{\theta}_1 = \theta + \rho' \cdot \frac{\nabla\hat{L}(\theta)}{\|\nabla\hat{L}(\theta)\|} = \theta + \rho' \cdot \frac{\mathbf{g}_0}{\|\mathbf{g}_0\|}. \quad (40)$$

We let $\tilde{\theta}_2 \triangleq \theta^{\text{adv}}$. Now applying Equation (38) to the calculation of $\nabla\|\nabla\hat{L}(\tilde{\theta}_2)\|$, we can get that

$$\begin{aligned} \nabla R_\rho^{(1)}(\theta) &\approx \rho \cdot \nabla\left\|\nabla\hat{L}(\tilde{\theta}_2)\right\| \approx \rho \cdot \frac{\nabla\hat{L}\left(\tilde{\theta}_2 + \rho' \cdot \frac{\nabla\hat{L}(\tilde{\theta}_2)}{\|\nabla\hat{L}(\tilde{\theta}_2)\|}\right) - \nabla\hat{L}(\tilde{\theta}_2)}{\rho'} \\ &= \frac{\rho}{\rho'} (\mathbf{g}_3 - \mathbf{g}_2), \end{aligned} \quad (41)$$

where

$$\mathbf{g}_2 = \nabla\hat{L}(\tilde{\theta}_2), \quad \mathbf{g}_3 = \nabla\hat{L}(\tilde{\theta}_3), \quad \text{and} \quad \tilde{\theta}_3 = \tilde{\theta}_2 + \rho' \cdot \frac{\nabla\hat{L}(\tilde{\theta}_2)}{\|\nabla\hat{L}(\tilde{\theta}_2)\|} = \tilde{\theta}_2 + \rho' \cdot \frac{\mathbf{g}_2}{\|\mathbf{g}_2\|}. \quad (42)$$

Algorithm 2 Accelerated GAM

```
1: Input: Batch size  $b$ , Learning rate  $\eta_t$ , Perturbation radius  $\rho_t, \rho'_t$ , Trade-off coefficient  $\alpha, \beta, \gamma$ , Small constant  $\xi$ 
2:  $t \leftarrow 0, \theta_0 \leftarrow$  initial parameters
3: while  $\theta_t$  not converged do
4:   Sample  $W_t$  from the training data with  $b$  instances
5:    $\mathbf{g}_{t,0} \leftarrow \nabla \hat{L}_{W_t}(\theta_t)$ 
6:    $\tilde{\theta}_{t,1} \leftarrow \theta_t + \rho'_t \cdot \mathbf{g}_{t,0} / (\|\mathbf{g}_{t,0}\| + \xi)$ 
7:    $\mathbf{g}_{t,1} \leftarrow \nabla \hat{L}_{W_t}(\tilde{\theta}_{t,1})$ 
8:    $\mathbf{h}_{t,0} \leftarrow \mathbf{g}_{t,1} - \mathbf{g}_{t,0}$ 
9:    $\tilde{\theta}_{t,2} \leftarrow \theta_t + \rho_t \cdot \mathbf{h}_{t,0} / (\|\mathbf{h}_{t,0}\| + \xi)$ 
10:   $\mathbf{g}_{t,2} \leftarrow \nabla \hat{L}_{W_t}(\tilde{\theta}_{t,2})$ 
11:   $\tilde{\theta}_{t,3} \leftarrow \tilde{\theta}_{t,2} + \rho'_t \cdot \mathbf{g}_{t,2} / (\|\mathbf{g}_{t,2}\| + \xi)$ 
12:   $\mathbf{g}_{t,3} \leftarrow \nabla \hat{L}_{W_t}(\tilde{\theta}_{t,3})$ 
13:   $\mathbf{h}_{t,+} \leftarrow \alpha \mathbf{g}_{t,1} + (1 - \alpha) \mathbf{g}_{t,3}$ 
14:   $\mathbf{h}_{t,-} \leftarrow \beta \mathbf{g}_{t,0} + (1 - \beta) \mathbf{g}_{t,2}$ 
15:   $\mathbf{h}_{t,-}^{\parallel}, \mathbf{h}_{t,-}^{\perp} \leftarrow \text{decompose}(\mathbf{h}_{t,-}; \mathbf{h}_{t,+})$   $\triangleright$  Decompose  $\mathbf{h}_{t,-}$  into components that are parallel or orthogonal to  $\mathbf{h}_{t,+}$ 
16:   $\theta_{t+1} \leftarrow \theta_t - \eta_t (\mathbf{h}_{t,+} - \gamma \mathbf{h}_{t,-}^{\perp})$ 
17:   $t \leftarrow t + 1$ 
18: end while
19: return  $\theta_t$ 
```

Accelerated GAM Based on the above approximations, we could obtain the accelerated version of GAM as shown in Algorithm 2. Besides them, the following modifications can further accelerate and suppress the effects of the above approximation.

1. We find that the gradient of the SAM regularization term, $\mathbf{g}_1 - \mathbf{g}_0$, is already calculated during the approximation steps. As a result, we directly optimize the target $\hat{L}(\theta) + \alpha' R_{\rho}^{(1)}(\theta) + \beta' R_{\rho'}^{(0)}(\theta)$ with hyper-parameter α', β' . The gradient of the target can be approximated as

$$\nabla \left(\hat{L}(\theta) + \alpha' R_{\rho}^{(1)}(\theta) + \beta' R_{\rho'}^{(0)}(\theta) \right) \approx \mathbf{g}_0 + \frac{\alpha' \rho}{\rho'} (\mathbf{g}_3 - \mathbf{g}_2) + \beta' (\mathbf{g}_1 - \mathbf{g}_0) = \beta' \mathbf{g}_1 + \frac{\alpha' \rho}{\rho'} \mathbf{g}_3 - (\beta' - 1) \mathbf{g}_0 - \frac{\alpha' \rho}{\rho'} \mathbf{g}_2, \quad (43)$$

which means that the gradient of our target is a linear combination of $\mathbf{g}_0, \mathbf{g}_1, \mathbf{g}_2, \mathbf{g}_3$. We further set three hyper-parameters to control the importance of these parts and preserve the signs of different parts, *i.e.*,

$$\nabla \left(\hat{L}(\theta) + \alpha' R_{\rho}^{(1)}(\theta) + \beta' R_{\rho'}^{(0)}(\theta) \right) \approx \alpha \mathbf{g}_1 + (1 - \alpha) \mathbf{g}_3 - \gamma (\beta \mathbf{g}_0 + (1 - \beta) \mathbf{g}_2) \quad (44)$$

with $1 \geq \alpha, \beta \geq 0, \gamma \geq 0$.

2. We find that the negative parts of Equation (44) may have side effects on the model's convergence and performance, and thus require fine-tuning of hyperparameters. Inspired by [69, 86], we decompose $\beta \mathbf{g}_0 + (1 - \beta) \mathbf{g}_2$ into components that are parallel and orthogonal to $\alpha \mathbf{g}_1 + (1 - \alpha) \mathbf{g}_3$. Specifically, we first let

$$\mathbf{h}_+ = \alpha \mathbf{g}_1 + (1 - \alpha) \mathbf{g}_3, \quad \mathbf{h}_- = \beta \mathbf{g}_0 + (1 - \beta) \mathbf{g}_2 \quad (45)$$

to denote the positive and negative parts in Equation (44), respectively. We then decompose the negative part \mathbf{h}_- into two components that are parallel or orthogonal to \mathbf{h}_+ and we get \mathbf{h}_-^{\parallel} and \mathbf{h}_-^{\perp} . As a result, the final gradient is given by $\mathbf{h}_+ - \gamma \mathbf{h}_-^{\perp}$.