

---

# Towards Theoretically Understanding Why SGD Generalizes Better Than ADAM in Deep Learning

---

Pan Zhou\*, Jiashi Feng†, Chao Ma‡, Caiming Xiong\*, Steven HOI\*, Weinan E‡

\*Salesforce Research, †National University of Singapore, ‡Princeton University

{pzhou,shoi,cxiong}@salesforce.com elefjia@nus.edu.sg {chaom@, weinan@math.}princeton.edu

## Abstract

It is not clear yet why ADAM-alike adaptive gradient algorithms suffer from worse generalization performance than SGD despite their faster training speed. This work aims to provide understandings on this generalization gap by analyzing their local convergence behaviors. Specifically, we observe the heavy tails of gradient noise in these algorithms. This motivates us to analyze these algorithms through their Lévy-driven stochastic differential equations (SDEs) because of the similar convergence behaviors of an algorithm and its SDE. Then we establish the escaping time of these SDEs from a local basin. The result shows that (1) the escaping time of both SGD and ADAM depends on the Radon measure of the basin positively and the heaviness of gradient noise negatively; (2) for the same basin, SGD enjoys smaller escaping time than ADAM, mainly because (a) the geometry adaptation in ADAM via adaptively scaling each gradient coordinate well diminishes the anisotropic structure in gradient noise and results in larger Radon measure of a basin; (b) the exponential gradient average in ADAM smooths its gradient and leads to lighter gradient noise tails than SGD. So SGD is more locally unstable than ADAM at sharp minima defined as the minima whose local basins have small Radon measure, and can better escape from them to flatter ones with larger Radon measure. As flat minima here which often refer to the minima at flat or asymmetric basins/valleys often generalize better than sharp ones [1, 2], our result explains the better generalization performance of SGD over ADAM. Finally, experimental results confirm our heavy-tailed gradient noise assumption and theoretical affirmation.

## 1 Introduction

Stochastic gradient descent (SGD) [3, 4] has become one of the most popular algorithms for training deep neural networks [5–11]. In spite of its simplicity and effectiveness, SGD uses one learning rate for all gradient coordinates and could suffer from unsatisfactory convergence performance, especially for ill-conditioned problems [12]. To avoid this issue, a variety of adaptive gradient algorithms have been developed that adjust learning rate for each gradient coordinate according to the current geometry curvature of the objective function [13–16]. These algorithms, especially for ADAM, have achieved much faster convergence speed than vanilla SGD in practice.

Despite their faster convergence behaviors, these adaptive gradient algorithms usually suffer from worse generalization performance than SGD [12, 17, 18]. Specifically, adaptive gradient algorithms often show faster progress in the training phase but their performance quickly reaches a plateaus on test data. Differently, SGD usually improves model performance slowly but could achieve higher test performance. One empirical explanation [1, 19–21] for this generalization gap is that adaptive gradient algorithms tend to converge to sharp minima whose local basin has large curvature and usually generalize poorly, while SGD prefers to find flat minima and thus generalizes better. However, recent evidence [2, 22] shows that (1) for deep neural networks, the minima at the asymmetric

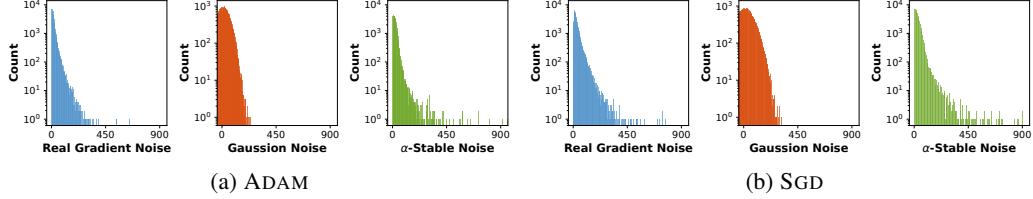


Figure 1: Illustration of gradient noise in ADAM and SGD on AlexNet trained with CIFAR10. (b) is produced under the same setting in [23]. By comparison, one can observe (1)  $\alpha$ -stable noise can better characterize real gradient noise and (2) SGD has heavier gradient noise tails than ADAM.

basins/valleys where both steep and flat directions exist also generalize well though they are sharp in terms of their local curvature, and (2) SGD often converges to these minima. So the argument of the conventional “flat” and “sharp” minima defined on curvature cannot explain these new results. Thus the reason for the generalization gap between adaptive gradient methods and SGD is still unclear.

In this work, we provide a new viewpoint for understanding the generalization performance gap. We first formulate ADAM and SGD as Lévy-driven stochastic differential equations (SDEs), since the SDE of an algorithm shares similar convergence behaviors of the algorithm and can be analyzed more easily than directly analyzing the algorithm. Then we analyze the escaping behaviors of these SDEs at local minima to investigate the generalization gap between ADAM and SGD, as escaping behaviors determine which basin that an algorithm finally converges to and thus affect the generalization performance of the algorithm. By analysis, we find that compared with ADAM, SGD is more locally unstable and is more likely to converge to the minima at the flat or asymmetric basins/valleys which often have better generalization performance over other type minima. So our results can explain the better generalization performance of SGD over ADAM. Our contributions are highlighted below.

Firstly, this work is the first one that adopts Lévy-driven SDE which better characterizes the algorithm gradient noise in practice, to analyze the adaptive gradient algorithms. Specifically, Fig. 1 shows that the gradient noise in ADAM and SGD, i.e. the difference between the full and stochastic gradients, has heavy tails and can be well characterized by symmetric  $\alpha$ -stable ( $S\alpha S$ ) distribution [24]. Based on this observation, we view ADAM and SGD as discretization of the continuous-time processes and formulate the processes as Lévy-driven SDEs to analyze their behaviors. Compared with Gaussian gradient noise assumption in SGD [25–27],  $S\alpha S$  distribution assumption can characterize the heavy-tailed gradient noise in practice more accurately as shown in Fig. 1, and also better explains the different generalization performance of SGD and ADAM as discussed in Sec. 3. This work extends [23, 28] from SGD on the over-simplified one-dimensional problems to much more complicated adaptive algorithms on high-dimensional problems. It also differs from [29], as [29] considers escaping behaviors of SGD along several fixed directions, while this work analyzes the dynamic underlying structures in gradient noise that plays an important role in the local escaping behaviors of both ADAM and SGD.

Next, we theoretically prove that for the Lévy-driven SDEs of ADAM and SGD, their escaping time  $\Gamma$  from a local basin  $\Omega$ , namely the least time for escaping from the inner of  $\Omega$  to its outside, is at the order of  $\mathcal{O}(\varepsilon^{-\alpha}/m(\mathcal{W}))$ , where the constant  $\varepsilon \in (0, 1)$  relies on the learning rate of algorithms and  $\alpha$  denotes the tail index of  $S\alpha S$  distribution. Here  $m(\mathcal{W})$  is a non-zero Radon measure on the escaping set  $\mathcal{W}$  of ADAM and SGD at the local basin  $\Omega$  (see Sec. 4.1), and actually negatively relies on the Radon measure of  $\Omega$ . So both ADAM and SGD have small escaping time at the “sharp” minima whose corresponding basins  $\Omega$  have small Radon measure. It means that ADAM and SGD are actually unstable at “sharp” minima and would escape them to “flatter” ones. Note, the Radon measure of  $\Omega$  positively depends on the volume of  $\Omega$ . So these results also well explain the observations in [1, 2, 20, 21] that the minima of deep networks found by SGD often locate at the flat or asymmetric valleys, as their corresponding basins have large volumes and thus large Radon measure.

Finally, our results can answer why SGD often converges to flatter minima than ADAM in terms of Radon measure, and thus explain the generalization gap between ADAM and SGD. Firstly, our analysis shows that even for the same basin  $\Omega$ , ADAM often has smaller Radon measure  $m(\mathcal{W})$  on the escaping set  $\mathcal{W}$  at  $\Omega$  than SGD, as the geometry adaptation in ADAM via adaptively scaling each gradient coordinate well diminishes underlying anisotropic structure in gradient noise and leads to smaller  $m(\mathcal{W})$ . Secondly, the empirical results in Sec. 5 and Fig. 1 show that SGD often has much

smaller tail index  $\alpha$  of gradient noise than ADAM for some optimization iterations and thus enjoys smaller factor  $\varepsilon^{-\alpha}$ . These results together show that SGD is more locally unstable and would like to converge to flatter minima with larger measure  $m(\mathcal{W})$  which often refer to the minima at the flat and asymmetric basins/valleys, according with empirical evidences in [12, 17, 30, 31]. Considering the observations in [1, 19–21] that the minima at the flat and asymmetric basins/valleys often generalize better, our results well explain the generalization gap between ADAM and SGD. Besides, our results also show that SGD benefits from its anisotropic gradient noise on its escaping behaviors, while ADAM does not.

## 2 Related Work

Adaptive gradient algorithms have become the default optimization tools in deep learning because of their fast convergence speed. But they often suffer from worse generalization performance than SGD [12, 17, 30, 31]. Subsequently, most works [12, 17, 18, 30, 31] empirically analyze this issue from the argument of flat and sharp minima defined on local curvature in [19] that flat minima often generalize better than sharp ones, as they observed that SGD often converges to flatter minima than adaptive gradient algorithms, *e.g.* ADAM. However, Sagun *et al.* [22] and He *et al.* [2] observed that the minima of modern deep networks at the asymmetric valleys where both steep and flat directions exist also generalize well, and SGD often converges to these minima. So the conventional flat and sharp argument cannot explain these new results. This work theoretically shows that SGD tends to converge to the minima whose local basin has larger Radon measure. It well explains the above new observations, as the minima with larger Radon measure often locate at the flat and asymmetric basins/valleys. Moreover, based on our results, exploring invariant Radon measure to parameter scaling in networks could resolve the issue in [32] that flat minima could become sharp via parameter scaling. See more details in Appendix B. Note, ADAM could achieve better performance than SGD when gradient clipping is required [33], *e.g.* attention models with gradient exploding issue, as adaptation in ADAM provides a clipping effect. This work considers a general non-gradient-exploding setting, as it is more practical across many important tasks, *e.g.* classification.

For theoretical generalization analysis, most works [25–27, 34] only focus on analyzing SGD. They formulated SGD into Brownian motion based SDE via assuming gradient noise to be Gaussian. For instance, Jastrzkebski *et al.* [26] proved that the larger ratio of learning rate to mini-batch size in SGD leads to flatter minima and better generalization. But Simsekli *et al.* [23] empirically found that the gradient noise has heavy tails and can be characterized by  $S\alpha S$  distribution instead of Gaussian distribution. Chaudhari *et al.* [27] also claimed that the trajectories of SGD in deep networks are not Brownian motion. Then Simsekli *et al.* [23] formulated SGD as a Lévy-driven SDE and adopted the results in [28] to show that SGD tends to converge to flat minima on one dimensional problems. Pavlyukevich *et al.* [29] extended the one-dimensional SDE in [28] and analyzed escaping behaviors of SGD along several fixed directions, differing from this work that analyzes dynamic underlying structures in gradient noise that greatly affect escaping behaviors of both ADAM and SGD.

The literature targeting theoretically understanding the generalization degeneration of adaptive gradient algorithms are limited mainly due to their more complex algorithms. Wilson *et al.* [17] constructed a binary classification problem and showed that ADAGRAD [13] tend to give undue influence to spurious features that have no effect on out-of-sample generalization. Unlike the above theoretical works that focus on analyzing SGD only or special problems, we target at revealing the different convergence behaviors of adaptive gradient algorithms and SGD and also analyzing their different generalization performance, which is of more practical interest especially in deep learning.

## 3 Lévy-driven SDEs of Algorithms in Deep Learning

In this section, we first briefly introduce SGD and ADAM, and formulate them as discretization of stochastic differential equations (SDEs) which is a popular approach to analyze algorithm behaviors. Suppose the objective function of  $n$  components in deep learning models is formulated as

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \mathbf{F}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n f_i(\boldsymbol{\theta}), \quad (1)$$

where  $f_i(\boldsymbol{\theta})$  is the loss of the  $i$ -th sample. Subsequently, we focus on analyzing SGD and ADAM. Note our analysis technique is applicable to other adaptive algorithms with similar results as ADAM.

### 3.1 SGD and ADAM

As one of the most effective algorithms, SGD [3] solves problem (1) by sampling a data mini-batch  $\mathcal{S}_t$  of size  $S$  and then running one gradient descent step:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \nabla f_{\mathcal{S}_t}(\boldsymbol{\theta}_t), \quad (2)$$

where  $\nabla f_{\mathcal{S}_t}(\boldsymbol{\theta}_t) = \frac{1}{S} \sum_{i \in \mathcal{S}_t} \nabla f_i(\boldsymbol{\theta}_t)$  denotes the gradient on mini-batch  $\mathcal{S}_t$ , and  $\eta$  is the learning rate. Recently, to improve the efficiency of SGD, adaptive gradient algorithms, such as ADAGRAD [13], RMSPROP [14] and ADAM [15], are developed which adjust the learning rate of each gradient coordinate according to the current geometric curvature. Among them, ADAM has become the default training algorithm in deep learning. Specifically, ADAM estimates the current gradient  $\nabla \mathbf{F}(\boldsymbol{\theta}_t)$  as

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \nabla f_{\mathcal{S}_t}(\boldsymbol{\theta}_t) \quad \text{with } \mathbf{m}_0 = \mathbf{0} \text{ and } \beta_1 \in (0, 1).$$

Then like natural gradient descent [35], ADAM adapts itself to the function geometry via a diagonal Fisher matrix approximation  $\text{diag}(\mathbf{v}_t)$  which serves as a preconditioner and is defined as

$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) [\nabla f_{\mathcal{S}_t}(\boldsymbol{\theta}_t)]^2 \quad \text{with } \mathbf{v}_0 = \mathbf{0} \text{ and } \beta_2 \in (0, 1).$$

Next ADAM preconditions the problem by scaling each gradient coordinate, and updates the variable

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \mathbf{m}_t / (1 - \beta_1^t) / \left( \sqrt{\mathbf{v}_t / (1 - \beta_2^t)} + \epsilon \right) \quad \text{with a small constant } \epsilon. \quad (3)$$

### 3.2 Lévy-driven SDEs

Let  $\mathbf{u}_t = \nabla \mathbf{F}(\boldsymbol{\theta}_t) - \nabla f_{\mathcal{S}_t}(\boldsymbol{\theta}_t)$  denote gradient noise. From Sec. 3.1, we can formulate SGD as follows

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \nabla \mathbf{F}(\boldsymbol{\theta}_t) + \eta \mathbf{u}_t.$$

To analyze behaviors of an algorithm, one effective approach is to obtain its SDE via making assumptions on  $\mathbf{u}_t$  and then analyze its SDE. For instance, to analyze SGD, most works [25–27, 34] assume that  $\mathbf{u}_t$  obeys a Gaussian distribution  $\mathcal{N}(\mathbf{0}, \Sigma_t)$  with covariance matrix

$$\Sigma_t = \frac{1}{S} \left[ \frac{1}{n} \sum_{i=1}^n \nabla f_i(\boldsymbol{\theta}_t) \nabla f_i(\boldsymbol{\theta}_t)^T - \nabla \mathbf{F}(\boldsymbol{\theta}_t) \nabla \mathbf{F}(\boldsymbol{\theta}_t)^T \right].$$

However, both recent work [23] and Fig. 1 show that the gradient noise  $\mathbf{u}_t$  has heavy tails and can be better characterized by  $\mathcal{S}\alpha\mathcal{S}$  distribution [24]. Moreover, the heavy-tail assumption can also better explain the behaviors of SGD than Gaussian noise assumption. Concretely, for the SDE of SGD on the one-dimensional problems, under Gaussian noise assumption its escaping time from a simple quadratic basin respectively exponentially and polynomially depends on the height and width of the basin [36], indicating that SGD gets stuck at deeper minima as opposed to wider/flatter minima. This contradicts with the observations in [1, 19–21] that SGD often converges to flat minima. By contrast, on the same problem, for Lévy-driven SDE, both [23] and this work show that SGD tends to converge to flat minima instead of deep minima, well explaining the convergence behaviors of SGD.

Following [23], we also assume  $\mathbf{u}_t$  obeys  $\mathcal{S}\alpha\mathcal{S}$  distribution but with a time-dependent covariance matrix  $\Sigma_t$  to better characterize the underlying structure in the gradient noise  $\mathbf{u}_t$ . In this way, when the learning rate  $\eta$  is small and  $\varepsilon = \eta^{(\alpha-1)/\alpha}$ , we can write the Lévy-driven SDE of SGD as

$$d\boldsymbol{\theta}_t = -\nabla \mathbf{F}(\boldsymbol{\theta}_t) + \varepsilon \Sigma_t dL_t. \quad (4)$$

Here the Lévy motion  $L_t \in \mathbb{R}^d$  is a random vector and its  $i$ -th entry  $L_{t,i}$  obeys the  $\mathcal{S}\alpha\mathcal{S}(1)$  distribution which is defined through the characteristic function  $\mathbb{E}[\exp(i\omega x)] = \exp(-\sigma^\alpha |\omega|^\alpha)$  if  $x \sim \mathcal{S}\alpha\mathcal{S}(\sigma)$ . Intuitively, the  $\mathcal{S}\alpha\mathcal{S}$  distribution is a heavy-tailed distribution with a decay density like  $1/|x|^{1+\alpha}$ . When the tail index  $\alpha$  is 2,  $\mathcal{S}\alpha\mathcal{S}(1)$  becomes a Gaussian distribution and thus has stronger data-fitting capacity over Gaussian distribution. In this sense, the SDE of SGD in [25–27, 34, 37] is actually a special case of the Lévy-driven SDE in this work. Moreover, Eqn. (4) extends the one-dimensional SDE of SGD in [23]. Note, Eqn. (4) differs from [29], since it considers dynamic covariance matrix  $\Sigma_t$  in gradient noise and shows great effects of its underlying structure to the escaping behaviors in both ADAM and SGD, while [29] analyzed escaping behaviors of SGD along several fixed directions.

Similarly, we can derive the SDE of ADAM. For brevity, we define  $\mathbf{m}'_t = \beta_1 \mathbf{m}'_{t-1} + (1 - \beta_1) \nabla \mathbf{F}(\boldsymbol{\theta}_t)$  with  $\mathbf{m}'_0 = \mathbf{0}$ . Then by the definitions of  $\mathbf{m}_t$  and  $\mathbf{m}'_t$ , we can compute

$$\mathbf{m}'_t - \mathbf{m}_t = (1 - \beta_1) \sum_{i=0}^t \beta_1^{t-i} [\nabla \mathbf{F}(\boldsymbol{\theta}_i) - \nabla f_{\mathcal{S}_i}(\boldsymbol{\theta}_i)] = (1 - \beta_1) \sum_{i=0}^t \beta_1^{t-i} \mathbf{u}_i.$$

As noise  $\mathbf{u}_t$  has heavy tails, their exponential average should have similar behaviors, which is also illustrated by Fig. 1. So we also assume  $\frac{1}{1-\beta_1^t}(\mathbf{m}'_t - \mathbf{m}_t)$  obeys  $\mathcal{S}\alpha\mathcal{S}(1)$  distribution with covariance matrix  $\Sigma_t$ . Meanwhile, we can write ADAM as

$$\theta_{t+1} = \theta_t - \eta \mathbf{m}'_t / z_t + \eta (\mathbf{m}'_t - \mathbf{m}_t) / z_t \quad \text{with} \quad z_t = (1 - \beta_1^t) \left( \sqrt{\mathbf{v}_t / (1 - \beta_2^t)} + \epsilon \right).$$

So we can derive the Lévy-driven SDE of ADAM:

$$d\theta_t = -\mu_t Q_t^{-1} \mathbf{m}_t + \varepsilon Q_t^{-1} \Sigma_t dL_t, \quad d\mathbf{m}_t = \beta_1 (\nabla F(\theta_t) - \mathbf{m}_t), \quad d\mathbf{v}_t = \beta_2 ([\nabla f_{S_t}(\theta_t)]^2 - \mathbf{v}_t), \quad (5)$$

where  $\varepsilon = \eta^{(\alpha-1)/\alpha}$ ,  $Q_t = \text{diag}(\sqrt{\omega_t} \mathbf{v}_t + \epsilon)$ ,  $\mu_t = 1/(1 - e^{-\beta_1 t})$  and  $\omega_t = 1/(1 - e^{-\beta_2 t})$  are two constants to correct the bias in  $\mathbf{m}_t$  and  $\mathbf{v}_t$ . Note, here we replace  $\mathbf{m}'_t$  with  $\mathbf{m}_t$  for brevity. Appendix A provides more construction details, randomness discussion and shows the fitting capacity of this SDE to ADAM. Subsequently, we will analyze escaping behaviors of the SDEs in Eqns. (4) and (5).

## 4 Analysis for Escaping Local Minima

Now we analyze the local stability of ADAM-alike adaptive algorithms and SGD. Suppose the process  $\theta_t$  in Eqns. (4) and (5) starts at a local basin  $\Omega$  with a minimum  $\theta^*$ , i.e.  $\theta_0 \in \Omega$ . Here we are particularly interested in the first escaping time  $\Gamma$  of  $\theta_t$  produced by an algorithm which reveals the convergence behaviors and generalization performance of the algorithm. Formally, let  $\Omega^{-\varepsilon^\gamma} = \{y \in \Omega \mid \text{dis}(\partial\Omega, y) \geq \varepsilon^\gamma\}$  denote the inner part of  $\Omega$ . Then we give two important definitions, i.e. (1) the escaping time  $\Gamma$  of the process  $\theta_t$  from the local basin  $\Omega$  and (2) the escaping set  $\mathcal{W}$  at  $\Omega$ , as

$$\Gamma = \inf\{t \geq 0 \mid \theta_t \notin \Omega^{-\varepsilon^\gamma}\} \quad \text{and} \quad \mathcal{W} = \{y \in \mathbb{R}^d \mid Q_{\theta^*}^{-1} \Sigma_{\theta^*} y \notin \Omega^{-\varepsilon^\gamma}\}, \quad (6)$$

where the constant  $\gamma > 0$  satisfies  $\lim_{\varepsilon \rightarrow 0} \varepsilon^\gamma = 0$ ,  $\Sigma_{\theta^*} = \lim_{\theta_t \rightarrow \theta^*} \Sigma_t$  for both SGD and ADAM, and  $Q_{\theta^*} = I$  for SGD and  $Q_{\theta^*} = \lim_{\theta_t \rightarrow \theta^*} Q_t$  for ADAM. Then we define Radon measure [38].

**Definition 1.** If a measure  $m(\mathcal{V})$  defined on Hausdorff topological space  $\mathcal{X}$  obeys (1) inner regular, i.e.  $m(\mathcal{V}) = \sup_{U \subseteq \mathcal{V}} m(U)$ , (2) outer regular, i.e.  $m(\mathcal{V}) = \inf_{U \supseteq \mathcal{V}} m(U)$ , and (3) local finiteness, i.e. every point of  $\mathcal{X}$  has a neighborhood  $U$  with finite  $m(U)$ , then  $m(\mathcal{V})$  is a Radon measure.

Then we define non-zero Radon measure which further obeys  $m(U) < m(\mathcal{V})$  if  $U \subset \mathcal{V}$ . Since larger set has larger volume,  $m(U)$  positively depends on the volume of the set  $U$ . Let  $m(\mathcal{W})$  be a non-zero Radon measure on the set  $\mathcal{W}$ . Then we first introduce two mild assumptions for analysis.

**Assumption 1.** For both ADAM and SGD, suppose the objective  $F(\theta)$  is a upper-bounded non-negative loss, and is locally  $\mu$ -strongly convex and  $\ell$ -smooth in the basin  $\Omega$ .

**Assumption 2.** For ADAM, suppose its process  $\theta_t$  satisfies  $\int_0^\Gamma \langle \frac{\nabla F(\theta_s)}{1+F(\theta_s)}, \mu_s Q_s^{-1} \mathbf{m}_s \rangle ds \geq 0$  almost surely, and its parameters  $\beta_1$  and  $\beta_2$  obey  $\beta_1 \leq \beta_2 \leq 2\beta_1$ . Moreover, for ADAM, we assume  $\|\mathbf{m}_t - \widehat{\mathbf{m}}_t\| \leq \tau_m \|\int_0^{t-} (\mathbf{m}_s - \widehat{\mathbf{m}}_s) ds\|$  and  $\|\widehat{\mathbf{m}}_t\| \geq \tau \|\nabla F(\widehat{\theta}_t)\|$  where  $\widehat{\mathbf{m}}_t$  and  $\widehat{\theta}_t$  are obtained by Eqn. (5) with  $\varepsilon = 0$ . Each coordinate  $v_{t,i}$  of  $\mathbf{v}_t$  in ADAM obeys  $v_{\min} \leq \sqrt{v_{t,i}} \leq v_{\max}$  ( $\forall i, t$ ).

Assumption 1 is very standard for analyzing stochastic optimization [39–44] and network analysis [45–48]. In Assumption 2, we indeed require similar directions of gradient estimate  $\mathbf{m}_t$  and full gradient  $\nabla F(\theta_t)$  in ADAM in most cases, as we assume their inner product is non-negative along the iteration trajectory. So this assumption can be satisfied in practice. To analyze the processes  $\theta_t$  and  $\widehat{\theta}_t$  in ADAM, we make an assumption on the distance between their corresponding gradient estimates  $\mathbf{m}_t$  and  $\widehat{\mathbf{m}}_t$  which can be easily fulfilled by their definitions. Then for ADAM, we mildly assume its estimated  $\mathbf{v}_t$  to be bounded. For  $v_{\min}$ , we indeed allow  $v_{\min} = 0$  because of the small constant  $\epsilon$ . The relation  $\beta_1 \leq \beta_2 \leq 2\beta_1$  is also satisfied under the default setting of ADAM. Actually, we also empirically investigate Assumption 2 on ADAM. In Fig. 2, we report the values of  $\rho_t = \frac{10}{t} \int_0^t \langle \frac{\nabla F(\theta_s)}{1+F(\theta_s)}, \mu_s Q_s^{-1} \mathbf{m}_s \rangle ds$ ,  $\tau'_m = \frac{\|\mathbf{m}_t - \widehat{\mathbf{m}}_t\|}{\|\int_0^{t-} (\mathbf{m}_s - \widehat{\mathbf{m}}_s) ds\|}$ ,  $\tau' = \frac{\|\widehat{\mathbf{m}}_t\|}{\|\nabla F(\widehat{\theta}_t)\|}$ ,  $v_{\min} = \min_i \sqrt{v_{t,i}}$ ,  $v_{\max} = \max_i \sqrt{v_{t,i}}$  in the SDE of ADAM on the 4-layered fully connected networks with width 20. Note that we scale some values of  $\rho_t$ ,  $\tau'_m$ ,  $\tau'$ ,  $v_{\min}$  and  $v_{\max}$  so that we can

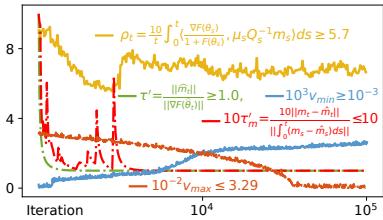


Figure 2: Empirical investigation of Assumption 2 on ADAM.

plot them in one figures. From Fig. 2, one can observe that  $\rho_t$ ,  $\tau'$  and  $v_{\min}$  are well lower bounded, and  $\tau'_m$  and  $v_{\max}$  are well upper bounded. These results demonstrate the validity of Assumption 2.

With these two assumptions, we analyze the escaping time  $\Gamma$  of process  $\theta_t$  and summarize the main results in Theorem 1. For brevity, we define a group of key constants for SGD and ADAM:  $\kappa_1 = \ell$  and  $\kappa_2 = 2\mu$  in SGD,  $\kappa_1 = \frac{c_1\ell}{(v_{\min}+\epsilon)|\tau_m-1|}$  and  $\kappa_2 = \frac{2\mu\tau}{\beta_1(v_{\max}+\epsilon)+\mu\tau}(\beta_1 - \frac{\beta_2}{4})$  in ADAM with a constant  $c_1$ .

**Theorem 1.** Suppose Assumptions 1 and 2 hold. Let  $\Theta(\varepsilon^{-1}) = \frac{2}{\alpha}\varepsilon^\alpha$ ,  $\rho_0 = \frac{1}{16(1+c_2\kappa_1)}$  and  $\ln\left(\frac{2\Delta}{\mu\varepsilon^{1/3}}\right) \leq \kappa_2\varepsilon^{-\frac{1}{3}}$  with  $\Delta = \mathbf{F}(\theta_0) - \mathbf{F}(\theta^*)$  and a constant  $c_2$ . Then for any  $\theta_0 \in \Omega^{-2\varepsilon^\gamma}$ ,  $u > -1$ ,  $\varepsilon \in (0, \varepsilon_0]$ ,  $\gamma \in (0, \gamma_0]$  and  $\rho \in (0, \rho_0]$  satisfying  $\varepsilon^\gamma \leq \rho_0$  and  $\lim_{\varepsilon \rightarrow 0} \rho = 0$ , SGD in (4) and ADAM in (5) obey

$$\frac{1-\rho}{1+u+\rho} \leq \mathbb{E} [\exp(-um(\mathcal{W})\Theta(\varepsilon^{-1})\Gamma)] \leq \frac{1+\rho}{1+u-\rho}.$$

See its proof in Appendix D.1. By setting  $\varepsilon$  small, Theorem 1 shows that for both ADAM and SGD, the upper and lower bounds of their expected escaping time  $\Gamma$  are at the order of  $\mathcal{O}(\frac{1}{m(\mathcal{W})\Theta(\varepsilon^{-1})})$ .

Note,  $m(\mathcal{W})$  has different values for SGD and ADAM due to their different  $Q_{\theta^*}$  in Eqn. (6). If the escaping time  $\Gamma$  is very large, it means that the algorithm cannot well escape from the basin  $\Omega$  and would get stuck in  $\Omega$ . Moreover, given the same basin  $\Omega$ , if one algorithm has smaller escaping time  $\Gamma$  than other algorithms, then it is more locally unstable and would faster escape from this basin to others. In the following sections, we discuss the effects of the geometry adaptation and the gradient noise structure of ADAM and SGD to the escaping time  $\Gamma$  which are respectively reflected by the factors  $m(\mathcal{W})$  and  $\Theta(\varepsilon^{-1})$ . Our results show that SGD has smaller escaping time than ADAM and can better escape from local basins with small Radon measure to those with larger Radon measure.

#### 4.1 Preference to Flat Minima

To interpret Theorem 1, we first define the “flat” minima in this work in terms of Radon measure.

**Definition 2.** A minimum  $\theta^* \in \Omega$  is said to be flat if its basin  $\Omega$  has large nonzero Radon measure.

Due to the factor  $m(\mathcal{W})$  in Theorem 1, both ADAM and SGD have large escaping time  $\Gamma$  at the “flat” minima. Specifically, if the basin  $\Omega$  has larger Radon measure, then the complementary set  $\mathcal{W}^c = \{y \in \mathbb{R}^d \mid Q_{\theta^*}^{-1}\Sigma_{\theta^*}y \in \Omega^{-\varepsilon^\gamma}\}$  of  $\mathcal{W}$  also has larger Radon measure. Meanwhile, the Radon measure on  $\mathcal{W}^c \cup \mathcal{W}$  is a constant, meaning the larger  $m(\mathcal{W}^c)$  the smaller  $m(\mathcal{W})$ . So ADAM and SGD have larger escaping time at “flat” minima. Thus, they would escape “sharp” minima due to their smaller escaping time, and tend to converge to “flat” ones. Since for basin  $\Omega$ , its Radon measure positively relies on its volume,  $m(\mathcal{W})$  negatively depends on the volume of  $\Omega$ . So ADAM and SGD are more stable at the minima with larger basin  $\Omega$  in terms of volume. This can be intuitively understood: for the process  $\theta_t$ , the volume of the basin determines the necessary jump size of the Lévy motion  $L_t$  in the SDEs to escape, which means the larger the basin the harder for an algorithm to escape.

Note the “flat” minima here is defined on Radon measure, and differ from the conventional flat ones whose local basins have no large curvature (no large eigenvalues in its Hessian matrix). In most cases, the flat minima here consist of the conventional flat ones and the minima at the asymmetric basins/valleys since local basins of these minima often have large volumes and thus larger Radon measures. Accordingly, our theoretical results can well explain the phenomena observed in many works [2, 12, 17, 22, 30, 31] that SGD often converges to the minima at the flat or asymmetric valleys which is interpreted by our theory to have larger Radon measure and attract SGD to stay at these places. In contrast, the conventional flat argument cannot explain asymmetric valleys, as asymmetric valleys means sharp minima under the conventional definition and should be avoided by SGD.

#### 4.2 Analysis of Generalization Gap between ADAM and SGD

Theorem 1 can also well explain the generalization gap between ADAM-alike adaptive algorithms and SGD. That is, compared with SGD, the minima given by ADAM often suffer from worse test performance [12, 17, 18, 30, 31]. On one hand, the observations in [1, 19–21] show that the minima at the flat or asymmetric basins/valleys often enjoy better generalization performance than others. On the other hand, Theorem 1 shows that ADAM and SGD can escape sharp minima to flat ones with larger Radon measure. As aforementioned, flat minima in terms of Radon measure often refer to the minima at the flat or asymmetric basins/valleys. This implies that if one algorithm can escape the

current minima faster, it is more likely for the algorithm to find flatter minima. These results together show the benefit of faster escaping behaviors of an algorithm to its generalization performance.

According to Theorem 1, two main factors, i.e. the gradient noise and geometry adaptation respectively reflected by the factors  $\Theta(\varepsilon^{-1}) = \frac{2}{\alpha} \varepsilon^{-\alpha}$  and  $m(\mathcal{W})$ , affects the escaping time  $\Gamma$  of both ADAM and SGD. We first look at the factor  $\Theta(\varepsilon^{-1})$  in the escaping time  $\Gamma$ . As illustrated in Fig. 3 in Sec. 5, the gradient noise in SGD enjoys very similar tail index  $\alpha$  with ADAM for most optimization iterations, but it has much smaller tail index  $\alpha$  than ADAM for some iterations, which means SGD has larger Lévy discontinuous jumps in these iterations and thus enjoys smaller escaping time  $\Gamma$ . This different tail property of gradient noise in these algorithms are caused by the following reason. SGD assumes the gradient noise  $\mathbf{u}_t = \nabla \mathbf{F}(\theta_t) - \nabla f_{S_t}(\theta_t)$  at one iteration has heavy tails, while ADAM considers the exponential gradient noise  $\frac{1-\beta_1}{1-\beta_1^t} \sum_{i=0}^t \beta_1^{t-i} \mathbf{u}_i$  which indeed smooths gradient noise over the iteration trajectory and prevents large occasional gradient noise. In this way, SGD reveals heavier tails of gradient noise than ADAM and thus has smaller tail index  $\alpha$  for some optimization iterations, helping escaping behaviors. Moreover, to guarantee convergence, ADAM needs to use smaller learning rate  $\eta$  than SGD due to the geometry adaptation in ADAM, e.g. default learning rate  $10^{-3}$  in ADAM and  $10^{-2}$  in SGD, leading to smaller  $\varepsilon = \eta^{(\alpha-1)/\alpha}$  and thus larger escaping time  $\Gamma$  in ADAM. Thus, compared with ADAM, SGD is more locally unstable and will converge to flatter minima which often locate at the flat or asymmetric basins/valleys and enjoy better generalization performance [1, 19–21].

Besides, the factor  $m(\mathcal{W})$  also plays an important role in the generalization degeneration phenomenon of ADAM. W.o.l.g., assume the minimizer  $\theta^* = \mathbf{0}$  in the basin  $\Omega$ . As the local basin  $\Omega$  is often small, following [34, 49] we adopt second-order Taylor expansion to approximate  $\Omega$  as a quadratic basin with center  $\theta^*$ , i.e.  $\Omega = \{\mathbf{y} \mid \mathbf{F}(\theta^*) + \frac{1}{2} \mathbf{y}^T \mathbf{H}(\theta^*) \mathbf{y} \leq h(\theta^*)\}$  with a basin height  $h(\theta^*)$  and Hessian matrix  $\mathbf{H}(\theta^*)$  at  $\theta^*$ . Then for SGD, since  $\mathbf{Q}_{\theta^*} = \mathbf{I}$  in Eqn. (6), its corresponding escaping set  $\mathcal{W}$  is

$$\mathcal{W}_{\text{SGD}} = \{\mathbf{y} \in \mathbb{R}^d \mid \mathbf{y}^T \Sigma_{\theta^*} \mathbf{H}(\theta^*) \Sigma_{\theta^*} \mathbf{y} \geq h_f^*\} \quad (7)$$

with  $h_f^* = 2(h(\theta^*) - \mathbf{F}(\theta^*))$ , while according to Eqn. (6), ADAM has escaping set

$$\mathcal{W}_{\text{ADAM}} = \{\mathbf{y} \in \mathbb{R}^d \mid \mathbf{y}^T \Sigma_{\theta^*} \mathbf{Q}_{\theta^*}^{-1} \mathbf{H}(\theta^*) \mathbf{Q}_{\theta^*}^{-1} \Sigma_{\theta^*} \mathbf{y} \geq h_f^*\}. \quad (8)$$

Then we prove that for most time interval except the jump time, the current variable  $\theta_t$  is indeed close to the minimum  $\theta^*$ . Specifically, we first decompose the Lévy process  $L_t$  into two components  $\xi_t$  and  $\zeta_t$ , i.e.  $L_t = \xi_t + \zeta_t$ , with the jump sizes  $\|\xi_t\| < \varepsilon^{-\delta}$  and  $\|\zeta_t\| \geq \varepsilon^{-\delta}$  ( $\delta \in (0, 1)$ ). In this way, the stochastic process  $\xi$  does not departure from  $\theta_t$  a lot due to its limited jump size. The process  $\zeta$  is a compound Poisson process with intensity  $\Theta(\varepsilon^{-\delta}) = \int_{\|\mathbf{y}\| \geq \varepsilon^{-\delta}} \nu(d\mathbf{y}) = \int_{\|\mathbf{y}\| \geq \varepsilon^{-\delta}} \frac{d\mathbf{y}}{\|\mathbf{y}\|^{1+\alpha}} = \frac{2}{\alpha} \varepsilon^{\alpha\delta}$  and jumps distributed according to the law of  $1/\Theta(\varepsilon^{-\delta})$ . Specifically, let  $0 = t_0 < t_1 < \dots < t_k < \dots$  denote the time of successive jumps of  $\zeta$ . Then the inner-jump time intervals  $\sigma_k = t_k - t_{k-1}$  are i.i.d. exponentially distributed random variables with mean value  $\mathbb{E}(\sigma_k) = \frac{1}{\Theta(\varepsilon^{-\delta})}$  and probability function  $\mathbb{P}(\sigma_k \geq x) = \exp(-x\Theta(\varepsilon^{-\delta}))$ . Based on this decomposition, we state our results in Theorem 2.

**Theorem 2.** Suppose Assumptions 1 and 2 hold. Assume the process  $\hat{\theta}_t$  is produced by setting  $\varepsilon = 0$  in the Lévy-driven SDEs of SGD and ADAM.

(1)  $\hat{\theta}_t$  exponentially converges to the minimizer  $\theta^*$  in  $\Omega$ . Specifically, by defining  $\Delta = \mathbf{F}(\theta_0) - \mathbf{F}(\theta^*)$ ,  $\kappa_3 = \frac{2\mu\tau}{\beta_1(v_{\max} + \epsilon) + \mu\tau} (\beta_1 - \frac{\beta_2}{4})$  in ADAM and  $\kappa_3 = 2\mu$  in SGD, for any  $\bar{\rho} > 0$ , it satisfies

$$\|\hat{\theta}_t - \theta^*\|_2^2 \leq \varepsilon^{\bar{\rho}} \quad \text{if } t \geq v_\varepsilon \triangleq \kappa_3^{-1} \ln(2\Delta\mu^{-1}\varepsilon^{-\bar{\rho}}).$$

(2) Assume  $\delta \in (0, 1)$ ,  $p = \min((\bar{\rho}(1 + c_3\kappa_1))/4, \bar{\rho})$ ,  $\bar{\rho} = \frac{1-\delta}{16(1+c_4\kappa_1)}$ , where  $\kappa_1$  (in Theorem 1),  $\bar{\rho}$ ,  $c_3$  and  $c_4$  are four positive constants. When  $\theta_t$  and  $\hat{\theta}_t$  have the same initialization  $\theta_0 = \hat{\theta}_0$ , we have

$$\sup_{\theta_0 \in \Omega} \mathbb{P}(\sup_{0 \leq t < \sigma_1} \|\theta_t - \hat{\theta}_t\|_2 \geq 2\varepsilon^{\bar{\rho}}) \leq 2 \exp(-\varepsilon^{-p}).$$

See its proof in Appendix D.2. By inspecting the first part of Theorem 2, one can observe that the gradient-noise-free processes  $\hat{\theta}_t$  produced by setting  $\varepsilon = 0$  in the Lévy-driven SDEs of SGD and ADAM locate in a very small neighborhood of the minimizer  $\theta^*$  in the local basin  $\Omega$  after a very small time interval  $v_\varepsilon = \kappa_3^{-1} \ln(2\Delta\mu^{-1}\varepsilon^{-\bar{\rho}})$ . The second part of Theorem 2 shows that before the first jump time  $t_1 = \sigma_1$  of the jump  $\zeta$  with size larger than  $\varepsilon^{-\delta}$  in Lévy motion  $L_t$ , the distance

between  $\theta_t$  and  $\widehat{\theta}_t$  is very small. So these two parts together guarantee small distance between  $\theta_t$  and  $\theta^*$  for the most time interval before the first big jump in the Lévy motion  $L_t$  since the mean jump time  $\mathbb{E}(\sigma_1) = \frac{\alpha}{2\varepsilon^{\alpha\delta}} = \mathcal{O}(\varepsilon^{-1})$  of the first big jump is much larger than  $v_\varepsilon = \mathcal{O}(\ln(\varepsilon^{-1}))$  when  $\varepsilon$  is small. Next after the first big jump, if  $\theta_t$  does not escape from the local basin  $\Omega$ , by using the first part of Theorem 2, after the time interval  $v_\varepsilon$ ,  $\theta_t$  becomes close to  $\theta^*$  again. This process will continue until the algorithm escapes from the basin. So for most time interval before escaping from  $\Omega$ , the stochastic process  $\theta_t$  locates in a very small neighborhood of the minimizer  $\theta^*$ .

The above analysis results on Theorem 2 hold for moderately ill-conditioned local basins (ICLBs). Specifically, the analysis requires  $v_\varepsilon \leq \sigma_1$  to guarantee small distance of current solution  $\theta_t$  to  $\theta^*$  before each big jump. So if  $\mu$  of ICLBs is larger than  $\mathcal{O}(\varepsilon^{\alpha\delta})$  which is very small as  $\varepsilon$  in SDE is often small to precisely mimic algorithm behaviors, The above analysis results 2 still hold. Moreover, to obtain the result (1) in Theorem 2, we assume the optimization trajectory goes along the eigenvector direction corresponding to  $\mu$  which is the worse case and leads to the worst convergence speed. As the measure of one/several eigenvector directions on high dimension is 0, optimization trajectory cannot always go along the eigenvector direction corresponding to  $\mu$ . So  $v_\varepsilon$  is actually much larger than  $\mathcal{O}(\frac{1}{\mu} \ln(\frac{1}{\mu\varepsilon^\delta}))$ , largely improving applicability of our theory. For extremely ICLBs ( $\mu \rightarrow 0$  or  $\mu = 0$ ), the above analysis does not hold which accords with the previous results that first-order gradient algorithms cannot escape from them provably [50]. Fortunately,  $\mu \rightarrow 0$  and  $\mu = 0$  give asymmetric basins which often generalize well [2, 22] and are not needed to escape.

By using the above results, we have  $\theta_t \approx \theta^*$  before escaping and thus  $v_t = \lim_{\theta_t \rightarrow \theta^*} [\nabla f_{\mathcal{S}_t}(\theta_t)]^2$ . Considering the randomness of the mini-batch  $\mathcal{S}_t$ ,  $\omega_t \approx 1$  and  $\epsilon \approx 0$ , we can approximate

$$\mathbb{E}[\mathbf{Q}_{\theta^*}] \approx \mathbb{E}[\lim_{\theta_t \rightarrow \theta^*} \text{diag}(\sqrt{\omega_t v_t})] \approx \text{diag}\left(\sqrt{\frac{1}{n} \sum_{i=1}^n [\nabla f_i(\theta^*)]^2}\right).$$

Meanwhile, since  $\Sigma_{\theta^*} = \frac{1}{S} \bar{\Sigma}_{\theta^*}$  because of  $\lim_{\theta_t \rightarrow \theta^*} \mathbf{F}(\theta_t) = \mathbf{0}$  where  $\bar{\Sigma}_{\theta^*} = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\theta^*) \nabla f_i(\theta^*)^T$ , one can approximately compute  $\mathbb{E}[\Sigma_{\theta^*} \mathbf{Q}_{\theta^*}^{-1}] \approx \frac{1}{S} \mathbf{I}$ . Plugging this result into the escaping set  $\mathcal{W}_{\text{ADAM}}$  yields

$$\mathcal{W}_{\text{ADAM}} \approx \left\{ \mathbf{y} \in \mathbb{R}^d \mid \mathbf{y}^T \mathbf{H}(\theta^*) \mathbf{y} \geq S^2 h_f^* \right\}.$$

Now we compare the escaping sets  $\mathcal{W}_{\text{SGD}}$  of SGD and  $\mathcal{W}_{\text{ADAM}}$  of ADAM. For clarity, we re-write  $\mathcal{W}_{\text{SGD}}$  in Eqn. (7) as

$$\mathcal{W}_{\text{SGD}} = \left\{ \mathbf{y} \in \mathbb{R}^d \mid \mathbf{y}^T \bar{\Sigma}_{\theta^*} \mathbf{H}(\theta^*) \bar{\Sigma}_{\theta^*} \mathbf{y} \geq S^2 h_f^* \right\}.$$

By comparison, one can observe that for ADAM, its gradient noise does not affect the escaping set  $\mathcal{W}_{\text{ADAM}}$  due to the geometry adaptation via scaling each gradient coordinate, while for SGD, its gradient noise plays an important role. Suppose  $\bar{\mathbf{H}}(\theta^*) = \bar{\Sigma}_{\theta^*} \mathbf{H}(\theta^*) \bar{\Sigma}_{\theta^*}$ , and the singular values of  $\mathbf{H}(\theta^*)$  and  $\bar{\Sigma}_{\theta^*}$  are respectively  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$  and  $\varsigma_1 \geq \varsigma_2 \geq \dots \geq \varsigma_d$ . Zhu *et al.* [34] proved that  $\bar{\Sigma}_{\theta^*}$  of SGD on deep neural networks well aligns the Hessian matrix  $\mathbf{H}(\theta^*)$ , namely the top eigenvectors associated with large eigenvalues in  $\Sigma_{\theta^*}$  have similar directions in those in  $\mathbf{H}(\theta^*)$ . Besides, for modern over-parameterized neural networks, both Hessian  $\mathbf{H}(\theta^*)$  and the gradient covariance matrix  $\Sigma_{\theta^*}$  are ill-conditioned and anisotropic near minima [22, 27]. Based on these results, we can approximate the singular values of  $\bar{\mathbf{H}}(\theta^*)$  as  $\lambda_1 \varsigma_1^2 \geq \lambda_2 \varsigma_2^2 \geq \dots \geq \lambda_d \varsigma_d^2$ , implying that  $\bar{\mathbf{H}}(\theta^*)$  becomes much more singular than  $\mathbf{H}(\theta^*)$ . Then the volume of the component set  $\mathcal{W}_{\text{ADAM}}^c$  of  $\mathcal{W}_{\text{ADAM}}$  is  $V(\mathcal{W}_{\text{ADAM}}^c) = \zeta \prod_{i=1}^d \lambda_i$  where  $\zeta = 2d^{-1} (\pi S/h_f^*)^{d/2} g^{-1}(d/2)$  with a gamma function  $g$ . Similarly, we can obtain the volume  $V(\mathcal{W}_{\text{SGD}}^c) = \zeta \prod_{i=1}^d \lambda_i \varsigma_i^2$  of the component set  $\mathcal{W}_{\text{SGD}}^c$  of  $\mathcal{W}_{\text{SGD}}$ . As aforementioned, covariance matrix  $\Sigma_{\theta^*}$  is ill-conditioned and anisotropic near minima and has only a few larger singular values [22, 27], indicating  $\prod_{i=1}^d \varsigma_i^2 \ll 1$ . So  $V(\mathcal{W}_{\text{SGD}}^c)$  is actually much smaller than  $V(\mathcal{W}_{\text{ADAM}}^c)$ . Hence  $\mathcal{W}_{\text{SGD}}$  has larger volume than  $\mathcal{W}_{\text{ADAM}}$  and thus has larger Radon measure  $m(\mathcal{W}_{\text{SGD}})$  than  $m(\mathcal{W}_{\text{ADAM}})$ . Accordingly, SGD has smaller escaping time at the local basin  $\Omega$  than ADAM. Thus, SGD would escape from  $\Omega$  and converges to flat minima whose local basins have large Radon measure, while ADAM will get stuck in  $\Omega$ . Since flat minima with large Radon measure usually locate at the flat or asymmetric basins/valleys and generalize better [12, 17, 30, 31, 51], SGD often enjoys better testing performance. From the above analysis, one can also observe that for SGD, the covariance matrix  $\Sigma_{\theta^*}$  helps increase Radon measure  $m(\mathcal{W}_{\text{SGD}})$  of  $\mathcal{W}_{\text{SGD}}$ . So anisotropic gradient noise helps SGD escape from the local basin but cannot help ADAM's escaping behaviors.

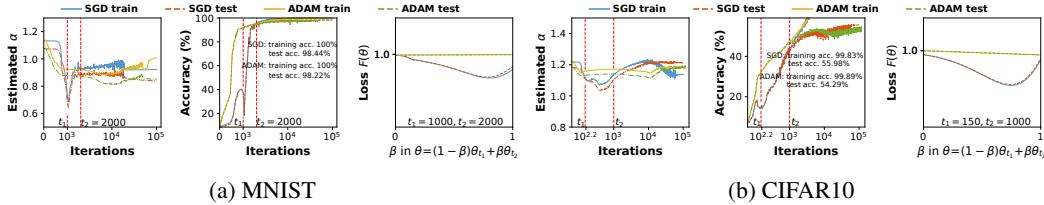


Figure 3: Behaviors illustration of SGD and ADAM on fully connected networks. In both (a) and (b), the left and middle figures respectively report the estimated tail index  $\alpha$  in  $\mathcal{S}\alpha\mathcal{S}$  distribution and classification accuracies; right figures show possible barriers between the solutions  $\theta_{1000}$  and  $\theta_{2000}$  on MNIST, and  $\theta_{150}$  and  $\theta_{1000}$  on CIFAR10, respectively. **Best viewed in  $\times 2$  sized color pdf file.**

## 5 Experiments

In this section, we first investigate the gradient noise in ADAM and SGD, and then show their iteration-based convergence behaviors to testify the implications of our escaping theory. The code is available at <https://panzhous.github.io>.

**Heavy Tails of Gradient Noise.** We respectively use SGD and ADAM to train AlexNet [52] on CIFAR10, and show the statistical behaviors of gradient noise on CIFAR10. To fit the noise via  $\mathcal{S}\alpha\mathcal{S}$  distribution, we consider covariance matrix  $\Sigma_t$  and use the approach in [23, 53] to estimate the tail index  $\alpha$ . Fig. 1 in Sec. 1 and Fig. 4 in Appendix A show that the gradient noise in both SGD and ADAM usually reveal the heavy tails and can be well characterized by  $\mathcal{S}\alpha\mathcal{S}$  distribution. This testifies the heavy tail assumption on the gradient noise in our theories.

**Escaping Behaviors.** We investigate the iteration-based convergence behaviors of SGD and ADAM, including their training and test accuracies and losses and tail index of their gradient noise. For MNIST [54] and CIFAR10 [55], we respectively use nine- and seven-layered fully-connected-networks. Each layer has 512 neurons and contains a linear layer and a ReLu layer. Firstly, the results in the middle figures show that SGD usually has better generalization performance than ADAM-alike adaptive algorithms which is consistent with the results in [12, 17, 18, 30].

Moreover, from the trajectories of the tail index  $\alpha$  and accuracy of SGD on MNIST and CIFAR10 in Fig. 3, one can observe two distinct phases. Specifically, for the first 1000 iterations in MNIST and 150 iterations in CIFAR10, both the training and test accuracies increase tardily, while the tail index parameter  $\alpha$  reduces quickly. This process continues until  $\alpha$  reaches its lowest value. When considering the barrier around inflection point (*e.g.* a barrier between  $\theta_{1000}$  and  $\theta_{2000}$  on MNIST), it seems that the process of SGD has a sudden jump from one basin to another one which leads to a sudden accuracy drop, and then gradually converges. Accordingly, the accuracies are improved quickly. In contrast, one cannot observe similar phenomenon in ADAM. This is because as our theory suggested, SGD is more locally unstable and converges to flatter minima than ADAM, which is caused by the geometry adaptation, exponential gradient average and smaller learning rate in ADAM. All these results are consistent with our theories and also explain the well observed evidences in [12, 17, 30, 31, 51] that SGD usually converges to flat minima which often locate at the flat or asymmetric basins/valleys, while ADAM does not. Because the empirical observations [1, 19–21] show that minima at the flat or asymmetric basins/valleys often generalize better than sharp ones, our empirical and theoretical results can well explain the generalization gap between ADAM-alike algorithms and SGD.

## 6 Conclusion

In this work, we analyzed the generalization performance degeneration of ADAM-alike adaptive algorithms over SGD. By looking into the local convergence behaviors of the Lévy-driven SDEs of these algorithms through analyzing their escaping time, we prove that for the same basin, SGD has smaller escaping time than ADAM and tends to converge to flatter minima whose local basins have larger Radon measure, explaining its better generalization performance. This result is also consistent with the widely observed convergence behaviors of SGD and ADAM in many literatures. Finally our experimental results testify the heavy gradient noise assumption and implications in our theory.

## Broader Impacts

This work theoretically analyzes a fundamental problem in deep learning field, namely the generalization gap between adaptive gradient algorithms and SGD, and reveals the essential reasons for the generalization degeneration of adaptive algorithms. The established theoretical understanding of these algorithms may inspire new algorithms with both fast convergence speed and good generalization performance, which alleviate the need for computational resource and achieve state-of-the-art results. Yet it still needs more efforts to provide more insights to design practical algorithms.

## References

- [1] N. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *Int'l Conf. Learning Representations*, 2017.
- [2] H. He, G. Huang, and Y. Yuan. Asymmetric valleys: Beyond sharp and flat local minima. In *Proc. Conf. Neural Information Processing Systems*, 2019.
- [3] H. Robbins and S. Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- [4] L. Bottou. Stochastic gradient learning in neural networks. *Proceedings of Neuro-Nimes*, 91(8):12, 1991.
- [5] Y. Bengio. Learning deep architectures for AI. *Foundations and trends® in Machine Learning*, 2(1):1–127, 2009.
- [6] G. Hinton, L. Deng, D. Yu, G. Dahl, A. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, and B. Kingsbury. Deep neural networks for acoustic modeling in speech recognition. *IEEE Signal processing magazine*, 29, 2012.
- [7] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
- [8] P. Zhou, Y. Hou, and J. Feng. Deep adversarial subspace clustering. In *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2018.
- [9] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [10] P. Zhou, X. Yuan, H. Xu, S. Yan, and J. Feng. Efficient meta learning via minibatch proximal update. In *Proc. Conf. Neural Information Processing Systems*, 2019.
- [11] P. Zhou, C. Xiong, R. Socher, and S. Hoi. Theory-inspired path-regularized differential network architecture search. In *Proc. Conf. Neural Information Processing Systems*, 2019.
- [12] N. Keskar and R. Socher. Improving generalization performance by switching from Adam to SGD. *arXiv preprint arXiv:1712.07628*, 2017.
- [13] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *J. of Machine Learning Research*, 12(Jul):2121–2159, 2011.
- [14] T. Tieleman and G. Hinton. Lecture 6.5—rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural Networks for Machine Learning*, 2012.
- [15] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *Int'l Conf. Learning Representations*, 2014.
- [16] S. Reddi, S. Kale, and S. Kumar. On the convergence of Adam and beyond. *arXiv preprint arXiv:1904.09237*, 2019.
- [17] A. Wilson, R. Roelofs, M. Stern, N. Srebro, and B. Recht. The marginal value of adaptive gradient methods in machine learning. In *Proc. Conf. Neural Information Processing Systems*, pages 4148–4158, 2017.
- [18] L. Luo, Y. Xiong, Y. Liu, and X. Sun. Adaptive gradient methods with dynamic bound of learning rate. In *Int'l Conf. Learning Representations*, 2019.
- [19] S. Hochreiter and J. Schmidhuber. Flat minima. *Neural Computation*, 9(1):1–42, 1997.
- [20] P. Izmailov, D. Podoprikin, T. Garipov, D. Vetrov, and A. Wilson. Averaging weights leads to wider optima and better generalization. *arXiv preprint arXiv:1803.05407*, 2018.

- [21] H. Li, Z. Xu, G. Taylor, C. Studer, and T. Goldstein. Visualizing the loss landscape of neural nets. In *Proc. Conf. Neural Information Processing Systems*, pages 6389–6399, 2018.
- [22] L. Sagun, U. Evci, V. Guney, Y. Dauphin, and L. Bottou. Empirical analysis of the hessian of over-parametrized neural networks. *arXiv preprint arXiv:1706.04454*, 2017.
- [23] U. Simsekli, L. Sagun, and M. Gurbuzbalaban. A tail-index analysis of stochastic gradient noise in deep neural networks. In *Proc. Int'l Conf. Machine Learning*, 2019.
- [24] P. Levy. Théorie de l'addition des variables aléatoires, gauthier-villars, paris, 1937. *LévyThéorie de l'addition des variables aléatoires1937*, 1954.
- [25] S. Mandt, M. Hoffman, and D. Blei. A variational analysis of stochastic gradient algorithms. In *Proc. Int'l Conf. Machine Learning*, pages 354–363, 2016.
- [26] S. Jastrzebski, Z. Kenton, D. Arpit, N. Ballas, A. Fischer, Y. Bengio, and A. Storkey. Three factors influencing minima in SGD. *arXiv preprint arXiv:1711.04623*, 2017.
- [27] P. Chaudhari and S. Soatto. Stochastic gradient descent performs variational inference, converges to limit cycles for deep networks. In *2018 Information Theory and Applications Workshop*, pages 1–10. IEEE, 2018.
- [28] I. Pavlyukevich. Cooling down Lévy flights. *Journal of Physics A: Mathematical and Theoretical*, 40(41), 2007.
- [29] I. Pavlyukevich. First exit times of solutions of stochastic differential equations driven by multiplicative lévy noise with heavy tails. *Stochastics and Dynamics*, 11(02n03):495–519, 2011.
- [30] S. Merity, N. Keskar, and R. Socher. Regularizing and optimizing LSTM language models. *arXiv preprint arXiv:1708.02182*, 2017.
- [31] I. Loshchilov and F. Hutter. SGDR: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- [32] L. Dinh, R. Pascanu, S. Bengio, and Y. Bengio. Sharp minima can generalize for deep nets. In *Proc. Int'l Conf. Machine Learning*, pages 1019–1028, 2017.
- [33] J. Zhang, S. Karimireddy, A. Veit, S. Kim, S. Reddi, S. Kumar, and S. Sra. Why adam beats sgd for attention models. *arXiv preprint arXiv:1912.03194*, 2019.
- [34] Z. Zhu, J. Wu, B. Yu, L. Wu, and J. Ma. The anisotropic noise in stochastic gradient descent: Its behavior of escaping from minima and regularization effects. In *Proc. Int'l Conf. Machine Learning*, 2019.
- [35] S. Amari. Natural gradient works efficiently in learning. *Neural computation*, 10(2):251–276, 1998.
- [36] P. Imkeller, I. Pavlyukevich, and T. Wetzel. The hierarchy of exit times of lévy-driven langevin equations. *The European Physical Journal Special Topics*, 191(1):211–222, 2010.
- [37] Q. Li, C. Tai, and W. E. Stochastic modified equations and adaptive stochastic gradient algorithms. In *Proc. Int'l Conf. Machine Learning*, pages 2101–2110, 2017.
- [38] L. Simon. *Lectures on geometric measure theory*. The Australian National University, Mathematical Sciences Institute, Centre ..., 1983.
- [39] S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- [40] Pan Zhou, Xiaotong Yuan, and Jiashi Feng. Efficient stochastic gradient hard thresholding. In *Proc. Conf. Neural Information Processing Systems*, 2018.
- [41] R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Proc. Conf. Neural Information Processing Systems*, pages 315–323, 2013.
- [42] P. Zhou, X. Yuan, and J. Feng. New insight into hybrid stochastic gradient descent: Beyond with-replacement sampling and convexity. In *Proc. Conf. Neural Information Processing Systems*, 2018.
- [43] P. Zhou, X. Yuan, and J. Feng. Faster first-order methods for stochastic non-convex optimization on riemannian manifolds. In *Int'l Conf. Artificial Intelligence and Statistics*, 2019.

- [44] P. Zhou and X. Tong. Hybrid stochastic-deterministic minibatch proximal gradient: Less-than-single-pass optimization with nearly optimal generalization. In *Proc. Int'l Conf. Machine Learning*, 2020.
- [45] S. Du, J. Lee, H. Li, L. Wang, and X. Zhai. Gradient descent finds global minima of deep neural networks. In *Proc. Int'l Conf. Machine Learning*, 2019.
- [46] Y. Tian. An analytical formula of population gradient for two-layered relu network and its applications in convergence and critical point analysis. In *Proc. Int'l Conf. Machine Learning*, pages 3404–3413, 2017.
- [47] P. Zhou and J. Feng. Understanding generalization and optimization performance of deep cnns. In *Proc. Int'l Conf. Machine Learning*, 2018.
- [48] P. Zhou and J. Feng. Empirical risk landscape analysis for understanding deep neural networks. In *Int'l Conf. Learning Representations*, 2018.
- [49] L. Wu, C. Ma, and W. E. How sgd selects the global minima in over-parameterized learning: A dynamical stability perspective. In *Proc. Conf. Neural Information Processing Systems*, pages 8279–8288, 2018.
- [50] A. Anandkumar and R. Ge. Efficient approaches for escaping higher order saddle points in non-convex optimization. In *Conf. on Learning Theory*, pages 81–102, 2016.
- [51] Y. Wu and K. He. Group normalization. In *Proc. European Conf. Computer Vision*, pages 3–19, 2018.
- [52] A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional neural networks. In *Proc. Conf. Neural Information Processing Systems*, pages 1097–1105, 2012.
- [53] M. Mohammadi, A. Mohammadpour, and H. Ogata. On estimating the tail index and the spectral measure of multivariate  $\alpha$ -stable distributions. *Metrika*, 78(5):549–561, 2015.
- [54] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient based learning applied to document recognition. *Proceedings of the IEEE*, page 2278–2324, 1998.
- [55] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. 2009.