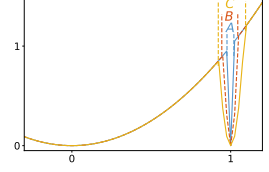


1 We thank all the reviewers for their insightful and encouraging comments, and will update revision to solve the issues.

2 **To Reviewer #2.** Our main goal is to theoretically show the stronger escaping ability of SGD over Adam at the same
 3 basin. For the by-product, i.e. relation between Radon measure and escaping time, we construct $f = \min(x^2, a(x-1)^2)$
 4 with a local basin at $x = 1$ shown in right side. By setting $a = 10^5, 500, 150$, we obtain three
 5 basins A, B and C , where their Radon measures obey $m(A) < m(B) < m(C)$. Then we run SDE
 6 of SGD with initialization $x_0 = 1$ for 2000 iterations, and repeat 1000 times. For A, B, C with
 7 same Lévy noise, their escaping probabilities are 100%, 65.6% and 10.1%, and their average
 8 iterations for successful escaping are 122, 457 and 1898. It confirms our theory: the larger
 9 Radon measure of the basin, the harder to escape. We will update it into revision.



10 The validity of flatness definition should be verified by existing observations. Our definition well explains recent
 11 observation, i.e. good generalization of minima at asymmetric valleys which cannot be explained by existing definitions.

12 **To Reviewer #3.** 1) Our theory also indicates that SGD with momentum (SGD-M) can generalize better than Adam.
 13 Specifically, as SGD-M does not adapt the geometry, it has the following Lévy SDE with $Q_t = I$:

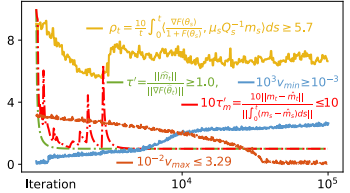
$$d\theta_t = -\mu_t Q_t^{-1} m_t + \varepsilon Q_t^{-1} \Sigma_t dL_t, \quad dm_t = \beta_1 (\nabla F(\theta_t) - m_t), \quad dv_t = \beta_2 ([\nabla f_{S_t}(\theta_t)]^2 - v_t). \quad (9)$$

14 Then we follow Eqn. (6) in manuscript and obtain escaping set $\mathcal{W} = \{y \in \mathbb{R}^d | Q_{\theta^*}^{-1} \Sigma_{\theta^*} y \notin \Omega^{-\varepsilon^\gamma}\}$ of SGD-M, where
 15 $Q_{\theta^*} = I$ and $\Sigma_{\theta^*} = \lim_{\theta_t \rightarrow \theta^*} \Sigma_t$. Since Adam has the same SDE (9) except $Q_t = \text{diag}(\sqrt{\omega_t v_t} + \varepsilon)$ and same escaping set \mathcal{W}
 16 except $Q_{\theta^*} = \lim_{\theta_t \rightarrow \theta^*} Q_t$, we can directly derive the escaping time $\Gamma = \mathcal{O}(\frac{1}{m(\mathcal{W})\Theta(\varepsilon^{-1})})$ of SGD-M with $\Theta(\varepsilon^{-1}) = \frac{2}{\alpha} \varepsilon^\alpha$.

17 As SGD-M and Adam use the same gradient estimation m_t , their gradient noise have the same tail index α and thus the
 18 same factor $\Theta(\varepsilon^{-1})$. For $m(\mathcal{W})$, due to different escaping sets $\mathcal{W}_{\text{SGD-M}}$ of SGD-M and $\mathcal{W}_{\text{Adam}}$ of Adam, $m(\mathcal{W}_{\text{SGD-M}})$ in
 19 SGD-M differs from $m(\mathcal{W}_{\text{Adam}})$ in Adam. By observation, $\mathcal{W}_{\text{SGD-M}}$ is as same as escaping set \mathcal{W}_{SGD} of SGD in Eqn. (6)
 20 in manuscript, as SGD(-M) have no geometry adaptation. Then Sec. 4.2 proves \mathcal{W}_{SGD} has much larger volume than
 21 $\mathcal{W}_{\text{Adam}}$. So $m(\mathcal{W}_{\text{SGD-M}})$ is much larger than $m(\mathcal{W}_{\text{Adam}})$. Thus, SGD-M has much smaller escaping time than Adam at the
 22 same basin, and can better escape sharp minima to flat ones for better generalization. We will update this into revision.

23 2) We follow [20] which analyzes behavior of SGD, and use standard tail index estimation method in [41] as mentioned
 24 in line 334. For learning rate (LR), many works analyze it and conclude: i) an initially large LR helps escape local
 25 minima and accelerates training; ii) decaying LR helps converge to local minima and avoid oscillation. This is testified
 26 by Fig.s 3 in [Jordan, arXiv:1908.01878; R(Kleinberg, arXiv:1802.06175)]. Moreover, Theorem 1 in [R] and analysis in
 27 [Lewkowycz, arXiv:2003.02218] show large LR in SGD help escape. Intuitively, with same basin and gradients, larger
 28 LR gives a larger step and escapes from the basin more easily. These results are consistent with [20] and ours where
 29 $\alpha > 1$ in most cases (see more investigations in [20]). We emphasize that one should focus on the overall variation trend
 30 of α instead of its exact value, as exact value is easily affected by estimation error but variation trend is more robust.

31 **To Reviewer #4.** 1) Assumption 2 often holds as explained in manuscript but is hard to theoretically prove. $\beta_1 \leq \beta_2 \leq 2\beta_1$
 32 holds under Adam's default setting. There always exist constants $v_{\min}, v_{\max}, \tau_m$ and τ
 33 such that $v_{\min} \leq \sqrt{v_t}, i \leq v_{\max}, \frac{\|m_t - \widehat{m}_t\|}{\|\int_0^t (m_s - \widehat{m}_s) ds\|} \leq \tau_m$ and $\frac{\|\widehat{m}_t\|}{\|\nabla F(\theta_t)\|} \geq \tau$ hold, as i) we allow
 34 $v_{\min} = 0$ due to constant ε , ii) $\|\int_0^t (m_s - \widehat{m}_s) ds\| \neq 0$ due to their different definitions and
 35 $\|\widehat{m}_t\| \neq 0$ almost sure when non-convergence. $\int_0^\Gamma \langle \frac{\nabla F(\theta_s)}{1 + F(\theta_s)}, \mu_s Q_s^{-1} m_s \rangle ds \geq 0$ generally
 36 holds, as $\nabla F(\theta_t)$ and its exponential average m_t often share similar directions. Right
 37 figure verifies the validity of these assumptions on 4-layered network (width 100).
 38



39 2) Our results hold for moderately ill-conditioned local basins (ICLBs). Theorem 2 shows that i) after time interval
 40 $v_\varepsilon = \mathcal{O}(\frac{1}{\mu} \ln(\frac{1}{\mu \varepsilon^\delta}))$, noise-free process $\widehat{\theta}_t$ ($\varepsilon = 0$ in SDEs) approaches the minimizer θ^* of a basin Ω , i.e. $\|\widehat{\theta}_t - \theta^*\| \leq \varepsilon^{-\delta}$;
 41 ii) time interval σ_1 between two big jumps ζ (size $\geq \varepsilon^{-\delta}$) is $\sigma_1 = \mathcal{O}(\frac{1}{\varepsilon \alpha \delta})$. Both Theorems 1 and 2 require $v_\varepsilon \leq \sigma_1$
 42 to guarantee small distance of current solution θ_t to θ^* before each big jump. So if μ of ICLBs is larger than $\mathcal{O}(\varepsilon^{\alpha \delta})$ which
 43 is very small as ε in SDE is often small to precisely mimic algorithm behaviors, our results still hold. Moreover, to
 44 obtain result i), we assume the optimization trajectory goes along the eigenvector direction corresponding to μ which is
 45 the worse case and leads to the worst convergence speed. As the measure of one/several eigenvector directions on high
 46 dimension is 0, optimization trajectory cannot always go along the eigenvector direction corresponding to μ . So v_ε
 47 is actually much larger than $\mathcal{O}(\frac{1}{\mu} \ln(\frac{1}{\mu \varepsilon^\delta}))$, largely improving applicability of our theory. We will update it into revision.

48 For extremely ICLBs ($\mu \rightarrow 0$ or $\mu = 0$), Anandkumar (arXiv:1602.05908v1) proved that first-order algorithms cannot
 49 escape from them, which is also the reasons why recent works (e.g. Jin Chi's works) on escaping saddle points do not
 50 discuss extremely ICLBs. Similarly, our theory also does not hold for this case, which accords with the previous works.
 51 Moreover, $\mu \rightarrow 0$ and $\mu = 0$ give asymmetric basins which often generalize well [2, 19] and are not needed to escape.

52 3) Loss around barrier \cap first decreases to the foot of \cap , then increases to climb \cap and finally decreases. Fig.2 shows the
 53 first two phases, as final loss is much smaller than the loss around \cap , indicating the third phase. We will re-plot it.