
Towards Theoretically Understanding Why SGD Generalizes Better Than ADAM in Deep Learning (Supplementary File)

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

This supplementary document contains the technical proofs of convergence results and some additional numerical results of the work entitled “Towards Theoretically Understanding Why SGD Generalizes Better Than ADAM in Deep Learning”. It is structured as follows. In Appendix A, we provide more construction details of the SDE for ADAM and also conduct experiments which show very similar convergence behaviors of ADAM (SGD) and its SDE. Appendix B compares our work with the related work [1, 2] in more details. Appendix C summarizes the notations throughout this document and also provides the auxiliary theories and lemmas for subsequent analysis whose proofs are deferred to Appendix E. Then Appendix D gives the proofs of the main results in Sec. 4, including Theorem 1 which analyzes the escaping time analysis of Lévy-driven SDEs and Theorem 2 which proves the processes with and without Lévy motion are close to each other. Finally, in Appendix E we present the proofs of auxiliary theories and lemmas in Appendix C, including Theorems 3 ~ 4 and Lemmas 1 ~ 3.

A More Discussion of SDE in ADAM

Here we provide more discussion and construction details for the SDE in ADAM. We first investigate the second order moment of the gradient noise in ADAM. Then we introduce the two types of randomness in the SDE of ADAM. Finally, we run experiments to investigate the validity of the constructed SDEs of ADAM and SGD.

A.1 $\mathcal{S}\alpha\mathcal{S}$ -distributed Gradient Noise in ADAM

In the manuscript, we have shown the gradient noise itself to be $\mathcal{S}\alpha\mathcal{S}$ -distributed. Here we further investigate the second-order moment of the gradient noise. From the bottom row of Figure 4, one can observe that (1) both the second-order moment of the gradient noise also reveals heavy tails; (2) compared with Gaussian distribution, $\mathcal{S}\alpha\mathcal{S}$ distribution can better characterize this kind of second-order moment of the gradient noise. All these results demonstrate that the gradient noise in both ADAM and SGD actually satisfies the $\mathcal{S}\alpha\mathcal{S}$ distribution. So the heavy-tailed gradient noise assumptions in our manuscript is very reasonable.

A.2 Randomness in SDE of ADAM

The SDE of ADAM approximates gradient noise m_t via the combination of full gradient and Lévy motion but does not approximate v_t . This SDE should be more accurate than the one which approximates both m_t and the coefficients v_t . So the randomness in the SDE of ADAM comes from the Lévy motion and also v_t caused by sampling a minibatch. But these two types of randomness actually do not depend on each other. Note that as shown in many literatures, *e.g.* [3, 4], SDE allows randomness in coefficients and also enjoys many good properties, such as stability and unique solution. This type of SDE is usually called “SDE with random coefficients”, and usually appears in

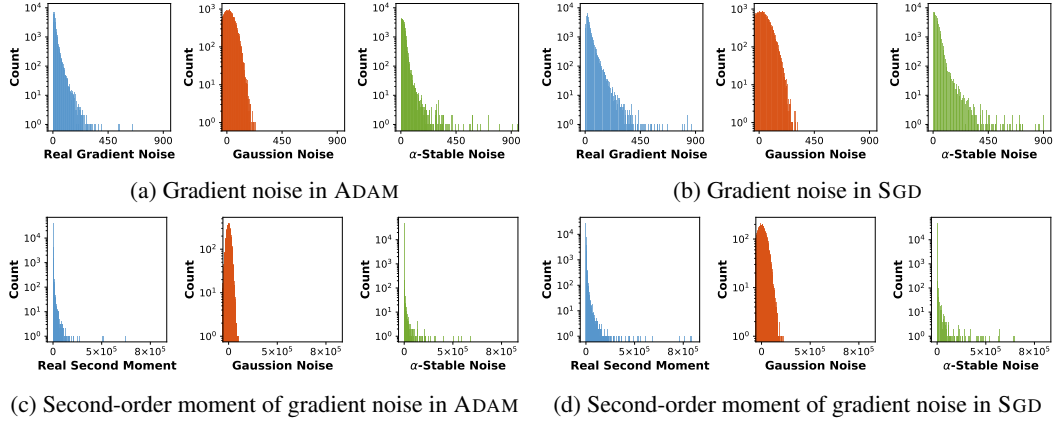


Figure 4: Illustration of gradient noise in ADAM and SGD. The left figures in (a) and (b) are the real gradient noise computed with AlexNet on CIFAR10. Similarly, the left figures in (c) and (d) are the second-order moment of gradient noise computed with AlexNet on CIFAR10. The middle and right figures in (a) ~ (d) are respectively the fitted Gaussian and systemic α -stable noise. By comparison, α -stable noise can better characterize real gradient noise in deep learning.

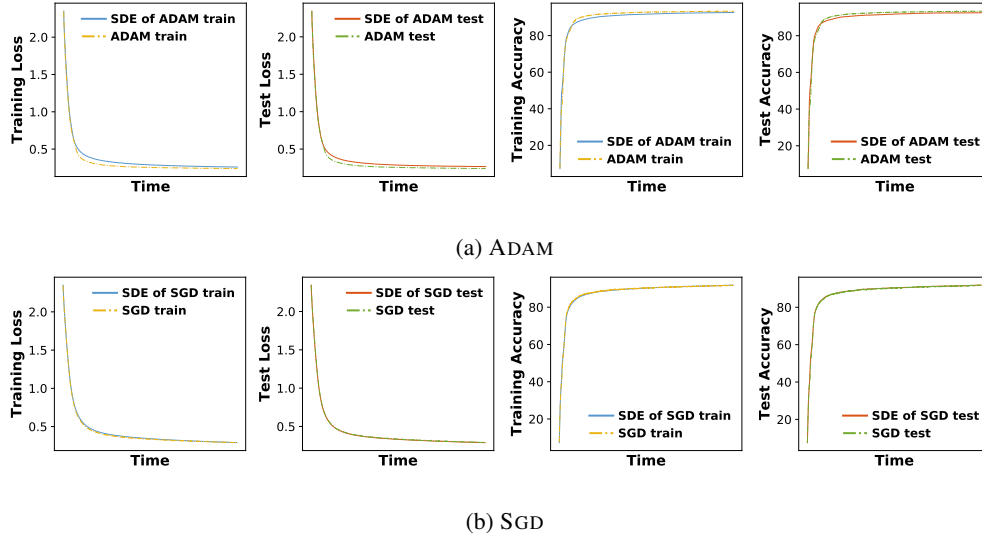


Figure 5: Illustration of convergence trajectories of ADAM, SGD and their SDEs. One can observe that for ADAM, its convergence trajectories are very similar to its SDE, which shows the validity of the SDE construction. Similarly, we can observe the same observations on SGD and its SDE.

stochastic jump systems [5], economics and finance [6, 7], biology [8, 9], mechanics and physics [10], etc. See more details of SDE with random coefficients in [3, 4].

A.3 Convergence Behavior Comparison between Algorithm and Its SDE

Here we conduct experiments on 784-10-10-sized networks and report the convergence behaviors of ADAM (SGD) and its SDE in Fig. 5. Note SDE actually equals to injecting heavy tailed noise into SGD and ADAM that use full gradients. We use a relatively small network since simulating high-dimensional gradient noise u_t and computing the huge covariance matrix Σ_t at each iteration are too computationally expensive to compute. From the convergence trajectories of both ADAM and its SDE in Fig. 5 (a), one can observe that they have very similar convergence behaviors. Similarly, in Fig. 5 (b) we can observe the same observations on SGD and its SDE. So injecting heavy tailed noise into SGD and ADAM that use full gradients leads to similar convergence behaviors to SGD and ADAM that use stochastic gradients. These results well demonstrate the validity of current SDE

construction. Note that here we do not observe jump behaviors, since the networks are very small and may have not very sharp minima. But these results as aforementioned can testify the validity of current SDE construction.

B Comparison to Related Works

Dinh et al. [1] showed flat minimum can become sharp by scaling two layers at the same time. But with this scaling, sharp minimum cannot be arbitrarily flat, as if the eigenvalues of two parameters in the same layer has large ratio, this scaling cannot change this ratio. So flat and sharp minimum are not totally equivalent. Combining the observation in many works that flat minima could achieve better generalization performance than sharper ones, one could conclude that flat minima can generalize well in most case, while sharp minima that can become flat one by linearly scaling two layers also can generalize but other sharp minima cannot. So analyzing the flat and sharp properties is still meaningful. Besides, the flatness in this work is defined on general non-zero Radon measure. If one finds an invariant measure to the scaling in [1], the flatness is also invariant, providing more insights to generalization. So it is promising to explore this invariant measure in the future.

C Notations and Auxiliary Lemmas

C.1 Notations

For analyzing the uniform Lévy-driven SDEs in Eqn. (4) and (5), we first decompose the Lévy process L_t into two components ξ_t and ζ_t , namely

$$L_t = \xi_t + \zeta_t \quad (7)$$

whose characteristic functions are respectively defined as

$$\begin{aligned} \mathbb{E}[e^{i\langle \lambda, \xi_t \rangle}] &= e^{t \int_{\mathbb{R}^d \setminus \{0\}} \zeta \mathbb{I}\{\|\mathbf{y}\| \leq \frac{1}{\varepsilon^\delta}\} \nu(d\mathbf{y})}, \\ \mathbb{E}[e^{i\langle \lambda, \zeta_t \rangle}] &= e^{t \int_{\mathbb{R}^d \setminus \{0\}} \zeta \mathbb{I}\{\|\mathbf{y}\| \geq \frac{1}{\varepsilon^\delta}\} \nu(d\mathbf{y})}, \end{aligned}$$

where $\zeta = e^{i\langle \lambda, \mathbf{y} \rangle} - 1 - i\langle \lambda, \mathbf{y} \rangle \mathbb{I}\{\|\mathbf{y}\| \leq 1\}$, ε (in Eqn. (4) and (5)) and δ are two small constants satisfying $\varepsilon^{-\delta} < 1$ and will be specified later. Define the Lévy measures ν as $\nu(d\mathbf{y}) = \frac{1}{\|\mathbf{y}\|^{1+\alpha}} d\mathbf{y}$. Accordingly, the Lévy measures ν of the stochastic processes ξ and ζ are

$$\nu_\xi = \nu(\mathbf{A} \cap \{0 < \|\mathbf{y}\| \leq \frac{1}{\varepsilon^\delta}\}), \quad \nu_\zeta = \nu(\mathbf{A} \cap \{\|\mathbf{y}\| \geq \frac{1}{\varepsilon^\delta}\}),$$

where $\mathbf{A} \in \mathcal{B}(\mathbb{R}^d)$. In this way, the stochastic process ξ has infinite Lévy measure with support $\{\mathbf{y} \mid 0 < \|\mathbf{y}\| \leq \varepsilon^{-\delta}\}$ and thus makes infinitely many jumps on any time interval. But the jump size does not exceed $\varepsilon^{-\delta}$ and thus is small which actually does not help escape the current local basin. In contrast, the Lévy measure $\nu_\zeta(\cdot)$ of ζ is finite and is computed as

$$\Theta(\varepsilon^{-\delta}) = \int_{\|\mathbf{y}\| \geq \varepsilon^{-\delta}} \nu(d\mathbf{y}) = \int_{\|\mathbf{y}\| \geq \varepsilon^{-\delta}} \frac{d\mathbf{y}}{\|\mathbf{y}\|^{1+\alpha}} = \frac{2}{\alpha} \varepsilon^{\alpha\delta}.$$

So the process ζ is a compound Poisson process with intensity $\Theta(\varepsilon^{-\delta})$ and jumps distributed according to the law of $1/\Theta(\varepsilon^{-\delta})$. Specifically, let $0 = t_1 < t_2 < \dots < t_k < \dots$ denote the times of successive jumps of ζ and J_k denote the jump size at the k -th jump. Then the inner-jump times $\sigma_k = t_k - t_{k-1}$ are i.i.d. exponentially distributed random variables with mean value $\mathbb{E}(\sigma_k) = \frac{1}{\Theta(\varepsilon^{-\delta})}$ and the probability distribution function $\mathbb{P}(\sigma_k \leq x) = 1 - \exp(-x\Theta(\varepsilon^{-\delta}))$. The probability law of J_k is also known explicitly in terms of the Lévy measure ν :

$$\mathbb{P}(J_k \in \mathbf{A}) = \frac{1}{\Theta(\varepsilon^{-\delta})} \nu(\mathbf{A} \cap \{\mathbf{y} \mid \|\mathbf{y}\| \geq \varepsilon^{-\delta}\}), \quad \mathbf{A} \in \mathcal{B}(\mathbb{R}^d).$$

So the main force for escaping the local basin comes from the big jumps in the process ζ which will be rigorous analyzed in the following sections.

Besides, for analysis, we usually need to consider affects of the Lévy motion (noise) L_t to the Lévy-driven SDEs of SGD and ADAM given in Eqn. (4) and (5). So here we define two Lévy-free SDEs which respectively correspond to Eqn. (4) and (5):

$$d\widehat{\boldsymbol{\theta}}_t = \nabla \mathbf{F}(\widehat{\boldsymbol{\theta}}_t), \quad (8)$$

and

$$\begin{cases} d\widehat{\boldsymbol{\theta}}_t = -\mu_t \widehat{\mathbf{Q}}_t^{-1} \widehat{\mathbf{m}}_t, \\ d\widehat{\mathbf{m}}_t = \beta_1 (\nabla \mathbf{F}(\widehat{\boldsymbol{\theta}}_t) - \widehat{\mathbf{m}}_t) \\ d\widehat{\mathbf{v}}_t = \beta_2 (\nabla f_{\mathcal{S}_t}(\widehat{\boldsymbol{\theta}}_t)^2 - \widehat{\mathbf{v}}_t). \end{cases} \quad (9)$$

where $\widehat{\mathbf{Q}}_t = \text{diag}(\sqrt{\widehat{\mathbf{v}}_t} + \epsilon)$. Then by analyzing the distance $\|\widehat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_t\|$ between the processes $\widehat{\boldsymbol{\theta}}_t$ without Lévy motion and $\boldsymbol{\theta}_t$ with Lévy motion, we can well know the effects of the Lévy motion to the escaping behaviors.

C.2 Auxiliary Theories and Lemmas

Theorem 3. *Suppose Assumptions 1 and 2 holds. Then for Lévy-driven SGD SDE (8) with $\widehat{\mathbf{Q}}_t = \mathbf{I}$ and $\beta_2 = 0$, the Lyapunov function $\mathcal{L}(t) = \mathbf{F}(\widehat{\boldsymbol{\theta}}_t) - \mathbf{F}(\boldsymbol{\theta}^*)$ obeys*

$$\mathcal{L}(t) \leq \Delta \exp(-2\mu t)$$

where $\Delta = \mathbf{F}(\widehat{\boldsymbol{\theta}}_0) - \mathbf{F}(\boldsymbol{\theta}^*)$ with the optimum solution $\boldsymbol{\theta}^*$ in the current local basin Ω . The sequence $\{\widehat{\boldsymbol{\theta}}_t\}$ produced by Eqn. (8) obeys

$$\|\widehat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^*\|_2^2 \leq \frac{2\Delta}{\mu} \exp(-2\mu t).$$

See its proof in Appendix E.1.

Theorem 4. *Suppose Assumptions 1 and 2 holds. Assume the sequence $\{(\widehat{\boldsymbol{\theta}}_t, \widehat{\mathbf{m}}_t, \widehat{\mathbf{v}}_t)\}$ are produced by Eqn. (9). Let $\widehat{\mathbf{s}}_t = \frac{h_t}{\mu_t} (\sqrt{\omega_t} \widehat{\mathbf{v}}_t + \epsilon)$ with $h_t = \beta_1$, $\mu_t = (1 - e^{-\beta_1 t})^{-1}$ and $\omega_t = (1 - e^{-\beta_2 t})^{-1}$. We define $\|\mathbf{x}\|_{\mathbf{y}}^2 = \sum_i \mathbf{y}_i \mathbf{x}_i^2$. Then for Lévy-driven ADAM SDEs in Eqn. (9), its Lyapunov function $\mathcal{L}(t) = \mathbf{F}(\widehat{\boldsymbol{\theta}}_t) - \mathbf{F}(\boldsymbol{\theta}^*) + \frac{1}{2} \|\widehat{\mathbf{m}}_t\|_{\widehat{\mathbf{s}}_t^{-1}}$ with the optimum solution $\boldsymbol{\theta}^*$ in the current local basin Ω obeys*

$$\mathcal{L}(t) \leq \Delta \exp\left(-\frac{2\mu\tau}{\beta_1(v_{\max} + \epsilon) + \mu\tau} \left(\beta_1 - \frac{\beta_2}{4}\right) t\right)$$

where $\Delta = \mathbf{F}(\widehat{\boldsymbol{\theta}}_0) - \mathbf{F}(\boldsymbol{\theta}^*)$ due to $\widehat{\mathbf{m}}_0 = \mathbf{0}$. The sequence $\{\widehat{\boldsymbol{\theta}}_t\}$ produced by Eqn. (9) obeys

$$\|\widehat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^*\|_2^2 \leq \frac{2\Delta}{\mu} \exp\left(-\frac{2\mu\tau}{\beta_1(v_{\max} + \epsilon) + \mu\tau} \left(\beta_1 - \frac{\beta_2}{4}\right) t\right).$$

See its proof in Appendix E.2.

Lemma 1. (1) *The process $\boldsymbol{\xi}$ in Eqn. (7) can be decomposed into two processes $\widehat{\boldsymbol{\xi}}$ and linear drift, namely,*

$$\boldsymbol{\xi}_t = \widehat{\boldsymbol{\xi}}_t + \mu_\epsilon t, \quad (10)$$

where $\widehat{\boldsymbol{\xi}}$ is a zero mean Lévy martingale with bounded jumps.

(2) *Let $\delta \in (0, 1)$, $\mu_\epsilon = \mathbb{E}[\boldsymbol{\xi}_1]$ and $T_\epsilon = \epsilon^{-\theta}$ for some $\theta > 0$, $\rho_0 = \rho_0(\delta) = \frac{1-\delta}{4} > 0$ and $\theta_0 = \theta_0(\delta) = \frac{1-\delta}{3} > 0$. Suppose ϵ is sufficient small such that $\Theta(1) \leq \epsilon^{-\frac{1-\delta}{6}}$ and $\epsilon^{-\rho} - 2(C + \Theta(1))\epsilon^{\frac{7}{6}(1-\delta) + \frac{\rho}{2}} \geq 1$ with a constant $C = \left| \int_{0 < u \leq 1} u^2 d\Theta(u) \right| \in (0, +\infty)$. Then for all $\delta \in (0, \delta_0)$, $\theta \in (0, \theta_0)$ there are $p_0 = p_0(\delta) = \frac{\delta}{2}$ and $\epsilon_0 = \epsilon_0(\delta, \rho)$ such that the estimates*

$$\|\epsilon \boldsymbol{\xi}_{T_\epsilon}\| = \epsilon \|\mu_\epsilon\| T_\epsilon < \epsilon^{2\rho} \quad \text{and} \quad \mathbb{P}([\epsilon \boldsymbol{\xi}]_{T_\epsilon}^d \geq \epsilon^\rho) \leq \exp(-\epsilon^{-\rho})$$

hold for all $p \in (0, p_0]$ and $\epsilon \in (0, \epsilon_0]$.

See its proof in Appendix E.3.

Lemma 2. Let $\delta \in (0, 1)$ and $\mathbf{g}_{t \geq 0}^t$ be a bounded adapted càdlàg stochastic process with values in \mathbb{R}^d , $T_\varepsilon = \varepsilon^{-\theta}$, $\theta > 0$. Suppose $\sup_{t \geq 0} \|\mathbf{g}^t\|$ is well bounded. Assume $\rho_0 = \rho_0(\delta) = \frac{1-\delta}{16} > 0$, $\theta_0 = \theta_0(\delta) = \frac{1-\delta}{3} > 0$, $p_0 = \frac{\rho}{2}$. For $\widehat{\boldsymbol{\xi}}_t$ in Eqn. (10), there is $\delta_0 = \delta_0(\delta) > 0$ such that for all $\rho \in (0, \rho_0)$ and $\theta \in (0, \theta_0)$, it holds

$$\mathbb{P} \left(\sup_{0 \leq t \leq T_\varepsilon} \varepsilon \left| \sum_{i=1}^d \int_0^t \mathbf{g}_{s-}^i d\widehat{\boldsymbol{\xi}}_s^i \right| \geq \varepsilon^\rho \right) \leq 2 \exp(-\varepsilon^{-p})$$

for all $p \in (0, p_0]$ and $0 < \varepsilon \leq \varepsilon_0$ with $\varepsilon_0 = \varepsilon_0(\rho)$, where $\widehat{\boldsymbol{\xi}}_s^i$ denotes the i -th entry in $\widehat{\boldsymbol{\xi}}_s$.

See its proof of Appendix E.4.

Lemma 3. Suppose Assumptions 1 and 2 holds. Assume $\delta \in (0, 1)$, $\rho_0 = \rho_0(\delta) = \frac{1-\delta}{16(1+c_1\kappa_1)} > 0$, $\theta_0 = \theta_0(\delta) = \frac{1-\delta}{3} > 0$, $p_0 = \min(\frac{\bar{\rho}(1+c_1\kappa_1)}{2}, p)$, $\frac{1}{c_2} \ln \left(\frac{2\Delta}{\mu\varepsilon^\rho} \right) \leq \varepsilon^{-\theta_0}$ where $\kappa_1 = \ell$ and $c_2 = 2\mu$ in SGD, $\kappa_1 = \frac{c_2\ell}{(v_{\min}+\varepsilon)|\tau_m-1|}$ and $c_3 = \frac{2\mu\tau}{\beta_1(v_{\max}+\varepsilon)+\mu\tau} \left(\beta_1 - \frac{\beta_2}{4} \right)$ in ADAM. Here $c_1 \sim c_3$ are positive constants. For all $\bar{\rho} \in (0, \rho_0)$, $p \in (0, p_0]$, $0 < \varepsilon \leq \varepsilon_0$ with $\varepsilon_0 = \varepsilon_0(\bar{\rho})$, and $\boldsymbol{\theta}_0 = \widehat{\boldsymbol{\theta}}_0$, we have

$$\sup_{\boldsymbol{\theta}_0 \in \Omega} \mathbb{P} \left(\sup_{0 \leq t < \sigma_1} \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}_t\| \geq 2\varepsilon^{\bar{\rho}} \right) \leq 2 \exp(-\varepsilon^{-p/2}), \quad (11)$$

where the sequences $\boldsymbol{\theta}_t$ and $\widehat{\boldsymbol{\theta}}_t$ are respectively produced by Eqn. (5) and (9) in Adam or Eqn. (4) and (8) in RMSPROP and SGD.

See its proof in Appendix E.5.

D Proof of Results in Sec. 4

D.1 Proof of Theorem 1

Proof. Here we first briefly introduce our proof idea. As we proved in Lemma 3, for any $\delta \in (0, 1)$, there exist ρ_0, p_0 and ε_0 such that for all $\bar{\rho} \in (0, \rho_0)$, $p \in (0, p_0]$ and $0 < \varepsilon \leq \varepsilon_0$, we have

$$\sup_{\boldsymbol{\theta}_0 \in \Omega} \mathbb{P} \left(\sup_{0 \leq t < \sigma_1} \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}_t\| \geq 2\varepsilon^{\bar{\rho}} \right) \leq 2 \exp(-\varepsilon^{-p/2}), \quad (12)$$

where the sequences $\boldsymbol{\theta}_t$ and $\widehat{\boldsymbol{\theta}}_t$ share the same initialization $\boldsymbol{\theta}_0 = \widehat{\boldsymbol{\theta}}_0$. Such a result holds for both SGD and Adam. Besides, from Theorems 3 and 4, we know that the sequence $\{\widehat{\boldsymbol{\theta}}_t\}$ produced by Eqn. (8) or (9) (namely, the dynamic systems of SGD and Adam) exponentially converges to the minimum $\boldsymbol{\theta}^*$ of the current local basin Ω . To escape the local basin Ω , there are two possible choices, the small jumps in the process $\boldsymbol{\xi}$ and the big jumps J_k in the process $\boldsymbol{\zeta}$. As the small jumps in the process $\boldsymbol{\xi}$ is well bounded, it is not very likely that these small jumps can help escape the local basin Ω which is verified by Eqn. (12). We will prove this more rigorously later. For the big jumps J , since the expectation jump time $\mathbb{E}(\sigma_1)$ is $1/\Theta(\varepsilon^{-\delta})$, such as $\mathbb{E}(\sigma_1) = \frac{2}{\varepsilon^{\alpha\delta}}$ in the α -stable ($\mathcal{S}\alpha\mathcal{S}$) distribution, $\mathbb{E}(\sigma_1)$ is usually much larger than the necessary time $t = \mathcal{O}(\ln(1/\varepsilon))$ to achieve $\|\widehat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^*\| \leq \varepsilon^\delta$. This means that before the jump time σ_1 the sequence $\widehat{\boldsymbol{\theta}}_t$ is very close to the optimum of Ω and thus $\boldsymbol{\theta}_t$ is very close to the minimum $\boldsymbol{\theta}^*$. In this way, the escaping time Γ of the sequence $\{\boldsymbol{\theta}_t\}$ most likely occurs at the time σ_1 if the big jump εJ_1 in the process $\boldsymbol{\zeta}$ is large. If the jump εJ_1 is small and $\boldsymbol{\theta}_{\sigma_1}$ does not escape Ω , then $\boldsymbol{\theta}_t$ will converge to the minimum $\boldsymbol{\theta}^*$ exponentially and stay in the small neighborhood of $\boldsymbol{\theta}^*$. Accordingly, before the second jump time $t_2 = t_1 + \sigma_2$, $\boldsymbol{\theta}_{t_2}$ will jump. This process will continue during the time interval $[0, t]$. Since for each jump time t_k , $\boldsymbol{\theta}_{t_k-}$ is very close to the optimum $\boldsymbol{\theta}^*$, the big jump size $\varepsilon \mathbf{Q}_{t_k}^{-1} \boldsymbol{\Sigma}_{t_k} J_{t_k} \approx \varepsilon \mathbf{Q}_{\boldsymbol{\theta}^*}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}^*} J_{t_k}$. So we can use $\varepsilon \mathbf{Q}_{t_k}^{-1} \boldsymbol{\Sigma}_{t_k} J_{t_k} \approx \varepsilon \mathbf{Q}_{\boldsymbol{\theta}^*}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}^*} J_{t_k} \notin \Omega$ to judge whether at time t_k , $\boldsymbol{\theta}_{t_k}$ escapes the local basin Ω . The events $\{\varepsilon J_1 \notin \mathcal{W}\} = \{\varepsilon \mathbf{Q}_{\boldsymbol{\theta}^*}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}^*} J_{t_k} \notin \Omega\}, \dots, \{\varepsilon J_{k-1} \notin \mathcal{W}\} = \{\varepsilon \mathbf{Q}_{\boldsymbol{\theta}^*}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}^*} J_{t_{k-1}} \notin \Omega\}, \{\varepsilon J_k \notin \mathcal{W}\} = \{\varepsilon \mathbf{Q}_{\boldsymbol{\theta}^*}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}^*} J_{t_k} \notin \Omega\}$ are independent.

Now we prove the desired results from two aspects, namely establishing upper bound and lower bound of $\mathbb{E} \left[\exp \left(-um(\mathcal{W})\Theta(\varepsilon^{-1})\Gamma \right) \right]$ for any $u > -1$. Before that, we first establish basic inequalities for lower and upper bounds.

Basic inequalities for lower and upper bounds. Since σ_1 is exponentially distributed with the parameter $\Theta(\varepsilon^{-\delta})$, we compute the Laplace transform of $m(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1$ as follows:

$$\begin{aligned} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \right] &= \mathbb{E} \left[\int_0^{+\infty} e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \cdot \Theta(\varepsilon^{-\delta})e^{-\Theta(\varepsilon^{-\delta})\sigma_1} d\sigma_1 \right] \\ &= \frac{\Theta(\varepsilon^{-\delta})}{\Theta(\varepsilon^{-\delta}) + um(\mathcal{W})\Theta(\varepsilon^{-1})} = \frac{1}{1 + ua_\varepsilon}, \end{aligned}$$

where $a_\varepsilon = m(\mathcal{W})\frac{\Theta(\varepsilon^{-1})}{\Theta(\varepsilon^{-\delta})}$ and $\Theta(\varepsilon^{-\delta}) = \Theta(-\varepsilon^\delta)$. Besides, for the probability law of the big jump we have

$$\mathbb{P} \left(\mathcal{Q}_{\theta^*}^{-1} \Sigma_{\theta^*} \varepsilon J_1 \notin \Omega \right) = \mathbb{P} \left(\varepsilon J_1 \in \mathcal{W} \right) = \frac{\nu(\mathcal{W}/\varepsilon)}{\Theta(\varepsilon^{-\delta})}.$$

Since for the Lévy measure, we have $m(\mathcal{W}) = \lim_{u \rightarrow +\infty} \frac{\nu(u\mathcal{W})}{\Theta(u)}$ according to [11]. So for any δ' , there always exists ε such that it holds

$$a_\varepsilon(1 - \delta') \leq \frac{\nu(\mathcal{W}/\varepsilon)}{\Theta(\varepsilon^{-\delta})} = \frac{\nu(\mathcal{W}/\varepsilon)}{\Theta(\varepsilon^{-1})} \frac{\Theta(\varepsilon^{-1})}{\Theta(\varepsilon^{-\delta})} \stackrel{\textcircled{1}}{\approx} m(\mathcal{W}) \frac{\Theta(\varepsilon^{-1})}{\Theta(\varepsilon^{-\delta})} = m(\mathcal{W}) \frac{\Theta(\varepsilon^{-1})}{\Theta(\varepsilon^{-\delta})} \leq a_\varepsilon(1 + \delta'). \quad (13)$$

where $\textcircled{1}$ holds since ε is enough small. Then with the help of the continuity of the function $(\theta, z) \rightarrow \mathcal{Q}_\theta^{-1} \Sigma_\theta z$ both in θ and z . Indeed, for any δ' we can choose $R > 0$ enough large such that for small ε we have

$$\mathbb{P} \left(\|\varepsilon J_1\| > R \right) \leq \frac{\delta' \Theta(\varepsilon^{-1})}{4 \Theta(\varepsilon^{-\delta})}.$$

Further, the function $(\theta, z) \rightarrow \mathcal{Q}_\theta^{-1} \Sigma_\theta z$ is uniformly continuous in z in the ball $\|z\| \leq R$ and is continuous in θ at the optimum θ^* . Following [11], by using the scaling property of the jump measure ν and the fact that the limiting measure m has no atoms we show that uniformly over $\|\theta - \theta^*\| \leq \varepsilon^\gamma$:

$$\begin{cases} \left| \mathbb{P} \left(\mathcal{Q}_\theta^{-1} \Sigma_\theta \varepsilon J_k \notin \Omega^{\pm \varepsilon^\gamma}, \|\varepsilon J_k\| \leq R \right) - \mathbb{P} \left(\mathcal{Q}_{\theta^*}^{-1} \Sigma_{\theta^*} \varepsilon J_k \notin \Omega, \|\varepsilon J_k\| \leq R \right) \right| \leq \frac{\delta' \Theta(\varepsilon^{-1})}{4 \Theta(\varepsilon^{-\delta})}, \\ \left| \mathbb{P} \left(\mathcal{Q}_\theta^{-1} \Sigma_\theta \varepsilon J_k \notin \Omega, \|\varepsilon J_k\| \leq R \right) - \mathbb{P} \left(\mathcal{Q}_{\theta^*}^{-1} \Sigma_{\theta^*} \varepsilon J_k \notin \Omega, \|\varepsilon J_k\| \leq R \right) \right| \leq \frac{\delta' \Theta(\varepsilon^{-1})}{4 \Theta(\varepsilon^{-\delta})}, \end{cases} \quad (14)$$

At the same time, we also can establish

$$\begin{aligned} &\mathbb{P} \left(\mathcal{Q}_{\theta^*}^{-1} \Sigma_{\theta^*} \varepsilon J_k \notin \Omega \right) - \mathbb{P} \left(\mathcal{Q}_{\theta^*}^{-1} \Sigma_{\theta^*} \varepsilon J_k \notin \Omega, \|\varepsilon J_k\| \leq R \right) \\ &= \mathbb{P} \left(\mathcal{Q}_{\theta^*}^{-1} \Sigma_{\theta^*} \varepsilon J_k \notin \Omega \right) - \mathbb{P} \left(\|\varepsilon J_k\| \leq R \mid \mathcal{Q}_{\theta^*}^{-1} \Sigma_{\theta^*} \varepsilon J_k \notin \Omega \right) \mathbb{P} \left(\mathcal{Q}_{\theta^*}^{-1} \Sigma_{\theta^*} \varepsilon J_k \notin \Omega \right) \\ &= \mathbb{P} \left(\mathcal{Q}_{\theta^*}^{-1} \Sigma_{\theta^*} \varepsilon J_k \notin \Omega \right) \left(1 - \mathbb{P} \left(\|\varepsilon J_k\| \leq R \mid \mathcal{Q}_{\theta^*}^{-1} \Sigma_{\theta^*} \varepsilon J_k \notin \Omega \right) \right) \\ &= \mathbb{P} \left(\mathcal{Q}_{\theta^*}^{-1} \Sigma_{\theta^*} \varepsilon J_k \notin \Omega \right) \mathbb{P} \left(\|\varepsilon J_k\| > R \mid \mathcal{Q}_{\theta^*}^{-1} \Sigma_{\theta^*} \varepsilon J_k \notin \Omega \right) \leq \mathbb{P} \left(\|\varepsilon J_k\| > R \right) \leq \frac{\delta' \Theta(\varepsilon^{-1})}{4 \Theta(\varepsilon^{-\delta})}. \end{aligned} \quad (15)$$

Upper bound of $\mathbb{E} \left[\exp \left(-um(\mathcal{W})\Theta(\varepsilon^{-1})\Gamma \right) \right]$. In this part, we consider both the big jumps in the process ζ and the small jumps in the process ξ which may escape the local minimum θ^* . Instead of estimate the escaping time Γ from Ω , we first estimate the escaping time $\tilde{\Xi}$ from $\Omega^{-\bar{\rho}}$. Here we define the inner part of Ω as $\Omega^{-\bar{\rho}} = \{y \in \Omega \mid \text{dis}(\partial\Omega, y) \geq \bar{\rho}\}$ and the outer $\bar{\rho}$ -neighborhood of Ω as $\Omega^{+\bar{\rho}} = \{y \mid \text{dis}(\partial\Omega, y) \geq \bar{\rho}\}$. Then by setting $\bar{\rho} \downarrow 0$, we can use $\tilde{\Xi}$ to estimate Γ well. Let $\bar{\rho} = \varepsilon^\gamma$ where γ is a constant such that the results of Lemmas 1~3 holds. Here we suppose the initial point $\theta_0 \in \Omega^{-2\varepsilon^\gamma}$.

Step 1. In this step we give the formulation of the upper bound of $\mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\Gamma} \right]$. For any $u > -1$, we can compute the formula of the total probability as follows

$$\mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\tilde{\Xi}} \right] \leq \sum_{k=1}^{+\infty} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})t_k} \mathbb{I} \left\{ \tilde{\Xi} = t_k \right\} + \text{Res}_k \right],$$

where

$$\text{Res}_k \leq \begin{cases} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})t_k} \mathbb{I} \left\{ \tilde{\Xi} \in (t_{k-1}, t_k) \right\} \right], & \text{if } u \in (-1, 0] \\ \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})t_{k-1}} \mathbb{I} \left\{ \tilde{\Xi} \in (t_{k-1}, t_k) \right\} \right], & \text{if } u \in (0, +\infty). \end{cases}$$

Step 2. In this step we specifically upper bounds the first term $\sum_{k=1}^{+\infty} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})t_k} \mathbb{I} \left\{ \tilde{\Xi} = t_k \right\} \right]$. For $k \geq 1$, we can use the strong Markov property and obtain

$$\begin{aligned} & \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})t_k} \mathbb{I} \left\{ \tilde{\Xi} = t_k \right\} \right] = \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})t_k} \mathbb{I} \left\{ \boldsymbol{\theta}_t \in \Omega^{-\varepsilon^\gamma}, t \in [0, t_k), \boldsymbol{\theta}_{t_k} \notin \Omega^{-\varepsilon^\gamma} \right\} \right] \\ &= \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_k} \mathbb{I} \left\{ \boldsymbol{\theta}_{t+t_{k-1}} \in \Omega^{-\varepsilon^\gamma}, t \in [0, \sigma_k) \right\} \right] \mathbb{I} \left\{ \boldsymbol{\theta}_{t_k} \notin \Omega^{-\varepsilon^\gamma} \right\} \\ & \quad \cdot \prod_{i=1}^{k-1} e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_i} \mathbb{I} \left\{ \boldsymbol{\theta}_{t+t_{i-1}} \in \Omega^{-\varepsilon^\gamma}, t \in [0, \sigma_i) \right\} \\ &\leq \sup_{\boldsymbol{\theta}_0 \in \Omega^{-2\varepsilon^\gamma}} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \mathbb{I} \left\{ \boldsymbol{\theta}_t \in \Omega^{-\varepsilon^\gamma}, t \in [0, \sigma_1) \right\} \right] \mathbb{I} \left\{ \boldsymbol{\theta}_{\sigma_1} \notin \Omega^{-\varepsilon^\gamma} \right\} \\ & \quad \cdot \sup_{\boldsymbol{\theta}_0 \in \Omega^{-2\varepsilon^\gamma}} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \mathbb{I} \left\{ \boldsymbol{\theta}_t \in \Omega^{-\varepsilon^\gamma}, t \in [0, \sigma_1) \right\} \right]^{k-1}. \end{aligned}$$

Recall $\bar{\rho} = \varepsilon^\gamma$ where γ is a constant such that the results of Lemmas 1~3 holds. The escaping from the basin $\Omega^{-\varepsilon^\gamma}$ with a big jump εJ_1 occurs when $\mathcal{Q}_{\sigma_1 - \Sigma_{\sigma_1 - \varepsilon} J_1}^{-1} \in \Omega^{-\varepsilon^\gamma}$. Furthermore, $\sup_{0 \leq t < \sigma_1} \|\boldsymbol{\theta}_t - \hat{\boldsymbol{\theta}}_t\| \leq \frac{1}{2}\varepsilon^\gamma$ with probability exponentially close to 1 (verified by Lemma 3). Meanwhile $\sigma_1 = \frac{2}{\alpha}\varepsilon^{\alpha\delta}$ in the α -stable ($\mathcal{S}\alpha\mathcal{S}$) distribution is much larger than $v_\varepsilon = \mathcal{O}(\ln(1/\varepsilon))$ with sufficient small ε , $\hat{\boldsymbol{\theta}}_t$ reaches a $\frac{1}{2}\varepsilon^\gamma$ -neighborhood of the optimum $\boldsymbol{\theta}^*$ which only requires time v_ε . So this actually means $\sup_{0 \leq t < \sigma_1} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}^*\| \leq \varepsilon^\gamma$. In this way, to obtain the final upper bound results, we only need to estimate the escaping probability $\mathbb{P}(\mathcal{Q}_{\boldsymbol{\theta}^*}^{-1} \Sigma_{\boldsymbol{\theta}^*} \varepsilon J_1 \in \Omega^{-\varepsilon^\gamma})$ and $\mathbb{P}(\mathcal{Q}_{\boldsymbol{\theta}^*}^{-1} \Sigma_{\boldsymbol{\theta}^*} \varepsilon J_1 \notin \Omega^{-\varepsilon^\gamma})$ uniformly over $\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| \leq \varepsilon^\gamma$. Then we first give two important inequalities which will used to bound each component later:

$$\begin{aligned} & \sup_{\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| \leq \varepsilon^\gamma} \mathbb{P} \left(\mathcal{Q}_{\boldsymbol{\theta}^*}^{-1} \Sigma_{\boldsymbol{\theta}^*} \varepsilon J_k \notin \Omega^{-\varepsilon^\gamma} \right) \\ &= \sup_{\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| \leq \varepsilon^\gamma} \mathbb{P} \left(\mathcal{Q}_{\boldsymbol{\theta}^*}^{-1} \Sigma_{\boldsymbol{\theta}^*} \varepsilon J_k \notin \Omega^{-\varepsilon^\gamma}, \|\varepsilon J_k\| \leq R \right) + \mathbb{P} \left(\mathcal{Q}_{\boldsymbol{\theta}^*}^{-1} \Sigma_{\boldsymbol{\theta}^*} \varepsilon J_k \notin \Omega^{-\varepsilon^\gamma}, \|\varepsilon J_k\| > R \right) \\ &\stackrel{\textcircled{1}}{\geq} \mathbb{P} \left(\mathcal{Q}_{\boldsymbol{\theta}^*}^{-1} \Sigma_{\boldsymbol{\theta}^*} \varepsilon J_k \notin \Omega^{-\varepsilon^\gamma}, \|\varepsilon J_k\| \leq R \right) - \frac{\delta' \Theta(\varepsilon^{-1})}{4 \Theta(\varepsilon^{-\delta})} + \mathbb{P} \left(\mathcal{Q}_{\boldsymbol{\theta}^*}^{-1} \Sigma_{\boldsymbol{\theta}^*} \varepsilon J_k \notin \Omega^{-\varepsilon^\gamma}, \|\varepsilon J_k\| > R \right) \\ &\geq \mathbb{P} \left(\mathcal{Q}_{\boldsymbol{\theta}^*}^{-1} \Sigma_{\boldsymbol{\theta}^*} \varepsilon J_k \notin \Omega^{-\varepsilon^\gamma}, \|\varepsilon J_k\| \leq R \right) - \frac{\delta' \Theta(\varepsilon^{-1})}{4 \Theta(\varepsilon^{-\delta})} \\ &\stackrel{\textcircled{2}}{\geq} \mathbb{P} \left(\mathcal{Q}_{\boldsymbol{\theta}^*}^{-1} \Sigma_{\boldsymbol{\theta}^*} \varepsilon J_k \notin \Omega^{-\varepsilon^\gamma} \right) - \frac{\delta' \Theta(\varepsilon^{-1})}{2 \Theta(\varepsilon^{-\delta})} \\ &\stackrel{\textcircled{3}}{\geq} m(\mathcal{W}) \left(1 - \delta' - \frac{\delta'}{2m(\mathcal{W})} \right) \frac{\Theta(\varepsilon^{-1})}{\Theta(\varepsilon^{-\delta})} \stackrel{\textcircled{4}}{\geq} m(\mathcal{W})(1 - \rho) \frac{\Theta(\varepsilon^{-1})}{\Theta(\varepsilon^{-\delta})}, \end{aligned}$$

where $\textcircled{1}$ uses the result in Eqn. (14), $\textcircled{2}$ uses Eqn. (15), $\textcircled{3}$ uses Eqn. (13), and in $\textcircled{4}$ we set δ' enough small such that $\rho \geq \delta' + \frac{\delta'}{2m(\mathcal{W})}$. So in this way, for any ρ we choose $\delta' > 0$ small enough to lower bound $\sup_{\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| \leq \varepsilon^\gamma} \mathbb{P}(\mathcal{Q}_{\boldsymbol{\theta}^*}^{-1} \Sigma_{\boldsymbol{\theta}^*} \varepsilon J_k \in \Omega^{-\varepsilon^\gamma})$ as follows:

$$\sup_{\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| \leq \varepsilon^\gamma} \mathbb{P} \left(\mathcal{Q}_{\boldsymbol{\theta}^*}^{-1} \Sigma_{\boldsymbol{\theta}^*} \varepsilon J_k \in \Omega^{-\varepsilon^\gamma} \right) = 1 - \sup_{\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| \leq \varepsilon^\gamma} \mathbb{P} \left(\mathcal{Q}_{\boldsymbol{\theta}^*}^{-1} \Sigma_{\boldsymbol{\theta}^*} \varepsilon J_k \notin \Omega^{-\varepsilon^\gamma} \right) \geq 1 - a_\varepsilon(1 - \rho).$$

Similarly, we only need to upper bound the remaining term $\sup_{\|\boldsymbol{\theta}-\boldsymbol{\theta}^*\|\leq\varepsilon^\gamma} \mathbb{P}(\mathbf{Q}_{\boldsymbol{\theta}}^{-1}\boldsymbol{\Sigma}_{\boldsymbol{\theta}}\varepsilon J_k \notin \boldsymbol{\Omega})$ as follows:

$$\begin{aligned}
& \sup_{\|\boldsymbol{\theta}-\boldsymbol{\theta}^*\|\leq\varepsilon^\gamma} \mathbb{P}(\mathbf{Q}_{\boldsymbol{\theta}}^{-1}\boldsymbol{\Sigma}_{\boldsymbol{\theta}}\varepsilon J_k \notin \boldsymbol{\Omega}^{-\varepsilon^\gamma}) \\
&= \sup_{\|\boldsymbol{\theta}-\boldsymbol{\theta}^*\|\leq\varepsilon^\gamma} \mathbb{P}(\mathbf{Q}_{\boldsymbol{\theta}}^{-1}\boldsymbol{\Sigma}_{\boldsymbol{\theta}}\varepsilon J_k \notin \boldsymbol{\Omega}^{-\varepsilon^\gamma}, \|\varepsilon J_k\| \leq R) + \mathbb{P}(\mathbf{Q}_{\boldsymbol{\theta}}^{-1}\boldsymbol{\Sigma}_{\boldsymbol{\theta}}\varepsilon J_k \notin \boldsymbol{\Omega}^{-\varepsilon^\gamma}, \|\varepsilon J_k\| > R) \\
&\stackrel{\textcircled{1}}{\leq} \sup_{\|\boldsymbol{\theta}-\boldsymbol{\theta}^*\|\leq\varepsilon^\gamma} \mathbb{P}(\mathbf{Q}_{\boldsymbol{\theta}}^{-1}\boldsymbol{\Sigma}_{\boldsymbol{\theta}}\varepsilon J_k \notin \boldsymbol{\Omega}^{-\varepsilon^\gamma}, \|\varepsilon J_k\| \leq R) + \frac{\delta' \Theta(\varepsilon^{-1})}{4 \Theta(\varepsilon^{-\delta})} \\
&\stackrel{\textcircled{2}}{\leq} \sup_{\|\boldsymbol{\theta}-\boldsymbol{\theta}^*\|\leq\varepsilon^\gamma} \mathbb{P}(\mathbf{Q}_{\boldsymbol{\theta}^*}^{-1}\boldsymbol{\Sigma}_{\boldsymbol{\theta}^*}\varepsilon J_k \notin \boldsymbol{\Omega}, \|\varepsilon J_k\| \leq R) + \frac{\delta' \Theta(\varepsilon^{-1})}{2 \Theta(\varepsilon^{-\delta})} \\
&\leq \sup_{\|\boldsymbol{\theta}-\boldsymbol{\theta}^*\|\leq\varepsilon^\gamma} \mathbb{P}(\mathbf{Q}_{\boldsymbol{\theta}^*}^{-1}\boldsymbol{\Sigma}_{\boldsymbol{\theta}^*}\varepsilon J_k \notin \boldsymbol{\Omega}) + \frac{\delta' \Theta(\varepsilon^{-1})}{2 \Theta(\varepsilon^{-\delta})} \\
&\stackrel{\textcircled{3}}{\leq} m(\mathcal{W}) \left(1 + \delta' + \frac{\delta'}{2m(\mathcal{W})}\right) \frac{\Theta(\varepsilon^{-1})}{\Theta(\varepsilon^{-\delta})} \leq m(\mathcal{W})(1 + \rho/3) \frac{\Theta(\varepsilon^{-1})}{\varepsilon^{-\delta}} = a_\varepsilon(1 + \rho/3),
\end{aligned}$$

where $\textcircled{1}$ uses $\mathbb{P}(\mathbf{Q}_{\boldsymbol{\theta}}^{-1}\boldsymbol{\Sigma}_{\boldsymbol{\theta}}\varepsilon J_k \notin \boldsymbol{\Omega}^{-\varepsilon^\gamma}, \|\varepsilon J_k\| > R) \leq \mathbb{P}(\|\varepsilon J_k\| > R) \leq \frac{\delta' \Theta(\varepsilon^{-1})}{4 \Theta(\varepsilon^{-\delta})}$, $\textcircled{2}$ uses the result in Eqn. (14), $\textcircled{3}$ uses Eqn. (15), and $\textcircled{3}$ uses Eqn. (13).

Next, for any $\rho > 0$ and ε we can obtain the Laplace transforms for any $u > -1$ as follows:

$$\begin{aligned}
& \sup_{\boldsymbol{\theta}_0 \in \boldsymbol{\Omega}^{-2\varepsilon^\gamma}} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \mathbb{I} \left\{ \boldsymbol{\theta}_t \in \boldsymbol{\Omega}^{-\varepsilon^\gamma}, t \in [0, \sigma_1] \right\} \right] \\
&\leq [1 - a_\varepsilon(1 + \rho)] \mathbb{E} \left[\int_0^{+\infty} e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \cdot \Theta(\varepsilon^{-\delta}) e^{-\Theta(\varepsilon^{-\delta})\sigma_1} \mathbf{d}\sigma_1 \right] \quad (16) \\
&= \frac{1 - a_\varepsilon(1 - \rho)}{1 + ua_\varepsilon}.
\end{aligned}$$

and

$$\begin{aligned}
& \sup_{\boldsymbol{\theta}_0 \in \boldsymbol{\Omega}^{-2\varepsilon^\gamma}} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \mathbb{I} \left\{ \boldsymbol{\theta}_t \in \boldsymbol{\Omega}^{-\varepsilon^\gamma}, t \in [0, \sigma_1] \right\} \mathbb{I} \left\{ \boldsymbol{\theta}_{\sigma_1} \notin \boldsymbol{\Omega}^{-\varepsilon^\gamma} \right\} \right] \\
&\leq a_\varepsilon \left(1 + \frac{\rho}{3}\right) \mathbb{E} \left[\int_0^{+\infty} e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \cdot \Theta(\varepsilon^{-\delta}) e^{-\Theta(\varepsilon^{-\delta})\sigma_1} \mathbf{d}\sigma_1 \right] = \frac{a_\varepsilon(1 - \rho/3)}{1 + ua_\varepsilon}.
\end{aligned}$$

Here we summarize the above results such that we can upper bound the first term $\sum_{k=1}^{+\infty} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})t_k} \mathbb{I} \left\{ \Gamma = t_k \right\} \right]$:

$$\begin{aligned}
\mathcal{R}_1 &= \sum_{k=1}^{+\infty} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})t_k} \mathbb{I} \left\{ \Gamma = t_k \right\} \right] \leq \frac{a_\varepsilon(1 + \rho/3)}{1 + ua_\varepsilon} \sum_{k=1}^{+\infty} \left(\frac{1 - a_\varepsilon(1 - \rho)}{1 + ua_\varepsilon} \right)^{k-1} \\
&\leq \frac{a_\varepsilon(1 + \rho/3)}{1 + ua_\varepsilon} \sum_{k=0}^{+\infty} \left(\frac{1 - a_\varepsilon(1 - \rho)}{1 + ua_\varepsilon} \right)^{k-1} = \frac{1 + \rho/3}{1 + u - \rho}.
\end{aligned}$$

Step 3. In this step we specifically upper bounds the second term $\sum_{k=1}^{+\infty} \mathbb{E} [\text{Res}_k]$. Specifically, we establish upper bound for each $\mathbb{E} [\text{Res}_k]$ as follows. We first consider the case where $k = 1$:

$$\begin{aligned}
\text{Res}_1 &\leq \begin{cases} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})t_1} \mathbb{I} \left\{ \Gamma \in (0, t_1) \right\} \right], & \text{if } u \in (-1, 0] \\ \mathbb{E} \left[\mathbb{I} \left\{ \Gamma \in (0, t_1) \right\} \right], & \text{if } u \in (0, +\infty). \end{cases} \\
&= \begin{cases} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \mathbb{I} \left\{ \exists t \in (0, \sigma_1) : \boldsymbol{\theta}_t \notin \boldsymbol{\Omega}^{-\varepsilon^\gamma} \right\} \right], & \text{if } u \in (-1, 0] \\ \mathbb{E} \left[\mathbb{I} \left\{ \exists t \in (0, \sigma_1) : \boldsymbol{\theta}_t \notin \boldsymbol{\Omega}^{-\varepsilon^\gamma} \right\} \right], & \text{if } u \in (0, +\infty). \end{cases} \\
&\leq \begin{cases} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \sup_{\boldsymbol{\theta}_0 \in \boldsymbol{\Omega}^{-2\varepsilon^\gamma}} \mathbb{I} \left\{ \exists t \in (0, \sigma_1) : \boldsymbol{\theta}_t \notin \boldsymbol{\Omega}^{-\varepsilon^\gamma} \right\} \right], & \text{if } u \in (-1, 0] \\ \mathbb{E} \left[\sup_{\boldsymbol{\theta}_0 \in \boldsymbol{\Omega}^{-2\varepsilon^\gamma}} \mathbb{I} \left\{ \exists t \in (0, \sigma_1) : \boldsymbol{\theta}_t \notin \boldsymbol{\Omega}^{-\varepsilon^\gamma} \right\} \right], & \text{if } u \in (0, +\infty). \end{cases}
\end{aligned}$$

For $k \geq 2$, it needs more efforts to be upper bounded:

$$\begin{aligned}
& \text{Res}_k \\
& \leq \begin{cases} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})t_k} \mathbb{I} \{ \Gamma \in (t_{k-1}, t_k) \} \right], & \text{if } u \in (-1, 0] \\ \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})t_{k-1}} \mathbb{I} \{ \Gamma \in (t_{k-1}, t_k) \} \right], & \text{if } u \in (0, +\infty). \end{cases} \\
& = \begin{cases} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})t_k} \mathbb{I} \{ t \in [0, t_{k-1}] : \boldsymbol{\theta}_t \in \Omega^{-\varepsilon^\gamma} \} \mathbb{I} \{ \exists t \in (t_{k-1}, t_k) : \boldsymbol{\theta}_t \notin \Omega^{-\varepsilon^\gamma} \} \right], & \text{if } u \in (-1, 0] \\ \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})t_{k-1}} \mathbb{I} \{ t \in [0, t_{k-1}] : \boldsymbol{\theta}_t \in \Omega^{-\varepsilon^\gamma} \} \mathbb{I} \{ \exists t \in (t_{k-1}, t_k) : \boldsymbol{\theta}_t \notin \Omega^{-\varepsilon^\gamma} \} \right], & \text{if } u \in (0, +\infty). \end{cases}
\end{aligned}$$

In this case, for all $u > 0$ we can upper bound Res_k as

$$\begin{aligned}
\text{Res}_k & \leq \left[\mathbb{E} \left[x \sup_{\boldsymbol{\theta}_0 \in \Omega^{-2\varepsilon^\gamma}} \mathbb{I} \{ t \in [0, \sigma_1] : \boldsymbol{\theta}_t \in \Omega^{-\varepsilon^\gamma} \} \right] \right]^{k-2} \\
& \quad \mathbb{E} \left[x \sup_{\boldsymbol{\theta}_0 \in \Omega^{-2\varepsilon^\gamma}} \mathbb{I} \{ t \in [0, \sigma_1] : \boldsymbol{\theta}_t \in \Omega^{-\varepsilon^\gamma} \} \mathbb{I} \{ \exists t \in (0, \sigma_1) : \boldsymbol{\theta}_t \notin \Omega^{-\varepsilon^\gamma} \} \right].
\end{aligned}$$

where $x = e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1}$. Let the event $E = \{\sup_{0 \leq t < \sigma_1} \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}_t\| \leq \varepsilon^\gamma\}$. Now we bound each term in the above inequalities:

$$\begin{aligned}
& \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \sup_{\boldsymbol{\theta}_0 \in \Omega^{-2\varepsilon^\gamma}} \mathbb{I} \{ \exists t \in (0, \sigma_1) : \boldsymbol{\theta}_t \notin \Omega^{-\varepsilon^\gamma} \} \right] \\
& = \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \sup_{\boldsymbol{\theta}_0 \in \Omega^{-2\varepsilon^\gamma}} \mathbb{I} \{ \exists t \in (0, \sigma_1) : \boldsymbol{\theta}_t \notin \Omega^{-\varepsilon^\gamma} \} (\mathbb{I}\{E\} + \mathbb{I}\{E^c\}) \right] \\
& \leq \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \sup_{\boldsymbol{\theta}_0 \in \Omega^{-2\varepsilon^\gamma}} \mathbb{I} \{ \exists t \in (0, \sigma_1) : \boldsymbol{\theta}_t \notin \Omega^{-\varepsilon^\gamma} \} \mathbb{I}\{E^c\} \right] \tag{17} \\
& \stackrel{\textcircled{1}}{\leq} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \exp(-\varepsilon^{-p}) \right] = \frac{\Theta(\varepsilon^{-\delta})}{\Theta(\varepsilon^{-\delta}) + um(\mathcal{W})\Theta(\varepsilon^{-1})} \cdot 2 \exp(-\varepsilon^{-p}) \\
& = \frac{1}{1 + ua_\varepsilon} \exp(-\varepsilon^{-p}) \stackrel{\textcircled{2}}{\leq} \frac{\rho/3}{1 + u - \rho},
\end{aligned}$$

where $\textcircled{1}$ uses the fact that $\sup_{\boldsymbol{\theta}_0 \in \Omega^{-2\varepsilon^\gamma}} \mathbb{I} \{ \exists t \in (0, \sigma_1) : \boldsymbol{\theta}_t \notin \Omega \} \leq 1$ and the sequence $\widehat{\boldsymbol{\theta}}_t$ obeys $\Omega^{-2\varepsilon^\gamma}$ due to $\boldsymbol{\theta}_0 \in \Omega^{-2\varepsilon^\gamma}$ and the results in Lemma 3:

$$\sup_{\boldsymbol{\theta}_0 \in \Omega} \mathbb{P} \left(\sup_{0 \leq t < \sigma_1} \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}_t\| \geq \varepsilon^\gamma \right) \leq 2 \exp(-\varepsilon^{-p}),$$

where the sequences $\boldsymbol{\theta}_t$ and $\widehat{\boldsymbol{\theta}}_t$ share the same initialization $\boldsymbol{\theta}_0 = \widehat{\boldsymbol{\theta}}_0$. In $\textcircled{2}$ we set ε small enough such that $2 \exp(-\varepsilon^{-p}) \leq \frac{\rho/3}{1+u-\rho}$. Similarly, we can upper bound

$$\begin{aligned}
\mathbb{E} \left[\sup_{\boldsymbol{\theta}_0 \in \Omega^{-2\varepsilon^\gamma}} \mathbb{I} \{ \exists t \in (0, t_1) : \boldsymbol{\theta}_t \notin \Omega^{-\varepsilon^\gamma} \} \right] & \leq \mathbb{E} \left[\sup_{\boldsymbol{\theta}_0 \in \Omega^{-2\varepsilon^\gamma}} \mathbb{I} \{ \exists t \in (0, t_1) : \boldsymbol{\theta}_t \notin \Omega^{-\varepsilon^\gamma} \} (\mathbb{I}\{E\} + \mathbb{I}\{E^c\}) \right] \\
& \leq \exp(-\varepsilon^{-p}) \leq \frac{\rho/3}{1 + u - \rho}.
\end{aligned} \tag{18}$$

Since p is much smaller than 1, then we have for $k = 2, \dots, k$

$$\begin{aligned}
\text{Res}_k & \leq \left[\mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \sup_{\boldsymbol{\theta}_0 \in \Omega^{-2\varepsilon^\gamma}} \mathbb{I} \{ t \in [0, \sigma_1] : \boldsymbol{\theta}_t \in \Omega^{-\varepsilon^\gamma} \} \right] \right]^{k-2} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \right. \\
& \quad \left. \sup_{\boldsymbol{\theta}_0 \in \Omega^{-2\varepsilon^\gamma}} \mathbb{I} \{ t \in [0, \sigma_1] : \boldsymbol{\theta}_t \in \Omega^{-\varepsilon^\gamma} \} \mathbb{I} \{ \exists t \in (0, \sigma_1) : \boldsymbol{\theta}_t \notin \Omega^{-\varepsilon^\gamma} \} \right] \leq \left[\frac{1 - a_\varepsilon(1 - \rho)}{1 + ua_\varepsilon} \right]^{k-2} \frac{a_\varepsilon(1 + \rho/3)}{1 + ua_\varepsilon}.
\end{aligned}$$

where we use the above results, namely, $\sup_{\theta_0 \in \Omega^{-2\varepsilon\gamma}} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \mathbb{I} \{ \theta_t \in \Omega^{-\varepsilon\gamma}, t \in [0, \sigma_1] \} \right] \leq \frac{1-a_\varepsilon(1-\rho)}{1+ua_\varepsilon}$ and $\sup_{\theta_0 \in \Omega^{-2\varepsilon\gamma}} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \mathbb{I} \{ \theta_t \in \Omega^{-\varepsilon\gamma}, t \in [0, \sigma_1] \} \mathbb{I} \{ \theta_{\sigma_1} \notin \Omega^{-\varepsilon\gamma} \} \right] \leq \frac{a_\varepsilon(1+\rho/3)}{1+ua_\varepsilon}$. So in this case, we have

$$\mathcal{R}_2 = \sum_{k=1}^{+\infty} \mathbb{E} [\text{Res}_k] \leq \frac{\rho/3}{1+u-\rho} + \sum_{k=2}^{+\infty} \left[\frac{1-a_\varepsilon(1-\rho)}{1+ua_\varepsilon} \right]^{k-2} \frac{a_\varepsilon(1+\rho/3)}{1+ua_\varepsilon} = \frac{1+2\rho/3}{1+u-\rho}$$

Therefore, for any $\theta_0 \in \Omega^{-2\varepsilon\gamma}$ we can upper bound

$$\mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\Gamma} \right] \leq \mathcal{R}_1 + \mathcal{R}_2 \leq \frac{1+\rho}{1+u-\rho},$$

where $\rho \downarrow 0$ as $\varepsilon \downarrow 0$.

Lower bound of $\mathbb{E} \left[\exp(-um(\mathcal{W})\Theta(\varepsilon^{-1})\Gamma) \right]$. In this part, we only consider the big jumps in the process ζ which may escape the local minimum θ^* , and ignore the possibility of the small jumps in the process ξ which may also help escape local minimum θ^* . Here we consider the result under $\theta_0 \in \Omega^{-\varepsilon\gamma}$ which is stronger than the results under $\theta_0 \in \Omega^{-2\varepsilon\gamma}$ due to $\Omega^{-2\varepsilon\gamma} \subset \Omega^{-\varepsilon\gamma}$.

Step 1. In this step we give the formulation of the lower bound of $\mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\Gamma} \right]$. For any $u > -1$, we can compute the formula of the total probability as follows

$$\mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\Gamma} \right] \geq \sum_{k=1}^{+\infty} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})t_k} \mathbb{I} \{ \Gamma = t_k \} \right].$$

This inequality holds, since we ignore the small jumps in the process ξ which may also help escape local minimum θ^* .

For any small $\bar{\rho} > 0$, we define the inner part of Ω as $\Omega^{-\bar{\rho}} = \{ \mathbf{y} \in \Omega \mid \text{dis}(\partial\Omega, \mathbf{y}) \geq \bar{\rho} \}$ and the outer $\bar{\rho}$ -neighborhood of Ω as $\Omega^{+\bar{\rho}} = \{ \mathbf{y} \mid \text{dis}(\partial\Omega, \mathbf{y}) \geq \bar{\rho} \}$. For $k \geq 1$, we can use the strong Markov property and obtain

$$\begin{aligned} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})t_k} \mathbb{I} \{ \Gamma = t_k \} \right] &= \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})t_k} \mathbb{I} \{ \theta_t \in \Omega, t \in [0, t_k], \theta_{t_k} \notin \Omega \} \right] \\ &= \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_k} \mathbb{I} \{ \theta_{t+\sigma_k} \in \Omega, t \in [0, \sigma_k] \} \mathbb{I} \{ \theta_{t_k} \notin \Omega \} \right. \\ &\quad \cdot \left. \prod_{i=1}^{k-1} e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_i} \mathbb{I} \{ \theta_{t+\sigma_i} \in \Omega, t \in [0, \sigma_i] \} \right] \\ &\geq \inf_{\theta_0 \in \Omega^{-\bar{\rho}}} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \mathbb{I} \{ \theta_t \in \Omega^{-\bar{\rho}}, t \in [0, \sigma_1] \} \mathbb{I} \{ \theta_{\sigma_1} \notin \Omega \} \right] \\ &\quad \cdot \inf_{\theta_0 \in \Omega^{-\bar{\rho}}} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \mathbb{I} \{ \theta_t \in \Omega^{-\rho}, t \in [0, \sigma_1] \} \right]^{k-1}. \end{aligned} \tag{19}$$

Step 2. In this step we specifically lower bounds each terms in the lower bound of $\mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\Gamma} \right]$. Recall $\bar{\rho} = \varepsilon^\gamma$ where γ is a constant such that the results of Lemmas 1~3 holds. The escaping from the basin Ω with a big jump εJ_1 occurs when $\mathcal{Q}_{\sigma_1}^{-1} \Sigma_{\sigma_1} \varepsilon J_1 \in \Omega$. Furthermore, $\sup_{0 \leq t < \sigma_1} \|\theta_t - \hat{\theta}_t\| \leq \frac{1}{2}\varepsilon^\gamma$ with probability exponentially close to 1 (verified by Lemma 3). Meanwhile $\sigma_1 = \frac{2}{\alpha} \varepsilon^{\alpha\delta}$ in the α -stable ($\mathcal{S}\alpha\mathcal{S}$) distribution is much larger than $v_\varepsilon = \mathcal{O}(\ln(1/\varepsilon))$ with sufficient small ε , $\hat{\theta}_t$ reaches a $\frac{1}{2}\varepsilon^\gamma$ -neighborhood of the optimum θ^* which only requires time v_ε . So this actually means $\sup_{0 \leq t < \sigma_1} \|\theta_t - \theta^*\| \leq \varepsilon^\gamma$. In this way, to obtain the final lower bound results, we only need to estimate the escaping probability $\mathbb{P}(\mathcal{Q}_{\theta}^{-1} \Sigma_{\theta} \varepsilon J_1 \in \Omega^{-\varepsilon\gamma})$ and $\mathbb{P}(\mathcal{Q}_{\theta}^{-1} \Sigma_{\theta} \varepsilon J_1 \notin \Omega)$ uniformly over $\|\theta - \theta^*\| \leq \varepsilon^\gamma$.

Based on the results in Eqn. (14) and (15) which provides the upper bound of $\mathbb{P}(\mathcal{Q}_{\theta^*}^{-1} \Sigma_{\theta^*} \varepsilon J_1 \notin \Omega)$ and some important inequalities, we first upper bound the term

$\inf_{\|\boldsymbol{\theta}-\boldsymbol{\theta}^*\| \leq \varepsilon^{-\gamma}} \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \varepsilon J_k \notin \boldsymbol{\Omega}^{-\varepsilon^{-\gamma}} \right)$ as follows:

$$\begin{aligned}
& \inf_{\|\boldsymbol{\theta}-\boldsymbol{\theta}^*\| \leq \varepsilon^{-\gamma}} \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \varepsilon J_k \notin \boldsymbol{\Omega}^{-\varepsilon^{-\gamma}} \right) \\
&= \inf_{\|\boldsymbol{\theta}-\boldsymbol{\theta}^*\| \leq \varepsilon^{-\gamma}} \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \varepsilon J_k \notin \boldsymbol{\Omega}^{-\varepsilon^{-\gamma}}, \|\varepsilon J_k\| \leq R \right) + \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \varepsilon J_k \notin \boldsymbol{\Omega}^{-\varepsilon^{-\gamma}}, \|\varepsilon J_k\| > R \right) \\
&\stackrel{\textcircled{1}}{\leq} \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}^*}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}^*} \varepsilon J_k \notin \boldsymbol{\Omega}, \|\varepsilon J_k\| \leq R \right) + \frac{\delta' \Theta(\varepsilon^{-1})}{4 \Theta(\varepsilon^{-\delta})} + \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \varepsilon J_k \notin \boldsymbol{\Omega}^{-\varepsilon^{-\gamma}}, \|\varepsilon J_k\| > R \right) \\
&\leq \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}^*}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}^*} \varepsilon J_k \notin \boldsymbol{\Omega} \right) + \frac{\delta' \Theta(\varepsilon^{-1})}{4 \Theta(\varepsilon^{-\delta})} + \mathbb{P} (\|\varepsilon J_k\| > R) \\
&\stackrel{\textcircled{2}}{\leq} m(\mathcal{W})(1 + \delta') \frac{\Theta(\varepsilon^{-1})}{\Theta(\varepsilon^{-\delta})} + \frac{\delta' \Theta(\varepsilon^{-1})}{4 \Theta(\varepsilon^{-\delta})} + \frac{\delta' \Theta(\varepsilon^{-1})}{4 \Theta(\varepsilon^{-\delta})} \\
&= m(\mathcal{W})(1 + \delta' + \frac{\delta'}{2m(\mathcal{W})}) \frac{\Theta(\varepsilon^{-1})}{\Theta(\varepsilon^{-\delta})} \stackrel{\textcircled{3}}{\leq} m(\mathcal{W})(1 + \rho) \frac{\Theta(\varepsilon^{-1})}{\Theta(\varepsilon^{-\delta})},
\end{aligned}$$

where $\textcircled{1}$ uses the result in Eqn. (14), $\textcircled{2}$ uses Eqn. (13), and in $\textcircled{3}$ we set δ' enough small via setting small ε such that $\rho \geq \delta' + \frac{\delta'}{2m(\mathcal{W})}$. So for any ρ we choose $\delta' > 0$ small enough to upper bound

$$\inf_{\|\boldsymbol{\theta}-\boldsymbol{\theta}^*\| \leq \varepsilon^{-\gamma}} \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \varepsilon J_k \in \boldsymbol{\Omega}^{-\varepsilon^{-\gamma}} \right) = 1 - \inf_{\|\boldsymbol{\theta}-\boldsymbol{\theta}^*\| \leq \varepsilon^{-\gamma}} \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \varepsilon J_k \notin \boldsymbol{\Omega}^{-\varepsilon^{-\gamma}} \right) \geq 1 - a_{\varepsilon}(1 + \rho).$$

Similarly, we only need to lower bound the remaining term $\inf_{\|\boldsymbol{\theta}-\boldsymbol{\theta}^*\| \leq \varepsilon^{-\gamma}} \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \varepsilon J_k \notin \boldsymbol{\Omega} \right)$ as follows:

$$\begin{aligned}
& \inf_{\|\boldsymbol{\theta}-\boldsymbol{\theta}^*\| \leq \varepsilon^{-\gamma}} \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \varepsilon J_k \notin \boldsymbol{\Omega} \right) \\
&= \inf_{\|\boldsymbol{\theta}-\boldsymbol{\theta}^*\| \leq \varepsilon^{-\gamma}} \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \varepsilon J_k \notin \boldsymbol{\Omega}, \|\varepsilon J_k\| \leq R \right) + \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \varepsilon J_k \notin \boldsymbol{\Omega}, \|\varepsilon J_k\| > R \right) \\
&\stackrel{\textcircled{1}}{\geq} \inf_{\|\boldsymbol{\theta}-\boldsymbol{\theta}^*\| \leq \varepsilon^{-\gamma}} \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}^*}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}^*} \varepsilon J_k \notin \boldsymbol{\Omega}^{-\varepsilon^{-\gamma}}, \|\varepsilon J_k\| \leq R \right) - \frac{\delta' \Theta(\varepsilon^{-1})}{4 \Theta(\varepsilon^{-\delta})} + \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \varepsilon J_k \in \boldsymbol{\Omega}, \|\varepsilon J_k\| > R \right) \\
&\geq \inf_{\|\boldsymbol{\theta}-\boldsymbol{\theta}^*\| \leq \varepsilon^{-\gamma}} \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}^*}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}^*} \varepsilon J_k \notin \boldsymbol{\Omega}, \|\varepsilon J_k\| \leq R \right) - \frac{\delta' \Theta(\varepsilon^{-1})}{4 \Theta(\varepsilon^{-\delta})} + \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}} \varepsilon J_k \in \boldsymbol{\Omega}, \|\varepsilon J_k\| > R \right) \\
&\geq \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}^*}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}^*} \varepsilon J_k \notin \boldsymbol{\Omega}, \|\varepsilon J_k\| \leq R \right) - \frac{\delta' \Theta(\varepsilon^{-1})}{4 \Theta(\varepsilon^{-\delta})} \\
&\stackrel{\textcircled{2}}{\geq} \mathbb{P} \left(\mathbf{Q}_{\boldsymbol{\theta}^*}^{-1} \boldsymbol{\Sigma}_{\boldsymbol{\theta}^*} \varepsilon J_k \notin \boldsymbol{\Omega} \right) - \frac{\delta' \Theta(\varepsilon^{-1})}{2 \Theta(\varepsilon^{-\delta})} \\
&\stackrel{\textcircled{3}}{\geq} m(\mathcal{W}) \left(1 - \delta' - \frac{\delta'}{2m(\mathcal{W})} \right) \frac{\Theta(\varepsilon^{-1})}{\Theta(\varepsilon^{-\delta})} \geq m(\mathcal{W})(1 - \rho) \frac{\Theta(\varepsilon^{-1})}{\varepsilon^{-\delta}} = a_{\varepsilon}(1 - \rho),
\end{aligned}$$

where $\textcircled{1}$ uses the result in Eqn. (14), $\textcircled{2}$ uses Eqn. (15), and $\textcircled{3}$ uses Eqn. (13).

Next, for any $\rho > 0$ and ε we can obtain Laplace transforms for any $u > -1$ as follows:

$$\begin{aligned}
& \inf_{\boldsymbol{\theta}_0 \in \boldsymbol{\Omega}^{-\varepsilon^{-\gamma}}} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \mathbb{I} \left\{ \boldsymbol{\theta}_t \in \boldsymbol{\Omega}^{-\varepsilon^{-\gamma}}, t \in [0, \sigma_1] \right\} \right] \\
&\geq [1 - a_{\varepsilon}(1 + \rho)] \mathbb{E} \left[\int_0^{+\infty} e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \cdot \Theta(\varepsilon^{-\delta}) e^{-\Theta(\varepsilon^{-\delta})\sigma_1} \mathbf{d}\sigma_1 \right] \\
&= \frac{1 - a_{\varepsilon}(1 + \rho)}{1 + ua_{\varepsilon}},
\end{aligned} \tag{20}$$

and

$$\begin{aligned}
& \inf_{\boldsymbol{\theta}_0 \in \boldsymbol{\Omega}^{-\varepsilon^{-\gamma}}} \mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \mathbb{I} \left\{ \boldsymbol{\theta}_t \in \boldsymbol{\Omega}^{-\varepsilon^{-\gamma}}, t \in [0, \sigma_1] \right\} \mathbb{I} \left\{ \boldsymbol{\theta}_{\sigma_1} \notin \boldsymbol{\Omega} \right\} \right] \\
&\geq [1 - a_{\varepsilon}(1 + \rho)] \mathbb{E} \left[\int_0^{+\infty} e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\sigma_1} \cdot \Theta(\varepsilon^{-\delta}) e^{-\Theta(\varepsilon^{-\delta})\sigma_1} \mathbf{d}\sigma_1 \right] = \frac{a_{\varepsilon}(1 - \rho)}{1 + ua_{\varepsilon}}.
\end{aligned}$$

Step 3. Here we summarize the results in Steps 1 and 2 such that we can lower bound the desired results $\mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\Gamma} \right]$. Specifically, from Eqn. (19), for any $\theta_0 \in \Omega^{-\varepsilon^\gamma}$ we can lower bound

$$\mathbb{E} \left[e^{-um(\mathcal{W})\Theta(\varepsilon^{-1})\Gamma} \right] \geq \frac{a_\varepsilon(1-\rho)}{1+ua_\varepsilon} \sum_{k=1}^{+\infty} \left(\frac{1-a_\varepsilon(1+\rho)}{1+ua_\varepsilon} \right)^{k-1} = \frac{1-\rho}{1+u+\rho},$$

where $\rho \downarrow 0$ as $\varepsilon \downarrow 0$. The proof is completed. \square

D.2 Proof of Theorem 2

Proof. In this step we prove the sequence $\{\widehat{\theta}_t\}$ produced by Eqn. (8) or (9) locates in a very small neighborhood of the optimum solution θ^* of the local basin Ω after a very small time interval.

Step 1. In this step, we prove the first part of Theorem 2. Since we assume the function is locally strongly convex, by using Lemmas 3 and 4, we know that the sequence $\{\widehat{\theta}_t\}$ produced by Eqn. (8) or (9) exponentially converges to the minimum θ^* at the current local basin Ω . So for any initialization $\theta_0 \in \Omega$, we have

$$\|\widehat{\theta}_t - \theta^*\|_2^2 \leq c_1 \exp(-c_2 t),$$

where $c_1 = \frac{2\Delta}{\mu}$ and $c_2 = \frac{2\mu\tau}{\beta_1(v_{\max}+\varepsilon)+\mu\tau} \left(\beta_1 - \frac{\beta_2}{4} \right)$ in Adam, $c_1 = \frac{2\Delta}{\mu}$ and $c_2 = 2\mu$ in SGD. Therefore, for any initialization $\theta_0 \in \Omega$ and sufficient small ε , we can obtain

$$\|\widehat{\theta}_t - \theta^*\|_2^2 \leq \varepsilon^{\bar{p}} \text{ if } t \geq v_\varepsilon = \frac{1}{c_3} \ln \left(\frac{2\Delta}{\mu\varepsilon^{\bar{p}}} \right).$$

where $c_3 = \frac{2\mu\tau}{\beta_1(v_{\max}+\varepsilon)+\mu\tau} \left(\beta_1 - \frac{\beta_2}{4} \right)$ in ADAM, $c_3 = 2\mu$ in SGD, and $\Delta = F(\theta_0) - F(\theta^*)$.

Step 2. In this step, we prove the second part of Theorem 2. By replacing p with $p/2$ in Lemma E.5, we can directly obtain the results. \square

E Proofs of Auxiliary Theories and Lemmas in Appendix C

Before analysis, we first introduce two useful lemmas which will be used in subsequent analysis.

Lemma 4 (Grönwall's Lemma [12]). *Suppose $g(s) : [0, t_0]$ is a non-negative continues function. If for almost $s \in [0, t_0]$*

$$g'(s) \leq q(s)g(s)$$

where $q(s)$ is a continuous function, then we have

$$g(t) \leq g(0) \exp \left(\int_0^t q(s) ds \right).$$

Lemma 5 (Theorem 5.3 in [13]). *Consider a set $\mathbf{A} \in \mathcal{B}(\mathbb{R} \setminus 0)$ with $0 \in \bar{\mathbf{A}}$ and a function $f : \mathbb{R} \rightarrow \mathbb{R}$ with Borel measurable and finite on \mathbf{A} . Then we have*

(1) *The process $(\int_0^t \int_{\mathbf{A}} f(x) \nu(ds, dx))_{0 \leq t \leq T}$ is a compound Poisson process with characteristic function*

$$\mathbb{E} \left(\exp \left(i\lambda \int_0^t \int_{\mathbf{A}} f(x) \mu^L(ds, dx) \right) \right) = \exp \left(t \int_{\mathbf{A}} (e^{i\lambda f(x)} - 1) \nu(dx) \right).$$

(2) *If $f \in L^1(\mathbf{A})$, then*

$$\mathbb{E} \left(\int_0^t \int_{\mathbf{A}} f(x) \mu^L(ds, dx) \right) = t \int_{\mathbf{A}} f(x) \nu(dx).$$

E.1 Proof of Theorem 3 for the Linear Convergence of Lévy-driven SGD SDE (8)

Proof. Step 1. In this step, we upper bound the gradient norm of the Lyapunov function $\mathcal{L}(t) = F(\hat{\boldsymbol{\theta}}_t) - F(\boldsymbol{\theta}^*)$ of (8) with $\hat{\mathbf{Q}}_t = \mathbf{I}$ and $\beta_2 = 0$. More specifically, we can upper bound $d\mathcal{L}(t)$ as follows:

$$d\mathcal{L}(t) = \langle \nabla F(\hat{\boldsymbol{\theta}}_t), d\hat{\boldsymbol{\theta}}_t \rangle = \langle \nabla F(\hat{\boldsymbol{\theta}}_t), -\nabla F(\hat{\boldsymbol{\theta}}_t) \rangle = -\|\nabla F(\hat{\boldsymbol{\theta}}_t)\|_2^2. \quad (21)$$

Step 2. Here we prove the linear convergence behavior of $\mathcal{L}(t) = F(\hat{\boldsymbol{\theta}}_t) - F(\boldsymbol{\theta}^*)$ by using the results in Step 1. Since $F(\boldsymbol{\theta})$ is locally μ -strongly convex, then we have

$$F(\mathbf{y}) \geq F(\boldsymbol{\theta}) + \langle \nabla F(\boldsymbol{\theta}), \mathbf{y} - \boldsymbol{\theta} \rangle + \frac{\mu}{2} \|\mathbf{y} - \boldsymbol{\theta}\|_2^2.$$

Next, by minimizing \mathbf{y} on both side ($\mathbf{y} = \boldsymbol{\theta}_*$ for the left side and $\mathbf{y} = \boldsymbol{\theta} - \frac{1}{\mu} \nabla F(\boldsymbol{\theta})$ for the right side), it yields

$$\|\nabla F(\boldsymbol{\theta})\|_2^2 \geq 2\mu(F(\boldsymbol{\theta}) - F(\boldsymbol{\theta}_*)). \quad (22)$$

Hence, plugging the above equation into Eqn. (21) gives

$$d\mathcal{L}(t) \leq -2\mu(F(\boldsymbol{\theta}) - F(\boldsymbol{\theta}_*)) = -2\mu\mathcal{L}(t).$$

In this way, by using the result in Lemma 4, we can easily obtain

$$\mathcal{L}(t) \leq \mathcal{L}(0) \exp\left(-\int_0^t 2\mu ds\right) \leq \Delta \exp(-2\mu t),$$

where we use $\mathcal{L}(0) = F(\hat{\boldsymbol{\theta}}_0) - F(\boldsymbol{\theta}^*) = \Delta$ where $\boldsymbol{\theta}^*$ is the optimum of the current basin.

Step 3. Finally, we explore the local strong-convexity of $F(\boldsymbol{\theta})$ to show the linear convergence of $\|\hat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^*\|_2^2$. Specifically, by using the strongly convex property of $F(\boldsymbol{\theta})$, we can obtain

$$F(\boldsymbol{\theta}) - F(\boldsymbol{\theta}^*) \geq \frac{\mu}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2^2.$$

So this gives

$$\|\hat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^*\|_2^2 \leq \frac{2\Delta}{\mu} \exp(-2\mu t).$$

The proof is completed. \square

E.2 Proof of Theorem 4 for the Linear Convergence of Lévy-driven ADAM SDE (9)

Proof. Step 1. In this step, we upper bound the gradient norm of the Lyapunov function of (9) defined as

$$\mathcal{L}(t) = F(\hat{\boldsymbol{\theta}}_t) - F(\boldsymbol{\theta}^*) + \frac{1}{2} \|\widehat{\mathbf{m}}_t\|_{\widehat{\mathbf{s}}_t^{-1}}^2, \quad (23)$$

where $\widehat{\mathbf{s}}_t = \frac{h_t}{\mu_t} (\sqrt{\omega_t \widehat{\mathbf{v}}_t} + \epsilon)$ with $h_t = \beta_1$, $\mu_t = (1 - e^{-\beta_1 t})^{-1}$ and $\omega_t = (1 - e^{-\beta_2 t})^{-1}$. Here we define $\|\mathbf{x}\|_{\mathbf{y}}^2 = \sum_i \mathbf{y}_i \mathbf{x}_i^2$. Then we can compute the derivative of Lyapunov function as

$$d\mathcal{L}(t) = \underbrace{\langle \nabla F(\hat{\boldsymbol{\theta}}_t), d\hat{\boldsymbol{\theta}}_t \rangle + \sum_{i=1}^d \frac{1}{\widehat{\mathbf{s}}_{t,i}} \widehat{\mathbf{m}}_{t,i} d\widehat{\mathbf{m}}_{t,i}}_{P_1} - \underbrace{\sum_{i=1}^d \frac{1}{2\widehat{\mathbf{s}}_{t,i}^2} \widehat{\mathbf{m}}_{t,i}^2 \nabla_{\widehat{\mathbf{v}}_t} \widehat{\mathbf{v}}_{t,i} d\widehat{\mathbf{v}}_{t,i} - \sum_{i=1}^d \frac{1}{2\widehat{\mathbf{s}}_{t,i}^2} \widehat{\mathbf{m}}_{t,i}^2 \nabla_t \widehat{\mathbf{s}}_{t,i}}_{P_2}, \quad (24)$$

where $\widehat{\mathbf{m}}_{t,i}$, $\widehat{\mathbf{v}}_{t,i}$ and $\widehat{\mathbf{s}}_{t,i}$ respectively denote the i -th entries of $\widehat{\mathbf{m}}_t$, $\widehat{\mathbf{v}}_t$ and $\widehat{\mathbf{s}}_t$.

We first consider Adam in which $h_t = \beta_1$, $\mu_t = (1 - e^{-\beta_1 t})^{-1}$, and $\omega_t = (1 - e^{-\beta_2 t})^{-1}$. We also assume $\beta_1 \leq \beta_2 \leq 2\beta_1$ which is consistent with the practical setting where $\beta_1 = 0.9$ and $\beta_2 = 0.999$.

Let $[\nabla f_{S_t}(\hat{\theta}_t)]_i$ denotes the i -th entry of the vector $\nabla f_{S_t}(\hat{\theta}_t)$. Under this setting, we can first upper bound the first term P_1 as follows:

$$\begin{aligned}
P_1 &= \langle \nabla F(\hat{\theta}_t), d\hat{\theta}_t \rangle + \sum_{i=1}^d \frac{1}{\hat{s}_{t,i}} \widehat{\mathbf{m}}_{t,i} d\widehat{\mathbf{m}}_{t,i} - \sum_{i=1}^d \frac{1}{2\widehat{s}_{t,i}^2} \widehat{\mathbf{m}}_{t,i}^2 \nabla_{\widehat{\mathbf{v}}_t \widehat{s}_{t,i}} d\widehat{\mathbf{v}}_{t,i} \\
&= \langle \nabla F(\hat{\theta}_t), -\frac{\mu_t \widehat{\mathbf{m}}_t}{\sqrt{\omega_t \widehat{\mathbf{v}}_t} + \epsilon} \rangle + \beta_1 \sum_{i=1}^d \frac{1}{\hat{s}_{t,i}} \widehat{\mathbf{m}}_{t,i} (\nabla F_i(\theta_t) - \widehat{\mathbf{m}}_{t,i}) - \sum_{i=1}^d \frac{\beta_2}{2\widehat{s}_{t,i}^2} \widehat{\mathbf{m}}_{t,i}^2 ([\nabla f_{S_t}(\hat{\theta}_t)]_i - \widehat{\mathbf{v}}_{t,i}) \nabla_{\widehat{\mathbf{v}}_t} \widehat{s}_{t,i} \\
&= -\beta_1 \left\langle \nabla F(\hat{\theta}_t), \frac{\widehat{\mathbf{m}}_t}{\widehat{\mathbf{s}}_t} \right\rangle + \beta_1 \left\langle \nabla F(\hat{\theta}_t), \frac{\widehat{\mathbf{m}}_t}{\widehat{\mathbf{s}}_t} \right\rangle - \beta_1 \sum_{i=1}^d \frac{1}{\hat{s}_{t,i}} \widehat{\mathbf{m}}_{t,i}^2 - \beta_2 \sum_{i=1}^d \frac{1}{2\widehat{s}_{t,i}^2} \widehat{\mathbf{m}}_{t,i}^2 ([\nabla f_{S_t}(\hat{\theta}_t)]_i - \widehat{\mathbf{v}}_{t,i}) \nabla_{\widehat{\mathbf{v}}_t} \widehat{s}_{t,i} \\
&= -\beta_1 \sum_{i=1}^d \frac{1}{\hat{s}_{t,i}} \widehat{\mathbf{m}}_{t,i}^2 - \beta_2 \sum_{i=1}^d \frac{1}{2\widehat{s}_{t,i}^2} \widehat{\mathbf{m}}_{t,i}^2 ([\nabla f_{S_t}(\hat{\theta}_t)]_i - \widehat{\mathbf{v}}_{t,i}) \nabla_{\widehat{\mathbf{v}}_t} \widehat{s}_{t,i}.
\end{aligned}$$

Next, we plug the specific formulation of $\nabla_{\widehat{\mathbf{v}}_t} \widehat{s}_{t,i} = \frac{\beta_1 \sqrt{\omega_t}}{2\mu_t \sqrt{\widehat{\mathbf{v}}_{t,i}}}$ into the above equation and obtain:

$$\begin{aligned}
P_1 &= -\beta_1 \sum_{i=1}^d \frac{1}{\hat{s}_{t,i}} \widehat{\mathbf{m}}_{t,i}^2 - \beta_2 \sum_{i=1}^d \frac{1}{2\widehat{s}_{t,i}^2} \widehat{\mathbf{m}}_{t,i}^2 ([\nabla f_{S_t}(\hat{\theta}_t)]_i - \widehat{\mathbf{v}}_{t,i}) \frac{\beta_1 \sqrt{\omega_t}}{2\mu_t \sqrt{\widehat{\mathbf{v}}_{t,i}}} \\
&= -\beta_1^2 \sum_{i=1}^d \frac{\widehat{\mathbf{m}}_{t,i}^2}{\mu_t \widehat{s}_{t,i}^2} \left(\epsilon + (1 - \frac{\beta_2}{4\beta_1}) \sqrt{\omega_t \widehat{\mathbf{v}}_{t,i}} + \frac{\beta_2}{4\beta_1} \frac{[\nabla f_{S_t}(\hat{\theta}_t)]_i \sqrt{\omega_t}}{\sqrt{\widehat{\mathbf{v}}_{t,i}}} \right) \\
&= -\beta_1 \sum_{i=1}^d \frac{\widehat{\mathbf{m}}_{t,i}^2}{\widehat{s}_{t,i}} \left(1 - \frac{\beta_2}{4\beta_1} + \frac{\beta_2 \epsilon}{4\beta_1 (\epsilon + \sqrt{\omega_t \widehat{\mathbf{v}}_{t,i}})} + \frac{\beta_2}{4\beta_1} \frac{[\nabla f_{S_t}(\hat{\theta}_t)]_i \sqrt{\omega_t}}{\sqrt{\widehat{\mathbf{v}}_{t,i} (\epsilon + \sqrt{\omega_t \widehat{\mathbf{v}}_{t,i}})}} \right) \\
&\leq -\left(\beta_1 - \frac{\beta_2}{4} \right) \sum_{i=1}^d \frac{\widehat{\mathbf{m}}_{t,i}^2}{\widehat{s}_{t,i}} = -\left(\beta_1 - \frac{\beta_2}{4} \right) \|\widehat{\mathbf{m}}_t\|_{\widehat{\mathbf{s}}_t^{-1}}^2.
\end{aligned}$$

Then we consider the second term P_2 under the setting $\widehat{\mathbf{s}}_t = \frac{\beta_1}{\mu_t} (\sqrt{\omega_t \widehat{\mathbf{v}}_t} + \epsilon)$ with $\mu_t = (1 - e^{-\beta_1 t})^{-1}$ and $\omega_t = (1 - e^{-\beta_2 t})^{-1}$. Similarly, we can upper bound P_2 as

$$\begin{aligned}
P_2 &= -\sum_{i=1}^d \frac{1}{2\widehat{s}_{t,i}^2} \widehat{\mathbf{m}}_{t,i}^2 \nabla_{\widehat{\mathbf{v}}_t} \widehat{s}_{t,i} \\
&= -\beta_1 \sum_{i=1}^d \frac{1}{2\widehat{s}_{t,i}^2} \widehat{\mathbf{m}}_{t,i}^2 \left(\beta_1 e^{-\beta_1 t} \left(\epsilon + \sqrt{\frac{\widehat{\mathbf{v}}_{t,i}}{1 - e^{-\beta_2 t}}} \right) - \frac{1}{2} \beta_2 e^{-\beta_2 t} \frac{1 - e^{-\beta_1 t}}{1 - e^{-\beta_2 t}} \sqrt{\frac{\widehat{\mathbf{v}}_{t,i}}{1 - e^{-\beta_2 t}}} \right) \\
&= -\frac{\beta_1^2}{2} \sum_{i=1}^d \frac{\widehat{\mathbf{m}}_{t,i}^2}{\mu_t \widehat{s}_{t,i}^2} \frac{e^{-\beta_1 t}}{1 - e^{-\beta_1 t}} \left(\epsilon + \left(1 - \frac{\beta_2 e^{-\beta_2 t} (1 - e^{-\beta_1 t})}{2\beta_1 e^{-\beta_1 t} (1 - e^{-\beta_2 t})} \right) \sqrt{\frac{\widehat{\mathbf{v}}_{t,i}}{1 - e^{-\beta_2 t}}} \right) \\
&\stackrel{\textcircled{1}}{\leq} -\frac{\beta_1^2}{2} \sum_{i=1}^d \frac{\widehat{\mathbf{m}}_{t,i}^2}{\mu_t \widehat{s}_{t,i}^2} \frac{e^{-\beta_1 t}}{1 - e^{-\beta_1 t}} \left(\epsilon + \left(1 - \frac{\beta_2}{2\beta_1} \right) \sqrt{\frac{\widehat{\mathbf{v}}_{t,i}}{1 - e^{-\beta_2 t}}} \right) \\
&= -\frac{\beta_1}{2} \sum_{i=1}^d \frac{\widehat{\mathbf{m}}_{t,i}^2}{\widehat{s}_{t,i}} \frac{e^{-\beta_1 t}}{1 - e^{-\beta_1 t}} \left(1 - \frac{\beta_2}{2\beta_1} + \frac{\beta_2 \epsilon}{2\beta_1 (\epsilon + \sqrt{\omega_t \widehat{\mathbf{v}}_{t,i}})} \right) \\
&\leq -\frac{1}{2} \left(\beta_1 - \frac{\beta_2}{2} \right) \frac{e^{-\beta_1 t}}{1 - e^{-\beta_1 t}} \sum_{i=1}^d \frac{\widehat{\mathbf{m}}_{t,i}^2}{\widehat{s}_{t,i}} = -\frac{1}{2} \left(\beta_1 - \frac{\beta_2}{2} \right) \frac{e^{-\beta_1 t}}{1 - e^{-\beta_1 t}} \|\widehat{\mathbf{m}}_t\|_{\widehat{\mathbf{s}}_t^{-1}}^2 \stackrel{\textcircled{2}}{\leq} 0,
\end{aligned}$$

where $\textcircled{1}$ uses $\frac{\beta_2 e^{-\beta_2 t} (1 - e^{-\beta_1 t})}{2\beta_1 e^{-\beta_1 t} (1 - e^{-\beta_2 t})} \leq \frac{\beta_2}{2\beta_1}$ since $\beta_2 \geq \beta_1$; in $\textcircled{2}$ we assume $\beta_1 - \frac{\beta_2}{2} > 0$. Therefore, by combining the upper bounds of P_1 and P_2 we can upper bound

$$d\mathcal{L}(t) \leq -\left[\beta_1 - \frac{\beta_2}{4} \right] \|\widehat{\mathbf{m}}_t\|_{\widehat{\mathbf{s}}_t^{-1}}^2. \quad (25)$$

On the other hand, noting $h_t = \beta_1$, $\mu_t = (1 - e^{-\beta_1 t})^{-1}$ and $\omega_t = (1 - e^{-\beta_2 t})^{-1}$, we have

$$\begin{aligned}\widehat{\mathbf{s}}_{t,i} &= \frac{h_t}{\mu_t} \left(\epsilon + \sqrt{\omega_t \widehat{\mathbf{v}}_{t,i}} \right) = \beta_1 (1 - e^{-\beta_1 t}) \left(\epsilon + \sqrt{\frac{\widehat{\mathbf{v}}_{t,i}}{1 - e^{-\beta_2 t}}} \right) \leq \beta_1 \left(\epsilon + \frac{1 - e^{-\beta_1 t}}{\sqrt{1 - e^{-\beta_2 t}}} \sqrt{\widehat{\mathbf{v}}_{t,i}} \right) \\ &\stackrel{\textcircled{1}}{\leq} \beta_1 \left(\epsilon + \frac{1 - e^{-\beta_1 t}}{1 - e^{-\beta_2 t/2}} \sqrt{\widehat{\mathbf{v}}_{t,i}} \right) \stackrel{\textcircled{2}}{\leq} \beta_1 (\epsilon + v_{\max}),\end{aligned}$$

where $\textcircled{1}$ uses $\sqrt{1-x} \geq 1 - \sqrt{x}$ for $0 \leq x \leq 1$ and $\textcircled{2}$ holds since $\sqrt{\widehat{\mathbf{v}}_{t,i}} \leq v_{\max}$. By using the assumption $\|\widehat{\mathbf{m}}_t\|^2 \geq \tau \|\nabla \mathbf{F}(\widehat{\boldsymbol{\theta}}_t)\|^2$, we can establish

$$\|\widehat{\mathbf{m}}_t\|_{\widehat{\mathbf{s}}_t}^2 \geq \frac{1}{\beta_1 (\epsilon + v_{\max})} \|\widehat{\mathbf{m}}_t\|_2^2 \geq \frac{\tau}{\beta_1 (\epsilon + v_{\max})} \|\nabla \mathbf{F}(\widehat{\boldsymbol{\theta}}_t)\|_2^2. \quad (26)$$

Then from the locally μ -strongly convex property Eqn. (22):

$$\|\nabla \mathbf{F}(\boldsymbol{\theta})\|_2^2 \geq 2\mu(\mathbf{F}(\boldsymbol{\theta}) - \mathbf{F}(\boldsymbol{\theta}^*)).$$

then we plug the above inequality into Eqn. (26) and establish

$$\|\widehat{\mathbf{m}}_t\|_{\widehat{\mathbf{s}}_t}^2 \geq \frac{1}{\beta_1 (\epsilon + v_{\max})} \|\widehat{\mathbf{m}}_t\|_2^2 \geq \frac{2\mu\tau}{\beta_1 (\epsilon + v_{\max})} (\mathbf{F}(\widehat{\boldsymbol{\theta}}_t) - \mathbf{F}(\boldsymbol{\theta}^*)).$$

Finally, we can write Eqn. (25) as

$$\begin{aligned}\mathbf{d}\mathcal{L}(t) &\leq -\frac{2\mu\tau}{\beta_1 (\epsilon + v_{\max}) + \mu\tau} \left[\beta_1 - \frac{\beta_2}{4} \right] \left(\frac{1}{2} + \frac{\beta_1 (\epsilon + v_{\max})}{2\mu\tau} \right) \|\widehat{\mathbf{m}}_t\|_{\widehat{\mathbf{s}}_t}^2 \\ &\leq -\frac{2\mu\tau}{\beta_1 (\epsilon + v_{\max}) + \mu\tau} \left[\beta_1 - \frac{\beta_2}{4} \right] \left(\mathbf{F}(\widehat{\boldsymbol{\theta}}_t) - \mathbf{F}(\boldsymbol{\theta}^*) + \frac{1}{2} \|\widehat{\mathbf{m}}_t\|_{\widehat{\mathbf{s}}_t}^2 \right) \\ &= -c_1 \mathcal{L}(t),\end{aligned}$$

where $c_1 = \frac{2\mu\tau}{\beta_1 (\epsilon + v_{\max}) + \mu\tau} \left[\beta_1 - \frac{\beta_2}{4} \right]$.

Step 2. Here we prove the linear convergence behavior of $\mathcal{L}(t) = \mathbf{F}(\widehat{\boldsymbol{\theta}}_t) - \mathbf{F}(\boldsymbol{\theta}^*)$ by using the results in Step 1. More specifically, by using the result in Lemma 4, we can easily obtain

$$\begin{aligned}\mathcal{L}(t) &\leq \mathcal{L}(0) \exp \left(\int_0^t c_1 \mathbf{d}s \right) = \mathcal{L}(0) \exp \left(-\frac{2\mu\tau}{\beta_1 (\epsilon + v_{\max}) + \mu\tau} \left(\beta_1 - \frac{\beta_2}{4} \right) t \right) \\ &\stackrel{\textcircled{1}}{\leq} (\mathbf{F}(\widehat{\boldsymbol{\theta}}_0) - \mathbf{F}(\boldsymbol{\theta}^*)) \exp \left(-\frac{2\mu\tau}{\beta_1 (\epsilon + v_{\max}) + \mu\tau} \left(\beta_1 - \frac{\beta_2}{4} \right) t \right),\end{aligned}$$

where $\textcircled{1}$ uses $\mathcal{L}(0) = \mathbf{F}(\widehat{\boldsymbol{\theta}}_0) - \mathbf{F}(\boldsymbol{\theta}^*) = \Delta$ due to $\widehat{\mathbf{m}}_0 = \mathbf{0}$.

Step 3. Finally, we explore the local strong-convexity of $\mathbf{F}(\boldsymbol{\theta})$ to show the linear convergence of $\|\widehat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^*\|_2^2$. Specifically, by using the strongly convex property of $\mathbf{F}(\boldsymbol{\theta})$, we can obtain

$$\mathbf{F}(\boldsymbol{\theta}) - \mathbf{F}(\boldsymbol{\theta}^*) \geq \frac{\mu}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2^2.$$

So this gives

$$\|\widehat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^*\|_2^2 \leq \frac{2\Delta}{\mu} \exp \left(-\frac{2\mu\tau}{\beta_1 (\epsilon + v_{\max}) + \mu\tau} \left(\beta_1 - \frac{\beta_2}{4} \right) t \right).$$

The proof is completed. \square

E.3 Proof of Lemma 1

Proof. To begin with, the process $\boldsymbol{\xi}$ is defined as $\boldsymbol{\xi}_t = \sum_{s \leq t} \Delta L_s \mathbb{I} \{ \|L_s\| \leq \varepsilon^{-\delta} \}$. Then by setting the set $\mathcal{A} = \{ \mathbf{y} \mid \|\mathbf{y}\| \leq \varepsilon^{-\delta} \}$ in Lemma 5 and noting $f(x) = x \in L^1(\mathcal{A})$, one can find

$\mathbb{E}[\xi_t] = t \int_{\mathbb{A}} f(x) \nu(\mathbf{d}x)$. Therefore, we can decompose the process ξ into two processes $\widehat{\xi}$ and linear drift, namely,

$$\xi_t = \widehat{\xi}_t + \mu_\varepsilon t,$$

where $\widehat{\xi}$ is a zero mean Lévy martingale with bounded jumps. Then we prove our results in two steps.

Step 1. We first estimate the value of μ_ε . Since ξ is a Lévy process, by Lévy-Itô decomposition theory [13, Theorem 6.1] its characteristic function is of form

$$\mathbb{E}[e^{i\langle \lambda, \xi_t \rangle}] = \exp \left(t \int_{\mathbb{R}^d \setminus \{0\}} \left(e^{i\langle \lambda, \mathbf{y} \rangle} - 1 - i\langle \lambda, \mathbf{y} \rangle \mathbb{I}\{\|\mathbf{y}\| \leq 1\} \right) \mathbb{I}\{\|\mathbf{y}\| \leq \varepsilon^{-\delta}\} \mathbf{d}\mathbf{y} \right),$$

which can be further split into two Lévy processes $\xi_{(1)}$ and $\xi_{(2)}$ with characteristic functions

$$\mathbb{E}[e^{i\langle \lambda, \xi_{(1)}, t \rangle}] = \exp \left(t \int_{0 < \|\mathbf{y}\| < 1} \left(e^{i\langle \lambda, \mathbf{y} \rangle} - 1 - i\langle \lambda, \mathbf{y} \rangle \right) \mathbf{d}\mathbf{y} \right)$$

and

$$\mathbb{E}[e^{i\langle \lambda, \xi_{(2)}, t \rangle}] = \exp \left(t \int_{1 \leq \|\mathbf{y}\| \leq \varepsilon^{-\delta}} \left(e^{i\langle \lambda, \mathbf{y} \rangle} - 1 \right) \mathbf{d}\mathbf{y} \right).$$

Let us consider ξ on the set $\{\mathbf{y} \mid 0 < \|\mathbf{y}\| \leq 1\}$. We construct a compensated compound Poisson process

$$L'_t = \sum_{s \leq t} \Delta L'_s \mathbb{I}\{1 > \|\Delta L_s\| > \varepsilon'\} - t \int_{1 > \|\mathbf{y}\| > \varepsilon'} \mathbf{y} \nu(\mathbf{d}\mathbf{y}) = \int_0^t \int_{1 > \|\mathbf{y}\| > \varepsilon'} \mathbf{y} \mu^L(\mathbf{d}\mathbf{y}, \mathbf{d}s) - t \int_{1 > \|\mathbf{y}\| > \varepsilon'} \mathbf{y} \nu(\mathbf{d}\mathbf{y}),$$

where ε' is a very small constant. By applying Lemma 5 on $\sum_{s \leq t} \Delta L'_s \mathbb{I}\{1 > \|\Delta L_s\| > \varepsilon'\}$, the characteristic function of L'_t is

$$\mathbb{E}[e^{i\langle \lambda, L'_t \rangle}] = \exp \left(t \int_{\varepsilon' < \|\mathbf{y}\| < 1} \left(e^{i\langle \lambda, \mathbf{y} \rangle} - 1 - i\langle \lambda, \mathbf{y} \rangle \right) \mathbf{d}\mathbf{y} \right).$$

This means that there exists a Lévy process L' which is a square integral martingale such that $L' \rightarrow \xi_{(1)}$ as $\varepsilon' \rightarrow 0$. As L' is a square integral martingale, we have $\mathbb{E}(\xi_{(1)}) = \mathbb{E}(L') = \mathbf{0}$, which means that μ_ε is only related to $\xi_{(2)}$. Therefore, we have

$$\begin{aligned} \mu_\varepsilon^i &= \mathbb{E}[\xi_{(2)}^i] = \int_{1 \leq \|\mathbf{y}\| \leq \varepsilon^{-\delta}} \mathbf{y}_i \nu(\mathbf{d}\mathbf{y}), \quad (i = 1, \dots, d) \\ \|\mu_\varepsilon\|^2 &= \int_{1 \leq \|\mathbf{y}\| \leq \varepsilon^{-\delta}} \|\mathbf{y}\|^2 \nu(\mathbf{d}\mathbf{y}) = - \int_1^{\varepsilon^{-\delta}} u^2 \mathbf{d}\Theta(u) = -u^2 \Theta(u) \Big|_1^{\varepsilon^{-\delta}} + 2 \int_1^{\varepsilon^{-\delta}} u \Theta(u) \mathbf{d}u \leq \varepsilon^{-2\delta} \Theta(1). \end{aligned}$$

Thus, we can bound $\|\mu_\varepsilon\| \leq \varepsilon^{-\delta} \sqrt{\Theta(1)}$. Finally, by setting $\theta_0 = (1 - \delta)/3$ and $\rho_0 = (1 - \delta)/4$ we can obtain $\varepsilon \|\mu_\varepsilon\| T_\varepsilon = \varepsilon^{1-\delta-\theta} \sqrt{\Theta(1)} \leq \varepsilon^{2\rho}$ by setting ε sufficient small such that $\Theta(1) \leq \frac{1}{\varepsilon^{1-2\rho-\delta-\theta}}$.

Step 2. Since the increment is non-negative, the quadratic variation process $[\varepsilon \widehat{\xi}]_t^d$ is a Lévy subordinator, namely,

$$[\varepsilon \widehat{\xi}]_t^d = \varepsilon^2 \sum_{s \leq t} \|\Delta \widehat{\xi}_s\|^2 = \varepsilon^2 \int_0^t \int_{0 < \|\mathbf{y}\| \leq \varepsilon^{-\delta}} \|\mathbf{y}\|^2 N(\mathbf{d}\mathbf{y}, \mathbf{d}s),$$

where $\Delta \widehat{\xi}_s = \widehat{\xi}_s - \widehat{\xi}_{s-}$ where $\widehat{\xi}_{s-} = \lim_{t \uparrow s} \widehat{\xi}_t$.

Since the jumps of $[\varepsilon \widehat{\xi}]^d$ are bounded, its Laplace transform is well-defined for all $\lambda \in \mathbb{R}$:

$$\mathbb{E} e^{\lambda [\varepsilon \widehat{\xi}]_t^d} = \exp \left(t \int_{0 < \|\mathbf{y}\| \leq \varepsilon^{-\delta}} (e^{\lambda \varepsilon^2 \|\mathbf{y}\|^2} - 1) \nu(\mathbf{d}\mathbf{y}) \right) = \exp \left(-t \int_{0 < u \leq \varepsilon^{-\delta}} (e^{\lambda \varepsilon^2 u^2} - 1) \mathbf{d}\Theta(u) \right).$$

For any $\lambda > 0$, the exponential Chebyshev inequality indicates

$$\begin{aligned} \mathbb{P}\left([\varepsilon \widehat{\boldsymbol{\xi}}]_{T_\varepsilon}^d > \varepsilon^\rho\right) &= \mathbb{P}\left(e^{\lambda[\varepsilon \widehat{\boldsymbol{\xi}}]_{T_\varepsilon}^d} > e^{\lambda \varepsilon^\rho}\right) \leq e^{-\lambda \varepsilon^\rho} \mathbb{E}[e^{\lambda[\varepsilon \widehat{\boldsymbol{\xi}}]_{T_\varepsilon}^d}] \\ &= \exp\left(-\lambda \varepsilon^\rho - T_\varepsilon \int_{0 < u \leq \varepsilon^{-\delta}} (e^{\lambda \varepsilon^2 u^2} - 1) d\Theta(u)\right). \end{aligned} \quad (27)$$

For $\lambda = \lambda_\varepsilon = \varepsilon^{-2\rho}$ with $0 < \rho < \rho_0 = (1 - \delta)/4$ we have $\max_{0 < u \leq \varepsilon^{-\delta}} \lambda \varepsilon^2 u^2 \leq \lambda_\varepsilon \varepsilon^{2(1-\delta)} \leq \varepsilon^{\frac{3}{2}(1-\delta)} \downarrow 0$ as $\varepsilon \downarrow 0$. With help of the elementary inequality $e^x - 1 \leq 2x$ for small positive x the second summand appearing in the exponent in right-hand side of (27) can be now established as

$$\begin{aligned} \left| T_\varepsilon \int_{0 < u \leq \varepsilon^{-\delta}} (e^{\lambda_\varepsilon \varepsilon^2 u^2} - 1) d\Theta(u) \right| &\leq \left| 2T_\varepsilon \lambda_\varepsilon \varepsilon^2 \left(\int_{0 < u \leq 1} + \int_{1 < u \leq \varepsilon^{-\delta}} \right) u^2 d\Theta(u) \right| \\ &\leq 2T_\varepsilon \lambda_\varepsilon \varepsilon^2 \left| \int_{0 < u \leq 1} u^2 d\Theta(u) \right| + 2T_\varepsilon \lambda_\varepsilon \varepsilon^{2(1-\delta)} \left| \int_{1 < u \leq \varepsilon^{-\delta}} d\Theta(u) \right| \\ &\leq 2CT_\varepsilon \lambda_\varepsilon \varepsilon^2 + 2\Theta(1)T_\varepsilon \lambda_\varepsilon \varepsilon^{2(1-\delta)} \end{aligned}$$

where $C = \left| \int_{0 < u \leq 1} u^2 d\Theta(u) \right| \in (0, +\infty)$ is a constant. Consequently, for all $0 < \rho \leq \rho_0$ and $0 < \theta < \theta_0$ we see that the exponential inequality

$$\mathbb{P}\left([\varepsilon \widehat{\boldsymbol{\xi}}]_{T_\varepsilon}^d > \varepsilon^\rho\right) \leq \exp\left(-\lambda_\varepsilon \varepsilon^\rho + 2CT_\varepsilon \lambda_\varepsilon \varepsilon^2 + 2\Theta(1)T_\varepsilon \lambda_\varepsilon \varepsilon^{2(1-\delta)}\right) \leq \exp(-\varepsilon^{-\rho/2})$$

holds for small enough ε with $p \in (0, \rho/2)$. This is because

$$\begin{aligned} -\lambda_\varepsilon \varepsilon^\rho + 2CT_\varepsilon \lambda_\varepsilon \varepsilon^2 + 2\Theta(1)T_\varepsilon \lambda_\varepsilon \varepsilon^{2(1-\delta)} &= -\varepsilon^{-\rho} + 2C\varepsilon^{2-\frac{1-\delta}{3}-\frac{1-\delta}{2}} + 2\Theta(1)\varepsilon^{2(1-\delta)-\frac{1-\delta}{3}-\frac{1-\delta}{2}} \\ &\leq -\varepsilon^{-\rho} + 2(C + \Theta(1))\varepsilon^{2(1-\delta)-\frac{1-\delta}{3}-\frac{1-\delta}{2}} \leq -\varepsilon^{-\rho} + 2(C + \Theta(1))\varepsilon^{\frac{7}{6}(1-\delta)} \stackrel{\textcircled{1}}{\leq} -\varepsilon^{-\rho/2}, \end{aligned}$$

where $\textcircled{1}$ holds by setting ε enough small such that $(\varepsilon^{-\rho} - 2(C + \Theta(1))\varepsilon^{\frac{7}{6}(1-\delta)})/\varepsilon^{-\rho/2} \geq \varepsilon^{-\rho/2} - 2(C + \Theta(1))\varepsilon^{\frac{7}{6}(1-\delta)+\frac{\rho}{2}} \geq 1$. The proof is completed. \square

E.4 Proof of Lemma 2

Proof. Step 1. Suppose $\sup_{t \geq 0} \|g^t\| \leq c_g$ for some constant $c_g > 0$. Then we consider the one-dimensional martingale

$$M_t = \sum_{i=1}^d \int_0^t g_s^i d\widehat{\boldsymbol{\xi}}_s^i.$$

We estimate the probability of a deviation of the size ε^ρ of εM_t from zero with help of the exponential inequality for martingales, see Theorem 26.17 (i) in [14]. Indeed for any $\rho > 0$ and $\theta > 0$, we have

$$\mathbb{P}\left(\sup_{t \leq T_\varepsilon} |\varepsilon M_t| \geq \varepsilon^\rho\right) \leq \mathbb{P}\left(\sup_{t \leq T_\varepsilon} |\varepsilon M_t| \geq \varepsilon^\rho \mid [\varepsilon M]_{T_\varepsilon} \leq \varepsilon^{4\rho}\right) + \mathbb{P}\left([\varepsilon M]_{T_\varepsilon} > \varepsilon^{4\rho}\right).$$

Inspecting the proofs of Lemma 26.19 and Theorem 26.17 (i) in [14] we get that for any $\lambda > 0$

$$\mathbb{P}\left(\sup_{t \leq T_\varepsilon} |\varepsilon M_t| \geq \varepsilon^\rho \mid [\varepsilon M]_{T_\varepsilon} \leq \varepsilon^{4\rho}\right) \leq \exp\left(-\lambda \varepsilon^\rho + \lambda^2 h(\lambda c_g \varepsilon^{1-\delta}) \varepsilon^{4\rho}\right),$$

where $h(x) = -(x + \ln(1-x))_+ x^{-2}$. For any $0 < \rho < \rho_1 = (1 - \delta)/2$ we set $\lambda = \lambda_\varepsilon = \varepsilon^{-2\rho}$ so that $h(\lambda_\varepsilon c_g \varepsilon^{1-\delta}) \rightarrow 1/2$ as $\varepsilon \rightarrow 0$ by using L'Hopital's rule. Hence we obtain the estimate

$$\mathbb{P}\left(\sup_{t \leq T_\varepsilon} |\varepsilon M_t| \geq \varepsilon^\rho \mid [\varepsilon M]_{T_\varepsilon} \leq \varepsilon^{4\rho}\right) \leq \exp\left(-\varepsilon^{-\rho} + \frac{1}{2}\right) \stackrel{\textcircled{1}}{\leq} \exp\left(-\varepsilon^{-\rho/2}\right) \leq \exp\left(-\varepsilon^{-p}\right),$$

which holds for small enough ε and $p \in (0, \rho/2]$. In $\textcircled{1}$, we set ε enough small such that $0 < \varepsilon^{-\rho/2} - \varepsilon^{-p} \leq 1$.

Step 2. Since $\|g^t\| \leq c_g$ is well bounded, then there is a constant c_1 with

$$[\varepsilon M]_t = \int_0^t g_s^2 d[\varepsilon \widehat{\boldsymbol{\xi}}]_s^d \leq c_1 [\varepsilon \widehat{\boldsymbol{\xi}}]_t^d.$$

Then we can use Lemma 1 to upper bound:

$$\mathbb{P}([\varepsilon M]_{T_\varepsilon} \geq \varepsilon^{4\rho}) \leq \mathbb{P}\left(c_1[\varepsilon \widehat{\boldsymbol{\xi}}]_t^d \geq \varepsilon^{4\rho}\right) \stackrel{\textcircled{1}}{\leq} \exp(-p),$$

where $\textcircled{1}$ uses $\rho < \rho_2 < \frac{\rho_0}{4}$ with $\rho_0 = \frac{1-\delta}{4}$ in Lemma 1 and sets ε sufficient small such that $\varepsilon^{\rho_0-4\rho} \leq c_1$. This is because if $\varepsilon^{\rho_0} \leq \frac{\varepsilon^{4\rho}}{c_1}$, then it yields $\mathbb{P}\left([\varepsilon \widehat{\boldsymbol{\xi}}]_t^d \geq \varepsilon^{4\rho}/c_1\right) \leq \exp(-p)$ due to $\mathbb{P}\left([\varepsilon \widehat{\boldsymbol{\xi}}]_t^d \geq \varepsilon^{\rho_0}\right) \leq \exp(-p)$. So the result in this lemma holds with $\rho_0 = \min(\rho_0 = \frac{1-\delta}{4}, \rho_1, \rho_2) = \frac{1-\delta}{16}$, $p_0 = \min(p_0 = \frac{\rho}{2}, p_1) = \frac{\delta}{2}$. The parameters ρ_0 and p_0 in the operator (\cdot) are from Lemma 1 as the results here is based on Lemma 1. Under this setting, we have

$$\mathbb{P}\left(\sup_{0 \leq t \leq T_\varepsilon} \varepsilon \left| \sum_{i=1}^d \int_0^t \mathbf{g}_{s-}^i d\widehat{\boldsymbol{\xi}}_s^i \right| \geq \varepsilon^\rho\right) \leq 2 \exp(-\varepsilon^{-p})$$

The proof is completed. \square

E.5 Proof of Lemma 3

Proof. Step 1. In this step we prove the sequence $\{\widehat{\boldsymbol{\theta}}_t\}$ produced by Eqn. (8) or (9) locates in a very small neighborhood of the optimum solution $\boldsymbol{\theta}^*$ of the local basin Ω after a very small time interval. Since we assume the function is locally strongly convex, by using Theorems 3 and 4, we know that the sequence $\{\widehat{\boldsymbol{\theta}}_t\}$ produced by Eqn. (8) or (9) exponentially converges to the minimum $\boldsymbol{\theta}^*$ at the current local basin Ω . So for any initialization $\boldsymbol{\theta}_0 \in \Omega$, we have

$$\|\widehat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^*\|_2^2 \leq c_1 \exp(-c_2 t),$$

where $c_1 = \frac{2\Delta}{\mu}$ and $c_2 = \frac{2\mu\tau}{\beta_1(v_{\max}+\varepsilon)+\mu\tau} \left(\beta_1 - \frac{\beta_2}{4}\right)$ in ADAM, $c_1 = \frac{2\Delta}{\mu}$ and $c_2 = 2\mu$ in SGD. Therefore, for any initialization $\boldsymbol{\theta}_0 \in \Omega$ and sufficient small ε , we can obtain

$$\|\widehat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^*\|_2^2 \leq \varepsilon^{\bar{p}} \text{ when } t \geq v_\varepsilon = \frac{1}{c_2} \ln\left(\frac{c_1}{\varepsilon^{\bar{p}}}\right).$$

Step 2. Here we prove that for the time $t \in [0, v_\varepsilon]$, the sequence $\{\boldsymbol{\theta}_t\}$ is always very close to the sequence $\{\widehat{\boldsymbol{\theta}}_t\}$ when they are with the same initialization $\boldsymbol{\theta}_0$ in the absence of the big jumps J_k in the stochastic process L .

To begin with, according to the updating rule in SGD, we have

$$\begin{aligned} \|\boldsymbol{\theta}_{t \wedge v_\varepsilon \wedge \sigma_1} - \widehat{\boldsymbol{\theta}}_{t \wedge v_\varepsilon \wedge \sigma_1}\| &= \left\| \int_0^{t \wedge v_\varepsilon \wedge \sigma_1} \left(-\nabla F(\boldsymbol{\theta}_s) + \nabla F(\widehat{\boldsymbol{\theta}}_s)\right) ds + \int_0^{t \wedge v_\varepsilon \wedge \sigma_1} \varepsilon \boldsymbol{\Sigma}_s dL_s \right\| \\ &\stackrel{\textcircled{1}}{\leq} \ell \int_0^{t \wedge v_\varepsilon \wedge \sigma_1} \|\boldsymbol{\theta}_s - \widehat{\boldsymbol{\theta}}_s\| ds + \varepsilon \left\| \int_0^{t \wedge v_\varepsilon \wedge \sigma_1} \boldsymbol{\Sigma}_s dL_s \right\|, \end{aligned} \tag{28}$$

where in $\textcircled{1}$, $F(\boldsymbol{\theta})$ is ℓ -smooth, namely $\|\nabla F(\boldsymbol{\theta}_1) - \nabla F(\boldsymbol{\theta}_2)\| \leq \ell \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|$ for any $\boldsymbol{\theta}_1$ and $\boldsymbol{\theta}_2$ in the local basin Ω .

Then we consider ADAM which needs more efforts. According to the dynamic system of ADAM, we can first establish

$$\mathbf{m}_t - \widehat{\mathbf{m}}_t = \int_0^t (\nabla F(\boldsymbol{\theta}_s) - \nabla F(\widehat{\boldsymbol{\theta}}_s)) ds - \int_0^t (\mathbf{m}_s - \widehat{\mathbf{m}}_s) ds.$$

Therefore, with the assumption $\|\mathbf{m}_t - \widehat{\mathbf{m}}_t\| \leq \tau_m \|\int_0^t (\mathbf{m}_s - \widehat{\mathbf{m}}_s) ds\|$, it yields

$$\begin{aligned} |1 - \tau_m| \cdot \left\| \int_0^t (\mathbf{m}_s - \widehat{\mathbf{m}}_s) ds \right\| &\leq \left\| \mathbf{m}_t - \widehat{\mathbf{m}}_t + \int_0^t (\mathbf{m}_s - \widehat{\mathbf{m}}_s) ds \right\| = \left\| \int_0^t (\nabla F(\boldsymbol{\theta}_s) - \nabla F(\widehat{\boldsymbol{\theta}}_s)) ds \right\| \\ &\leq \ell \int_0^t \|\boldsymbol{\theta}_s - \widehat{\boldsymbol{\theta}}_s\| ds. \end{aligned}$$

Moreover, we can upper bound $\frac{\mu_s}{\sqrt{\omega_s \mathbf{v}_s + \epsilon}} = \frac{\sqrt{1-e^{-\beta_2 t}}}{1-e^{-\beta_1 t}} \cdot \frac{1}{1+\epsilon \sqrt{1-e^{-\beta_2 t}}}$. Then let $q(x) = \frac{\sqrt{1-e^{-\beta_2 t}}}{1-e^{-\beta_1 t}} \leq c_4 = \min(q(0), q(+\infty), q(t^*))$, where t^* is a time such that $q'(t^*) = 0$. Since $q(0) = \frac{\beta_2}{2\beta_1}$ by L'Hopital's rule, $q(+\infty) = 1$ and $q(t^*) < \infty$ is a constant, $c_4 < \infty$ is a constant. So there exists a constant c_5 such that $\frac{\mu_s}{\sqrt{\omega_s \mathbf{v}_s + \epsilon}} \leq \frac{c_5}{v_{\min} + \epsilon}$. Then similarly, in ADAM, we also can establish

$$\begin{aligned} \|\boldsymbol{\theta}_{t \wedge v_\epsilon \wedge \sigma_1 -} - \widehat{\boldsymbol{\theta}}_{t \wedge v_\epsilon \wedge \sigma_1 -}\| &= \left\| \int_0^{t \wedge v_\epsilon \wedge \sigma_1 -} \left(-\frac{\mu_s \mathbf{m}_s}{\sqrt{\omega_s \mathbf{v}_s + \epsilon}} + \frac{\mu_s \widehat{\mathbf{m}}_s}{\sqrt{\omega_s \widehat{\mathbf{v}}_s + \epsilon}} \right) \mathbf{d}s + \int_0^{t \wedge v_\epsilon \wedge \sigma_1 -} \varepsilon \mathbf{Q}_s^{-1} \boldsymbol{\Sigma}_s \mathbf{d}L_s \right\| \\ &\stackrel{\textcircled{1}}{\leq} \frac{c_5 \ell}{(v_{\min} + \epsilon) |\tau_m - 1|} \int_0^{t \wedge v_\epsilon \wedge \sigma_1 -} \|\boldsymbol{\theta}_s - \widehat{\boldsymbol{\theta}}_s\| \mathbf{d}s + \varepsilon \left\| \int_0^{t \wedge v_\epsilon \wedge \sigma_1 -} \mathbf{Q}_s^{-1} \boldsymbol{\Sigma}_s \mathbf{d}L_s \right\|. \end{aligned}$$

Next, we can employ Gronwall's to estimate

$$\sup_{0 \leq t \leq \sigma_1 \wedge v_\epsilon} \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}_t\| \leq \exp(\kappa_1 v_\epsilon) \sup_{0 \leq t \leq v_\epsilon} \varepsilon \left\| \int_0^t \mathbf{Q}_s^{-1} \boldsymbol{\Sigma}_s \mathbf{d}\boldsymbol{\xi}_s \right\|,$$

where $\kappa_1 = \ell$ in SGD, and $\kappa_1 = \frac{c_5 \ell}{(v_{\min} + \epsilon) |\tau_m - 1|}$ in ADAM. Since when ε is small enough, $v_\epsilon = \frac{1}{c_2} \ln\left(\frac{c_1}{\varepsilon^{\bar{\rho}}}\right)$ is much smaller than $T_\varepsilon = \varepsilon^{-\theta}$ when ε is sufficient small. It yields

$$\begin{aligned} \mathbb{P} \left(\sup_{0 \leq t \leq \sigma_1 \wedge v_\epsilon} \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}_t\| \geq \varepsilon^{\bar{\rho}} \right) &\leq \mathbb{P} \left(\exp(\kappa_1 v_\epsilon) \sup_{0 \leq t \leq v_\epsilon} \varepsilon \left\| \int_0^t \mathbf{Q}_s^{-1} \boldsymbol{\Sigma}_s \mathbf{d}\boldsymbol{\xi}_s \right\| \geq \varepsilon^{\bar{\rho}} \right) \\ &\stackrel{\textcircled{1}}{\leq} \mathbb{P} \left(\sup_{0 \leq t \leq v_\epsilon} \varepsilon \left\| \int_0^t \mathbf{Q}_s^{-1} \boldsymbol{\Sigma}_s \mathbf{d}\widehat{\boldsymbol{\xi}}_s \right\| + \varepsilon \|\mu_\varepsilon\| T_\varepsilon \geq \varepsilon^{\bar{\rho} + c_3 \kappa_1 \bar{\rho}} \right) \\ &= \mathbb{P} \left(\sup_{0 \leq t \leq v_\epsilon} \varepsilon \left\| \int_0^t \mathbf{Q}_s^{-1} \boldsymbol{\Sigma}_s \mathbf{d}\widehat{\boldsymbol{\xi}}_s \right\| \geq \varepsilon^\rho (\varepsilon^{\bar{\rho}(1+c_3 \kappa_1)} - \varepsilon^\rho) \right) \\ &\stackrel{\textcircled{2}}{\leq} \exp(-p), \end{aligned}$$

where $\textcircled{1}$ uses Lemma 1: (1) the process $\boldsymbol{\xi}$ can be decomposed into two processes $\widehat{\boldsymbol{\xi}}$ and linear drift, namely, $\boldsymbol{\xi}_t = \widehat{\boldsymbol{\xi}}_t + \mu_\varepsilon t$, where $\widehat{\boldsymbol{\xi}}$ is a zero mean Lévy martingale with bounded jumps; (2) $\|\varepsilon \boldsymbol{\xi}_{T_\varepsilon}\| = \varepsilon \|\mu_\varepsilon\| T_\varepsilon < \varepsilon^{2\rho}$. In $\textcircled{2}$, (1) we set $\bar{\rho}(1 + c_3 \kappa_1) < \rho$ and also set ε sufficient small such that $\varepsilon^{\bar{\rho}(1+c_3 \kappa_1)} - \varepsilon^\rho \geq 1$; (2) by assume $\rho_0 = \rho_0(\delta) = \frac{1-\delta}{16} > 0$, $\theta_0 = \theta_0(\delta) = \frac{1-\delta}{3} > 0$ and $p_0 = p_0(\rho) = \frac{\rho}{2}$, we use Lemma 2 by setting $\mathbf{g}^t = \mathbf{Q}_t^{-1} \boldsymbol{\Sigma}_t$ and obtain $\mathbb{P} \left(\sup_{0 \leq t \leq T_\varepsilon} \varepsilon \left\| \sum_{i=1}^d \int_0^t \mathbf{g}_s^i - \mathbf{d}\widehat{\boldsymbol{\xi}}_s^i \right\| \geq \varepsilon^\rho \right) \leq 2 \exp(-\varepsilon^{-p})$ for all $p \in (0, p_0]$ and $0 < \varepsilon \leq \varepsilon_0$ with $\varepsilon_0 = \varepsilon_0(\rho)$.

Step 3. In the first step, we have analyzed that the sequence $\{\widehat{\boldsymbol{\theta}}_t\}$ will converge to the optimum $\boldsymbol{\theta}^*$ of the basin Ω . Moreover, in the second step, we prove that $\boldsymbol{\theta}_t$ is very close to $\widehat{\boldsymbol{\theta}}_t$. In this step, we show that in absence of the big jumps of the driving process L the sequence $\boldsymbol{\theta}_t$ is close to $\boldsymbol{\theta}^*$. For brevity, we set $\boldsymbol{\theta}^* = \mathbf{0}$. Then we define a function $h(\boldsymbol{\theta}) = \ln(1 + \mathbf{F}(\boldsymbol{\theta})) \geq 0$. Since for a small local convex basin Ω , the function $\mathbf{F}(\boldsymbol{\theta})$ can be well approximated by a quadratic function. In this way, for small $\boldsymbol{\theta}$ one can always estimate $c_6 \|\boldsymbol{\theta}\|^2 \leq h(\boldsymbol{\theta}) \leq c_7 \|\boldsymbol{\theta}\|^2$ for some positive constants c_6 and c_7 . Furthermore, the derivatives $\partial_i h(\boldsymbol{\theta}) = \frac{\partial_i \mathbf{F}(\boldsymbol{\theta})}{1 + \mathbf{F}(\boldsymbol{\theta})}$ and $\partial_i \partial_j h(\boldsymbol{\theta}) = \frac{\partial_{ij} \mathbf{F}(\boldsymbol{\theta})(1 + \mathbf{F}(\boldsymbol{\theta})) - \partial_i \mathbf{F}(\boldsymbol{\theta}) \partial_j \mathbf{F}(\boldsymbol{\theta})}{(1 + \mathbf{F}(\boldsymbol{\theta}))^2}$ are bounded since the assumptions on the function $\mathbf{F}(\boldsymbol{\theta})$, namely $\mathbf{F}(\boldsymbol{\theta})$ being upper bounded, ℓ -smooth. Next we can apply the Itô formulation to the process $h(\boldsymbol{\theta}_t)$:

$$\begin{aligned} 0 \leq h(\boldsymbol{\theta}_{t \wedge T_\varepsilon \wedge \sigma_1 -}) &= h(\boldsymbol{\theta}) + \sum_{i=1}^d \int_0^{t \wedge T_\varepsilon \wedge \sigma_1 -} \partial_i h(\boldsymbol{\theta}_{s-}) \mathbf{d}\boldsymbol{\theta}_{s-}^i + \frac{1}{2} \sum_{i,j=1}^d \int_0^{t \wedge T_\varepsilon \wedge \sigma_1 -} \partial_i \partial_j h(\boldsymbol{\theta}_{s-}) \mathbf{d}[\boldsymbol{\theta}^i, \boldsymbol{\theta}^j]_s^c \\ &\quad + \sum_{s < t \wedge T_\varepsilon \wedge \sigma_1} \left(h(\boldsymbol{\theta}_s) - h(\boldsymbol{\theta}_{s-}) - \sum_{i=1}^d \partial_i h(\boldsymbol{\theta}_{s-}) \Delta \boldsymbol{\theta}_s^i \right) \\ &\stackrel{\textcircled{1}}{\leq} h(\boldsymbol{\theta}) - \int_0^{t \wedge T_\varepsilon \wedge \sigma_1 -} \left\langle \frac{\nabla \mathbf{F}(\boldsymbol{\theta}_{s-})}{1 + \mathbf{F}(\boldsymbol{\theta}_{s-})}, \frac{\mu_t \mathbf{m}_s}{\epsilon + \sqrt{\omega_s - \mathbf{v}_s}} \right\rangle \mathbf{d}s + \int_0^{t \wedge T_\varepsilon \wedge \sigma_1 -} \frac{\varepsilon (\nabla \mathbf{F}(\boldsymbol{\theta}_{s-}))^T \mathbf{Q}_{s-}^{-1} \boldsymbol{\Sigma}_{s-}}{1 + \mathbf{F}(\boldsymbol{\theta}_{s-})} \mathbf{d}L_s \\ &\quad + \sum_{s < t \wedge T_\varepsilon \wedge \sigma_1} \left(h(\boldsymbol{\theta}_s) - h(\boldsymbol{\theta}_{s-}) - \sum_{i=1}^d \partial_i h(\boldsymbol{\theta}_{s-}) \Delta \boldsymbol{\theta}_s^i \right), \end{aligned}$$

where ① uses $\mathbf{d}\theta_s = -\frac{\mu_s \mathbf{m}_s}{\epsilon + \sqrt{\omega_s \mathbf{v}_s}} + \epsilon \mathbf{Q}_s^{-1} \Sigma_s \mathbf{d}L_s$ and the path-by-path continuous part $[\theta^i, \theta^j]_s^c = 0$ of the quadratic covariation of θ^i and θ^j . Since in Adam by assumption $\int_0^{t \wedge T_\epsilon \wedge \sigma_1} \left\langle \frac{\nabla \mathbf{F}(\theta_{s-})}{1 + \mathbf{F}(\theta_{s-})}, \frac{\mu_t \mathbf{m}_{s-}}{\epsilon + \sqrt{\omega_{s-} \mathbf{v}_{s-}}} \right\rangle \mathbf{d}s \geq 0$, the second term is non-negative due to $\mathbf{F}(\theta) \geq 0$. Note in SGD, $\mathbf{m}_s = \nabla \mathbf{F}(\theta_s)$. So in SGD we do not make the assumption $\langle \nabla \mathbf{F}(\theta_t), \mathbf{m}_t \rangle \geq 0$. In SGD, $\epsilon + \sqrt{\omega_{s-} \mathbf{v}_{s-}}$ equals to one. In this way, we can estimate the last term as

$$\begin{aligned} & \sum_{s < t \wedge T_\epsilon \wedge \sigma_1} \left| h(\theta_s) - h(\theta_{s-}) - \sum_{i=1}^d \partial_i h(\theta_{s-}) \Delta \theta_s^i \right| \\ & \leq \frac{1}{2} \sum_{i,j=1}^d \sum_{s < t \wedge T_\epsilon \wedge \sigma_1} \left| \int_0^1 (1-v) \partial_i \partial_j h(\theta_{s-} + v \Delta \theta_s) \mathbf{d}v \right| \cdot |\Delta \theta_s^i \Delta \theta_s^j| \leq c_8 \sum_{s \leq t} \|\Delta \theta_s\|^2 = c_8 [\theta]_t^d, \end{aligned}$$

holds with some $c_8 > 0$. Furthermore, since \mathbf{v}_t and Σ_t are assumed to be bounded, then we can upper bound $[\theta]_t^d$ as follows:

$$[\theta]_t^d \leq c_9 [\varepsilon L]_t^d \stackrel{\text{①}}{=} c_9 [\varepsilon \xi]_t^d$$

hold for some constant c_9 for all $t \leq \sigma_1$. ① holds since we assume there is no big jump during $t \leq \sigma_1$. Then by combining all the results and letting $\mathbf{g}^s = \frac{(\nabla \mathbf{F}(\theta_s))^T \mathbf{Q}_s^{-1} \Sigma_s}{1 + \mathbf{F}(\theta_s)}$ and considering $\mathbf{F}(\theta) \leq c_7 \|\theta\|$, we can obtain the following results when $\|\theta\| = \|\theta_0\| \leq \varepsilon^{\bar{\rho}}$ with enough small ε :

$$0 \leq \|\theta_{t \wedge T_\epsilon \wedge \sigma_1}\|^2 \leq \frac{1}{c_6} h(\theta_{t \wedge T_\epsilon \wedge \sigma_1}) \leq c_{10} \left(\varepsilon^{2\bar{\rho}} + \varepsilon \sup_{0 \leq t \leq T_\epsilon} \left| \int_0^t \mathbf{g}^{s-} \widehat{\mathbf{d}}\xi_s \right| + \varepsilon \|\mu_\varepsilon\|_{T_\epsilon} + \varepsilon^2 [\xi]_{T_\epsilon}^d \right).$$

where c_{10} is a certain constant. Combining the above results gives

$$\begin{aligned} \mathbb{P} \left(\sup_{0 \leq t \leq T_\epsilon \wedge \sigma_1} \|\theta_t\| \geq \varepsilon^{\bar{\rho}} \right) & \leq \mathbb{P} \left(\varepsilon^{2\bar{\rho}} \geq \frac{\varepsilon^{\bar{\rho}}}{4c_{10}} \right) + \mathbb{P} \left(\varepsilon \sup_{0 \leq t \leq T_\epsilon} \left| \int_0^t \mathbf{g}^{s-} \widehat{\mathbf{d}}\xi_s \right| \geq \frac{\varepsilon^{\bar{\rho}}}{4c_{10}} \right) \\ & + \mathbb{P} \left(\varepsilon \|\mu_\varepsilon\|_{T_\epsilon} \geq \frac{\varepsilon^{\bar{\rho}}}{4c_{10}} \right) + \mathbb{P} \left(\varepsilon^2 [\xi]_{T_\epsilon}^d \geq \frac{\varepsilon^{\bar{\rho}}}{4c_{10}} \right). \end{aligned}$$

Then by setting $\bar{\rho} < \rho$ and sufficient small ε such that $\frac{\varepsilon^{\bar{\rho}-\rho}}{4c_{10}} \geq 1$ giving $\frac{\varepsilon^{\bar{\rho}}}{4c_{10}} \geq \varepsilon^\rho$. Then let the results in Lemma 1 and 2 hold simultaneously by setting $\rho_0 = \rho_0(\delta) = \frac{1-\delta}{16} > 0$, $\theta_0 = \theta_0(\delta) = \frac{1-\delta}{3} > 0$, $p_0 = \frac{\rho}{2}$, and small enough ε , we have $\|\varepsilon \xi_{T_\epsilon}\| = \varepsilon \|\mu_\varepsilon\|_{T_\epsilon} < \varepsilon^{2\rho}$ and $\mathbb{P}([\varepsilon \xi]_{T_\epsilon}^d \geq \varepsilon^\rho) \leq \exp(-\varepsilon^{-p})$ in Lemma 1, and $\mathbb{P}(\sup_{0 \leq t \leq T_\epsilon} \varepsilon \left| \sum_{i=1}^d \int_0^t \mathbf{g}_s^i \widehat{\mathbf{d}}\xi_s^i \right| \geq \varepsilon^\rho) \leq 2 \exp(-\varepsilon^{-p})$ in Lemma 2. By using these results, we have

$$\mathbb{P} \left(\sup_{0 \leq t \leq T_\epsilon \wedge \sigma_1} \|\theta_t\| \geq \varepsilon^{\bar{\rho}} \right) \leq 4 \exp(-\varepsilon^{-p}).$$

for all $p \in (0, p_0]$ and $0 < \varepsilon \leq \varepsilon_0$ with $\varepsilon_0 = \varepsilon_0(\rho)$.

Step 4. In Steps 1 and 2, we guarantee $\mathbb{P}(\sup_{0 \leq t \leq v_\varepsilon \wedge \sigma_1} \|\theta_t - \hat{\theta}_t\| \geq \varepsilon^{\bar{\rho}}) \leq 4 \exp(-\varepsilon^{-p})$. Then after v_ε time, we have $\|\theta_t\| \leq \varepsilon^{\bar{\rho}}$ for all $t \geq v_\varepsilon$. In this way, the result in Step 4 holds. So in this step, we combine the results in Steps 1, 2 and 3 and extend the initialization in Step 3 to all possible parameter in $\theta_0 \in \Omega$:

$$\mathbb{P} \left(\sup_{0 \leq t \leq v_\varepsilon \wedge \sigma_1} \|\theta_t - \hat{\theta}_t\| \geq \varepsilon^{\bar{\rho}} \right) \leq 4 \exp(-\varepsilon^{-p}),$$

for all $p \in (0, p_0]$ and $0 < \varepsilon \leq \varepsilon_0$ with $\varepsilon_0 = \varepsilon_0(\rho)$ by setting $\rho_0 = \rho_0(\delta) = \frac{1-\delta}{16} > 0$, $\theta_0 = \theta_0(\delta) = \frac{1-\delta}{3} > 0$, $p_0 = \frac{\rho}{2}$, $\bar{\rho}(1 + c_3 \kappa_1) < \rho$ and small enough ε . Note here we can remove the extra factor ρ by setting $\varepsilon_0 = \varepsilon_0(\bar{\rho})$, $\rho_0 = \rho_0(\delta) = \frac{1-\delta}{16(1+c_3 \kappa_1)} > 0$, $\theta_0 = \theta_0(\delta) = \frac{1-\delta}{3} > 0$, $p_0 = \frac{\bar{\rho}(1+c_3 \kappa_1)}{2}$, $p \in (0, p_0]$.

Step 5. In this step, we extend the result in Step 4 from the time interval $[0, T_\epsilon \wedge \sigma_1]$ to the time interval $[0, \sigma_1]$.

Let $\boldsymbol{\theta}_t^\xi$ denote the sequence produced by SGD (4) or Adam (5) driven by the process $\boldsymbol{\xi}$. Then it is easy to check that for any $t < \sigma_1$, we have $\boldsymbol{\theta}_t^\xi = \boldsymbol{\theta}_t$, since there are no big jumps in $\boldsymbol{\theta}_t$. Then consider any $\boldsymbol{\theta}_0 \in \Omega$ and $k \geq 1$, we have for any $\bar{\rho} > 0$ and $\theta > 0$

$$\begin{aligned} & \mathbb{P} \left(\sup_{0 \leq t < \sigma_1} \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}_t\| \geq \varepsilon^{\bar{\rho}} \right) \\ & \leq \mathbb{P} \left(\sup_{0 \leq t < kT_\varepsilon \wedge \sigma_1} \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}_t\| \geq \varepsilon^{\bar{\rho}} \mid kT_\varepsilon < \sigma_1 \right) + \mathbb{P} \left(\sup_{0 \leq t < \sigma_1} \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}_t\| \geq \varepsilon^{\bar{\rho}} \mid kT_\varepsilon \geq \sigma_1 \right) \\ & \leq \mathbb{P} \left(\sup_{0 \leq t < kT_\varepsilon \wedge \sigma_1} \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}_t\| \geq \varepsilon^{\bar{\rho}} \right) + \mathbb{P}(kT_\varepsilon \geq \sigma_1) \\ & \leq \mathbb{P} \left(\sup_{0 \leq t < kT_\varepsilon} \|\boldsymbol{\theta}_t - \widehat{\boldsymbol{\theta}}_t\| \geq \varepsilon^{\bar{\rho}} \right) + \mathbb{P}(kT_\varepsilon \geq \sigma_1). \end{aligned}$$

Besides, by using the linear convergence results of $\widehat{\boldsymbol{\theta}}_t$ to the optimum solution $\boldsymbol{\theta}^* = \mathbf{0}$ in the local basin Ω , for enough small ε we have $\|\widehat{\boldsymbol{\theta}}_{T_\varepsilon}\| \leq \varepsilon^{2\bar{\rho}}$ with initialization $\boldsymbol{\theta}_0 \in \Omega$. Then we let $\widehat{\boldsymbol{\theta}}_t(\boldsymbol{\theta})$ denote the sequence $\widehat{\boldsymbol{\theta}}_t$ but with initialization $\boldsymbol{\theta}$ and define

$$\mathbf{E}_i = \left\{ \sup_{t \in [iT_\varepsilon, (i+1)T_\varepsilon]} \|\boldsymbol{\theta}_t^\xi - \widehat{\boldsymbol{\theta}}_{t-iT_\varepsilon}(\boldsymbol{\theta}_{iT_\varepsilon}^\xi)\| < \varepsilon^{\bar{\rho}} \right\}, \quad 0 \leq i \leq k-1.$$

Note that the probability of $\mathbf{E}_0^c = \left\{ \sup_{t \in [0, T_\varepsilon]} \|\boldsymbol{\theta}_t^\xi - \widehat{\boldsymbol{\theta}}_t(\boldsymbol{\theta}_0^\xi)\| \geq \varepsilon^{\bar{\rho}} \right\}$ is given in Step 4 where $\boldsymbol{\theta}_0^\xi = \boldsymbol{\theta}_0$. Furthermore for any $k \geq 1$, we have

$$\bigcap_{i=0}^{k-1} \mathbf{E}_i \subseteq \left\{ \sup_{t \in [0, kT_\varepsilon]} \|\boldsymbol{\theta}_t^\xi - \widehat{\boldsymbol{\theta}}_t\| < 2\varepsilon^{\bar{\rho}} \right\}.$$

As a result, we can obtain

$$\begin{aligned} \mathbb{P} \left(\sup_{t \in [0, kT_\varepsilon]} \|\boldsymbol{\theta}_t^\xi - \widehat{\boldsymbol{\theta}}_t\| \geq 2\varepsilon^{\bar{\rho}} \right) & \leq \mathbb{P} \left(\bigcup_{i=0}^{k-1} \mathbf{E}_i^c \right) = \mathbb{P} \left(\mathbf{E}_0^c \cup (\mathbf{E}_0 \mathbf{E}_1^c) \cup (\mathbf{E}_0 \mathbf{E}_1 \mathbf{E}_2^c) \cup \dots \cup \left(\bigcup_{i=0}^{k-2} \mathbf{E}_i \mathbf{E}_{k-1}^c \right) \right) \\ & \leq \sum_{i=0}^{k-1} \mathbb{P} \left(\mathbf{E}_i^c, \boldsymbol{\theta}_{iT_\varepsilon}^\xi \in \Omega \right) \leq k \sup_{\boldsymbol{\theta}_0 \in \Omega} \mathbb{P}(\mathbf{E}_0^c). \end{aligned}$$

For $k = k_\varepsilon = \varepsilon^{-2r}$ and any $\theta > 0$ we have

$$\mathbb{P}(\sigma_1 \geq k_\varepsilon T_\varepsilon) = \exp(-k_\varepsilon T_\varepsilon \Theta(\varepsilon^{-\delta})) \leq \exp(-\varepsilon^{r\delta - \theta - 2r} \Theta(\varepsilon^{-\delta})) \leq \exp(-\varepsilon^{-p})$$

for all $0 < p \leq (2 - \delta)r$ with enough small ε . On the other hand, we have

$$\mathbb{P} \left(\sup_{t \in [0, kT_\varepsilon]} \|\boldsymbol{\theta}_t^\xi - \widehat{\boldsymbol{\theta}}_t\| \geq 2\varepsilon^{\bar{\rho}} \right) \leq k \sup_{\boldsymbol{\theta}_0 \in \Omega} \mathbb{P}(\mathbf{E}_0^c) \leq \varepsilon^{-2r} \exp(-\varepsilon^{-p}) \leq \exp(-\varepsilon^{-p/2})$$

for any $p \leq \frac{2 \log(r \log(\varepsilon))}{\log(\varepsilon)}$. Therefore, the result in this lemma holds

$$\begin{aligned} & \mathbb{P} \left(\sup_{t \in [0, \sigma_1]} \|\boldsymbol{\theta}_t^\xi - \widehat{\boldsymbol{\theta}}_t\| \geq 2\varepsilon^{\bar{\rho}} \right) \\ & = \mathbb{P} \left(\sup_{t \in [0, \sigma_1]} \|\boldsymbol{\theta}_t^\xi - \widehat{\boldsymbol{\theta}}_t\| \geq 2\varepsilon^{\bar{\rho}}, \sigma_1 < kT_\varepsilon \right) + \mathbb{P} \left(\sup_{t \in [0, \sigma_1]} \|\boldsymbol{\theta}_t^\xi - \widehat{\boldsymbol{\theta}}_t\| \geq 2\varepsilon^{\bar{\rho}}, \sigma_1 \geq kT_\varepsilon \right) \\ & \leq \mathbb{P} \left(\sup_{t \in [0, kT_\varepsilon]} \|\boldsymbol{\theta}_t^\xi - \widehat{\boldsymbol{\theta}}_t\| \geq 2\varepsilon^{\bar{\rho}} \right) + \mathbb{P}(\sigma_1 \geq kT_\varepsilon) \leq 2 \exp(-\varepsilon^{-p/2}). \end{aligned}$$

for all $p \in (0, p_0]$ and $0 < \varepsilon \leq \varepsilon_0$ with $\varepsilon_0 = \varepsilon_0(\bar{\rho})$ by setting $\rho_0 = \rho_0(\delta) = \frac{1-\delta}{16(1+c_3\kappa_1)} > 0$, $\theta_0 = \theta_0(\delta) = \frac{1-\delta}{3} > 0$, $p_0 = \min\left(\frac{\bar{\rho}(1+c_3\kappa_1)}{2}, p\right)$ with $p > 0$ and small enough ε . Besides, we also require $v_\varepsilon = \frac{1}{c_2} \ln\left(\frac{c_1}{\varepsilon^{\bar{\rho}}}\right) = \frac{1}{c_2} \ln\left(\frac{2\Delta}{\mu\varepsilon^{\bar{\rho}}}\right) \leq \varepsilon^{-\theta_0}$ where $c_1 = \frac{2\Delta}{\mu}$ and $c_2 = \frac{2\mu\tau}{\beta_1(v_{\max} + \varepsilon) + \mu\tau} \left(\beta_1 - \frac{\beta_2}{4}\right)$ in ADAM, $c_1 = \frac{2\Delta}{\mu}$ and $c_2 = 2\mu$ in SGD. That is, The proof is completed. \square

References

- [1] 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. 1, 3
- [2] 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. 1
- [3] A. Bishop and P. Del Moral. Stability properties of systems of linear stochastic differential equations with random coefficients. *SIAM Journal on Control and Optimization*, 57(2):1023–1042, 2019. 1, 2
- [4] A. Kohatsu-Higa, J. León, and D. Nualart. Stochastic differential equations with random coefficients. *Bernoulli*, 3(2):233–245, 1997. 1, 2
- [5] Y. Fang and K. Loparo. Stabilization of continuous-time jump linear systems. *IEEE Transactions on Automatic Control*, 47(10):1590–1603, 2002. 2
- [6] Andrew EB Lim and Xun Yu Zhou. Mean-variance portfolio selection with random parameters in a complete market. *Mathematics of Operations Research*, 27(1):101–120, 2002. 2
- [7] Stephen J Turnovsky. Optimal stabilization policies for deterministic and stochastic linear economic systems. *The Review of Economic Studies*, 40(1):79–95, 1973. 2
- [8] Jawahar Lal Tiwari and John E Hobbie. Random differential equations as models of ecosystems: Monte carlo simulation approach. *Mathematical Biosciences*, 28(1-2):25–44, 1976. 2
- [9] Chris P Tsokos and William J Padgett. *Random integral equations with applications to life sciences and engineering*. Academic Press, 1974. 2
- [10] Brad A Finney, David S Bowles, and Michael P Windham. Random differential equations in river water quality modeling. *Water resources research*, 18(1):122–134, 1982. 2
- [11] 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. 6
- [12] T. Gronwall. Note on the derivatives with respect to a parameter of the solutions of a system of differential equations. *Annals of Mathematics*, pages 292–296, 1919. 12
- [13] A. Papapantoleon. An introduction to lévy vy processes with applications in finance. *arXiv preprint arXiv:0804.0482*, 2008. 12, 16
- [14] O. Kallenberg. *Foundations of modern probability*. Springer Science & Business Media, 2006. 17