

ON GENERALIZATION ERROR BOUNDS OF NOISY GRADIENT METHODS FOR NON-CONVEX LEARNING

Jian Li*
Tsinghua University

Xuanyuan Luo†
Tsinghua University

Mingda Qiao‡
Stanford University

ABSTRACT

Generalization error (also known as the out-of-sample error) measures how well the hypothesis learned from training data generalizes to previously unseen data. Proving tight generalization error bounds is a central question in statistical learning theory. In this paper, we obtain generalization error bounds for learning general non-convex objectives, which has attracted significant attention in recent years. We develop a new framework, termed *Bayes-Stability*, for proving *algorithm-dependent* generalization error bounds. The new framework combines ideas from both the PAC-Bayesian theory and the notion of algorithmic stability. Applying the Bayes-Stability method, we obtain new data-dependent generalization bounds for stochastic gradient Langevin dynamics (SGLD) and several other noisy gradient methods (e.g., with momentum, mini-batch and acceleration, Entropy-SGD). Our result recovers (and is typically tighter than) a recent result in Mou et al. (2018) and improves upon the results in Pensia et al. (2018). Our experiments demonstrate that our data-dependent bounds can distinguish randomly labelled data from normal data, which provides an explanation to the intriguing phenomena observed in Zhang et al. (2017a). We also study the setting where the total loss is the sum of a bounded loss and an additional ℓ_2 regularization term. We obtain new generalization bounds for the continuous Langevin dynamic in this setting by developing a new Log-Sobolev inequality for the parameter distribution at any time. Our new bounds are more desirable when the noise level of the process is not very small, and do not become vacuous even when T tends to infinity.

1 INTRODUCTION

Non-convex stochastic optimization is the major workhorse of modern machine learning. For instance, the standard supervised learning on a model class parametrized by \mathbb{R}^d can be formulated as the following optimization problem:

$$\min_{w \in \mathbb{R}^d} \mathbb{E}_{z \sim \mathcal{D}} [F(w, z)],$$

where w denotes the model parameter, \mathcal{D} is an unknown data distribution over the instance space \mathcal{Z} , and $F : \mathbb{R}^d \times \mathcal{Z} \rightarrow \mathbb{R}$ is a given objective function which may be non-convex. A learning algorithm takes as input a sequence $S = (z_1, z_2, \dots, z_n)$ of n data points sampled i.i.d. from \mathcal{D} , and outputs a (possibly randomized) parameter configuration $\hat{w} \in \mathbb{R}^d$.

A fundamental problem in learning theory is to understand the *generalization performance* of learning algorithms—is the algorithm guaranteed to output a model that generalizes well to the data distribution \mathcal{D} ? Specifically, we aim to prove upper bounds on the *generalization error* $\text{err}_{\text{gen}}(S) = \mathcal{L}(\hat{w}, \mathcal{D}) - \mathcal{L}(\hat{w}, S)$, where $\mathcal{L}(\hat{w}, \mathcal{D}) = \mathbb{E}_{z \sim \mathcal{D}}[\mathcal{L}(\hat{w}, z)]$ and $\mathcal{L}(\hat{w}, S) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\hat{w}, z_i)$ are the population and empirical losses, respectively. We note that the loss function \mathcal{L} (e.g., the 0/1 loss) could be different from the objective function F (e.g., the cross-entropy loss) used in the training process (which serves as a surrogate for the loss \mathcal{L}).

*lijian83@mail.tsinghua.edu.cn

†luo-xy19@mails.tsinghua.edu.cn

‡mqiao@stanford.edu

Classical learning theory relates the generalization error to various complexity measures (e.g., the VC-dimension and Rademacher complexity) of the model class. Directly applying these classical complexity measures, however, often fails to explain the recent success of over-parametrized neural networks, where the model complexity significantly exceeds the amount of available training data (see e.g., Zhang et al. (2017a)). By incorporating certain data-dependent quantities such as margin and compressibility into the classical framework, some recent work (e.g., Bartlett et al. (2017); Arora et al. (2018); Wei & Ma (2019)) obtains more meaningful generalization bounds in the deep learning context.

An alternative approach to generalization is to provide algorithm-dependent bounds. One celebrated example along this line is the algorithmic stability framework initiated by Bousquet & Elisseeff (2002). Roughly speaking, the generalization error can be bounded by the stability of the algorithm (see Section 2 for the details). Using this framework, Hardt et al. (2016) study the stability (hence the generalization) of stochastic gradient descent (SGD) for both convex and non-convex functions. Their work motivates recent study of the generalization performance of several other gradient-based optimization methods: Kuzborskij & Lampert (2018); London (2016); Chaudhari et al. (2017); Raginsky et al. (2017); Mou et al. (2018); Pensia et al. (2018); Chen et al. (2018).

In this paper, we study the algorithmic stability and generalization performance of various iterative gradient-based method, with certain continuous noise injected in each iteration, in a non-convex setting. As a concrete example, we consider the stochastic gradient Langevin dynamics (SGLD) (see Raginsky et al. (2017); Mou et al. (2018); Pensia et al. (2018)). Viewed as a variant of SGD, SGLD adds an isotropic Gaussian noise at every update step:

$$W_t = W_{t-1} - \eta g_t(W_{t-1}) + \sqrt{\frac{\eta}{2}} N(0; I_d); \quad (1)$$

where $g_t(W_{t-1})$ denotes either the full gradient or the gradient over a mini-batch sampled from training dataset. We also study a continuous version of (1), which is the dynamic defined by the following stochastic differential equation (SDE):

$$dW_t = -\eta F(W_t) dt + \sqrt{\frac{\eta}{2}} dB_t; \quad (2)$$

where B_t is the standard Brownian motion.

1.1 RELATED WORK

Most related to our work is the study of algorithm-dependent generalization bounds of stochastic gradient methods. Hardt et al. (2016) first study the generalization performance of SGD via algorithmic stability. They prove a generalization bound that scales linearly with the number of iterations, when the loss function is convex, but their results for general non-convex optimization are more restricted. London (2017) and Rivasplata et al. (2018) also combine ideas from both PAC-Bayesian and algorithm stability. However, these works are essentially different from ours. In London (2017), the prior and posterior are distributions on the hyperparameter space instead of distributions on the hypothesis space. Rivasplata et al. (2018) study the hypothesis stability measured by the distance on the hypothesis space in a setting where the returned hypothesis (model parameter) is perturbed by a Gaussian noise. Our work is a follow-up of the recent work by Mou et al. (2018), in which they provide generalization bounds for SGLD from both stability and PAC-Bayesian perspectives. Another closely related work by Pensia et al. (2018) derives similar bounds for noisy stochastic gradient methods, based on the information theoretic framework of Xu & Raginsky (2017). However, their bounds scale as $\frac{1}{\sqrt{n}}$ (n is the size of the training dataset) and are sub-optimal even for SGLD.

We acknowledge that besides the algorithm-dependent approach that we follow, recent advances in learning theory aim to explain the generalization performance of neural networks from many other perspectives. Some of the most prominent ideas include bounding the network capacity by the norms of weight matrices Neyshabur et al. (2015); Liang et al. (2019), margin theory Bartlett et al. (2017); Wei et al. (2019), PAC-Bayesian theory Dziugaite & Roy (2017); Neyshabur et al. (2018); Dziugaite & Roy (2018), network compressibility Arora et al. (2018), and over-parametrization Du et al. (2019); Allen-Zhu et al. (2019); Zou et al. (2018); Chizat et al. (2019). Most of these results are stated in the context of neural networks (some are tailored to networks with specific architecture), whereas our work addresses generalization in non-convex stochastic optimization in general. We

also note that some recent work provides explanations for the phenomenon reported in Zhang et al. (2017a) from a variety of different perspectives (e.g., Bartlett et al. (2017); Arora et al. (2018; 2019)).

Welling & Teh (2011) first consider stochastic gradient Langevin dynamics (SGLD) as a sampling algorithm in the Bayesian inference context. Raginsky et al. (2017) give a non-asymptotic analysis and establish the finite-time convergence guarantee of SGLD to an approximate global minimum. Zhang et al. (2017b) analyze the hitting time of SGLD and prove that SGLD converges to an approximate local minimum. These results are further improved and generalized to a family of Langevin dynamics based algorithms by the subsequent work of Xu et al. (2018).

1.2 OVERVIEW OF OUR RESULTS

In this paper, we provide generalization guarantees for the noisy variants of several popular stochastic gradient methods.

The Bayes-Stability method and data-dependent generalization bounds. We develop a new method for proving generalization bounds, termed as Bayes-Stability, by incorporating ideas from the PAC-Bayesian theory into the stability framework. In particular, assuming the loss takes value in $[0, C]$, our method shows that the generalization error is bounded by $2C\sqrt{\frac{1}{n} \mathbb{E}_P[\sum_{t=1}^n \mathbb{E}_{S_t} \|\mathbf{g}_e(t)\|^2]}$ and $2C \mathbb{E}_z[\sqrt{2\text{KL}(Q_z; P)}]$, where P is a prior distribution independent of the training set and Q_z is the expected posterior distribution conditioned on z (i.e., the last training data z). The formal definition and the results can be found in Definition 5 and Theorem 7.

Inspired by Lever et al. (2013), instead of using a fixed prior distribution, we bound the KL-divergence from the posterior to a distribution-dependent prior. This enables us to derive the following generalization error bound that depends on the expected norm of the gradient along the optimization path:

$$\text{err}_{\text{gen}} = O\left(\frac{C}{n} \sqrt{\mathbb{E}_S \sum_{t=1}^n \mathbb{E}_{P_t} \|\mathbf{g}_e(t)\|^2}\right) \quad (3)$$

Here S is the dataset and $\mathbf{g}_e(t) = \mathbb{E}_{W_{t-1}} \left[\frac{1}{n} \sum_{i=1}^n \nabla F(W_{t-1}; z_i) \right]^2$ is the expected squared gradient norm at step t ; see Theorem 11 for the details.

Compared with the previous $\frac{LC}{n} \sqrt{\frac{1}{n} \sum_{t=1}^n \mathbb{E}_P \|\mathbf{g}_e(t)\|^2}$ bound in (Mou et al., 2018, Theorem 1), where L is the global Lipschitz constant of the loss, our new bound (3) depends on the data distribution and is typically tighter (as the gradient norm is at most L). In modern deep neural networks, the worst-case Lipschitz constant can be quite large, and typically much larger than the expected empirical gradient norm along the optimization trajectory. Specifically, in the later stage of the training, the expected empirical gradient is small (see Figure 1(d) for the details). Hence, our generalization bound does not grow much even if we train longer at this stage.

Our new bound also offers an explanation to the difference between training on correct and random labels observed by Zhang et al. (2017a). In particular, we show empirically that the sum of expected squared gradient norm (along the optimization path) is significantly higher when the training labels are replaced with random labels (Section 3.1, Figure 1, Appendix C.2).

We would also like to mention the PAC-Bayesian bound (for SGLD with regularization) proposed by Mou et al. (2018). (This bound is different from what we mentioned before; see Theorem 2 in their paper.) Their bound scales as $\frac{1}{n}$ and the numerator of their bound has a similar sum of gradient norms (with a decaying weight if the regularization coefficient is not 0). Their bound is based on the PAC-Bayesian approach and holds with high probability, while our bound only holds in expectation.

Extensions. We remark that our technique allows for an arguably simpler proof of (Mou et al., 2018, Theorem 1); the original proof is based on SDE and Fokker-Planck equation. More importantly, our technique can be easily extended to handle mini-batches and a variety of general settings as follows.

1. Extension to other gradient-based methods Our results naturally extends to other noisy stochastic gradient methods including momentum due to Polyak (1964) (Theorem 26), Nes-

terov's accelerated gradient method in Nesterov (1983) (Theorem 26), and Entropy-SGD proposed by Chaudhari et al. (2017) (Theorem 27).

2. Extension to general noises. The proof of the generalization bound in Mou et al. (2018) relies heavily on that the noise is Gaussian, which makes it difficult to generalize to other noise distributions such as the Laplace distribution. In contrast, our analysis easily carries over to the class of log-Lipschitz noises (i.e., noises drawn from distributions with Lipschitz log densities).
3. Pathwise stability. In practice, it is also natural to output a certain function of the entire optimization path, e.g., the one with the smallest empirical risk or a weighted average. We show that the same generalization bound holds for all such variants (Remark 12). We note that the analysis in an independent work of Pensia et al. (2018) also satisfies this property, yet their bound is $O(C^2 L^2 n^{-1} \sum_{t=1}^T \frac{1}{t^2})$ (see Corollary 1 in their work), which scales at a slow $O(1/\sqrt{n})$ rate (instead of $O(1/n)$) when dealing with C -bounded loss².

Generalization bounds with ℓ_2 regularization via Log-Sobolev inequalities. We also study the setting where the total objective function is the sum of a C -bounded differentiable objective F_0 and an additional ℓ_2 regularization term $\frac{\lambda}{2} \|w\|_2^2$. In this case F can be treated as a perturbation of a quadratic function, and the continuous Langevin dynamics (CLD) is well understood for quadratic functions. We obtain two generalization bounds for CLD, both via the technique of Log-Sobolev inequalities, a powerful tool for proving the convergence rate of CLD. One of our bounds is as follows (Theorem 15):

$$\text{err}_{\text{gen}} \leq \frac{2e^4 C^2 L^2}{n} \frac{1}{1 - \exp\left(-\frac{T}{e^8 C}\right)} \quad (4)$$

The above bound has the following advantages:

1. Applying $e^{-x} \leq 1 - x$, one can see that our bound is at most $O(1/n)$, which matches the previous bound in (Mou et al., 2018, Proposition 8)
2. As time T grows, the bound is upper bounded by and approaches $C L^2/n$ (unlike the previous $O(1/\sqrt{n})$ bound that goes to infinity as $n \rightarrow \infty$).
3. If the noise level is not so small (i.e., is not very large), the generalization bound is quite desirable.

Our analysis is based on a Log-Sobolev inequality (LSI) for the parameter distribution at time t , whereas most known LSIs only hold for the stationary distribution of the Markov process. We prove the new LSI by exploiting the variational formulation of the entropy formula.

2 PRELIMINARIES

Notations. We use D to denote the data distribution. The training dataset $S = (z_1, \dots, z_n)$ is a sequence of independent samples drawn from D ; $S^0 \subset Z^n$ are called neighboring datasets if and only if they differ at exactly one data point (we could assume without loss of generality that $z_n \in z_n^0$). Let $F(w; z)$ and $L(w; z)$ be the objective and the loss functions, respectively, where $w \in \mathbb{R}^d$ denotes a model parameter and z is a data point. Define $F(w; S) = \frac{1}{n} \sum_{i=1}^n F(w; z_i)$ and $F(w; D) = \mathbb{E}_{z \sim D} [F(w; z)]$; $L(w; S)$ and $L(w; D)$ are defined similarly. A learning algorithm A takes as input a dataset S and outputs a parameter $w \in \mathbb{R}^d$ randomly. Let \mathcal{G} be the set of all possible mini-batches $G_n = \{B \subset G : n \in B\}$ denotes the collection of mini-batches that contain the n -th data point, while $\overline{\mathcal{G}}_n = \mathcal{G} \setminus G_n$. Let $\text{diam}(A) = \sup_{x, y \in A} \|x - y\|_2$ denote the diameter of a set A .

¹In particular, their proof leverages the Fokker-Planck equation, which describes the time evolution of the density function associated with the Langevin dynamics and can only handle Gaussian noise.

²They assume the loss is sub-Gaussian. By Hoeffding's lemma, bounded random variables are sub-Gaussian with parameter σ .

³The proof of their $O(1/\sqrt{n})$ bound can be easily extended to our setting with regularization.

Definition 1 (L-lipschitz). A function $F : \mathbb{R}^d \times \mathcal{Z} \rightarrow \mathbb{R}$ is L-lipschitz if and only if $|F(w_1; z) - F(w_2; z)| \leq L \|w_1 - w_2\|_2$ holds for any $w_1, w_2 \in \mathbb{R}^d$ and $z \in \mathcal{Z}$.

Definition 2 (Expected generalization error) The expected generalization error of a learning algorithm A is defined as

$$\text{err}_{\text{gen}} := \mathbb{E}_{S \sim D^n} [\text{err}_{\text{gen}}(S)] = \mathbb{E}_{S \sim D^n; A} [L(A(S); D) - L(A(S); S)]:$$

Algorithmic Stability. Intuitively, a learning algorithm that is stable (i.e., a small perturbation of the training data does not affect its output too much) can generalize well. In the seminal work of Bousquet & Elisseeff (2002) (see also Hardt et al. (2016)), the authors formally defined algorithmic stability and established a close connection between the stability of a learning algorithm and its generalization performance.

Definition 3 (Uniform stability). (Bousquet & Elisseeff (2002); Elisseeff et al. (2005)) A randomized algorithm A is ϵ -uniformly stable w.r.t. loss l , if for all neighboring sets $S, S^0 \in \mathcal{Z}^n$, it holds that

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_A [L(w_S; z)] - \mathbb{E}_A [L(w_{S^0}; z)] \leq \epsilon;$$

where w_S and w_{S^0} denote the outputs of A on S and S^0 respectively.

Lemma 4 (Generalization in expectation) (Hardt et al. (2016)) Suppose a randomized algorithm A is ϵ -uniformly stable. Then $\mathbb{E}[\text{err}_{\text{gen}}] \leq \epsilon$.

3 BAYES-STABILITY METHOD

In this section, we incorporate ideas from the PAC-Bayesian theory (see e.g., Lever et al. (2013)) into the algorithmic stability framework. Combined with the technical tools introduced in previous sections, the new framework enables us to prove tighter data-dependent generalization bounds.

First, we define the posterior of a dataset and the posterior of a single data point.

Definition 5 (Single-point posterior) Let Q_S be the posterior distribution of the parameter for a given training dataset $S = (z_1; \dots; z_n)$. In other words, it is the probability distribution of the output of the learning algorithm on dataset S (e.g., for T iterations of SGLD in (1), Q_S is the pdf of W_T). This single-point posterior $Q_{(i; z)}$ is defined as

$$Q_{(i; z)} = \mathbb{E}_{(z_1; \dots; z_{i-1}; z_{i+1}; \dots; z_n)} Q_{(z_1; \dots; z_{i-1}; z; z_{i+1}; \dots; z_n)} :$$

For convenience, we make the following natural assumption on the learning algorithm:

Assumption 6 (Order-independent) For any fixed dataset $S = (z_1; \dots; z_n)$ and any permutation p , Q_S is the same as Q_{S^p} , where $S^p = (z_{p_1}; \dots; z_{p_n})$.

Assumption 6 implies $Q_{(1; z)} = \dots = Q_{(n; z)}$, so we use Q_z as a shorthand for $Q_{(i; z)}$ in the following. Note that this assumption can be easily satisfied by letting the learning algorithm randomly permute the training data at the beginning. It is also easy to verify that both SGD and SGLD satisfy the order-independent assumption.

Now, we state our new Bayes-stability framework, which holds for any prior distribution over the parameter space that is independent of the training dataset.

Theorem 7 (Bayes-Stability) Suppose the loss function $l(w; z)$ is C -bounded and the learning algorithm is order-independent (Assumption 6). Then for any prior distribution P not depending on S , the generalization error is bounded by $2C \mathbb{E}_z \sqrt{2\text{KL}(P; Q_z)}$ and $2C \mathbb{E}_z \sqrt{2\text{KL}(Q_z; P)}$.

Remark 8. Our Bayes-Stability framework originates from the algorithmic stability framework, and hence is similar to the notions of uniform stability and leave-one-out error (see Elisseeff et al. (2003)). However, there are important differences. Uniform stability is a distribution-independent property, while Bayes-Stability can incorporate the information of the data distribution (through the prior P). Leave-one-out error measures the loss of a learned model on an unseen data point, yet Bayes-Stability focuses on the extent to which a single data point affects the outcome of the learning algorithm (compared to the prior).

To establish an intuition, we first apply this framework to obtain an expectation generalization bound for (full) gradient Langevin dynamics (GLD), which is a special case of SGLD in (1) (i.e., GLD uses the full gradient $\nabla_w F(W_{t-1}; S)$ as $g(W_{t-1})$).

Theorem 9. Suppose that the loss function is C -bounded. Then we have the following expected generalization bound for T iterations of GLD:

$$\text{err}_{\text{gen}} \leq \frac{2^p C^p}{n} \mathbb{E}_{S \sim D^n} \left[\sum_{t=1}^T \frac{1}{t} \mathbb{E}_{W_{t-1}} \left[\sum_{i=1}^n \text{kr} F(w_{t-1}; z_i) k_2^2 \right] \right];$$

where $g_e(t) = \mathbb{E}_{W_{t-1}} \left[\frac{1}{n} \sum_{i=1}^n \text{kr} F(w_{t-1}; z_i) k_2^2 \right]$ is the empirical squared gradient norm, and μ is the parameter at step t of GLD.

Proof The proof builds upon the following technical lemma, which we prove in Appendix A.2.

Lemma 10. Let (W_0, \dots, W_T) and (W_0^0, \dots, W_T^0) be two independent sequences of random variables such that for each $t \in \{0, \dots, T\}$, W_t and W_t^0 have the same support. Suppose W_0 and W_0^0 follow the same distribution. Then,

$$\text{KL}(W_T; W_T^0) = \sum_{t=1}^T \mathbb{E}_{W_{<t}} \left[\text{KL}(W_t | W_{<t} = w_{<t}; W_t^0 | W_{<t}^0 = w_{<t}) \right];$$

where $W_{<t}$ denotes (W_0, \dots, W_{t-1}) and $W_{<t}^0$ denotes $(W_0^0, \dots, W_{t-1}^0)$.

Define $P = \mathbb{E}_{S \sim D^n} [Q_{(\bar{S}; 0)}]$, where 0 denotes the zero data point (i.e., $\mathbb{E}(w; 0) = 0$ for any w). Theorem 7 shows that

$$\text{err}_{\text{gen}} \leq 2C \mathbb{E}_Z \left[\frac{1}{n} \sum_{i=1}^n \text{KL}(Q_{(\bar{S}; z)}; P) \right]; \quad (5)$$

By the convexity of KL-divergence, for a fixed $z \in Z$, we have

$$\text{KL}(Q_Z; P) = \mathbb{E}_S \left[\text{KL}(Q_{(\bar{S}; z)}; \mathbb{E}_S [Q_{(\bar{S}; 0)}]) \right] \leq \mathbb{E}_S \left[\text{KL}(Q_{(\bar{S}; z)}; Q_{(\bar{S}; 0)}) \right]; \quad (6)$$

Let $(W_t)_{t=0}^T$ and $(W_t^0)_{t=0}^T$ be the training process of GLD for $S = (\bar{S}; z)$ and $S^0 = (\bar{S}; 0)$, respectively. Note that for a fixed $w_{<t}$, both $W_t | W_{<t} = w_{<t}$ and $W_t^0 | W_{<t}^0 = w_{<t}$ are Gaussian distributions. Since $\text{KL}(N(\mu_1; \Sigma_1); N(\mu_2; \Sigma_2)) = \frac{1}{2} \text{kr} \left(\frac{\Sigma_1^{-1} \mu_2 - \mu_1}{\Sigma_1} \right)^2 + \frac{1}{2} \text{tr} \left(\frac{\Sigma_2}{\Sigma_1} \right) - \frac{1}{2}$ (see Lemma 18 in Appendix A.2).

$$\text{KL}(W_t | W_{<t} = w_{<t}; W_t^0 | W_{<t}^0 = w_{<t}) = \frac{1}{2} \text{kr} F(w_{t-1}; z) k_2^2.$$

Applying Lemma 10 and $\text{KL}(W_T; W_T^0) = \text{KL}(W_T; W_T^0)$ gives

$$\text{KL}(Q_S; Q_{S^0}) \leq \frac{1}{n^2} \sum_{t=1}^T \frac{1}{t} \mathbb{E}_{W_{t-1}} \left[\text{kr} F(w_{t-1}; z) k_2^2 \right];$$

Recall that W_{t-1} is the parameter at step $t-1$ using $S = (\bar{S}; z)$ as dataset. In this case, we can rewrite z as z_n since it is the n -th data point of S . Note that SGLD satisfies the order-independent assumption, we can rewrite z as z_i for all $i \in [n]$. Together with (5), (6), and using $\frac{1}{n} \sum_{i=1}^n x_i$, we can prove this theorem. ■

More generally, we give the following bound for SGLD. The proof is similar to that of Theorem 9; the difference is that we need to bound the KL-divergence between two Gaussian mixtures instead of two Gaussians. This proof is more technical and deferred to Appendix A.3.

Theorem 11. Suppose that the loss function is C -bounded and the objective function is L -Lipschitz. Assume that the following conditions hold:

1. Batch size $\phi = n=2$.

2. Learning rate $\eta = \frac{1}{20L}$.

Then, the following expected generalization error bound holds for iterations of SGLD(1):

$$\text{err}_{\text{gen}} \leq \frac{8.12C}{n} \sqrt{\frac{1}{t} \sum_{s=1}^t \mathbb{E}_{S^s} \left[\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{z_i} \left[\mathbb{E}_{w_i} \left[\sum_{t=1}^T \mathbb{E}_{z_i} \left[\mathbb{E}_{w_i} \left[\left\| \nabla F(w_i; z_i) \right\|_2^2 \right] \right] \right] \right] \right]} \right]} ; \quad (\text{empirical norm})$$

where $g_e(t) = \mathbb{E}_{w_i} \left[\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{z_i} \left[\mathbb{E}_{w_i} \left[\left\| \nabla F(w_i; z_i) \right\|_2^2 \right] \right] \right]$ is the empirical squared gradient norm, and w_i is the parameter at step t of SGLD.

Furthermore, based on essentially the same proof, we can obtain the following bound that depends on the population gradient norm

$$\text{err}_{\text{gen}} \leq \frac{8.12C}{n} \sqrt{\frac{1}{t} \sum_{s=1}^t \mathbb{E}_{S^s} \left[\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{z_i} \left[\mathbb{E}_{w_i} \left[\sum_{t=1}^T \mathbb{E}_{z_i} \left[\mathbb{E}_{w_i} \left[\left\| \nabla F(w_i; z_i) \right\|_2^2 \right] \right] \right] \right] \right]} \right]} ;$$

The full proofs of the above results are postponed to Appendix A, and we provide some remarks about the new bounds.

Remark 12. In fact, our proof establishes that the above upper bound holds for the two sequences W_T and W_T^0 : $\text{KL}(W_T; W_T^0) \leq \frac{8.12C}{n^2} \sum_{t=1}^T \mathbb{E}_{S^t} \left[\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{z_i} \left[\mathbb{E}_{w_i} \left[\sum_{t=1}^T \mathbb{E}_{z_i} \left[\mathbb{E}_{w_i} \left[\left\| \nabla F(w_i; z_i) \right\|_2^2 \right] \right] \right] \right] \right]$. Hence, our bound holds for any sufficiently regular function over the parameter sequences $(f(W_T); f(W_T^0))$. In particular, our generalization error bound automatically extends to several variants of SGLD, such as outputting the average of the trajectory, the average of the suffix of certain length, or the exponential moving average.

Remark 13 (High-probability bounds) By relaxing the expected squared gradient norm term and using the uniform stability framework, our proof can be adapted to recover the bound in (Mou et al., 2018, Theorem 1). Then, we can apply the recent results of Feldman & Vondrak (2019) to provide a generalization error bound $\mathcal{O}(\frac{1}{\sqrt{n}})$ that holds with high probability. (Here \mathcal{O} hides poly-logarithmic factors.) When n is at least linear in $\ln n$, the additional $\frac{1}{\sqrt{n}}$ term is not dominating.

3.1 EXPERIMENT

Distinguish random from normal. Inspired by Zhang et al. (2017a), we run both GLD (Figure 1) and SGLD (Appendix C.2) to fit both normal data and randomly labeled data (see Appendix C for more experiment details). As shown in Figure 1 and Figure 3 in Appendix C.2, a larger random label portion leads to both much higher generalization error and much larger generalization error bound. Moreover, the shapes of the curves of our bounds look quite similar to those of the generalization error curves.

Note that in (b) and (c) of Figure 1, the scales in the axis are different. We list some possible reasons that may explain why our bound is larger than the actual generalization error. (1) as we explained in Remark 12, our bounds (Theorem 9 and 11) hold for any trajectory-based output, and are much stronger than upper bounds for the last point on the trajectory. (2) The constant we can prove in Lemma 21 may not be very tight. (3) The variance of Gaussian noise is not large enough in our experiment. However, if we choose a larger variance, fitting the random labeled training data becomes quite slow. Hence, we use a small data size (10000) for the above reason. We also run an extra experiment for GLD on the full MNIST dataset (60000) without label corruption (see Figure 2 in the Appendix C). We can see that our bound is vacuous since GLD—which computes the full gradients—took a long time to converge, we stopped when we achieved 90% training accuracy⁴.

⁴ We highlight another difficulty in proving non-vacuous generalization error bounds when the data are randomly labeled. Consider a 10-class classification setting where all the labels are random. For any sufficiently small data size, there is always a deep neural network that perfectly fits the dataset. Thus, the training error is zero while the population error is 90%. In this case, any valid generalization error bound should be larger than 0.9. Then, the theoretical bound would still be vacuous even if it is only loose by a factor of

(a) (b) (c) (d)

Figure 1: Training MLP with GLD ($\epsilon_t = 0:2^{\frac{p}{2} - t}$) on a smaller version of MNIST with different random label portions. (a) shows the training accuracy. (b) shows the generalization error, i.e., the gap between the 0/1 loss⁵ on the training data and on the test data. (c) plots our bound in Theorem 9. (d) shows that for $\epsilon = 0$, the gradient norms become much smaller at later stages of training.

Relax the step size constraint. The condition on the step size in Theorem 11 may seem restrictive in the practical use. We provide several ways to relax this constraint:

1. The proof of Theorem 11 still goes through if we replace $\max_{i \in [n]} \|\nabla F(W_{t-1}; z_i)\|_2$ with $\max_{i \in [n]} \|\nabla F(W_{t-1}; z_i)\|_2$ in the constraint.
2. The maximum gradient norm can be controlled by gradient clipping, i.e., multiplying $\frac{\min(C_L; \|\nabla F(W_{t-1}; z_i)\|_2)}{\|\nabla F(W_{t-1}; z_i)\|_2}$ to each $\nabla F(W_{t-1}; z_i)$.
3. Replacing the constant 2 with 2 in this constraint will only increase the constant of our bound from $8:12$ to $84:4$.

We also provide an experiment combining the above ideas to make our Theorem 11 applicable in the practical use (see Figure 4 in Appendix C).

4 GENERALIZATION OF CLD AND GLD WITH ℓ_2 REGULARIZATION

In this section, we study the generalization error of Continuous Langevin Dynamics (CLD) with regularization. Throughout this section, we assume that the objective function over training set is defined as $F(w; S) = F_0(w; S) + \frac{\lambda}{2} \|w\|_2^2$, and moreover, the following assumption holds.

Assumption 14. The loss function F and the original objective F_0 are C -bounded. Moreover, F_0 is differentiable and L -lipschitz.

The Continuous Langevin Dynamics is defined by the following SDE:

$$dW_t = -\nabla F(W_t; S) dt + \sqrt{\frac{p}{2}} dB_t; \quad W_0 \sim \mu_0; \quad (\text{CLD})$$

where $(B_t)_{t \geq 0}$ is the standard Brownian motion in \mathbb{R}^d and the initial distribution μ_0 is the centered Gaussian distribution in \mathbb{R}^d with covariance $\frac{1}{2} I_d$. We show that the generalization error of CLD is upper bounded by $O(e^{4C} n^{-1} \frac{p}{2})$, which is independent of the training time (Theorem 15).

Furthermore, as λ goes to infinity, we have a tighter generalization error bound $O(C^2 n^{-1})$ (Theorem 39 in Appendix B). We also study the generalization of Gradient Langevin Dynamics (GLD), which is the discretization of CLD:

$$W_{k+1} = W_k - \eta \nabla F(W_k; S) + \sqrt{\frac{p}{2}} \epsilon_k; \quad (\text{GLD})$$

where ϵ_k is the standard Gaussian random vector in \mathbb{R}^d . By leveraging a result developed in Raginsky et al. (2017), we show that, as λ tends to zero, GLD has the same generalization as CLD (see Theorems 15 and 39). We first formally state our first main result in this section.

⁵The condition $\epsilon_t = O(\epsilon_{t=L})$ is also required in (Mou et al., 2018, Theorem 1)

Theorem 15. Under Assumption 14, CLD (with initial probability measure $\mu_0 = \frac{1}{Z} e^{-\frac{k\|w\|^2}{2}} dw$) has the following expected generalization error bound:

$$\text{err}_{\text{gen}} \leq \frac{2e^{4C} CL}{n} - 1 \exp\left(-\frac{T}{e^{8C}}\right) : \quad (7)$$

In addition, if L is M -smooth and non-negative, by setting $\eta > 0$ and $\alpha > 0; 1 \wedge \frac{\eta}{8M^2}$, GLD (running K iterations with the same μ_0 as CLD) has the expected generalization error bound:

$$\text{err}_{\text{gen}} \leq 2C^p \frac{1}{2KC_1^2 + \frac{2CLE^{4C}}{n}} - 1 \exp\left(-\frac{K}{e^{8C}}\right) ; \quad (8)$$

where C_1 is a constant that only depends on η, α, b, L and d .

The following lemma is crucial for establishing the above generalization bound for CLD. In particular, we need to establish a Log-Sobolev inequality for the parameter distribution at time t for every time step $t > 0$. In contrast, most known LSIs only characterize the stationary distribution of the Markov process. The proof of the lemma can be found in Appendix B.

Lemma 16. Under Assumption 14, let μ_t be the probability measure μ_t in CLD (with $\mu_0 = \frac{1}{Z} e^{-\frac{k\|w\|^2}{2}} dw$). Let ν be a probability measure that is absolutely continuous with respect to μ_t . Suppose $d\nu = \nu(w) dw$ and $d\mu_t = \mu_t(w) dw$. Then, it holds that

$$\text{KL}(\nu; \mu_t) \leq \frac{\exp(8C)}{2} \int_{\mathbb{R}^d} r \log \frac{\nu(w)}{\mu_t(w)}^2 (w) dw:$$

We sketch the proof of Theorem 15, and the complete proof is relegated to Appendix B.

Proof Sketch of Theorem 15 Suppose S and S^0 are two neighboring datasets. $(W_t)_{t=0}$ and $(W_t^0)_{t=0}$ be the process of CLD running on S and S^0 , respectively. Let μ_t and μ_t^0 be the pdf of W_t^0 and W_t . Let $F_S(w)$ denote $F(w; S)$. We have

$$\begin{aligned} \frac{d}{dt} \text{KL}(\mu_t; \mu_t^0) &= \frac{1}{2} \int_{\mathbb{R}^d} \mu_t^0 r \log \frac{\mu_t}{\mu_t^0}^2 dw + \int_{\mathbb{R}^d} \mu_t r \log \frac{\mu_t}{\mu_t^0}; r F_S - r F_{S^0} dw \\ &= \frac{1}{2} \int_{\mathbb{R}^d} \mu_t^0 r \log \frac{\mu_t}{\mu_t^0}^2 dw + \frac{1}{2} \int_{\mathbb{R}^d} \mu_t k r F_S - r F_{S^0} k^2 dw: \\ &= \frac{1}{e^{8C}} \text{KL}(\mu_t; \mu_t^0) + \frac{2L^2}{n^2} \end{aligned} \quad (\text{Lemma 16})$$

Solving this inequality gives $\text{KL}(\mu_t; \mu_t^0) \leq \frac{1}{2} \frac{2L^2 e^{8C}}{n^2} (1 - e^{-t e^{8C}})$. Hence the generalization error of CLD can be bounded by $\frac{1}{2} \text{KL}(\mu_T; \mu_T^0)$, which proves the first part. The second part of the theorem follows from Lemma 36 in Appendix B. ■

Our second generalization bound for CLD (Theorem 39 in Appendix B) is

$$\text{err}_{\text{gen}} \leq \frac{8C^2}{n} + 4C \exp\left(-\frac{T}{e^{4C}}\right) \frac{p}{C}:$$

The high level idea to prove this bound is very similar to that in Raginsky et al. (2017). We first observe that the (stationary) Gibbs distribution has a small generalization error. Then, we bound the distance from μ_t to μ . In our setting, we can use the Holley-Stroock perturbation lemma which allows us to bound the Logarithmic Sobolev constant, and we can thus bound the above distance easily.

5 FUTURE DIRECTIONS

In this paper, we prove new generalization bounds for a variety of noisy gradient-based methods. Our current techniques can only handle continuous noises for which we can bound the KL-divergence.

One future direction is to study the discrete noise introduced in SGD (in this case the KL-divergence may not be well defined). For either SGLD or CLD, if the noise level is small (i.e., large), it may take a long time for the diffusion process to reach the stable distribution. Hence, another interesting future direction is to consider the local behavior and generalization of the diffusion process in finite time through the techniques developed in the studies of metastability (see e.g., Bovier et al. (2005); Bovier & den Hollander (2006); Tzen et al. (2018)). In particular, the technique may be helpful for further improving the bounds in Theorems 15 and 39 (when ϵ is not very large).

6 ACKNOWLEDGEMENT

We would like to thank Liwei Wang for several helpful discussions during various stages of the work. The research is supported in part by the National Natural Science Foundation of China Grant 61822203, 61772297, 61632016, 61761146003, and the Zhongguancun Haihua Institute for Frontier Information Technology and Turing AI Institute of Nanjing.

REFERENCES

- Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameterized neural networks, going beyond two layers. *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 6155–6166, 2019.
- Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger generalization bounds for deep nets via a compression approach. *International Conference on Machine Learning (ICML)*, pp. 254–263, 2018.
- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 8139–8148, 2019.
- Dominique Bakry, Ivan Gentil, and Michel Ledoux. *Analysis and geometry of Markov diffusion operators*, volume 348. Springer Science & Business Media, 2013.
- Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds for neural networks. *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 6240–6249, 2017.
- Olivier Bousquet and Andrzej Elisseeff. Stability and generalization. *Journal of Machine Learning Research (JMLR)*, pp. 499–526, 2002.
- Anton Bovier and Frank den Hollander. Metastability: a potential theoretic approach. *International Congress of Mathematicians*, volume 3, pp. 499–518. Eur. Math. Society, 2006.
- Anton Bovier, Véronique Gayrard, and Markus Klein. Metastability in reversible diffusion processes ii: Precise asymptotics for small eigenvalues. *Journal of the European Mathematical Society*, (1):69–99, 2005.
- Pratik Chaudhari, Anna Choromanska, Stefano Soatto, Yann LeCun, Carlo Baldassi, Christian Borgs, Jennifer Chayes, Levent Sagun, and Riccardo Zecchina. Entropy-SGD: Biasing gradient descent into wide valleys. *International Conference on Learning Representations (ICLR)*, 2017.
- Yuansi Chen, Chi Jin, and Bin Yu. Stability and convergence trade-off of iterative optimization algorithms. *arXiv preprint arXiv:1804.01619*, 2018.
- Lenaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 2933–2943, 2019.
- Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. *International Conference on Machine Learning (ICML)*, pp. 1675–1685, 2019.
- John Duchi. *Derivations for linear algebra and optimization*. Berkeley, California, 2007.

- Gintare Karolina Dziugaite and Daniel Roy. Entropy-SGD optimizes the prior of a PAC-Bayes bound: Generalization properties of entropy-SGD and data-dependent priors. *International Conference on Machine Learning (ICML)*, pp. 1377–1386, 2018.
- Gintare Karolina Dziugaite and Daniel M Roy. Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. *Uncertainty in Artificial Intelligence (UAI)*, 2017.
- André Elisseeff, Massimiliano Pontil, et al. Leave-one-out error and stability of learning algorithms with applications. *NATO science series sub series iii computer and systems sciences*, pp. 111–130, 2003.
- Andre Elisseeff, Theodoros Evgeniou, and Massimiliano Pontil. Stability of randomized learning algorithms. *Journal of Machine Learning Research (JMLR)*, (Jan):55–79, 2005.
- Vitaly Feldman and Jan Vondrak. High probability generalization bounds for uniformly stable algorithms with nearly optimal rate. *Conference on Learning Theory (COLT)*, pp. 1270–1279, 2019.
- Moritz Hardt, Benjamin Recht, and Yoram Singer. Train faster, generalize better: stability of stochastic gradient descent. *International Conference on Machine Learning (ICML)*, pp. 1225–1234, 2016.
- Richard Holley and Daniel Stroock. Logarithmic sobolev inequalities and stochastic ising models. *Journal of statistical physics*, 46(5):1159–1194, 1987.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 1097–1105, 2012.
- Ilya Kuzborskij and Christoph Lampert. Data-dependent stability of stochastic gradient descent. In *International Conference on Machine Learning (ICML)*, pp. 2815–2824, 2018.
- Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 6(11):2278–2324, 1998.
- Guy Lever, François Laviolette, and John Shawe-Taylor. Tighter pac-bayes bounds through distribution-dependent priors. *Theoretical Computer Science*, 473:4–28, 2013.
- Tengyuan Liang, Tomaso Poggio, Alexander Rakhlin, and James Stokes. Fisher-rao metric, geometry, and complexity of neural networks. *International Conference on Artificial Intelligence and Statistics (AISTATS)*, pp. 888–896, 2019.
- Ben London. Generalization bounds for randomized learning with application to stochastic gradient descent. In *NIPS Workshop on Optimizing the Optimization*, 2016.
- Ben London. A pac-bayesian analysis of randomized learning with application to stochastic gradient descent. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 2931–2940, 2017.
- Georg Menz, André Schlichting, et al. Poincaré and logarithmic sobolev inequalities by decomposition of the energy landscape. *The Annals of Probability*, 42(5):1809–1884, 2014.
- Wenlong Mou, Liwei Wang, Xiyu Zhai, and Kai Zheng. Generalization bounds of SGLD for non-convex learning: Two theoretical viewpoints. *Conference on Learning Theory (COLT)*, pp. 605–638, 2018.
- Yurii E Nesterov. A method for solving the convex programming problem with convergence rate $O(1/k^2)$. In *Dokl. Akad. Nauk SSSR*, volume 269, pp. 543–547, 1983.
- Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. Norm-based capacity control in neural networks. In *Conference on Learning Theory (COLT)*, pp. 1376–1401, 2015.

- Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nathan Srebro. A PAC-Bayesian approach to spectrally-normalized margin bounds for neural networks. *International Conference on Learning Representations (ICLR)*, 2018.
- Grigorios A Pavliotis. *Stochastic processes and applications: diffusion processes, the Fokker-Planck and Langevin equations*, volume 60. Springer, 2014.
- Ankit Pensia, Varun Jog, and Po-Ling Loh. Generalization error bounds for noisy, iterative algorithms. In *International Symposium on Information Theory (ISIT)*. 546–550, 2018.
- Boris T Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*(5):1–17, 1964.
- Maxim Raginsky, Alexander Rakhlin, and Matus Telgarsky. Non-convex learning via stochastic gradient langevin dynamics: a nonasymptotic analysis. *Conference on Learning Theory (COLT)*, pp. 1674–1703, 2017.
- Hannes Risken. Fokker-planck equation. *The Fokker-Planck Equation*, pp. 63–95. Springer, 1996.
- Omar Rivasplata, Csaba Szepesvári, John S Shawe-Taylor, Emilio Parrado-Hernandez, and Shiliang Sun. Pac-bayes bounds for stable algorithms with instance-dependent priors. *Advances in Neural Information Processing Systems (NeurIPS)* 9214–9224, 2018.
- Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. *International Conference on Machine Learning (ICML)*, pp. 1139–1147, 2013.
- Flemming Topsoe. Some inequalities for information divergence and related measures of discrimination. *IEEE Transactions on Information Theory*(46(4)):1602–1609, 2000.
- Belinda Tzen, Tengyuan Liang, and Maxim Raginsky. Local optimality and generalization guarantees for the langevin algorithm via empirical metastability. *Conference On Learning Theory (COLT)*, pp. 857–875, 2018.
- Colin Wei and Tengyu Ma. Data-dependent sample complexity of deep neural networks via lipschitz augmentation. *Advances in Neural Information Processing Systems (NeurIPS)* 9722–9733, 2019.
- Colin Wei, Jason D Lee, Qiang Liu, and Tengyu Ma. Regularization matters: Generalization and optimization of neural nets vs their induced kernel. *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 9709–9721, 2019.
- Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *International Conference on Machine Learning (ICML)*, pp. 681–688, 2011.
- Aolin Xu and Maxim Raginsky. Information-theoretic analysis of generalization capability of learning algorithms. In *Advances in Neural Information Processing Systems (NeurIPS)* 2524–2533, 2017.
- Pan Xu, Jinghui Chen, Difan Zou, and Quanquan Gu. Global convergence of langevin dynamics based algorithms for nonconvex optimization. *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 3126–3137, 2018.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *International Conference on Learning Representations (ICLR)* 2017a.
- Yuchen Zhang, Percy Liang, and Moses Charikar. A hitting time analysis of stochastic gradient langevin dynamics. *Conference on Learning Theory (COLT)*, pp. 1980–2022, 2017b.
- Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes over-parameterized deep relu networks. *arXiv preprint arXiv:1811.08882*, 2018.

A PROOFS IN SECTION 3

A.1 BAYES-STABILITY FRAMEWORK

Lemma 17. Under Assumption 6, for any prior distribution P not depending on the dataset $S = (z_1, \dots, z_n)$, the generalization error is upper bounded by

$$\mathbb{E}_Z \mathbb{E}_w \mathbb{E}_P L(w; z) - \mathbb{E}_w \mathbb{E}_{Q_z} L(w; z) + \mathbb{E}_Z \mathbb{E}_w \mathbb{E}_P L(w) - \mathbb{E}_w \mathbb{E}_{Q_z} L(w) ;$$

where $L(w)$ denotes the population loss $L(w; D)$.

Proof of Lemma 17 Let $\text{err}_{\text{train}} = \mathbb{E}_S \mathbb{E}_w \mathbb{E}_{Q_S} L(w; S)$ and $\text{err}_{\text{test}} = \mathbb{E}_S \mathbb{E}_w \mathbb{E}_{Q_S} L(w)$. We can rewrite generalization error $\text{err}_{\text{gen}} = \text{err}_{\text{test}} - \text{err}_{\text{train}}$, where

$$\begin{aligned} \text{err}_{\text{test}} &= \mathbb{E}_Z \mathbb{E}_w \mathbb{E}_{Q_{(1; z)}} L(w) = \mathbb{E}_Z \mathbb{E}_w \mathbb{E}_{Q_z} L(w) && \text{(Assumption 6)} \\ &= \mathbb{E}_Z \int_{\mathbb{R}^d} (Q_z(w) - P(w)) L(w) dw + \int_{\mathbb{R}^d} P(w) L(w) dw \end{aligned}$$

and

$$\begin{aligned} \text{err}_{\text{train}} &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_S \mathbb{E}_w \mathbb{E}_{Q_S} L(w; z_i) \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Z \mathbb{E}_w \mathbb{E}_{Q_{(i; z)}} L(w; z) = \mathbb{E}_Z \mathbb{E}_w \mathbb{E}_{Q_z} L(w; z) && \text{(Assumption 6)} \\ &= \mathbb{E}_Z \int_{\mathbb{R}^d} (Q_z(w) - P(w)) L(w; z) dw + \int_{\mathbb{R}^d} P(w) \mathbb{E}_z L(w; z) dw && \text{(P is a prior)} \\ &= \mathbb{E}_Z \int_{\mathbb{R}^d} (Q_z(w) - P(w)) L(w; z) dw + \int_{\mathbb{R}^d} P(w) L(w) dw && \text{(definition of } f(w)) \end{aligned}$$

Thus, we have

$$\begin{aligned} |\text{err}_{\text{gen}}| &= |\text{err}_{\text{test}} - \text{err}_{\text{train}}| \\ &= \mathbb{E}_Z \int_{\mathbb{R}^d} (Q_z(w) - P(w)) L(w) dw - \mathbb{E}_Z \int_{\mathbb{R}^d} (Q_z(w) - P(w)) L(w; z) dw \\ &= \mathbb{E}_Z \mathbb{E}_w \mathbb{E}_{Q_z} L(w; z) - \mathbb{E}_w \mathbb{E}_P L(w; z) + \mathbb{E}_Z \mathbb{E}_w \mathbb{E}_{Q_z} L(w) - \mathbb{E}_w \mathbb{E}_P L(w) ; \end{aligned}$$

Now we are ready to prove Theorem 7, which we restate in the following.

Theorem 7 (Bayes-Stability) Suppose the loss function $L(w; z)$ is C -bounded and the learning algorithm is order-independent (Assumption 6), then for any prior distribution P not depending on S , the generalization error is bounded by $2C \mathbb{E}_z \sqrt{\frac{1}{n} 2\text{KL}(P; Q_z)}$ and $2C \mathbb{E}_z \sqrt{\frac{1}{n} 2\text{KL}(Q_z; P)}$.

Proof By Lemma 17,

$$\begin{aligned} \text{err}_{\text{gen}} &= \mathbb{E}_Z \mathbb{E}_w \mathbb{E}_P L(w; z) - \mathbb{E}_w \mathbb{E}_{Q_z} L(w; z) + \mathbb{E}_Z \mathbb{E}_w \mathbb{E}_P L(w) - \mathbb{E}_w \mathbb{E}_{Q_z} L(w) \\ &= \mathbb{E}_Z [2C \text{TV}(P; Q_z) + 2C \text{TV}(P; Q_z)] && \text{(C-boundedness)} \\ &= 4C \mathbb{E}_Z \sqrt{\frac{1}{2} \text{KL}(P; Q_z)} && \text{(Pinsker's inequality)} \end{aligned}$$

The other bound follows from a similar argument. ■

A.2 TECHNICAL LEMMAS

Now we turn to the proof of Theorem 11. The following lemma allows us to reduce the proof of algorithmic stability to the analysis of a single update step.

Lemma 10. Let $(W_0; \dots; W_T)$ and $(W_0^0; \dots; W_T^0)$ be two independent sequences of random variables such that for each $t \in \{0; \dots; T\}$, W_t and W_t^0 have the same support. Suppose W_0 and W_0^0 follow the same distribution. Then,

$$KL(W_T; W_T^0) = \sum_{t=1}^T \mathbb{E}_{W_{<t}} [KL(W_t | W_{<t} = w_{<t}; W_t^0 | W_{<t}^0 = w_{<t})];$$

where $W_{<t}$ denotes $(W_0; \dots; W_{t-1})$ and $W_{<t}^0$ denotes $(W_0^0; \dots; W_{t-1}^0)$.

Proof By the chain rule of the KL-divergence,

$$KL(W_t; W_t^0) = KL(W_{<t}; W_{<t}^0) + \mathbb{E}_{W_{<t}} [KL(W_t | W_{<t} = w_{<t}; W_t^0 | W_{<t}^0 = w_{<t})];$$

The lemma follows from a summation over $t \in \{1; \dots; T\}$. ■

The following lemma (see e.g., (Duchi, 2007, Section 9)) gives a closed-form formula for the KL-divergence between two Gaussian distributions.

Lemma 18. Suppose that $P = N(\mu_1; \Sigma_1)$ and $Q = N(\mu_2; \Sigma_2)$ are two Gaussian distributions on \mathbb{R}^d . Then,

$$KL(P; Q) = \frac{1}{2} \left(\text{tr}(\Sigma_2^{-1} \Sigma_1) + (\mu_2 - \mu_1)^T \Sigma_2^{-1} (\mu_2 - \mu_1) \right) + \ln \frac{\det(\Sigma_2)}{\det(\Sigma_1)};$$

The following lemma (Topsoe, 2000, Theorem 3) helps us to upper bound the KL-divergence.

Definition 19. Let P and Q be two probability distributions on \mathbb{R}^d . The directional triangular discrimination from P to Q is defined as

$$(P; Q) = \sum_{k=0}^{\infty} \frac{1}{2^k} \int_{\mathbb{R}^d} \frac{P(w) - Q(w)}{P(w) + Q(w)} dw;$$

where

$$(P; Q) = \int_{\mathbb{R}^d} \frac{P(w) - Q(w)}{P(w) + Q(w)} dw;$$

Lemma 20. For any two probability distributions P and Q on \mathbb{R}^d ,

$$KL(P; Q) \leq \ln 2 \cdot (P; Q);$$

Recall that \mathcal{G} is the set of all possible mini-batches $G = \{B \subseteq \mathcal{X} : |B| = n\}$ and $\mathcal{G}_n = \{B \subseteq \mathcal{X} : |B| = n, B \subseteq G\}$ denotes the collection of mini-batches that contain n points from G , while $\overline{\mathcal{G}}_n = \mathcal{G} \setminus \mathcal{G}_n$. $\text{diam}(A) = \sup_{x, y \in A} \|x - y\|$ denotes the diameter of a set A . The following technical lemma upper bounds the KL-divergence between two Gaussian mixtures induced by sampling a mini-batch from neighbouring datasets.

Lemma 21. Suppose that batch size $n \geq 2$. $f_B : \mathcal{X} \rightarrow \mathbb{R}^d$ and $f_B^0 : \mathcal{X} \rightarrow \mathbb{R}^d$ are two collections of points in \mathbb{R}^d labeled by mini-batches of size n that satisfy the following conditions for some constant $\epsilon \in [0, 1]$:

- $\|f_B - f_B^0\| \leq \epsilon$ for $B \in \mathcal{G}_n$ and $\|f_B - f_B^0\| \leq \epsilon$ for $B \in \overline{\mathcal{G}}_n$.
- $\text{diam}(f_B : B \in \mathcal{G}_n \cup f_B^0 : B \in \mathcal{G}_n) \leq \epsilon$.

Let p_B denote the Gaussian distribution $N(f_B; \frac{\sigma^2}{n} I_d)$. Let $P = \frac{1}{|\mathcal{G}|} \sum_{B \in \mathcal{G}} p_B$ and $P^0 = \frac{1}{|\mathcal{G}|} \sum_{B \in \mathcal{G}} p_B^0$ be two mixture distributions over all mini-batches. Then,

$$KL(P; P^0) \leq \frac{8 \cdot 23 \epsilon^2}{2n^2};$$

Proof of Lemma 21 By Lemma 20, $KL(P; P^0)$ is bounded by

$$\begin{aligned} \ln 2 \quad (P; P^0) &= \ln 2 \int_{\mathbb{R}^d} \frac{4^{-k} (P(w) - P^0(w))^2}{2^{-k} P(w) + (2^{-k}) P^0(w)} dw \\ &= \ln 2 \int_{\mathbb{R}^d} \frac{4^{-k} (P(w) - P^0(w))^2}{2^{-k} P(w) + (2^{-k}) P^0(w)} dw \end{aligned}$$

The numerator of the above integrand is upper bounded by

$$\begin{aligned} 4^{-k} (P - P^0)^2 &= 4^{-k} \frac{1}{|G_j|} \int_{B_{2G}} (p_{B;}; p_{B^0;})^2 \\ &= \frac{4^{-k} |G_n|^2}{|G_j|^2} \frac{1}{|G_n|} \int_{B_{2G_n}} (p_{B;}; p_{B^0;})^2 \\ &= \frac{4^{-k} b^2}{n^2} \frac{1}{|G_n|} \int_{B_{2G_n}} (p_{B;}; p_{B^0;})^2; \end{aligned} \quad (9)$$

while the denominator can be lower bounded as follows:

$$\begin{aligned} 2^{-k} P + (2^{-k}) P^0 &= \frac{2^{-k}}{|G_j|} \int_{B_{2G_n}} p_{B;}; + \frac{2^{-k}}{|G_j|} \int_{B_{2G_n}} p_{B^0;}; \\ &= \frac{2^{-k}}{|G_j|} \int_{B_{2G_n}} p_{B;}; \quad (B = B^0 \text{ for } B \in \overline{G_n}) \\ &= \frac{1}{|G_n|} \frac{2(n-b)}{n} \int_{B_{2G_n}} p_{B;}; \\ &= \frac{1}{|G_n|} \int_{B_{2G_n}} p_{B;}; \quad (b \leq n=2) \end{aligned}$$

which implies, by the convexity of $f(x) = \frac{1}{x}$, that

$$\frac{1}{2^{-k} P + (2^{-k}) P^0} \leq \frac{1}{\frac{1}{|G_n|} \int_{B_{2G_n}} p_{B;};} \leq \frac{1}{|G_n|} \int_{B_{2G_n}} \frac{1}{p_{B;}}; \quad (10)$$

Inequalities (9) and (10) together imply

$$2^{-k} P + (1 - 2^{-k}) P^0 \leq \frac{4^{-k} b^2}{n^2 |G_n| |G_n|} \int_{A_{2G_n}} \int_{B_{2G_n}} \int_{\mathbb{R}^d} \frac{(p_{B;}(w) - p_{B^0;}(w))^2}{p_{A;}(w)} dw \quad (11)$$

Now we bound the right-hand side of (11) for $x \in A$ and B . By applying a translation and a rotation, we can assume without loss of generality that $x = 0$, and the last $d-2$ coordinates of B and B^0 are all zero. Note that the integral is unchanged when we project the space to the two-dimensional subspace corresponding to the first two coordinates. Thus, it suffices to prove a bound for $d = 2$. We rewrite (11) as

$$(2^{-k} P + (1 - 2^{-k}) P^0) \leq \frac{4^{-k} b^2}{n^2 |G_n| |G_n|} \int_{A_{2G_n}} \int_{B_{2G_n}} \frac{1}{2} \int_{\mathbb{R}^2} \frac{e^{-k \|w - B\|^2} e^{-\frac{w \cdot B^0}{2}}}{e^{-k \|w\|^2}} dw \quad (12)$$

Let I be the integral in the right-hand side of (12). Note that $\frac{0}{B} < \frac{0}{B} < 0.1$ and $\frac{B}{B} = 1$. Let $\frac{0}{B} = \frac{w}{B}$ and $r = \frac{w}{B}$. Our goal is to bound $\max_{y \in [0; 0.1]} (I - \frac{1}{2})$. Let $(x)^+ = \max(x, 0)$. Since

$$I = \int_0^1 \frac{\max_{y \in [2^{-(r-0.1)^+}; r+0.1]} (e^{-y^2} - e^{-(y+)^2})^2}{e^{-r^2}} 2r \, dr:$$

We have

$$\max_{y \in [0; 0.1]} \frac{1}{2} \int_0^1 e^{-r^2} 2r \max_{y \in [2^{-(r-0.1)^+}; r+0.1]} \max_{y \in [2^{-(r-0.1)^+}; r+0.1]} \frac{(e^{-y^2} - e^{-(y+)^2})^2}{2} \, dr:$$

Let $(y; r) = (e^{-y^2} - e^{-(y+)^2})^2$, we make two claims which we will prove later:

1. For all $y; r \in [0, 1]$, $(y; r) \leq \frac{2}{e}$.
2. For all $r \in [\frac{1}{2}, 1]$, $(y; r)$ is non-increasing in y .

The above claims imply that:

1. For any $r \in [0; \frac{1}{2} + 0.1]$, $\max_{y \in [2^{-(r-0.1)^+}; r+0.1]} (y; r) \leq \frac{2}{e}$.
2. For any $r \in [\frac{1}{2} + 0.1; 1]$, we have

$$\begin{aligned} \max_{y \in [2^{-(r-0.1)^+}; r+0.1]} (y; r) &= \max_{y \in [2^{-(r-0.1)^+}; r+0.1]} \lim_{t \rightarrow 0} [(y; r)] \\ &= \max_{y \in [2^{-(r-0.1)^+}; r+0.1]} 4y^2 e^{-2y^2} \\ &= 4(r-0.1)^2 e^{-2(r-0.1)^2}. \end{aligned}$$

The last step holds since $y^2 e^{-2y^2}$ is decreasing on $[\frac{1}{2}; 1]$.

Thus we have

$$\begin{aligned} \max_{y \in [0; 0.1]} \frac{1}{2} \int_0^1 e^{-r^2} 2r \, dr + \int_{\frac{1}{2} + 0.1}^1 e^{-r^2} 2r \cdot 4(r-0.1)^2 e^{-2(r-0.1)^2} \, dr \\ = 18.6487 \frac{1}{2}: \end{aligned}$$

Plugging the above into (12) gives

$$(2^{-k} P + (1 - 2^{-k}) P^0; P^0) \frac{4^{-k} b^2}{n^2} \max_{y \in [0; 0.1]} (I - \frac{1}{2}) \frac{4^{-k} b^2}{n^2} 18.6487 \frac{1}{2}:$$

We conclude that

$$\begin{aligned} KL(P; P^0) &\leq \sum_{k=0}^{\infty} 2^k (2^{-k} P + (1 - 2^{-k}) P^0; P^0) \\ &\leq \frac{37.2974 b^2 \ln 2}{n^2} + \frac{8.233 b^2}{n^2}: \end{aligned}$$

Finally, we prove the two claims used above:

1. For all $y; r \in [0, 1]$, let $h(x) = e^{-x^2}$, we have $y^2 - e^{-(y+)^2} = \int_0^y h^0(t) \, dt$. Since $h^0(t) = -2te^{-2t^2}$, we have $h^0(t) \leq 0$, we have $t = 1 - \frac{1}{2}$. Thus, $h^0(t) \leq e^{-2(\frac{1}{2})^2} = \frac{1}{2}$ and $y^2 - e^{-(y+)^2} \leq \frac{1}{2}$.

2. Suppose $r \in [\frac{1}{2}, 1]$, we have

$$\begin{aligned} \frac{\partial}{\partial y} \frac{e^{-y^2} - e^{-(y+)^2}}{2} &= \frac{e^{-(y+)^2} (2y + 2^2 + 1) - e^{-y^2}}{2} \\ &= \frac{e^{-y^2}}{2} [e^{-2y^2} (2y + 2^2 + 1) - 1]: \end{aligned}$$

Let $g(y; r) = e^{-2y^2} (2y + 2^2 + 1) - 1$. Note that

$$(a) \lim_{\eta \rightarrow 0} \frac{\partial}{\partial \eta} \left(\frac{e^{-y^2} - e^{-(y+\eta)^2}}{\eta} \right) = 0.$$

$$(b) \frac{\partial}{\partial \eta} g(\eta) = 4e^{-(y+\eta)^2} - 2e^{-(y+\eta)^2} < 0 \text{ and } g(0) = 0.$$

It implies that $\frac{\partial}{\partial \eta} \left(\frac{e^{-(1-\eta)^2} - e^{-(1-\eta+\eta)^2}}{\eta} \right) > 0$ for $\eta > 0$. Since $\frac{\partial}{\partial \eta} g(\eta) = 4e^{-(y+\eta)^2} - 2e^{-(y+\eta)^2} < 0$ for $\eta > 0$, we conclude that for any $\eta > 0$:

$$\frac{\partial}{\partial \eta} \frac{e^{-y^2} - e^{-(y+\eta)^2}}{\eta} > \frac{\partial}{\partial \eta} \frac{e^{-(1-\eta)^2} - e^{-(1-\eta+\eta)^2}}{\eta} > 0.$$

■

A.3 MAIN THEOREM

Recall that SGLD on data \mathcal{S} is defined as

$$W_t = W_{t-1} - \eta \nabla_w F(W_{t-1}; S_{B_t}) + \sqrt{\frac{\eta}{2}} N(0; I_d).$$

Here η is the step size. $B_t = \{i_1, \dots, i_b\}$ is a subset of $\{1, \dots, n\}$ of size b , and $S_{B_t} = (z_{i_1}, \dots, z_{i_b})$ is the mini-batch indexed by B_t . Recall that $F(w; S)$ denotes $\frac{1}{|S|} \sum_{i \in S} F(w; z_i)$. We restate and prove Theorem 11 in the following.

Theorem 11. Suppose that the loss function is C -bounded and the objective function is L -Lipschitz. Assume that the following conditions hold:

1. Batch size $b = n/2$.
2. Learning rate $\eta = \frac{1}{20L}$.

Then, the following expected generalization error bound holds for iterations of SGLD(1):

$$\text{err}_{\text{gen}} \leq \frac{8:12C}{n} \sqrt{\frac{\sum_{t=1}^T \mathbb{E}_{S_D} \left[\frac{1}{|S_t|} \sum_{i \in S_t} \|\nabla F(w_t; z_i)\|_2^2 \right]}{t}}; \quad (\text{empirical norm})$$

where $g_e(t) = \mathbb{E}_w \sum_{i=1}^n \|\nabla F(w; z_i)\|_2^2$ is the empirical squared gradient norm, and η is the parameter at step t of SGLD.

Proof By Theorem 7, we have

$$\text{err}_{\text{gen}} \leq 2C \mathbb{E}_z^P \sqrt{2\text{KL}(Q_z; P)} \quad (13)$$

for any prior distribution P . In particular, we define the prior $\mathcal{P}(w) = \mathbb{E}_{S \sim \mathcal{D}} [P_S(w)]$, where $P_S(w) = Q_{(\bar{S}; 0)}$. By the convexity of KL-divergence,

$$\text{KL}(Q_z; P) = \text{KL} \left(\mathbb{E}_S [Q_{(\bar{S}; z)}]; \mathbb{E}_S [Q_{(\bar{S}; 0)}] \right) \leq \mathbb{E}_S \text{KL}(Q_{(\bar{S}; z)}; Q_{(\bar{S}; 0)}); \quad (14)$$

Fix a data point $z \in Z$. Let $(W_t)_{t=0}^\infty$ and $(W_t^0)_{t=0}^\infty$ be the training process of SGLD for $\bar{S} = (\bar{S}; z)$ and $\bar{S}^0 = (\bar{S}; 0)$, respectively. Fix a time step t and $w_{<t} = (w_0, \dots, w_{t-1})$. Let P_t and P_t^0 denote the distribution of W_t and W_t^0 conditioned on $w_{<t} = w_{<t}$ and $w_{<t}^0 = w_{<t}$, respectively. By the definition of SGLD, we have $P_t = \frac{1}{|G|} \int_{B \in G} P_B$ and $P_t^0 = \frac{1}{|G|} \int_{B \in G} P_B^0$, where $B = w_{t-1} - \eta \nabla_w F(w_{t-1}; S_B)$, $P_B^0 = w_{t-1} - \eta \nabla_w F(w_{t-1}; S_B^0)$, and p denotes the Gaussian distribution $N(\cdot; \frac{\eta}{2} I_d)$. We note that:

1. $k_B^0 = k_B$ for $B \in G_n$ and $k_B = 0$ for $B \in \bar{G}_n$.

$$2. \text{diam}(f : B \rightarrow B) \leq 2L \quad t=10.$$

By applying Lemma 21 with $\epsilon = \frac{1}{b} \text{kr} F(w_{t-1}; z) k_2$ and $\delta = \epsilon$,

$$\text{KL}(P_t; P_t^0) \leq \frac{8 \cdot 23 \epsilon^2 \text{kr} F(w_{t-1}; z) k_2^2}{t n^2}.$$

By Lemma 10,

$$\text{KL}(W_T; W_T^0) = \sum_{t=1}^T \mathbb{E}_{W_{<t}} [\text{KL}(P_t; P_t^0)] \leq \sum_{t=1}^T \mathbb{E}_{W_{t-1}} \frac{8 \cdot 23 \epsilon^2 \text{kr} F(w; z) k_2^2}{t n^2};$$

which implies that

$$\text{KL}(Q_S; Q_S^0) = \text{KL}(W_T; W_T^0) \leq \sum_{t=1}^T \mathbb{E}_{W_{t-1}} \frac{8 \cdot 23 \epsilon^2 \text{kr} F(w; z) k_2^2}{t n^2};$$

Together with (13) and (14), we have

$$\text{err}_{\text{gen}} \leq 2C \mathbb{E}_z \frac{\sum_{t=1}^T \mathbb{E}_{W_{t-1}} \frac{8 \cdot 23 \epsilon^2 \text{kr} F(w; z) k_2^2}{t n^2}}{2 \mathbb{E}_S} \leq 2C \mathbb{E}_z \frac{\sum_{t=1}^T \mathbb{E}_{W_{t-1}} \frac{8 \cdot 23 \epsilon^2 \text{kr} F(w; z_n) k_2^2}{t n^2}}{2 \mathbb{E}_S} \quad (\text{concavity of } \mathbb{E} \bar{x})$$

Since SGLD is order-independent, we can replace $F(w; z)$ with $F(w; z_i)$ for any $i \in [n]$ in the right-hand side of the above bound. Our theorem then follows from the concavity of $\mathbb{E} \bar{x}$. Furthermore, if we bound $\text{KL}(P; Q_z)$ instead of $\text{KL}(Q_z; P)$ in the above proof, we obtain the following bound that depends on the population squared gradient norm

$$\text{err}_{\text{gen}} \leq \frac{8 \cdot 12C}{n} \sum_{t=1}^T \mathbb{E}_S \frac{\sum_{t=1}^T \mathbb{E}_{W_{t-1}} \mathbb{E}_{z \sim D} \text{kr} F(w; z) k_2^2}{t n^2};$$

■

A.4 EXTENSION TO GENERAL NOISES

We can extend the generalization bounds in previous sections, which require the noise to be Gaussian, to other general noises, namely the family of log-lipschitz noises.

Definition 22 (Log-Lipschitz Noises) A probability distribution on \mathbb{R}^d with density p is L -log-lipschitz if and only if $\text{kr} \ln p(w) k \leq L$ holds for any $w \in \mathbb{R}^d$. A random variable is called an L -log-lipschitz noise if and only if it is drawn from an L -log-lipschitz distribution.

The analog of SGLD, noisy momentum method (Definition 24), and noisy NAG (Definition 25) can be naturally defined by replacing the Gaussian noise at each iteration with an independent L -log-lipschitz noise in the definition.

The following lemma is an analog of Lemma 21 under log-lipschitz noises. Recall that \mathcal{B} denotes a collection of mini-batches of size b . Lemma 23 readily implies the analogs of Theorems 11, 26 and 27 under more general noise distributions.

Lemma 23. Suppose that batch size $n=2$ and N is an L_{noise} -log-lipschitz distribution on \mathbb{R}^d . $f_B : B \rightarrow \mathbb{G}$ and $f_B^0 : B \rightarrow \mathbb{G}$ are two collections of points in \mathbb{R}^d that satisfy the following conditions for some constant $0 < \frac{1}{L_{\text{noise}}}$:

1. $\|k_B - \frac{0}{B} k\|$ for $B \in \mathbb{G}_n$ and $\|B - \frac{0}{B}\|$ for $B \in \overline{\mathbb{G}_n}$.
2. $\text{diam}(f_B : B \in \mathbb{G}_n \cup [f_B^0 : B \in \mathbb{G}_n]) \leq 1$.

For $w \in \mathbb{R}^d$, let p denote the distribution of w when w is drawn from N . Let $P = \frac{1}{|\mathbb{G}_n|} \sum_{B \in \mathbb{G}_n} p_B$ and $P^0 = \frac{1}{|\mathbb{G}_n|} \sum_{B \in \mathbb{G}_n} p_B^0$ be mixture distributions over all mini-batches. Then,

$$\text{KL}(P; P^0) \leq \frac{C_0 b^2}{n^2}$$

for some constant C_0 that only depends on L_{noise} .

Proof of Lemma 23 Following the same argument as in the proof of Lemma 21, we have

$$\text{KL}(P; P^0) \leq \sum_{k=0}^{\infty} \frac{1}{2^k} \left(\sum_{B \in \mathbb{G}_n} p_B + \sum_{B \in \overline{\mathbb{G}_n}} p_B^0 \right) \quad (15)$$

where

$$\left(\sum_{B \in \mathbb{G}_n} p_B + \sum_{B \in \overline{\mathbb{G}_n}} p_B^0 \right) \leq \frac{4 b^2}{n^2 |\mathbb{G}_n|} \sum_{A \in \mathbb{G}_n} \sum_{B \in \overline{\mathbb{G}_n}} \int_{\mathbb{R}^d} \frac{(p_B(w) - p_B^0(w))^2}{p_A(w)} dw \quad (16)$$

Fix $A \in \mathbb{G}_n$ and $B \in \overline{\mathbb{G}_n}$. Let p_{noise} denote the density of the noise distribution. Since $\|k_A - k_B\| \leq 1$ and p_{noise} is L_{noise} -log-lipschitz, we have

$$p_B(w) = p_{\text{noise}}(w - k_B) - p_{\text{noise}}(w - k_A) \leq e^{-L_{\text{noise}} \|k_A - k_B\|} p_A(w)$$

Similarly, since $\|k_B - \frac{0}{B} k\| \leq 1$, we have

$$e^{-L_{\text{noise}} \|k_B - \frac{0}{B} k\|} p_B(w) \leq p_B^0(w) \leq e^{L_{\text{noise}} \|k_B - \frac{0}{B} k\|} p_B(w)$$

Then, it follows from $L_{\text{noise}} \leq 1$ that

$$(p_B(w) - p_B^0(w))^2 \leq (e^{L_{\text{noise}}} - 1)^2 p_B(w)^2 \leq 2 L_{\text{noise}}^2 p_B(w)^2$$

Therefore, the integral on the righthand side of (16) can be upper bounded as follows:

$$\begin{aligned} \int_{\mathbb{R}^d} \frac{(p_B(w) - p_B^0(w))^2}{p_A(w)} dw &\leq 2 L_{\text{noise}}^2 \int_{\mathbb{R}^d} \frac{p_B(w)^2}{p_A(w)} dw \\ &\leq 2 L_{\text{noise}}^2 \int_{\mathbb{R}^d} p_B(w) e^{-L_{\text{noise}} \|k_A - k_B\|} dw \\ &= 2 L_{\text{noise}}^2 e^{-L_{\text{noise}}}. \end{aligned}$$

Plugging the above inequality into (15) and (16) gives

$$\left(\sum_{B \in \mathbb{G}_n} p_B + \sum_{B \in \overline{\mathbb{G}_n}} p_B^0 \right) \leq \frac{4 b^2}{n^2 |\mathbb{G}_n|} \sum_{A \in \mathbb{G}_n} \sum_{B \in \overline{\mathbb{G}_n}} 2 L_{\text{noise}}^2 e^{-L_{\text{noise}}} = L_{\text{noise}}^2 e^{-L_{\text{noise}}} \frac{4 b^2}{n^2}$$

and

$$\text{KL}(P; P^0) \leq \sum_{k=0}^{\infty} \frac{1}{2^k} L_{\text{noise}}^2 e^{-L_{\text{noise}}} \frac{4 b^2}{n^2} = 2 \ln 2 L_{\text{noise}}^2 e^{-L_{\text{noise}}} \frac{b^2}{n^2}.$$

■

A.5 EXTENSION TO OTHER GRADIENT-BASED METHODS

A.5.1 STABILITY BOUND FOR MOMENTUM AND NESTEROV'S ACCELERATED GRADIENT

We adopt the formulation of Classical Momentum and Nesterov's Accelerated Gradient (NAG) methods in Sutskever et al. (2013) and consider the noisy versions of them.

Definition 24 (Noisy Momentum Method) Noisy Momentum Method on objective function $F(w; z)$ and dataset \mathcal{S} is defined as

$$\begin{aligned} V_t &= V_{t-1} + \eta \nabla_w F(W_{t-1}; \mathcal{S}_{B_t}) + \xi_t \\ W_t &= W_{t-1} + V_t \end{aligned}$$

Definition 25 (Noisy Nesterov's Accelerated Gradient) Noisy Nesterov's Accelerated Gradient (NAG) on objective function $F(w; z)$ and dataset \mathcal{S} is defined as

$$\begin{aligned} V_t &= V_{t-1} + \eta \nabla_w F(W_{t-1} + V_{t-1}; \mathcal{S}_{B_t}) + \xi_t; \\ W_t &= W_{t-1} + V_t; \end{aligned}$$

In both definitions, η is the step size, mini-batch B_t is drawn uniformly from \mathcal{G} , ξ_t is a Gaussian noise drawn from $\mathcal{N}(0; \frac{\sigma^2}{2} I_d)$, and $\beta \in [0, 1]$ is the momentum coefficient.

Theorem 26. Under the same assumptions on the loss function, objective function, batch size and learning rate as in Theorem 11, the generalization bounds in Theorem 11 still hold for noisy momentum method and noisy NAG.

Proof of Theorem 26 For any time step t and $w_{<t} = (w_0; w_1; \dots; w_{t-1})$, let P_t and P_t^0 denote the distribution of W_t and W_t^0 conditioned on $W_{<t} = w_{<t}$ and $W_{<t}^0 = w_{<t}$, respectively. By definition, we have $P_t = \frac{1}{|\mathcal{G}|} \prod_{B \in \mathcal{G}} P_B$ and $P_t^0 = \frac{1}{|\mathcal{G}|} \prod_{B \in \mathcal{G}} P_B^0$:

If $t = 1$, for both noisy momentum method and noisy NAG, we have

$$\begin{aligned} W_B &= W_{t-1} + \eta \nabla_w F(W_{t-1}; \mathcal{S}_B); \\ W_B^0 &= W_{t-1} + \eta \nabla_w F(W_{t-1}; \mathcal{S}_B^0); \end{aligned}$$

For $t > 1$, if noisy momentum method is used, we have

$$\begin{aligned} W_B &= W_{t-1} + (W_{t-1} - W_{t-2}) + \eta \nabla_w F(W_{t-1}; \mathcal{S}_B); \\ W_B^0 &= W_{t-1} + (W_{t-1} - W_{t-2}) + \eta \nabla_w F(W_{t-1}; \mathcal{S}_B^0); \end{aligned}$$

Similarly, the following holds under noisy NAG:

$$\begin{aligned} W_B &= W_{t-1} + (W_{t-1} - W_{t-2}) + \eta \nabla_w F(W_{t-1} + (W_{t-1} - W_{t-2}); \mathcal{S}_B); \\ W_B^0 &= W_{t-1} + (W_{t-1} - W_{t-2}) + \eta \nabla_w F(W_{t-1} + (W_{t-1} - W_{t-2}); \mathcal{S}_B^0); \end{aligned}$$

In either case, it can be verified that the conditions of Lemma 21 hold for $\frac{2\eta L}{b}$ and $\sigma = \frac{\sigma}{\sqrt{2}}$. The rest of the proof is the same as the proof of Theorem 11. \blacksquare

A.5.2 STABILITY BOUND FOR ENTROPY-SGD

In the Entropy-SGD algorithm due to Chaudhari et al. (2017), instead of directly optimizing the original objective $F(w)$, we minimize the negative local entropy defined as follows:

$$E(w; \eta) = -\log \int_{\mathcal{R}^d} \exp(-F(w^0)) \frac{1}{2} \eta \|w - w^0\|_2^2 dw \quad (17)$$

Intuitively, a wider local minimum has a lower loss (i.e. $E(w; \eta)$) than sharper local minima. See Chaudhari et al. (2017) for more details. The Entropy-SGD algorithm invokes standard SGD to minimize the negative local entropy. However, the gradient of negative local entropy

$$\nabla_w E(w; \eta) = -\int_{\mathcal{R}^d} \frac{w - w^0}{\eta} \exp(-F(w^0)) \frac{1}{2} \eta \|w - w^0\|_2^2 dw \quad (18)$$

is hard to compute. Thus, the algorithm uses exponential averaging to estimate the gradient in the SGLD loop; see Algorithm 1 for more details.

We have the following generalization bound for Entropy-SGD.

Algorithm 1: Entropy-SGD

Input: Training set $S = (z_1; \dots; z_n)$ and loss function $\ell(w; z)$.
Hyper-parameters: Scope \mathcal{S} , SGD learning rate η , SGLD step size σ^2 and batch size b .

```

1 for t = 1 to T do
2   //SGD iteration
3    $W_{t,0}; \dots; W_{t,1;K+1}$ ;
4   for k = 0 to K - 1 do
5     //SGLD iteration
6      $B_{t,k}$  mini-batch with size  $b$ ;
7      $W_{t,k+1} = W_{t,k} - \eta \nabla_w g(W_{t,k}; S_{B_{t,k}}) + \sigma \sqrt{\eta} N(0; \frac{1}{2} I_d)$ ;
8      $W_{t,k+1} = (1 - \eta \eta_{t,k}) W_{t,k} + W_{t,k}$ ;
9   end
10   $W_{t,K+1} = W_{t,K} - \eta_{t,K} \nabla_w \ell(W_{t,K}; S)$ ;
11 end
12 return  $W_{T,K+1}$ ;

```

Theorem 27. Suppose that the loss function is C -bounded and the objective function is L -lipschitz. If batch size $b = n/2$ and $\sigma^2 = (20L)^2$, the following expected generalization error bound holds for Entropy-SGD:

$$\text{err}_{\text{gen}} \leq \frac{8 \cdot 12 C^2 \sigma^2}{n} \sum_{t=1}^T \sum_{k=0}^{K-1} \mathbb{E}_S \left[\frac{1}{n} \sum_{i=1}^n \|\nabla_w \ell(w_{t,k}; z_i)\|_2^2 \right] \quad (\text{empirical norm})$$

where $g_e(t; k) = \mathbb{E}_w \mathbb{E}_{W_{t,k}} \left[\frac{1}{n} \sum_{i=1}^n \|\nabla_w \ell(w; z_i)\|_2^2 \right]$ is the empirical squared gradient norm, and $W_{t,k}$ denotes the training process with respect to S .

Since $g_e(t; k)$ is at most L^2 , it further implies the generalization error of Entropy-SGD is bounded by $O\left(\frac{C^2 \sigma^2}{n} T K\right)$.

Proof of Theorem 27. Denote the history before time step (t, k) as follows:

$$W_{(t,k)} = (W_{0,0}; \dots; W_{0,K+1}; \dots; W_{t,1,0}; \dots; W_{t,1,K+1}; W_{t,0}; \dots; W_{t,k}): \quad (19)$$

Since $W_{(t,k)}$ is only determined by W , we only need to focus on W . This proof is similar to the proof of Theorem 11. By setting $\mathcal{S} = \mathbb{E}_S[\mathcal{Q}(S; 0)]$. Suppose $\mathcal{S} = (S; z)$ and $S^0 = (S; 0)$ are fixed, let W and W^0 denote their training process, respectively. Considering the following 3 cases:

1. $W_{t,0} = W_{t,1;K+1}$: In this case, for a fixed $W_{(t,1;K+1)}$, we have

$$\text{KL}(W_{t,0}; W_{(t,1;K+1)}) = \text{KL}(W_{(t,1;K+1)}; W_{(t,1;K+1)}^0) = 0;$$

2. $W_{t,k+1} = W_{t,k} - \eta \nabla_w g(W_{t,k}; S_{B_{t,k}}) + \sigma \sqrt{\eta} N(0; \frac{1}{2} I_d)$: In this case, $W_{(t,k)}$ is fixed, applying Lemma 21 gives

$$\text{KL}(W_{t,k+1}; W_{(t,k)}) = \text{KL}(W_{(t,k)}; W_{(t,k)}^0) = \frac{8 \cdot 23 \sigma^2 \mathbb{E}_S \left[\sum_{i=1}^n \|\nabla_w \ell(w_{t,k}; z_i)\|_2^2 \right]}{2n^2};$$

3. $W_{t,K+1} = W_{t,K} - \eta_{t,K} \nabla_w \ell(W_{t,K}; S)$: In this case, for a fixed $W_{(t,K)}$, we have

$$\text{KL}(W_{t,K+1}; W_{(t,K)}) = \text{KL}(W_{(t,K)}; W_{(t,K)}^0) = 0;$$

By applying Lemma 10, we have

$$\text{KL}(W_{T,K+1}; W_{T,K+1}^0) \leq \frac{8 \cdot 23 \sigma^2 \sum_{t=1}^T \sum_{k=0}^{K-1} \mathbb{E}_S \left[\sum_{i=1}^n \|\nabla_w \ell(w_{t,k}; z_i)\|_2^2 \right]}{2n^2};$$

and where $g_e(t; k)$ is the empirical squared gradient norm of the k -th SGLD iteration in the t -th SGD iteration, respectively. The rest of the proof is the same as the proof of Theorem 11. ■

B PROOFS IN SECTION 4

B.1 MARKOV SEMIGROUP AND LOG-SOBOLEV INEQUALITY

The continuous version of the noisy gradient descent method is the Langevin dynamics, described by the following stochastic differential equation:

$$dW_t = -\nabla F(W_t) dt + \sqrt{2} dB_t; \quad W_0 = w_0; \quad (20)$$

where B_t is the standard Brownian motion. To analyze the above Langevin dynamics, we need some preliminary knowledge about Log-Sobolev inequalities.

Let $p_t(w; y)$ denote the probability density function (i.e., probability kernel) describing the distribution of W_t starting from w . For a given SDE such as (20), we can define the associated diffusion semigroup P_t :

Definition 28 (Diffusion Semigroup) (see e.g., (Bakry et al., 2013, p. 39)) Given a stochastic differential equation (SDE), the associated diffusion semigroup $(P_t)_{t \geq 0}$ is a family of operators that satisfy for every $t \geq 0$, P_t is a linear operator sending any real-valued bounded measurable function f on \mathbb{R}^d to

$$P_t f(w) = \mathbb{E}[f(W_t) | W_0 = w] = \int_{\mathbb{R}^d} f(y) p_t(w; dy)$$

The semigroup property $P_{t+s} = P_t \circ P_s$ holds for every $t, s \geq 0$. Another useful property of P_t is that it maps a nonnegative function to a nonnegative function. The **Carathéodory** operator of this diffusion semigroup (w.r.t (20)) is (Bakry et al., 2013, p. 42)

$$\mathcal{L}f = \frac{1}{2} \text{tr}(\nabla^2 f) - \nabla f \cdot \nabla F$$

We use the shorthand notation $\text{Ent}(f) = \int_{\mathbb{R}^d} f \log f \, d\mu = \int_{\mathbb{R}^d} f \log f \, d\mu - \int_{\mathbb{R}^d} f \log \int_{\mathbb{R}^d} f \, d\mu \, d\mu$, and define (with the convention that $0 \log 0 = 0$)

$$\text{Ent}(f) = \int_{\mathbb{R}^d} f \log f \, d\mu - \int_{\mathbb{R}^d} f \log \int_{\mathbb{R}^d} f \, d\mu \, d\mu$$

Definition 29 (Logarithmic Sobolev Inequality) (see e.g., (Bakry et al., 2013, p. 237)) A probability measure μ is said to satisfy a logarithmic Sobolev inequality $\text{LS}(\rho)$ (with respect to μ), if for all functions $f: \mathbb{R}^d \rightarrow \mathbb{R}^+$ in the Dirichlet domain $\mathcal{D}(E)$,

$$\text{Ent}(f) \leq \frac{\rho}{2} \int_{\mathbb{R}^d} \frac{(\mathcal{L}f)^2}{f} \, d\mu$$

$\mathcal{D}(E)$ is the set of functions $f \in L^2(\mu)$ for which the quantity $\int_{\mathbb{R}^d} \frac{(\mathcal{L}f)^2}{f} \, d\mu$ has a finite (decreasing) limit as ρ decreases to 0.

A well-known Logarithmic Sobolev Inequality is the following result for Gaussian measures.

Lemma 30 (Logarithmic Sobolev Inequality for Gaussian measure) (Bakry et al., 2013, p. 258)

Let μ be the centered Gaussian measure on \mathbb{R}^d with covariance matrix $2^{-1}I_d$. Then μ satisfies the following LSI:

$$\text{Ent}(f) \leq \frac{1}{2} \int_{\mathbb{R}^d} \frac{(\mathcal{L}f)^2}{f} \, d\mu$$

Lemma 30 states that the centered Gaussian measure with covariance matrix Σ satisfies $\text{LS}(\frac{1}{2})$ (with respect to μ), where $\mathcal{L}f = \frac{1}{2} \text{tr}(\nabla^2 f) - \nabla f \cdot \nabla F$ is the **Carathéodory** operator of the diffusion semigroup defined above.

Before proving our results, we need some known results from Markov diffusion process. It is well known that the invariant measure (Bakry et al., 2013, p. 10) of the above CLD is the Gibbs measure $\mu = \frac{1}{Z} \exp(-F(w)) \, dw$ (Menz et al., 2014, (1.3)). In other words, $\mu = \int_{\mathbb{R}^d} P_t f \, d\mu = \int_{\mathbb{R}^d} f \, d\mu$ for every bounded positive measurable function f on \mathbb{R}^d where P_t is the Markov semigroup in Definition 28. The following lemma by Holley and Stroock (Holley & Stroock (1987) (see also (Bakry et al., 2013, p. 240))) allows us to determine the Logarithmic Sobolev constant of the invariant measure μ .

Lemma 31 (Bounded perturbation) Assume that the probability measure μ satisfies LS(σ) (with respect to ν). Let μ' be a probability measure such that $\mu' \ll \mu$ and $\frac{d\mu'}{d\mu} \leq C$ for some constant $C > 1$. Then μ' satisfies LS($\frac{\sigma}{C}$) (with respect to ν).

In fact, Lemma 31 is a simple consequence of the following variational formula in the special case that $\nu(x) = x \log x$, which we will also need in our proof:

Lemma 32 (Variational formula) (see .g., (Bakry et al., 2013, p. 240)) Let $I \subset \mathbb{R}$ on some open interval \mathbb{R} be convex of class \mathcal{C}^2 . For every (bounded or suitably integrable) measurable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ with values in I ,

$$\int_{\mathbb{R}^d} f(x) d\mu(x) - \int_{\mathbb{R}^d} f(x) d\nu(x) = \inf_{r \in I} \int_{\mathbb{R}^d} [f(x) - (r) - \frac{\sigma}{2} (f(x) - r)^2] d\nu(x) \quad (21)$$

It is worth noting the integrand of the right-hand side is nonnegative due to the convexity of

B.2 LOGARITHMIC SOBOLEV INEQUALITY FOR CLD

Recall that $F_S(w) = F(w; S) := F_0(w; S) + \frac{\sigma}{2} \|w\|_2^2$ is the sum of the empirical original objective $F_0(w; S)$ and $\frac{\sigma}{2}$ regularization. Let $\mu = \frac{1}{Z} \exp(-F_S(w)) dw$ be the invariant (Gibbs) measure of CLD, and $\nu = \frac{1}{Z} \exp(-\frac{\sigma}{2} \|w\|_2^2) dw$. Invoking Lemma 30 with $\sigma = \frac{\sigma}{2}$ shows that μ satisfies LS($\frac{\sigma}{4}$) (with respect to ν). Consider the density $h(w) = \frac{d\mu}{d\nu} = \frac{1}{Z} \exp(-F_0(w; S))$. If the original objective function F_0 is C -bounded, we have $\exp(-2C) \leq h(w) \leq \exp(2C)$. By applying Lemma 31 with $\mu' = \mu$, we have the following lemma.

Lemma 33. Under Assumption 14, let $(f; g) = \langle \cdot, \cdot \rangle_{\mu}$ be the carré du champ operator of the diffusion semigroup associated to CLD, and μ the invariant measure of the SDE. Then, μ satisfies LS($\frac{\sigma}{4} C$) with respect to ν .

Let μ_t be the probability measure μ_t . By definition of P_t , for any real-valued bounded measurable function f on \mathbb{R}^d and any $s; t \geq 0$,

$$\int_{\mathbb{R}^d} f(w) d\mu_{t+s}(w) = \int_{\mathbb{R}^d} [P_t f](w) d\mu_s(w) \quad (22)$$

In particular, if the invariant measure $\mu = \mu_1$ exists, we have

$$\int_{\mathbb{R}^d} f(w) d\mu(w) = \int_{\mathbb{R}^d} [P_1 f](w) d\mu(w) = \int_{\mathbb{R}^d} f(w) d\mu(w) = \int_{\mathbb{R}^d} [P_t f](w) d\mu(w) \quad (23)$$

The following lemma is crucial for establishing the first generalization bound for CLD. In fact, we establish a Log-Sobolev inequality for μ_t , the parameter distribution at time t for any time $t > 0$. Note that our choice of the initial distribution μ_0 is important for the proof⁶.

Lemma 34. Under Assumption 14, let μ_t be the probability measure μ_t in (CLD) with initial probability measure $\mu_0 = \frac{1}{Z} e^{-\frac{\sigma}{2} \|w\|_2^2} dw$. Let $(f; g)$ be the carré du champ operator of diffusion semigroup associated to (CLD). Then, for any $f : \mathbb{R}^d \rightarrow \mathbb{R}^+$ in $D(E)$:

$$\text{Ent}_{\mu_t}(f) \leq \frac{e^{2C}}{2} \int_{\mathbb{R}^d} \frac{(f; g)}{f} d\mu_t$$

Proof Let μ be the invariant measure of CLD. By Lemma 33 and Definition 29,

$$\text{Ent}_{\mu_t}(f) \leq \frac{e^{4C}}{2} \int_{\mathbb{R}^d} \frac{(f; g)}{f} d\mu \quad (24)$$

⁶ For arbitrary initial distribution, it is impossible to prove similar inequality for any 0 (unless the loss is strongly convex).

By applying Lemma 32 with $(x) = x \log x$, we rewrite the left-hand side as

$$\begin{aligned} \text{Ent}(f) &:= \int_{\mathbb{R}^d} f \log f \, d\mu - \int_{\mathbb{R}^d} f \, d\mu \log \int_{\mathbb{R}^d} f \, d\mu \\ &= \inf_{r \geq 1} \int_{\mathbb{R}^d} [P_t(f)(r) - \int_{\mathbb{R}^d} P_t(f)(r)] \, d\mu \\ &= \inf_{r \geq 1} \int_{\mathbb{R}^d} [P_t(f)(r) - \int_{\mathbb{R}^d} P_t(f)(r)] \, d\mu : \end{aligned}$$

where the last equation holds by the definition of invariant measure $\mu = \int_{\mathbb{R}^d} f \, d\mu$. Thus, we have

$$\inf_{r \geq 1} \int_{\mathbb{R}^d} [P_t(f)(r) - \int_{\mathbb{R}^d} P_t(f)(r)] \, d\mu = \text{Ent}(f) \frac{e^{4C}}{2} \int_{\mathbb{R}^d} \frac{kr f k_2^2}{f} \, d\mu; \quad (25)$$

Let ν_t be the probability measure ν_t . Lemma 32 and (22) together imply that

$$\begin{aligned} \text{Ent}_{\nu_t}(f) &= \inf_{r \geq 1} \int_{\mathbb{R}^d} [P_t(f)(r) - \int_{\mathbb{R}^d} P_t(f)(r)] \, d\nu_t \\ &= \inf_{r \geq 1} \int_{\mathbb{R}^d} [P_t(f)(r) - \int_{\mathbb{R}^d} P_t(f)(r)] \, d\mu_0 \end{aligned} \quad (26)$$

Since $P_t(f)(r) - \int_{\mathbb{R}^d} P_t(f)(r) \geq 0$ and $\frac{d\nu_t}{d\mu_0} \leq \exp(2C)$, we have

$$\begin{aligned} \text{Ent}_{\nu_t}(f) &= \inf_{r \geq 1} \int_{\mathbb{R}^d} [P_t(f)(r) - \int_{\mathbb{R}^d} P_t(f)(r)] \frac{d\nu_t}{d\mu_0} \, d\mu_0 \\ &\leq \exp(2C) \text{Ent}(f) \frac{e^{6C}}{2} \int_{\mathbb{R}^d} \frac{kr f k_2^2}{f} \, d\mu; \end{aligned} \quad (27)$$

Since $\frac{d\nu_t}{d\mu_0} \leq \exp(2C)$ and μ is the invariant measure, we conclude that

$$\begin{aligned} \text{Ent}_{\nu_t}(f) &\leq \frac{e^{6C}}{2} \int_{\mathbb{R}^d} \frac{kr f k_2^2}{f} \, d\mu = \frac{e^{6C}}{2} \int_{\mathbb{R}^d} P_t \left(\frac{kr f k_2^2}{f} \right) \, d\mu \\ &= \frac{e^{6C}}{2} \int_{\mathbb{R}^d} P_t \left(\frac{kr f k_2^2}{f} \right) \frac{d\nu_t}{d\mu_0} \, d\mu_0 \\ &\leq \frac{e^{8C}}{2} \int_{\mathbb{R}^d} P_t \left(\frac{kr f k_2^2}{f} \right) \, d\mu_0 \\ &= \frac{e^{8C}}{2} \int_{\mathbb{R}^d} \frac{1}{f} \, d\nu_t = \frac{e^{8C}}{2} \int_{\mathbb{R}^d} \frac{(f)}{f} \, d\nu_t \end{aligned} \quad (28)$$

■

Lemma 16. Under Assumption 14, let ν_t be the probability measure ν_t in CLD (with $d\nu_0 = \frac{1}{Z} e^{-\frac{kwk^2}{2}} dw$). Let μ be a probability measure that is absolutely continuous with respect to ν_t . Suppose $d\mu = \mu(w) dw$ and $d\nu_t = \nu_t(w) dw$. Then it holds that:

$$\text{KL}(\mu; \nu_t) \leq \frac{\exp(8C)}{2} \int_{\mathbb{R}^d} \mu(w) \log \frac{\mu(w)}{\nu_t(w)} \, d\nu_t; \quad (29)$$

Proof Let $f(w) = \mu(w) - \nu_t(w)$, by Lemma 34 and $\int_{\mathbb{R}^d} f \, d\nu_t = 1$, we have

$$\int_{\mathbb{R}^d} f \log f \, d\nu_t \leq \frac{e^{8C}}{2} \int_{\mathbb{R}^d} \frac{kr f k_2^2}{f} \, d\nu_t \quad (30)$$

⁷ This is because \log is convex and P_t is a positive operator.

We can see that the left-hand side is equal to $\int_{\mathbb{R}^d} \frac{e^{\beta C}}{2} \int_{\mathbb{R}^d} \frac{r \frac{(w)}{t(w)}^2}{(w)=t(w)} t(w) dw$; t)⁸, and the right-hand side is equal to

$$\frac{e^{\beta C}}{2} \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \frac{r \frac{(w)}{t(w)}^2}{(w)=t(w)} t(w) dw = \frac{e^{\beta C}}{2} \int_{\mathbb{R}^d} r \log \frac{(w)}{t(w)}^2 (w) dw:$$

This concludes the proof. ■

B.3 THE DISCRETIZATION LEMMA FROM RAGINSKY ET AL. (2017)

Let $h(w; z) = F_0(w; z) + \frac{kwk_2^2}{2}$. We can rewrite $F_S(w) = \frac{1}{n} \sum_{i=1}^n h(w; z_i)$. Define $\nu_{s;k}$ and $\nu_{s;t}$ as the probability measure ν_k (in GLD) and ν_t (in CLD), respectively. Raginsky et al. (2017) provided a bound on $\text{KL}(\nu_{s;k}; \nu_{s;t})$ under Assumption 35. This bound enables us to derive a generalization error bound for the discrete GLD from the bound for the continuous CLD. We use the assumption from Raginsky et al. (2017). Their work considers the following SGLD:

$$W_{k+1} = W_k - g_S(W_k) + \frac{1}{2} \frac{1}{k} \epsilon_k$$

Where $g_S(w_k)$ is a conditionally unbiased estimate of the gradient $\nabla F_S(w_k)$. In our GLD setting, $g_S(W_k)$ is equal to $\nabla F_S(w_k)$.

Assumption 35. Let $F_S(w) = \frac{1}{n} \sum_{i=1}^n h(w; z_i) = F_0(w; S) + \frac{1}{2} kwk_2^2$.

1. The function h takes non-negative real values, and there exist constants $A, B \geq 0$, such that

$$|h(0; z)| \leq A \quad \text{and} \quad |h(0; z)| \leq B \quad \forall z \in \mathbb{Z}^d$$

2. For each $z \in \mathbb{Z}^d$, the function $h(\cdot; z)$ is M -smooth: for some $M > 0$,

$$|h(w; z) - h(v; z)| \leq M \|w - v\|_2; \quad \forall w, v \in \mathbb{R}^d$$

3. For each $z \in \mathbb{Z}^d$, the function $h(\cdot; z)$ is (m, b) -dissipative: for some $m > 0$ and $b \geq 0$,

$$|hw; r h(w; z)| \leq m kwk_2^2 + b; \quad \forall w \in \mathbb{R}^d$$

4. There exists a constant $\alpha \in [0, 1)$, such that, for each $S \in \mathbb{Z}^d$,

$$E[kg_S(w) - \nabla F_S(w)]_2^2 \leq M^2 kwk_2^2 + B^2; \quad \forall w \in \mathbb{R}^d$$

5. The probability law ν_0 of the initial hypothesis W_0 has a bounded and strictly positive density ν_0 with respect to the Lebesgue measure ν , and

$$\nu_0 := \int_{\mathbb{R}^d} e^{-kwk_2^2} \nu_0(w) dw < 1$$

Lemma 36. (Raginsky et al., 2017, Lemma 7) Suppose that Assumption 35 holds and let $\nu_{s;0} = \nu_0$. Then, for any $k \in \mathbb{N}$ and any $\beta \in (0, 1 \wedge \frac{m}{4M^2})$, the following inequality holds

$$\text{KL}(\nu_{s;k}; \nu_{s;k}) \leq (C_0 + C_1) k;$$

where C_0 and C_1 are constants that only depend on ν_0, m, b, α, B and d .

B.4 PROOFS FOR MAIN THEOREMS

Theorem 15. Under Assumption 14, CLD (with initial probability measure $\nu_0 = \frac{1}{Z} e^{-\frac{kwk_2^2}{2}} dw$) has the following expected generalization error bound:

$$\text{err}_{\text{gen}} \leq \frac{2e^{\beta C} \text{CL}}{n} - 1 \exp \left(-\frac{T}{e^{\beta C}} \right) \quad (31)$$

⁸Indeed, $\int_{\mathbb{R}^d} \log \frac{f}{g} d\nu = \int_{\mathbb{R}^d} \log \left(\frac{f}{g} \right) d\nu = \text{KL}(f; g)$

In addition, if F_0 is also M -smooth and non-negative, by setting $\gamma > 2$, $\beta > 0$ and $\alpha \in [0; 1 \wedge \frac{1}{8M^2})$, the GLD (running K iterations with the same θ_0 as CLD) has the expected generalization error bound:

$$\text{err}_{\text{gen}} \leq 2C \frac{p}{2KC_1^2 + \frac{2CLE^4 C}{n}} - 1 \exp \frac{K}{e^{8C}}; \quad (32)$$

where C_1 is a constant that only depends on m, β, b, L and d .

Proof of Theorem 15 We apply the uniform stability framework. Suppose S and S^0 are two neighboring datasets that differ on exactly one data point. $(W_t)_{t=0}$ and $(W_t^0)_{t=0}$ be the process of CLD running on S and S^0 , respectively. Let μ_t and μ_t^0 be the pdf of W_t and W_t^0 . We have

$$\begin{aligned} \frac{d}{dt} \text{KL}(\mu_t; \mu_t^0) &= \frac{d}{dt} \int_{\mathbb{R}^d} \mu_t \log \frac{\mu_t}{\mu_t^0} dw \\ &= \int_{\mathbb{R}^d} \frac{d\mu_t}{dt} \log \frac{\mu_t}{\mu_t^0} + \mu_t \frac{\frac{d\mu_t}{dt} \mu_t^0 - \mu_t \frac{d\mu_t^0}{dt}}{\mu_t^2} dw \\ &= \int_{\mathbb{R}^d} \frac{d\mu_t}{dt} \log \frac{\mu_t}{\mu_t^0} dw - \int_{\mathbb{R}^d} \frac{\mu_t}{\mu_t^0} \frac{d\mu_t}{dt} dw \end{aligned} \quad (33)$$

According to Fokker-Planck equation (see Risken (1996)) for CLD, we know that

$$\frac{\partial \mu_t}{\partial t} = \frac{1}{t+r} \nabla \cdot (\mu_t r F_{S^0}); \quad \frac{\partial \mu_t^0}{\partial t} = \frac{1}{t+r} \nabla \cdot (\mu_t^0 r F_S);$$

It follows that

$$\begin{aligned} I &:= \int_{\mathbb{R}^d} \frac{d\mu_t}{dt} \log \frac{\mu_t}{\mu_t^0} dw \\ &= \int_{\mathbb{R}^d} \frac{1}{t+r} \nabla \cdot (\mu_t r F_{S^0}) \log \frac{\mu_t}{\mu_t^0} dw \\ &= - \int_{\mathbb{R}^d} \text{tr} \log \frac{\mu_t}{\mu_t^0}; r \nabla \mu_t dw + \int_{\mathbb{R}^d} \text{tr} \log \frac{\mu_t}{\mu_t^0}; r F_{S^0} dw; \end{aligned} \quad (\text{integration by parts})$$

and

$$\begin{aligned} J &:= \int_{\mathbb{R}^d} \frac{\mu_t}{\mu_t^0} \frac{d\mu_t}{dt} dw \\ &= \int_{\mathbb{R}^d} \frac{\mu_t}{\mu_t^0} \frac{1}{t+r} \nabla \cdot (\mu_t^0 r F_S) dw \\ &= - \int_{\mathbb{R}^d} \text{tr} \frac{\mu_t}{\mu_t^0}; r \nabla \mu_t^0 dw + \int_{\mathbb{R}^d} \text{tr} \frac{\mu_t}{\mu_t^0}; r F_S dw; \end{aligned} \quad (\text{integration by parts})$$

Together with (33), we have

$$\begin{aligned} \frac{d}{dt} \text{KL}(\mu_t; \mu_t^0) &= I - J \\ &= \int_{\mathbb{R}^d} \text{tr} \log \frac{\mu_t}{\mu_t^0}; r \nabla \mu_t dw - \int_{\mathbb{R}^d} \text{tr} \log \frac{\mu_t}{\mu_t^0}; r F_{S^0} dw \\ &\quad - \int_{\mathbb{R}^d} \text{tr} \log \frac{\mu_t}{\mu_t^0}; r F_S dw + \int_{\mathbb{R}^d} \text{tr} \log \frac{\mu_t}{\mu_t^0}; r F_{S^0} dw \\ &= \int_{\mathbb{R}^d} \text{tr} \log \frac{\mu_t}{\mu_t^0}; r \nabla \mu_t dw + \int_{\mathbb{R}^d} \text{tr} \log \frac{\mu_t}{\mu_t^0}; r F_S - r F_{S^0} dw \\ &= \int_{\mathbb{R}^d} \text{tr} \log \frac{\mu_t}{\mu_t^0}; r \nabla \mu_t dw + \int_{\mathbb{R}^d} \text{tr} \log \frac{\mu_t}{\mu_t^0}; r F_S - r F_{S^0} dw; \end{aligned}$$

The last step holds because $\frac{1}{2} \leq \frac{1}{2} + \frac{1}{2} = 1$. Since $F_S = r F_{S^0} k_2^2 = \frac{4L^2}{n^2}$, by Lemma 16, we have

$$KL(W_t; W_{t'}) = \frac{e^{8C}}{2} \int_{R^d} \log \frac{t}{t'} dw;$$

which implies

$$\frac{1}{e^{8C}} KL(W_t; W_{t'}) = \frac{1}{2} \int_{R^d} \log \frac{t}{t'} dw;$$

Hence,

$$\frac{d}{dt} KL(W_t; W_{t'}) = \frac{1}{e^{8C}} KL(W_t; W_{t'}) + \frac{2L^2}{n^2}; \quad \text{with } KL(W_0; W_0) = 0: \quad (34)$$

Solving this differential inequality gives

$$KL(W_t; W_{t'}) = \frac{2L^2 e^{8C} (1 - e^{-t e^{8C}})}{n^2}. \quad (35)$$

By Pinsker's inequality, we can finally see that

$$\sup_z \mathbb{E} |L(W_T^0; z) - L(W_T; z)| \leq 2C \sqrt{\frac{1}{2} KL(W_T; W_T)} = \frac{2e^{4C} CL}{n} \sqrt{\frac{1 - e^{-T e^{8C}}}{n^2}};$$

By Lemma 4, the generalization error of CLD is bounded by the right-hand side of the above inequality.

Now, we prove the second part of the theorem. $(W_k)_{k=0}$ and $(W_k^0)_{k=0}$ be the (discrete) GLD processes training \mathcal{D} and \mathcal{S}^0 , respectively. Then for any Z :

$$\begin{aligned} & \mathbb{E} |L(W_k; z) - \mathbb{E} L(W_k^0; z)| \\ & \leq 2C \text{TV}(S_{\cdot; k}; S^0_{\cdot; k}) \quad (\text{C-boundedness}) \\ & \leq 2C (\text{TV}(S_{\cdot; k}; S_{\cdot; k}) + \text{TV}(S_{\cdot; k}; S^0_{\cdot; k}) + \text{TV}(S^0_{\cdot; k}; S^0_{\cdot; k})): \end{aligned}$$

Since $\gamma > 2$ and $\gamma > \frac{1}{2}$, Assumption 35 holds with $A = C, B = L, m = \frac{L^2}{2}, b = \frac{L^2}{2}, \alpha = 0$ and $\beta = \frac{d}{2} \log(1 + \frac{2}{\gamma})$. By applying Pinsker's inequality and Lemma 36, we have

$$\text{TV}(S_{\cdot; k}; S_{\cdot; k}) \leq \frac{1}{2} KL(S_{\cdot; k}; S_{\cdot; k}) \leq \frac{1}{2} KC_1^2 \quad (36)$$

and

$$\text{TV}(S^0_{\cdot; k}; S^0_{\cdot; k}) \leq \frac{1}{2} KL(S^0_{\cdot; k}; S^0_{\cdot; k}) \leq \frac{1}{2} KC_1^2; \quad (37)$$

From (35), we have

$$\text{TV}(S_{\cdot; k}; S^0_{\cdot; k}) \leq \frac{1}{2} KL(S_{\cdot; k}; S^0_{\cdot; k}) \leq \frac{L^2 e^{8C} (1 - e^{-\frac{k}{e^{8C}}})}{n^2} \quad (38)$$

Combining (36), (37) and (38), we have

$$\mathbb{E} |L(W_k; z) - \mathbb{E} L(W_k^0; z)| \leq 2C \left(\frac{1}{2} KC_1^2 + \frac{2CLe^{4C}}{n} \sqrt{\frac{1 - e^{-\frac{k}{e^{8C}}}}{n^2}} \right) := \eta_n;$$

By Definition 3, GLD is η_n -uniformly stable. Applying Lemma 4 gives the generalization bound of GLD. \blacksquare

Lemma 37 (Exponential decay in entropy) (Bakry et al., 2013, Theorem 5.2.1) The logarithmic Sobolev inequality (LSI) for the probability measure is equivalent to saying that for every positive function f in $L^1(\mu)$ (with finite entropy),

$$\text{Ent}(P_t f) \leq e^{-2t} \text{Ent}(f)$$

for every $t \geq 0$.

The following Lemma shows that $\mathbb{P}_t(\frac{d}{d^0}) = \pi_t$ in our diffusion process.

Lemma 38. Let P denote the diffusion semigroup of CLD. Let π denote the invariant measure of P and let π_t denote the probability measure \mathbb{P}_t . Then $\mathbb{P}_t(\frac{d}{d^0}) = \pi_t$.

Proof Let $d = \int (x) dx$ and $d^0 = \int (x) dx$. As shown in (Pavliotis, 2014, page 118), our diffusion process (Smoluchowski dynamics) is reversible, which means $\pi_t(x; y) = \pi_t(y; x)$. Thus for any $g(x)$, we have

$$\begin{aligned} \mathbb{E}_{x \sim \mathbb{P}_t(\frac{d}{d^0})} [g(x)] &= \int g(x) \pi_t(x) (\mathbb{P}_t(d^0=d))(x) dx \\ &= \int g(x) \pi_t(x) dx \int \pi_t(y) \pi_t(x; y) dy \\ &= \int g(x) \pi_t(x) dx \int \pi_t(y) \pi_t(y; x) dy \\ &= \int g(x) \pi_t(x) \pi_t(x) dx = \int g(x) \pi_t(x) dx = \mathbb{E}_{x \sim \pi_t} [g(x)]. \end{aligned}$$

Since g is arbitrary, $\mathbb{P}_t(\frac{d}{d^0})$ and π_t must be the same. ■

Theorem 39. Suppose that $t > 8C$. Under Assumption 14, CLD (with initial distribution $\pi_0 = \frac{1}{Z} e^{-\frac{k_w k^2}{2}} dw$) has the following expected generalization error bound:

$$\text{err}_{\text{gen}} \leq \frac{8C^2}{n} + 4C \exp\left(-\frac{T}{e^{4C}}\right) \frac{P}{C};$$

In addition, if F_0 is also M -smooth and non-negative, by setting $\gamma > 2$, $\beta > \frac{1}{2}$ and $\alpha \in [0; 1 \wedge \frac{2}{8M^2}]$, the GLD process (running k iterations with the same π_0 as CLD) has the expected generalization error bound:

$$\text{err}_{\text{gen}} \leq 2C^{\frac{P}{2KC_1^2}} + \frac{8C^2}{n} + 4C \exp\left(-\frac{K}{e^{4C}}\right) \frac{P}{C};$$

where C_1 is a constant that only depends on γ, β, b, L and d .

Proof of Theorem 39 Suppose S and S^0 are two datasets that differ on exactly one data point. Let $(W_t)_t$ and $(W_t^0)_t$ be their processes, respectively. Let $\pi_t = \pi_t(w) dw$ and $\pi_t^0 = \pi_t^0(w) dw$ be the probability measure \mathbb{P}_t and \mathbb{P}_t^0 , respectively. The invariant measure of CLD π and π^0 are denoted as π and π^0 , respectively. Recall that

$$d = \frac{1}{Z} e^{-F_S(w)} dw; \quad d^0 = \frac{1}{Z^0} e^{-F_{S^0}(w)} dw;$$

The total variation distance of π and π^0 is

$$\begin{aligned} \text{TV}(\pi; \pi^0) &= \frac{1}{2} \int_{\mathbb{R}^d} \left| \frac{d}{d^0} - 1 \right| d \\ &= \frac{1}{2} \int_{\mathbb{R}^d} \left| \frac{Z}{Z^0} \exp(-F_{S^0}(w) + F_S(w)) - 1 \right| \frac{1}{Z} e^{-F_S(w)} dw. \end{aligned} \tag{39}$$

Since $\frac{Z}{Z^0} \exp(-F_{S^0}(w) + F_S(w)) \leq e^{\frac{4C}{n}}; e^{-\frac{4C}{n}}$ and $\frac{4C}{n} < 1=2$, we have

$$\text{TV}(\pi; \pi^0) \leq \max\left\{ \frac{1}{2} \left(1 - e^{-\frac{4C}{n}}\right); \frac{1}{2} \left(e^{\frac{4C}{n}} - 1\right) \right\} \frac{4C}{n}. \tag{40}$$

Since π and π^0 satisfy $LS(e^{4C})$ (Lemma 33), applying Lemma 37 with $\pi = \frac{d}{d^0}$ and $\pi^0 = \frac{d^0}{d^0}$ and Lemma 38 yields:

$$\text{KL}(\pi_t; \pi) \leq \exp\left(-\frac{2t}{e^{4C}}\right) \text{KL}(\pi_0; \pi); \quad \text{KL}(\pi_t^0; \pi^0) \leq \exp\left(-\frac{2t}{e^{4C}}\right) \text{KL}(\pi_0^0; \pi^0); \tag{41}$$

Since $\text{KL}(\theta_0; \theta)$ and $\text{KL}(\theta_0; \theta')$ are upper bounded by C , Pinsker's inequality implies that $\text{TV}(\theta; \theta')$ and $\text{TV}(\theta_0; \theta)$ are upper bounded by $\exp\left(\frac{2t}{e^4 C}\right) C$. Combining with (40) and note that $\text{TV}(\theta; \theta') = \text{TV}(\theta; \theta_0) + \text{TV}(\theta_0; \theta) + \text{TV}(\theta_0; \theta')$, we have

$$\sup_{z \sim \mathcal{A}} \mathbb{E}[|L(W_T; z) - L(W_T^0; z)|] \leq 2C \text{TV}(\theta; \theta') \leq 4C \exp\left(\frac{2t}{e^4 C}\right) C + \frac{8C^2}{n}.$$

By Lemma 4, the generalization error of CLD is bounded by the right-hand side.

The proof for GLD proceeds in the same way as the second part of the proof of Theorem 5.

C EXPERIMENT DETAILS

We first present the general setup of our experiments:

Dataset: We use MNIST (LeCun et al., 1998) and CIFAR10 (Krizhevsky & Hinton, 2009) in our experiments.

Neural network: In our experiments, we test two different neural networks: a smaller version of AlexNet (Krizhevsky et al., 2012) and MLP. The structures of the networks are similar to what are used in Zhang et al. (2017a).

Small AlexNet: k is the kernel size, d is the depth of a convolution layer, fn is the fully-connected layer that has n neurons. The ReLU activation are used in the first 6 layers.

1	2	3	4	5	6	7
conv(k:5,d:64)	pool(k:3)	conv(k:5,d:192)	pool(k:3)	fc(384)	fc(192)	fc(10)

MLP: The MLP used in our experiment has 3 hidden layers, each having width 512. We also use ReLU as the activation function in MLP.

Objective function: For a data point $\mathbf{x} = (x; y)$ in MNIST, the objective function is

$$F(W; z) = -\ln(\text{softmax}(\text{net}_W(\mathbf{x}))[\mathbf{y}]);$$

where $\text{softmax}(\mathbf{a})[i] = \frac{e^{\mathbf{a}[i]}}{\sum_{j=1}^{10} e^{\mathbf{a}[j]}}$, and $\text{net}_W(\mathbf{x})$ is the output of the neural network (10 dimensional vector). Note that the objective function is exactly the cross-entropy loss.

0/1 loss: The 0-1 loss L^{01} is defined as:

$$L^{01}(W; (\mathbf{x}; y)) = \begin{cases} 1 & (\arg \max_i \text{net}_W(\mathbf{x})[i]) \neq y; \\ 0 & \text{otherwise} \end{cases} \quad (42)$$

Random labels: Suppose the dataset contains n datapoint, and the corruption portion is p . We randomly select pn data points, and replace their labels with random labels, as in Zhang et al. (2017a).

C.1 EXPERIMENTAL RESULTS FOR GLD

The result of this experiment (see Figure 1) is discussed in Section 3.1. Here we present our implementation details.

We repeat our experiment 5 times. At every individual run, we first randomly sample 10000 data points from the complete MNIST training data. The initial learning rate is 0.003. It decays 0.995 after every 60 steps, and it stops decaying when it is lower than 0.0005. During the training, we keep $\eta = 0.2 \sqrt{t}$. Recall that the empirical squared gradient $\sigma_g(t)$ in our bound (Theorem 9)

is $E_{W_{t-1}}[\frac{1}{n} \sum_{i=1}^n \text{kr f}(W_{t-1}; z_i)k^2]$. Since it is time-consuming to compute the exact $g_e(t)$, in our experiment, we use an unbiased estimation instead. At every step, we randomly sample a mini-batch B with batch size 200 from the training data, and use $\frac{1}{|B|} \sum_{z_i \in B} \text{kr f}(W_{t-1}; z_i)k^2$ as $g_e(t)$ to compute our bound in Figure 1. The estimation $g_e(t)$ at every step is shown in Figure 1(d). Since $g_e(t)$ is not very stable, in our figure, we plot its moving average over a window of size 100 to make the curve smoother (i.e. $g_{\text{avg}}(t) = \frac{1}{100} \sum_{s=t}^{t+100} g_e(s)$).

C.2 EXPERIMENTAL RESULTS FOR SGLD

In this subsection, we present some experiment results for running SGLD on both MNIST and CIFAR10 datasets, to demonstrate that our bound (see Theorem 11), in particular the sum of the empirical squared gradient norms along the training path, can distinguish normal dataset from dataset that contains random labels. As shown in Figure 3, the curves of our bounds look quite similar to the generalization curves. Due to the sub-optimal constants in our bound, the bound is currently greater than 1, and hence we omit the numbers on the y-axis.

We note that in our experiments presented in Figure 3, the learning rate that we choose is larger than that required by the second condition of Theorem 11. This is because the global Lipschitz constant L is hard to estimate and the model is not able to fit training data under a very large noise. As discussed in Section 3.1, we can relax $(20L) \eta$ to $\eta (2 \max_{i \in [n]} \text{kr F}(W_{t-1}; z_i)k)$. By applying gradient clipping trick, we can further relax this condition to

$$\eta \leq \min\{C_L; (2 \max_{i \in [n]} \text{kr F}(W_{t-1}; z_i)k)g^{-1}(\eta)\} \quad (43)$$

where C_L is defined in Section 3.1. In order to show that our observation (“random normal”) still holds when the step size satisfies the requirement of our theory, we run an experiment that using gradient clipping trick with $C_L = 1$. The model is trained on a small subset of MNIST as fitting the original data set with random labels under such a large Gaussian noise is extremely slow. As shown in Figure 4, the experimental results remain unchanged when all the conditions of our bound are met.

These experiments indicate that the sum of squared empirical gradient norms is highly related to the generalization performance, and we believe by further optimizing the constants in our bound, it is possible to achieve a generalization bound that is much closer to the real generalization error.

Figure 2: Training MLP with GLD ($\eta = 0.2 \eta$) on the full MNIST dataset without label corruption. Learning rate $\eta = 0.01 \cdot 0.95^{60c}$. Note that in the early stage, the testing accuracy is even higher than the training accuracy, thus we plot the absolute value of generalization error. As shown in this figure, even when the training accuracy approaches 90%, our bound is still relatively small.

