

STABILITY BASED GENERALIZATION BOUNDS FOR EXPONENTIAL FAMILY LANGEVIN DYNAMICS

Anonymous authors

Paper under double-blind review

ABSTRACT

We study the generalization of noisy stochastic mini-batch based iterative algorithms based on the notion of stability. Recent years have seen key advances in data-dependent generalization bounds for noisy iterative learning algorithms such as stochastic gradient Langevin dynamics (SGLD) based on (Mou et al., 2018; Li et al., 2020) and related approaches (Negrea et al., 2019; Haghifam et al., 2020). In this paper, we unify and substantially generalize stability based generalization bounds and make three technical advances. First, we bound the generalization error of general noisy stochastic iterative algorithms (not necessarily gradient descent) in terms of expected stability, which in turn can be bounded by the expected Le Cam Style Divergence (LSD). Such bounds have a $O(1/n)$ sample dependence unlike many existing bounds with $O(1/\sqrt{n})$ dependence. Second, we introduce Exponential Family Langevin Dynamics (EFLD) which is a substantial generalization of SGLD and which allows exponential family noise to be used with gradient descent. We establish data-dependent expected stability based generalization bounds for general EFLD. Third, we consider an important new special case of EFLD: noisy sign-SGD, which extends sign-SGD by using Bernoulli noise over $\{-1, +1\}$, and we establish optimization guarantees for the algorithm. Further, we present empirical results on benchmark datasets to illustrate the our bounds are non-vacuous and quantitatively much sharper than existing bounds.

1 INTRODUCTION

Recent years have seen renewed interest and advances in characterizing generalization performance of learning algorithms in terms of stability, which considers change in performance of a learning algorithm based on change of a single training point (Hardt et al., 2016; Bousquet & Elisseeff, 2002; Li et al., 2020; Mou et al., 2018). For stochastic gradient descent (SGD), Hardt et al. (2016) established generalization bounds based on uniform stability (Hardt et al., 2016; Bousquet & Elisseeff, 2002), although the analysis needed rather small step sizes $\eta_t = 1/t$ which is not useful in practice. While improving the analysis for SGD has remained a challenge, advances have been made on noisy SGD algorithms, especially stochastic gradient Langevin dynamics (SGLD) (Welling & Teh, 2011; Mou et al., 2018; Li et al., 2020), which adds Gaussian noise to the stochastic gradients of marginal variance σ_t^2 . In parallel, there has been key developments on related information-theoretic generalization bounds applicable to SGLD type algorithms (Negrea et al., 2019; Haghifam et al., 2020; Xu & Raginsky, 2017; Russo & Zou, 2016; Pensia et al., 2018).

While these developments have led to major advances in analyzing generalization of noisy SGD algorithms, they each have certain limitations which leave room for further improvements. Using uniform stability, Mou et al. (2018) established a bound for SGLD of the form $\frac{K}{n} \sqrt{\sum_t \eta_t^2 / \sigma_t^2}$ which depends on K , the global Lipschitz constant for the loss, and with step size $\eta_t \leq \sigma_t \ln 2/K$. The bound has a desirable dependency of $O(1/n)$ on the samples, but has an undesirable dependence on K , and the step sizes, bounded by σ_t/K , are too small. Mou et al. (2018) also presented another bound which addresses some of these issues, but gets an undesirable $O(1/\sqrt{n})$ sample dependence. By building on the developments of Russo & Zou (2016); Xu & Raginsky (2017); Pensia et al. (2018), Negrea et al. (2019) made advances from the information theoretic perspective and established bounds for SGLD which have the desirable dependence on the norm of gradient incoherence, i.e., difference in gradients over different mini-batches, avoids dependence on Lipschitz constant K , and is applicable to unbounded sub-Gaussian losses, but have an undesirable $O(1/\sqrt{n})$ sample

dependence. Haghifam et al. (2020) made further advances on the problem from an information theoretic perspective based on the conditional mutual information framework of Steinke & Zakyntinou (2020) and obtained generalization bounds based on gradient incoherence with $O(1/n)$ sample dependence, but their analysis holds for full batch Langevin dynamics, not mini-batch SGLD. Li et al. (2020) made advances on such bounds based on the notion of Bayes-stability, by combining ideas from PAC-Bayes bounds into stability, and established a bound of the form $\frac{c}{n} \sqrt{\sum_t \eta_t^2 \mathbf{g}_e(t) / \sigma_t^2}$ for bounded losses, where $\mathbf{g}_e(t)$ is the expected gradient norm square at step t . While the bound avoids dependency on the Lipschitz constant K , the dependence on the gradient norm makes such bounds much weaker than the information theoretic bounds of Negrea et al. (2019); Haghifam et al. (2020) which depend on the norm of gradient incoherence, which are typically orders of magnitude smaller. Further, the analysis of Li et al. (2020) still needs small step sizes, bounded by σ_t/K .

In this paper, we build on the core strengths of such existing approaches, most notably the $O(1/n)$ sample dependence of stability based bounds (Mou et al., 2018; Li et al., 2020) and the dependence on gradient incoherence for information theoretic bounds (Negrea et al., 2019; Haghifam et al., 2020), and develop a framework (Section 2) for developing generalization bounds for noisy stochastic iterative (NSI) algorithms. Our framework considers generalization based on the concept of *expected stability*, rather than uniform stability (Hardt et al., 2016; Bousquet & Elisseeff, 2002; Bousquet et al., 2020; Mou et al., 2018), which yields distribution dependent generalization bounds and avoids the worst-case setting of uniform stability. Building on Li et al. (2020), we show that expected stability of general NSI algorithms can be bounded by the expected Le Cam Style Divergence with dependence on parameter distributions from mini-batches differing by one sample. In Section 3, we introduce Exponential Family Langevin Dynamics (EFLD), a family of noisy gradient descent algorithms based on exponential family noise. Special cases of EFLD include SGLD and noisy versions of Sign-SGD or quantized SGD algorithms widely used in practice (Bernstein et al., 2018a;b; Jin et al., 2020; Alistarh et al., 2017). Our main result provides an expected stability based generalization bound applicable to any EFLD algorithm with a $O(1/n)$ sample dependence and a dependence on gradient incoherence, rather than gradient norms. Existing generalization bounds for SGLD (Li et al., 2020; Negrea et al., 2019) usually use properties of the Gaussian distribution, and our analysis on EFLD illustrates that this was unnecessary. We also consider optimization guarantees for EFLD and establish such results for noisy Sign-SGD and SGLD. Through experiments on benchmark datasets (Section 4), we illustrate that our bounds are non-vacuous and quantitatively much sharper than existing bounds (Li et al., 2020; Negrea et al., 2019).

Related work. Uniform stability has been a classical approach for bounding generalization error (Bousquet & Elisseeff, 2002; Bousquet et al., 2020; Feldman & Vondrak, 2018; 2019), pioneered by Rogers & Wagner (1978); Devroye & Wagner (1979). Beyond the aforementioned work, there has been recent work on differential privacy that analyzes the uniform stability of differentially private SGD (DP-SGD) (Hardt et al., 2016; Bassily et al., 2020). Beyond uniform stability, information-theoretic approaches (Russo & Zou, 2016; Xu & Raginsky, 2017) that bounds the generalization error by the mutual information between the algorithm input S and the algorithm output \mathbf{w} , have been used for deriving generalization bounds for noisy iterative algorithms (Pensia et al., 2018; Bu et al., 2019). Along this line of literature, Negrea et al. (2019); Haghifam et al. (2020); Rodríguez-Gálvez et al. (2021) prove data-dependent generalization bounds dropping dependence on the Lipschitz constant. Further, tighter bounds (Haghifam et al., 2020; Zhou et al., 2021; Rodríguez-Gálvez et al., 2021; Neu, 2021; Hellström & Durisi, 2021) are proposed based on conditional mutual information (Steinke & Zakyntinou, 2020; Grünwald et al., 2021; Hellström & Durisi, 2020). Due to space limitations, an extended discussion of the related work is deferred to Appendix A.

2 GENERALIZATION BOUNDS WITH EXPECTED STABILITY

In the setting of statistical learning, there is an instance space \mathcal{Z} , a hypothesis space \mathcal{W} , and a loss function $\ell : \mathcal{W} \times \mathcal{Z} \mapsto \mathbb{R}_+$. Let D be an unknown distribution of \mathcal{Z} and let $S \sim D^n$ be n i.i.d. draws from D . For any specific hypothesis $\mathbf{w} \in \mathcal{W}$, the population and empirical loss are respectively given by $L_D(\mathbf{w}) \triangleq \mathbb{E}_{z \sim D}[\ell(\mathbf{w}, z)]$, and $L_S(\mathbf{w}) \triangleq \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}, z_i)$. For any distribution P over the hypothesis space, we respectively denote the expected population and empirical loss as

$$L_D(P) \triangleq \mathbb{E}_{z \sim D} \mathbb{E}_{\mathbf{w} \sim P}[\ell(\mathbf{w}, z)], \quad \text{and} \quad L_S(P) \triangleq \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathbf{w} \sim P}[\ell(\mathbf{w}, z_i)]. \quad (1)$$

Consider a randomized algorithm A which works with $S = \{z_1, \dots, z_n\} \sim D^n$ and cre-

ates a distribution over the hypothesis space \mathcal{W} . For convenience, we will denote the distribution as $A(S)$. The focus of the analysis is to bound the generalization error of A defined as: $\text{gen}(A(S)) \triangleq L_D(A(S)) - L_S(A(S))$. We will assume A is permutation invariant, i.e., the ordering of samples in S do not modify $A(S)$, an assumption satisfied by most learning algorithms. We will focus on developing bounds for the expectation $\mathbb{E}_S[L_D(A(S)) - L_S(A(S))]$, and discuss high-probability bounds in the Appendix B.

2.1 BOUNDS BASED ON EXPECTED STABILITY

We start our analysis by noting that the expected generalization error can be upper bounded by *expected stability* based on the Hellinger divergence (Sason & Verdu, 2016; Li et al., 2020): $H^2(P\|P') = \frac{1}{2} \int_{\mathbf{w}} (\sqrt{p(\mathbf{w})} - \sqrt{p'(\mathbf{w})})^2 d\mathbf{w}$.

Proposition 1. *Let $S_n \sim D^n$ and let S'_n be a dataset obtained by replacing $z_n \in S_n$ with $z'_n \sim D$. Let $A(S_n), A(S'_n)$ respectively denote the distributions over the hypothesis space \mathcal{W} obtained by running randomized algorithm A on S_n, S'_n . Assume that for $S_n \sim D^n, \forall z \in \mathcal{Z}$, $\mathbb{E}_{W \sim A(S_n)}[\ell^2(W, z)] \leq c_0/2, c_0 > 0$. With $H(\cdot, \cdot)$ denoting the Hellinger divergence, we have*

$$|\mathbb{E}_{S_n \sim D^n}[L_D(A(S_n)) - L_S(A(S_n))]| \leq c_0 \mathbb{E}_{S_n \sim D^n} \mathbb{E}_{z'_n \sim D} \sqrt{2H^2(A(S_n), A(S'_n))}. \quad (2)$$

Remark 2.1. Proposition 1 does not need bounded losses. Just the *second moment* of $\ell(W, z), W \sim A(S_n), S_n \sim D^n, \forall z \in \mathcal{Z}$ need to be bounded. The assumption is satisfied by bounded losses. It is instructive to compare the assumption to that in recent information theoretic bounds (Haghifam et al., 2020; Xu & Raginsky, 2017), where one assumes $\ell(\mathbf{w}, Z), Z \sim D, \forall \mathbf{w} \in \mathcal{W}$ to be *sub-Gaussian*.

Remark 2.2. The bound in Proposition 1 is in terms of *expected stability* where we consider $\mathbb{E}_{S \sim D^n} \mathbb{E}_{z'_n \sim D}[\cdot \cdot \cdot]$, an important departure from bounds based on *uniform stability* (Elisseeff et al., 2005; Bousquet & Elisseeff, 2002; Mou et al., 2018; Bousquet et al., 2020; Feldman & Vondrak, 2018; 2019) where one considers $\sup_{S, S' \in \mathcal{Z}^n, |S \setminus S'|=1}[\cdot \cdot \cdot]$. Replacing sup by \mathbb{E} makes the bounds distribution dependent, and arguably leads to quantitatively tighter bounds.

Note that the Hellinger divergence can be bounded by the KL divergence.

Proposition 2. *For any distributions P and P' , $2H^2(P, P') \leq \min \{KL(P, P'), \sqrt{\frac{1}{2}KL(P, P')}\}$.*

2.2 EXPECTED STABILITY OF NOISY STOCHASTIC ITERATIVE ALGORITHMS

We consider a general family of noisy stochastic iterative (NSI) algorithms. Given $S \sim D^n$, such iterative algorithms have two additional sources of randomness in each iteration t : (a) a stochastic mini-batch of samples S_{B_t} , with $|S_{B_t}| = b$, drawn uniformly at random with replacement from S ; and (b) noise ξ_t suitably included in the iterative update. Given a trajectory of past iterates $W_{0:(t-1)} = \mathbf{w}_{0:(t-1)}$, the new iterate W_t is drawn from a distribution $P_{B_t, \xi_t | \mathbf{w}_{0:(t-1)}}$ over \mathcal{W} :

$$W_t \sim P_{B_t, \xi_t | \mathbf{w}_{0:(t-1)}}(W). \quad (3)$$

We will drop the conditioning $\mathbf{w}_{0:(t-1)}$ to avoid clutter in the sequel. Let P_T, P'_T denote the marginal distributions over hypotheses $\mathbf{w} \in \mathcal{W}$ after T steps of the algorithm based on S_n, S'_n respectively. Further, let $P_{0:(t-1)}$ denote the joint distribution over $W_{0:(t-1)} = (W_0, \dots, W_{t-1})$, and let $P_t \equiv P_{B_t, \xi_t | \mathbf{w}_{0:(t-1)}}$ compactly denote the conditional distribution on W_t conditioned on the trajectory $W_{0:(t-1)} = \mathbf{w}_{0:t-1}$. Following (Negrea et al., 2019; Haghifam et al., 2020; Pensia et al., 2018), we use the following chain rule for KL-divergence: $KL(P_T \| P'_T) \leq KL(P_{0:T} \| P'_{0:T}) = \sum_{t=1}^T \mathbb{E}_{P_{0:(t-1)}}[KL(P_t \| P'_t)]$. Let $\bar{S} \sim D^{n+1}$, and let S_n, S'_n be size n subsets of \bar{S} such that $S_n = \{Z_1, \dots, Z_{n-1}, Z_n\}$ and $S'_n = \{Z_1, \dots, Z_{n-1}, Z'_n\}$, where $Z'_n = Z_{n+1}$. Let $S_0 = \{Z_1, \dots, Z_{n-1}\}$. The algorithms we consider use a mini-batch of size b in each iteration uniformly sampled from n samples. Let the set of all mini-batch index sets be denoted by G . Let the set of all mini-batch index sets A drawn from S_0 be denoted by G_0 . Note that $|G_0| = \binom{n-1}{b}$. Let G_1 denote the set of all mini-batch index sets B which includes the last sample, viz. z_n for S with mini-batches and z'_n for S'_n . Note that $|G_1| = \binom{n-1}{b-1}$. Also note that $|G_0| + |G_1| = \binom{n-1}{b} + \binom{n-1}{b-1} = \binom{n}{b} = |G|$.

Following Li et al. (2020), we can bound their conditional KL-divergences $KL(P_t \| P'_t)$ in terms of a Le Cam Style Divergence (LSD). While the classical Le Cam divergence (Sason & Verdu,

2016) is $LCD(P||P') \triangleq \frac{1}{2} \int \frac{(dP-dP')^2}{dP+dP'}$ (where dP denotes the density), our bounds in terms of $LSD(P_t||P'_t) \triangleq \int \frac{(dP_{B_t, \xi_t} - dP'_{B_t, \xi_t})^2}{dP_{A_t, \xi_t}}$, $B_t \in G_1, A_t \in G_0$. Note that P_{B_t, ξ_t} and P'_{B_t, ξ_t} represent the distribution of W_t for S_n and S'_n respectively since the mini-batch S_{B_t} of S_n and S'_n differs in the n -th sample. Putting everything together, we have the following LSD based bound.

Lemma 1. Consider a noisy stochastic iterative algorithms of the form (3) with mini-batch size $b \leq n/2$. Then, with $c_1 = \sqrt{2}c_0$ (with c_0 as in Proposition 1), we have

$$|\mathbb{E}_{S_n}[L_D(A(S_n)) - L_{S_n}(A(S_n))]| \leq c_1 \frac{b}{n} \mathbb{E}_{S_n} \mathbb{E}_{z'_n} \sqrt{\sum_{t=1}^T \mathbb{E}_{W_{0:(t-1)}} \mathbb{E}_{B_t \in G_1} \mathbb{E}_{A_t \in G_0} \left[\int_{\xi_t} \frac{(dP_{B_t, \xi_t} - dP'_{B_t, \xi_t})^2}{dP_{A_t, \xi_t}} d\xi_t \right]}. \quad (4)$$

Remark 2.3. Li et al. (2020) essentially has this result for SGLD and inspired our work. Our proofs are significantly simpler and, more importantly, illustrates applicability to general noisy iterative algorithms of the form (3), not just SGLD with Gaussian noise as in Li et al. (2020).

Remark 2.4. Note that the bound does not assume the loss to be bounded, depends on expectations over samples S_n, z'_n , trajectories $w_{0:(t-1)}$, and mini-batches B_t, A_t . Further, the bound depends on the distribution discrepancy as captured by the expected LSD.

Remark 2.5. The bound seems to worsen with b , the size of the mini-batch. As we shown in Section 3, the expected LSD term has a $\frac{1}{b^2}$ dependence for the Exponential Family Langevin dynamics (EFLD) models we introduce, so the leading b is neutralized.

3 EXPONENTIAL FAMILY LANGEVIN DYNAMICS

In recent years, considerable advances have been made in establishing generalization bounds for stochastic gradient Langevin dynamics (SGLD) (Li et al., 2020; Pensia et al., 2018; Negrea et al., 2019; Haghifam et al., 2020). As an example of NSI algorithms of the form (3), SGLD adds an isotropic Gaussian noise at every step of SGD: $w_{t+1} = w_t - \eta_t \nabla \ell(w_t, S_{B_t}) + \mathcal{N}(0, \sigma_t^2 \mathbb{I}_d)$, where $\nabla \ell(w_t, S_{B_t})$ is the stochastic gradient on mini-batch B_t , η_t is the step size, and σ_t^2 is noise variance.

In this paper, we introduce a substantial generalization of SGLD called Exponential Family Langevin Dynamics (EFLD) which uses general exponential family noise in noisy iterative updates of the form (3). In addition to being a mathematical generalization of the popular SGLD, the proposed EFLD provides flexibility to use noise gradient algorithms with different representation of the gradient, e.g., Bernoulli noise for Sign-SGD, discrete distribution for quantized or finite precision SGD, etc. (Canonne et al., 2020; Alistarh et al., 2017; Jiang & Agrawal, 2018; Yang et al., 2019).

3.1 EXPONENTIAL FAMILY LANGEVIN DYNAMICS (EFLD)

Exponential families (Barndorff-Nielsen, 2014; Brown, 1986; Wainwright & Jordan, 2008) constitute a large family of parametric distributions which include Gaussian, Bernoulli, gamma, Poisson, Dirichlet, etc., as special cases. Exponential families are typically represented in terms of natural parameters $\theta_\alpha = \theta/\alpha$ with scaling $\alpha > 0$, i.e., $p_\psi(\xi, \theta_\alpha) = \exp(\langle \xi, \theta_\alpha \rangle - \psi(\theta_\alpha)) \pi_0(\xi) = \prod_{j=1}^p \exp(\xi_j \theta_{j\alpha} - \psi_j(\theta_{j\alpha})) \pi_0(\xi_j)$, where ξ is the sufficient statistic, $\psi(\theta_\alpha) = \sum_{j=1}^p \psi_j(\theta_{j\alpha})$ is the log-partition function, and $\pi_0(\xi) = \prod_{j=1}^p \pi_0(\xi_j)$ is the base measure. Note that $\alpha = 1$ gives the canonical form of the exponential family distributions. For general scaling $\alpha > 0$, for some cases the base measure π_0 may depend on the scaling, i.e., $\pi_{0,\alpha}$. A scaling $\alpha > 0$ is valid as long as $\exp(\langle \xi, \theta_\alpha \rangle)$ is integrable, i.e., $\int_{\xi} \exp(\langle \xi, \theta_\alpha \rangle) \pi_0(\xi) d\xi < \infty$. Further, ψ is a smooth function by construction (Barndorff-Nielsen, 2014; Banerjee et al., 2005; Wainwright & Jordan, 2008) and the smoothness of ψ implies $\nabla^2 \psi(\theta_\alpha) \leq c_2 \mathbb{I}$.

Exponential family Langevin dynamics (EFLD) uses noisy stochastic gradient updates similar to SGLD, but using exponential family noise rather than Gaussian noise as in SGLD. In particular, for mini-batch S_{B_t} , EFLD updates are as follows: with step size $\rho_t > 0$

$$w_t = w_{t-1} - \rho_t \xi_t, \quad \xi_t \sim p_\psi(\xi; \theta_{B_t, \alpha_t}), \quad (5)$$

where

$$p_\psi(\xi; \theta_{B_t, \alpha_t}) = \exp(\langle \xi, \theta_{B_t, \alpha_t} \rangle - \psi(\theta_{B_t, \alpha_t})) \pi_0(\xi), \quad \theta_{B_t, \alpha_t} \triangleq \frac{\theta_{B_t}}{\alpha_t} = \frac{\nabla \ell(w_{t-1}, S_{B_t})}{\alpha_t}. \quad (6)$$

For EFLD, the natural parameter θ_{B_t, α_t} at step t is simply a scaled version of the mini-batch gradient $\nabla \ell(\mathbf{w}_{t-1}, S_{B_t})$. We first show that EFLD becomes SGLD when the exponential family is Gaussian, and becomes a noisy version of sign-SGD (Bernstein et al., 2018a;b) when the exponential family is Bernoulli over $\{-1, +1\}$. More details and examples are in Appendix C.1.

Example 3.1 (SGLD). SGLD uses scaled Gaussian noise with $\psi(\theta) = \|\theta\|_2^2/2$, $\alpha_t = \sqrt{\sigma_t/\eta_t}$, so that $p_\psi(\xi; \theta_{B_t, \alpha_t}) = \mathcal{N}(\theta_{B_t}, \alpha_t^2 \mathbb{I}_d)$. By taking $\rho_t = \sqrt{\eta_t \sigma_t}$, the update (5) based on $\rho_t \xi_t$ is distributed as $\mathcal{N}(\rho_t \theta_{B_t}, \rho_t^2 \alpha_t^2 \mathbb{I}_d) = \mathcal{N}(\eta_t \nabla \ell(\mathbf{w}_{t-1}, S_{B_t}), \sigma_t^2 \mathbb{I}_d)$. Thus the EFLD update reduces to the SGLD update: $\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \nabla \ell(\mathbf{w}_{t-1}, S_{B_t}) + \mathcal{N}(0, \sigma_t^2 \mathbb{I}_d)$. \square

Example 3.2 (Noisy Sign-SGD). By taking $\rho_t = \eta_t$ and component-wise $\xi_j \in \{-1, 1\}$, $\pi_0(\xi_j) = 1$, $\psi(\theta) = \log(\exp(-\theta) + \exp(\theta))$ in exponential family update equation (5), the j -th component of exponential family distribution $p_\psi(\xi; \theta_{B_t, \alpha_t})$ becomes $p_{\theta_{B_t, \alpha_t, j}}(\xi_j) = \frac{\exp(\xi_j \theta_{B_t, \alpha_t, j})}{\exp(-\theta_{B_t, \alpha_t, j}) + \exp(\theta_{B_t, \alpha_t, j})}$. Thus, the EFLD update reduces to a noisy version of Sign-SGD: $\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \xi_t$, $\xi_{t,j} \sim p_{\theta_{B_t, \alpha_t, j}}(\xi_j)$, $j \in [d]$, where $\theta_{B_t, \alpha_t} = \nabla \ell(\mathbf{w}_{t-1}, S_{B_t})/\alpha_t$ is the scaled mini-batch gradient. \square

3.2 EXPECTED STABILITY OF EXPONENTIAL FAMILY LANGEVIN DYNAMICS

From Lemma 1, conditioned on a trajectory $\mathbf{w}_{0:(t-1)}$, mini-batches S_{B_t}, S_{A_t} , we can get generalization bound by suitably bounding the Le Cam Style Divergence (LSD) given by: $I_{A_t, B_t} = \int_{\xi_t} \frac{(dP_{B_t, \xi_t} - dP'_{B_t, \xi_t})^2}{dP_{A_t, \xi_t}} d\xi_t$. For EFLD, the density functions dP_{B_t, ξ_t} are exponential family densities $p_\psi(\xi; \theta_{B_t, \alpha_t})$ as in (5)-(6), and we have the following bound on the per step LSD:

Theorem 1. For a given set $\bar{S} \sim D^{n+1}$ and \mathbf{w}_{t-1} at iteration $(t-1)$, let $\Delta_{t|\mathbf{w}_{t-1}}(\bar{S}) = \max_{z, z' \in \bar{S}} \|\nabla \ell(\mathbf{w}_{t-1}, z) - \nabla \ell(\mathbf{w}_{t-1}, z')\|_2$. Further, for a c_2 -smooth log-partition function ψ , let the scaling $\alpha_{t|\mathbf{w}_{t-1}}$ be data-dependent such that $\alpha_{t|\mathbf{w}_{t-1}}^2 \geq 8c_2 \Delta_{t|\mathbf{w}_{t-1}}^2(S_{n+1})$. Then, we have

$$I_{A_t, B_t} \leq 5c_2 \|\theta_{B_t, \alpha_t} - \theta'_{B_t, \alpha_t}\|_2^2 = \frac{5c_2}{2\alpha_{t|\mathbf{w}_{t-1}}^2} \left[\|\nabla \ell(\mathbf{w}_{t-1}, S_{B_t}) - \nabla \ell(\mathbf{w}_{t-1}, S'_{B_t})\|_2^2 \right], \quad (7)$$

Note that S_{B_t} and S'_{B_t} only differ at samples z_n and z'_n . The above bound can now be directly applied to Lemma 1 to get expected stability based generalization bounds for any EFLD algorithm.

Theorem 2. Consider an exponential family Langevin dynamics (EFLD) algorithm of the form (5)-(6) with a c_2 -smooth log-partition function ψ . Then, for mini-batch size $b \leq n/2$, with $c = c_0 \sqrt{5c_2}$ (with c_0 as in Lemma 1) and $\alpha_{t|}^2 \geq 8c_2 \Delta_{t|}^2(S_{n+1})$ (as in Theorem 1, with the conditioning on \mathbf{w}_{t-1} hidden to avoid clutter), we have

$$\mathbb{E}_S [L_D(A(S)) - L_S(A(S))] \leq c \frac{1}{n} \mathbb{E}_{S_{n+1}} \sqrt{\sum_{t=1}^T \mathbb{E}_{W_{0:(t-1)}} \frac{1}{\alpha_{t|}^2} \left[\|\nabla \ell(\mathbf{w}_{t-1}, z_n) - \nabla \ell(\mathbf{w}_{t-1}, z'_n)\|_2^2 \right]}. \quad (8)$$

Remark 3.1. Theorem 2 captures the generalization error of SGLD, which is a special case of EFLD. Our bound has the same dependence on n , T , step size η_t as the bound in Li et al. (2020). However, our bound is numerically sharper because we replace the *gradient norms*, i.e., $\frac{1}{n} \sum_{z \in S} \|\ell(\mathbf{w}_t, z)\|$ in Li et al. (2020) and with gradient discrepancy $\|\nabla \ell(\mathbf{w}_t, z) - \nabla \ell(\mathbf{w}_t, z')\|$, which is quantitatively smaller than gradient norms as we show in the experiment section. The bound in Negrea et al. (2019) depends on *gradient incoherence* which is empirically smaller than gradient discrepancy as observed in the experiment section, their bound depends on $1/\sqrt{n}$, which is worse than the $1/n$ dependence in our bound.

Remark 3.2. EFLD can be extended to work with anisotropic noise by using $\theta_{B_t, \alpha_t} = \nabla \ell(\mathbf{w}_{t-1}, S_{B_t}) \oslash \alpha_t$ in (6) where $\alpha_t \in \mathbb{R}^p$ determines different scaling for each dimension and \oslash denotes Hadamard division. Theorems 1 and 2 can be extended to such anisotropic noise by using α -scaled norms for the gradient discrepancy, i.e., $\|\mathbf{g} - \mathbf{g}'\|_{2, \alpha}^2 = \sum_j (g_j - g'_j)^2 / \alpha_j^2$. \square

Remark 3.3. The condition on α_t is a data-dependent quantity, which can be computed along the training process. It gives much more benign condition of the step size compared to those in the related work (Mou et al., 2018; Li et al., 2020, Hardt et al. 2016), which require step size being bounded by σ_t/L . However, the step sizes in Theorem 2 need to be bounded by $\sigma_t/\Delta_t(\bar{S})$, which is considerably more relaxed since $\Delta_t(\bar{S})$ is much smaller than Lipschitz constant L , which is a uniform bound over the whole parameter space. Also, usually one would expect $\Delta_t(\bar{S})$ to decrease as training proceeds since the gradients shrink as the loss function being minimized. Thus, the constraint on step size does not require the step sizes to be as small as σ_t/L .

3.3 PROOF SKETCHES OF MAIN RESULTS: THEOREMS 1 AND 2

We focus on Theorem 1. To avoid clutter, we drop the subscript t for the analysis and note that the analysis holds for any step t . When the density $dP_{B,\xi} = p_\psi(\xi; \theta_{B,\alpha})$, by mean-value theorem, for each ξ , we have $p_\psi(\xi; \theta_{B,\alpha}) - p_\psi(\xi; \theta_{B',\alpha}) = \langle \theta_{B,\alpha} - \theta_{B',\alpha}, \nabla_{\tilde{\theta}_{B,\alpha}} p_\psi(\xi; \theta_{B,\alpha}) \rangle$, for some $\tilde{\theta}_{B,\alpha} = \gamma\xi\theta_{B,\alpha} + (1 - \gamma\xi)\theta'_{B,\alpha}$ where $\gamma\xi \in [0, 1]$. Then,

$$I_{A,B} = \int_{\xi} \frac{(p_\psi(\xi; \theta_{B,\alpha}) - p_\psi(\xi; \theta_{B',\alpha}))^2}{p_\psi(\xi; \theta_{A,\alpha})} d\xi = \int_{\xi} \frac{\langle \theta_{B,\alpha} - \theta'_{B,\alpha}, \xi - \nabla_{\tilde{\theta}_{B,\alpha}} \psi(\xi; \tilde{\theta}_{B,\alpha}) \rangle^2 p_\psi^2(\xi; \tilde{\theta}_{B,\alpha})}{p_\psi(\xi; \theta_{A,\alpha})} d\xi,$$

since $p_\psi(\xi; \tilde{\theta}_{B,\alpha}) = \exp(\langle \xi, \tilde{\theta}_{B,\alpha} \rangle - \psi(\tilde{\theta}_{B,\alpha}))\pi_0(\xi)$.

Handling Distributional Dependence of $\tilde{\theta}_B$. Note that we cannot proceed with the analysis with the density term depending on $\tilde{\theta}_{B,\alpha}$ since $\tilde{\theta}_{B,\alpha}$ depends on ξ . So, we first bound the density term depending on $\tilde{\theta}_{B,\alpha}$ in terms of exponential family densities with parameters $\theta_{B,\alpha}$ and $\theta_{B',\alpha}$ using c_2 -smoothness of ψ .

Lemma 2. For some $\gamma\xi \in [0, 1]$, $\tilde{\theta}_{B,\alpha} = \gamma\xi\theta_{B,\alpha} + (1 - \gamma\xi)\theta'_{B,\alpha}$, we have

$$\frac{\exp[\langle \xi, \tilde{\theta}_{B,\alpha} \rangle - \psi(\tilde{\theta}_{B,\alpha})]}{\max(\exp[\langle \xi, \theta_{B,\alpha} \rangle - \psi(\theta_{B,\alpha})], \exp[\langle \xi, \theta_{B',\alpha} \rangle - \psi(\theta_{B',\alpha})])} \leq \exp[c_2\|\theta_{B,\alpha} - \theta_{B',\alpha}\|_2^2].$$

Bounding the Density Ratio. Next we focus on the density ratio $p_\psi^2(\xi; \tilde{\theta}_{B,\alpha})/p_\psi(\xi; \theta_{A,\alpha})$. By Lemma 2, it suffices to focus on $p_\psi^2(\xi; \theta_{B,\alpha})/p_\psi(\xi; \theta_{A,\alpha})$ or the equivalent term for $\theta_{B',\alpha}$. We show that the density ratio can be bounded by another exponential family with parameters $(2\theta_{B,\alpha} - \theta_{A,\alpha})$.

Lemma 3. For any ξ , we have

$$\frac{\exp[\langle \xi, 2\theta_{B,\alpha} \rangle - 2\psi(\theta_{B,\alpha})]}{\exp[\langle \xi, \theta_{A,\alpha} \rangle - \psi(\theta_{A,\alpha})]} \leq \exp[2c_2\|\theta_{B,\alpha} - \theta_{A,\alpha}\|_2^2] \exp[\langle \xi, (2\theta_{B,\alpha} - \theta_{A,\alpha}) - \psi(2\theta_{B,\alpha} - \theta_{A,\alpha}) \rangle].$$

The analysis for the term $p_\psi^2(\xi; \theta_{B',\alpha})/p_\psi(\xi; \theta_{A,\alpha})$ is exactly the same.

Bounding the Integral. Ignoring multiplicative terms which do not depend on ξ for the moment, the analysis needs to bound an integral term of the form $\int_{\xi} \langle \theta_{B,\alpha} - \theta'_{B,\alpha}, \xi - \nabla\psi(\xi; \tilde{\theta}_{B,\alpha}) \rangle^2 p_\psi(\xi; 2\theta_{B,\alpha} - \theta_{A,\alpha}) d\xi$, and a similar term with $p_\psi^2(\xi; 2\theta_{B',\alpha} - \theta_{A,\alpha})$. First, note that $\nabla\psi(\xi; \tilde{\theta}_{B,\alpha}) = \tilde{\mu}_{B,\alpha}$, the expectation parameter for $p_\psi(\xi; \tilde{\theta}_{B,\alpha})$ (Wainwright & Jordan (2008); Banerjee et al. (2005)). The integral, however, is with respect to $p_\psi(\xi; 2\theta_{B,\alpha} - \theta_{A,\alpha})$. We handle this discrepancy by using $\xi - \nabla\psi(\xi; \tilde{\theta}_{B,\alpha}) = (\xi - \mathbb{E}[\xi]) + (\mathbb{E}[\xi] - \nabla\psi(\xi; \tilde{\theta}_{B,\alpha}))$, and decomposing as sum-of-squares. Quadratic form for the first term yields the covariance $\mathbb{E}[(\xi - \mathbb{E}[\xi])(\xi - \mathbb{E}[\xi])^T] = \nabla^2\psi(\theta_{2\theta_{B,\alpha} - \theta_{A,\alpha}}) \leq c_2\mathbb{I}$, by smoothness. The second term depends on the difference of gradients $\nabla\psi(2\theta_{B,\alpha} - \theta_{A,\alpha}) - \nabla\psi(\tilde{\theta}_{B,\alpha})$ which, using smoothness and additional analysis, can be bounded by the norm of $(\theta_{B,\alpha} - \theta_{A,\alpha})$. All the pieces can be put together to get the bound in Theorem 1, which when used in Lemma 1 yields Theorem 2.

3.4 OPTIMIZATION GUARANTEES FOR EFLD

We now establish optimization guarantees for two examples of EFLD, i.e., Noisy Sign-SGD with Bernoulli noise over $\{-1, +1\}$ and SGLD with Gaussian noise.

Noisy Sign-SGD. For mini-batch B_t and scaling α_t , mini-batch Noisy Sign-SGD updates the parameters as $\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \xi_t$, where each component $j \in [d]$

$$\xi_{t,j} \sim p_{\theta_{B_t, \alpha_t, j}}(x) = \frac{\exp(x\theta_{B_t, \alpha_t, j})}{\exp(-\theta_{B_t, \alpha_t, j}) + \exp(\theta_{B_t, \alpha_t, j})}, \quad x \in \{-1, +1\} \quad (9)$$

where $\theta_{B_t, \alpha_t} = \nabla\ell(\mathbf{w}_{t-1}, S_{B_t})/\alpha_t$ is the scaled mini-batch gradient. The full-batch version uses parameters $\mathbb{E}_{B_t}[\theta_{B_t, \alpha_t}] = \nabla L_S(\mathbf{w}_{t-1})$. For the optimization analysis, we assume that the loss is smooth and mini-batch gradients are unbiased, symmetric, and sub-Gaussian.

Assumption 1. The loss function L_S satisfies: for all \mathbf{w} and \mathbf{w}' , for some non-negative constant $\vec{K} := [K_1, \dots, K_d]$, we have $L_S(\mathbf{w}) \leq L_S(\mathbf{w}') + \nabla L_S(\mathbf{w}')^T (\mathbf{w} - \mathbf{w}') + \frac{1}{2} \sum_i K_i (\mathbf{w}_i - \mathbf{w}'_i)^2$.

Assumption 2. Given \mathbf{w}_{t-1} , the mini-batch gradient $\nabla \ell(\mathbf{w}_{t-1}, S_{B_t})$ is (a) unbiased, i.e., $\mathbb{E}_{B_t | \mathbf{w}_{t-1}} \nabla \ell(\mathbf{w}_{t-1}, S_{B_t}) = \nabla L_S(\mathbf{w}_{t-1})$; (b) symmetric, i.e., the density $p(x) \equiv \nabla \ell(\mathbf{w}_{t-1}, S_{B_t})$ is symmetric around its expectation $L_S(\mathbf{w}_{t-1})$: $p(x) = p(2\nabla L_S(\mathbf{w}_{t-1}) - x)$ and (c) sub-Gaussian, i.e., for any $\lambda > 0$, any \mathbf{v} s.t. $\|\mathbf{v}\|_2 = 1$, $\mathbb{E}_{B_t | \mathbf{w}_{t-1}} \exp \lambda \langle \mathbf{v}, \nabla \ell(\mathbf{w}_{t-1}, S_{B_t}) - \nabla L_S(\mathbf{w}_{t-1}) \rangle \leq \exp(\lambda^2 \kappa_t^2 / 2)$ for some constant $\kappa_t > 0$.

Based on the assumptions, we have the following optimization guarantee for mini-batch noisy Sign-SGD. We defer the optimization guarantee for full-batch noisy Sign-SGD to Appendix D.

Theorem 3. Under Assumption 1 and 2, for mini-batch noisy Sign-SGD with step size $\eta_t = 1/\sqrt{T}$, α_t satisfying $c \geq \alpha_t \geq \max[\sqrt{2}\kappa_t, 4\|\nabla L_S(\mathbf{w}_{t-1})\|_\infty]$, we have for any S and any initialization \mathbf{w}_0

$$\mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \|\nabla L_S(\mathbf{w}_t)\|_2^2 \right] \leq \frac{4c}{\sqrt{T}} \left(L_S(\mathbf{w}_0) - L_S(\mathbf{w}^*) + \frac{1}{2} \|\vec{K}\|_1 \right), \quad (10)$$

where the expectation is taken over the randomness of algorithm.

SGLD. We acknowledge that the following optimization result of SGLD exists in various forms, as noisy gradient descent algorithms have been studied in literature such as differential privacy, where SGLD can be viewed as DP-SGD (Bassily et al., 2014; Wang & Xu, 2019) and the proof technique boils down to bounding the stochastic variance of the noisy gradient (Shamir & Zhang, 2013).

Theorem 4. Under Assumption 1 and 2, with $K_i = K, \forall i \in [d]$, for any S , SGLD (EFLD with step size $\rho_t = \sqrt{\eta_t \sigma_t}$, $\alpha_t = \sqrt{\sigma_t / \eta_t}$, $|B_t| = b$, and $\eta_t = \frac{1}{\sqrt{T}}$, can achieve

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla L_S(\mathbf{w}_t)\|^2 \leq O\left(\frac{1}{\sqrt{T}}\right) + O\left(K \frac{p \sum_{t=1}^T \alpha_t^4 + \log T}{\sqrt{T}}\right), \quad (11)$$

where the expectation is over the randomness of the algorithm.

The error rate of SGLD depends on the noise variance α_t . One can choose a decaying noise variance such as $\alpha_t = 1/\sqrt[4]{t}$ to guarantee the convergence. Then the rate will become $O(\log T / \sqrt{T})$. We note that similar to the optimization guarantees of DP-SGD, the convergence rate depends on the dimension of the gradient p due to the isotropic Gaussian noise. Special noise structures such as anisotropic noise that aligned with the gradient structure can reduce the dependence on dimension (Kairouz et al., 2020; Zhang et al., 2021; Asi et al., 2021; Zhou et al., 2020).

4 EXPERIMENTS

In this section, we conduct a series of experiments to evaluate our generalization error bounds. For SGLD, we aim to compare the proposed bound in Theorem 2 with existing bounds in Li et al. (2020), Negrea et al. (2019), and Rodríguez-Gálvez et al. (2021) for various datasets. Note that the bound presented in Rodríguez-Gálvez et al. (2021) is an extension of that in Haghifam et al. (2020) from full-batch setting to mini-batch setting. We also evaluate the optimization performance of proposed Noisy Sign-SGD by comparing it with the original sign-SGD (Bernstein et al., 2018a) and present the corresponding generalization bound in Theorem 2.

The details of our model architectures, learning rate scheduling, hyper-parameter selections and additional experimental results can be found in Appendix E. We acknowledge that we did not achieve the state-of-the-art predictive performance, mainly due to the simplicity of our model architectures. With more complex model and further tuning, the prediction results could be improved.

4.1 STOCHASTIC GRADIENT LANGEVIN DYNAMICS

Comparison with existing work. We have derived theoretical generalization error bounds that depend on the data-dependent quantity *gradient discrepancy*, i.e., $\|\nabla \ell(\mathbf{w}_t, z_n) - \nabla \ell(\mathbf{w}_t, z'_n)\|_2^2$. Existing bounds in Li et al. (2020) and Negrea et al. (2019) have also improved the Lipschitz constant in Mou et al. (2018) to a data-dependent quantity. As shown in Figure 1 (a)-(d), by combining with the empirical training error, all four generalization error bounds can be used to bound the empirical test error, but our bound is able to generate a much tighter upper bound. Such difference is mainly due to the fact that we replace the squared *gradient norm* in Li et al. (2020), the squared norm of

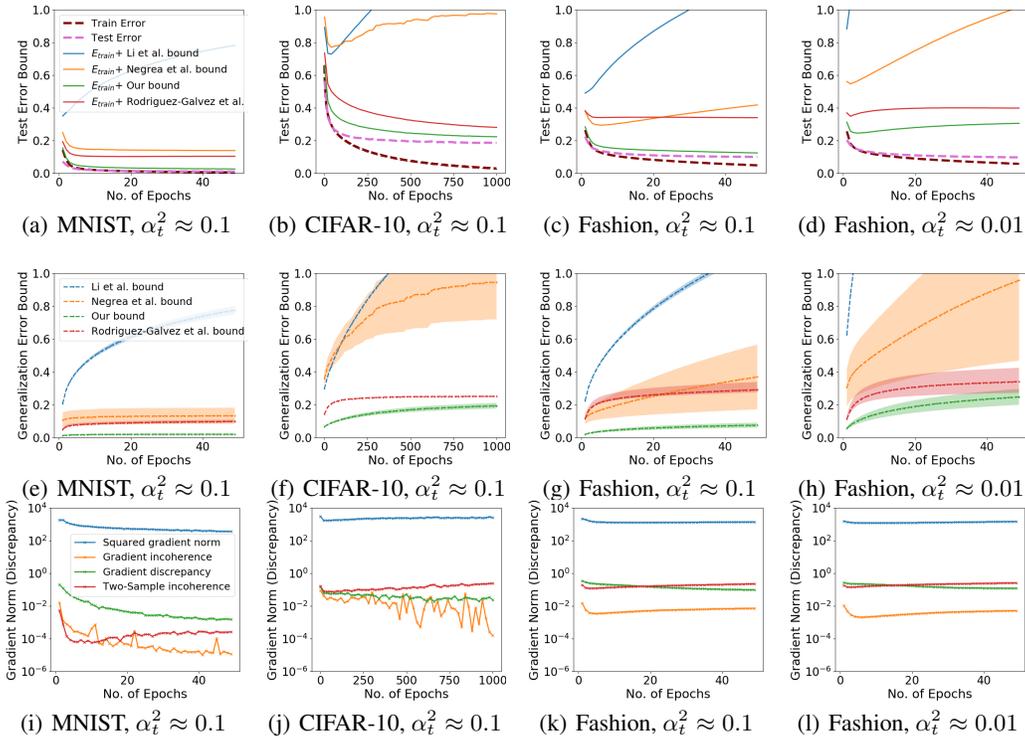


Figure 1: Numerical results for training CNN using SGLD ($\sigma_t = \sqrt{2\eta_t/\beta_t}$) on MNIST, Fashion-MNIST and CIFAR-10. X-axis shows the number of training epochs. (a)-(d) shows our bound is non-vacuous and can be used to bound the empirical test error. (e)-(h) compare our bound with the existing bounds and show the effect on α_t^2 . (i)-(l) show the key factors in each bound, i.e., the squared gradient norm in Li et al. (2020), the gradient incoherence in Negrea et al. (2019), the two-sample incoherence in Rodríguez-Gálvez et al. (2021), and the gradient discrepancy in our bound. Our bounds are numerically sharper than existing bounds, and larger α_t^2 leads to tighter generalization bounds which is consistent with the theoretical analysis.

gradient incoherence in Negrea et al. (2019), and that of *two-sample incoherence* in Rodríguez-Gálvez et al. (2021) with the gradient discrepancy. Results in Figure 1 (e)-(h) show that our bounds are much sharper than those of Li et al. (2020) because our gradient discrepancy (Figure 1 (i)-(l)) is usually 2-4 order of magnitude smaller than the squared gradient norms appeared in Li et al. (2020). Our bounds are also sharper than those of Negrea et al. (2019) and Rodríguez-Gálvez et al. (2021) due to an improved dependence on n from an order of $1/\sqrt{n}$ to $1/n$. Note that, even though the gradient incoherence in Negrea et al. (2019) is about 1 to 2 order of magnitude smaller than the gradient discrepancy for simple problems such as MNIST and Fashion-MNIST, the difference between the gradient incoherence and our gradient discrepancy reduces as the problem becomes harder (see results for CIFAR-10 in Figure 1(j)).

Effect of Randomness. Motivated by Zhang et al. (2017), we train CNN with SGLD on a smaller subset of MNIST dataset ($n = 10000$) with randomly corrupted labels. The corruption fraction varies from 0% (without label corruption) to 60%. As shown in Figure 2 (d), for long enough training time, all experiments with different level of label randomness can achieve almost zero training error. However, the one with higher level of randomness has higher generalization/test error (Figure 2 (a) dashed lines). Our generalization bound also becomes larger as the randomness increases since the corresponding gradient discrepancy increases.

4.2 NOISY SIGN-SGD

Optimization. Figure 3 (a)-(d) show the training dynamics of Noisy Sign-SGD under various selections of α_t . As $\alpha_t \rightarrow 0$, Noisy Sign-SGD matches both the optimization trajectory as well as the final test accuracy of the original Sign-SGD (Bernstein et al., 2018a). However, as α_t increases, the probability of getting 1 approaches 0.5, and ξ_t approximates a uniform distribution. As a result, the corresponding Noisy Sign-SGD still converges, but the generalization performance is much worse.

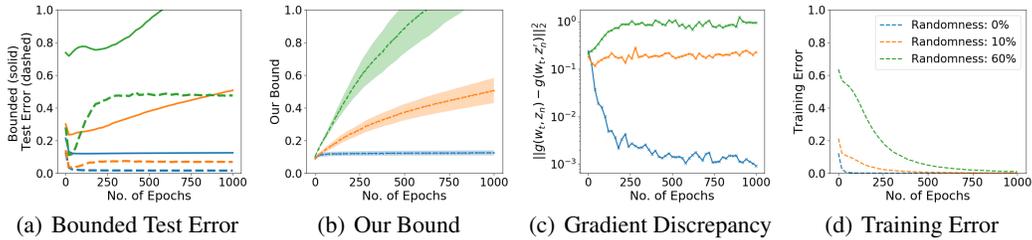


Figure 2: Numerical results for training CNN using SGLD ($\sigma_t = 0.2\eta_t$) on a subset of MNIST ($n = 10000$) with different randomness on labels. (a) demonstrates that, as the randomness increases, the empirical test error (dashed lines) increases but still can be bounded by our generalization bound by combining the empirical training error (solid lines). (b) presents our bound in Theorem 2. (c) shows the gradient discrepancy $\|\nabla \ell(\mathbf{w}_t, z_n) - \nabla \ell(\mathbf{w}_t, z'_n)\|_2$. (d) plots the training error. The gradient discrepancy increases as randomness increases, so does our generalization bound.

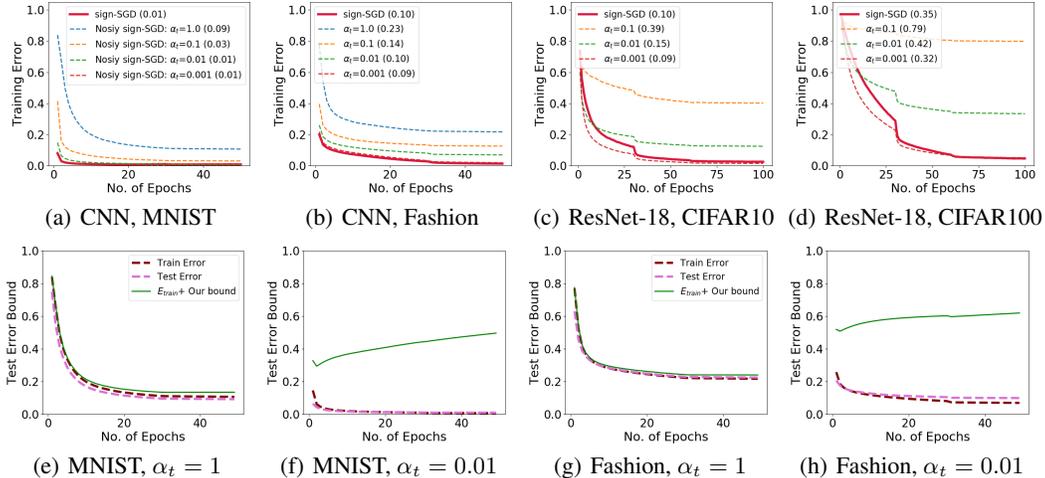


Figure 3: (a)-(d) show the training dynamics of CNN on MNIST and Fashion-MNIST, and ResNet-18 on CIFAR-10 and CIFAR-100 using noisy sign-SGD with different scaling α_t . Legends indicate the choice of α_t and the numbers in brackets are test errors at convergence. As $\alpha_t \rightarrow 0$, Noisy sign-SGD matches both the optimization trajectory as well as the final test accuracy of the original sign-SGD (Bernstein et al., 2018a). (e)-(f) show that empirical test error can be bounded by our bound and the corresponding training error. The larger α_t is the sharper our bound is.

Generalization Bound. Figure 3(e)-(f) show that our bound successfully bounds the empirical test error. The larger α_t is the sharper the upper bound is. However, larger α_t would slow down and adversely affect the optimization, e.g., Figure 3 (a)-(d) blue and orange lines. In practice, one needs to balance the optimization error and generalization by choosing a suitable scaling α_t .

5 CONCLUSIONS

Inspired by recent advances in stability based and information theoretic approaches to generalization bounds (Mou et al., 2018; Pensia et al., 2018; Negrea et al., 2019; Li et al., 2020; Haghifam et al., 2020), we have presented a framework for developing such bounds based on expected stability for noisy stochastic iterative (NSI) learning algorithms. We have also introduced Exponential Family Langevin Dynamics (EFLD), a family of noisy gradient descent algorithms based on exponential family noise, including SGLD and Noisy Sign-SGD as two special cases. We have developed an expected stability based generalization bound applicable to any EFLD algorithm with a $O(1/n)$ sample dependence and a dependence on gradient incoherence, rather than gradient norms. Further, we have provided optimization guarantees for EFLD and establish such results for Noisy Sign-SGD and SGLD. Our experiments on various benchmarks illustrate that our bounds are non-vacuous and quantitatively much sharper than existing bounds (Li et al., 2020; Negrea et al., 2019).

REFERENCES

- Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. Qsgd: Communication-efficient sgd via gradient quantization and encoding. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 30*, pp. 1709–1720. Curran Associates, Inc., 2017.
- Hilal Asi, John Duchi, Alireza Fallah, Omid Javidbakht, and Kunal Talwar. Private adaptive gradient methods for convex optimization. In *International Conference on Machine Learning*, pp. 383–392. PMLR, 2021.
- Arindam Banerjee, Srujana Merugu, Inderjit S Dhillon, and Joydeep Ghosh. Clustering with bregman divergences. *Journal of machine learning research*, 6(10), 2005.
- Ole Barndorff-Nielsen. *Information and exponential families: in statistical theory*. John Wiley & Sons, 2014.
- Raef Bassily, Adam Smith, and Abhradeep Thakurta. Private empirical risk minimization: Efficient algorithms and tight error bounds. In *2014 IEEE 55th Annual Symposium on Foundations of Computer Science*, pp. 464–473. IEEE, 2014.
- Raef Bassily, Vitaly Feldman, Kunal Talwar, and Abhradeep Guha Thakurta. Private stochastic convex optimization with optimal rates. *Advances in neural information processing systems*, 2019.
- Raef Bassily, Vitaly Feldman, Cristóbal Guzmán, and Kunal Talwar. Stability of stochastic gradient descent on nonsmooth convex losses. *Advances in Neural Information Processing Systems*, 33, 2020.
- Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signsgd: Compressed optimisation for non-convex problems. In *International Conference on Machine Learning*, pp. 560–569. PMLR, 2018a.
- Jeremy Bernstein, Jiawei Zhao, Kamyar Azizzadenesheli, and Anima Anandkumar. signsgd with majority vote is communication efficient and fault tolerant. In *International Conference on Learning Representations*, 2018b.
- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- Olivier Bousquet and André Elisseeff. Stability and generalization. *Journal of Machine Learning Research*, 2:499–526, 2002.
- Olivier Bousquet, Yegor Klochkov, and Nikita Zhivotovskiy. Sharper bounds for uniformly stable algorithms. In *Conference on Learning Theory*, pp. 610–626. PMLR, 2020.
- Lawrence D Brown. *Fundamentals of statistical exponential families: with applications in statistical decision theory*. Ims, 1986.
- Yuheng Bu, Shaofeng Zou, and Venugopal V Veeravalli. Tightening mutual information based bounds on generalization error. In *2019 IEEE International Symposium on Information Theory (ISIT)*, pp. 587–591. IEEE, 2019.
- Mark Bun, Cynthia Dwork, Guy N Rothblum, and Thomas Steinke. Composable and versatile privacy via truncated cdp. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pp. 74–86, 2018.
- Clément L Canonne, Gautam Kamath, and Thomas Steinke. The discrete gaussian for differential privacy. In *NeurIPS*, 2020.
- Xiangyi Chen, Tiancong Chen, Haoran Sun, Zhiwei Steven Wu, and Mingyi Hong. Distributed training with heterogeneous data: Bridging median-and mean-based algorithms. *arXiv preprint arXiv:1906.01736*, 2019.

- Luc Devroye and Terry Wagner. Distribution-free inequalities for the deleted and holdout error estimates. *IEEE Transactions on Information Theory*, 25(2):202–207, 1979.
- Andre Elisseeff, Theodoros Evgeniou, Massimiliano Pontil, and Leslie Pack Kaelbling. Stability of randomized learning algorithms. *Journal of Machine Learning Research*, 6(1), 2005.
- Vitaly Feldman and Jan Vondrak. Generalization bounds for uniformly stable algorithms. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pp. 9770–9780, 2018.
- Vitaly Feldman and Jan Vondrak. High probability generalization bounds for uniformly stable algorithms with nearly optimal rate. In *Conference on Learning Theory*, pp. 1270–1279. PMLR, 2019.
- Peter Grünwald, Thomas Steinke, and Lydia Zakyntinou. Pac-bayes, mac-bayes and conditional mutual information: Fast rate bounds that handle general vc classes. *arXiv preprint arXiv:2106.09683*, 2021.
- Mahdi Haghifam, Jeffrey Negrea, Ashish Khisti, Daniel M Roy, and Gintare Karolina Dziugaite. Sharpened generalization bounds based on conditional mutual information and an application to noisy, iterative algorithms. *Advances in Neural Information Processing Systems*, 2020.
- Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International Conference on Machine Learning*, pp. 1225–1234, 2016.
- Fredrik Hellström and Giuseppe Durisi. Generalization bounds via information density and conditional information density. *IEEE Journal on Selected Areas in Information Theory*, 1(3):824–839, 2020.
- Fredrik Hellström and Giuseppe Durisi. Fast-rate loss bounds via conditional information measures with applications to neural networks. In *2021 IEEE International Symposium on Information Theory (ISIT)*, pp. 952–957. IEEE, 2021.
- Peng Jiang and Gagan Agrawal. A linear speedup analysis of distributed deep learning with sparse and quantized communication. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 31*, pp. 2525–2536. Curran Associates, Inc., 2018.
- Chi Jin, Rong Ge, Praneeth Netrapalli, Sham M Kakade, and Michael I Jordan. How to escape saddle points efficiently. In *International Conference on Machine Learning*, pp. 1724–1732, 2017.
- Chi Jin, Praneeth Netrapalli, Rong Ge, Sham M Kakade, and Michael I Jordan. On nonconvex optimization for machine learning: Gradients, stochasticity, and saddle points. *arXiv preprint arXiv:1902.04811*, 2019.
- Richeng Jin, Yufan Huang, Xiaofan He, Tianfu Wu, and Huaiyu Dai. Stochastic-sign sgd for federated learning with theoretical guarantees. *arXiv preprint arXiv:2002.10940*, 2020.
- Peter Kairouz, Mónica Ribero, Keith Rush, and Abhradeep Thakurta. Dimension independence in unconstrained private erm via adaptive preconditioning. *arXiv preprint arXiv:2008.06570*, 2020.
- Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images. Technical Report Vol. 1. No. 4., University of Toronto, 2009.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Jian Li, Xuanyuan Luo, and Mingda Qiao. On generalization error bounds of noisy gradient methods for non-convex learning. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=SkxxtgHKPS>.
- Wenlong Mou, Liwei Wang, Xiyu Zhai, and Kai Zheng. Generalization bounds of sgd for non-convex learning: Two theoretical viewpoints. In *Conference on Learning Theory*, pp. 605–638. PMLR, 2018.

- Jeffrey Negrea, Mahdi Haghifam, Gintare Karolina Dziugaite, Ashish Khisti, and Daniel M Roy. Information-theoretic generalization bounds for sgld via data-dependent estimates. In *Advances in Neural Information Processing Systems*, 2019.
- Gergely Neu. Information-theoretic generalization bounds for stochastic gradient descent. *arXiv preprint arXiv:2102.00931*, 2021.
- Ankit Pensia, Varun Jog, and Po-Ling Loh. Generalization error bounds for noisy, iterative algorithms. In *2018 IEEE International Symposium on Information Theory (ISIT)*, pp. 546–550. IEEE, 2018.
- David Pollard. *A user’s guide to measure theoretic probability*. Number 8. Cambridge University Press, 2002.
- Borja Rodríguez-Gálvez, Germán Bassi, Ragnar Thobaben, and Mikael Skoglund. On random subset generalization error bounds and the stochastic gradient langevin dynamics algorithm. In *2020 IEEE Information Theory Workshop (ITW)*, pp. 1–5. IEEE, 2021.
- William H Rogers and Terry J Wagner. A finite sample distribution-free performance bound for local discrimination rules. *The Annals of Statistics*, pp. 506–514, 1978.
- Daniel Russo and James Zou. Controlling bias in adaptive data analysis using information theory. In *Artificial Intelligence and Statistics*, pp. 1232–1240. PMLR, 2016.
- Igal Sason and Sergio Verdu. f -divergence inequalities. *IEEE Transactions on Information Theory*, 62, 2016.
- Ohad Shamir and Tong Zhang. Stochastic gradient descent for non-smooth optimization: Convergence results and optimal averaging schemes. In *International conference on machine learning*, pp. 71–79. PMLR, 2013.
- Thomas Steinke and Lydia Zakyntinou. Reasoning about generalization via conditional mutual information. In *Conference on Learning Theory*, pp. 3437–3452. PMLR, 2020.
- Alexandre B Tsybakov. *Introduction to nonparametric estimation*. Springer Science & Business Media, 2008.
- Martin J Wainwright and Michael Irwin Jordan. *Graphical models, exponential families, and variational inference*. Now Publishers Inc, 2008.
- Di Wang and Jinhui Xu. Differentially private empirical risk minimization with smooth non-convex loss functions: A non-stationary view. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 1182–1189, 2019.
- Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *International Conference on Machine Learning, ICML ’11*, pp. 681–688, 2011.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms, 2017.
- Aolin Xu and Maxim Raginsky. Information-theoretic analysis of generalization capability of learning algorithms. *Advances in Neural Information Processing Systems*, 2017:2525–2534, 2017.
- Guandao Yang, Tianyi Zhang, Polina Kirichenko, Junwen Bai, Andrew Gordon Wilson, and Christopher De Sa. Swalp: Stochastic weight averaging in low-precision training. *36th International Conference on Machine Learning (ICML)*, 2019.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL <https://openreview.net/forum?id=Sy8gdB9xx>.
- Huanyu Zhang, Ilya Mironov, and Meisam Hejazinia. Wide network learning with differential privacy. *arXiv preprint arXiv:2103.01294*, 2021.

Ruida Zhou, Chao Tian, and Tie Liu. Individually conditional individual mutual information bound on generalization error. In *2021 IEEE International Symposium on Information Theory (ISIT)*, pp. 670–675. IEEE, 2021.

Yingxue Zhou, Steven Wu, and Arindam Banerjee. Bypassing the ambient dimension: Private sgd with gradient subspace identification. In *International Conference on Learning Representations*, 2020.