

# Two Facets of SDE Under an Information-Theoretic Lens: Generalization of SGD via Training Trajectories and via Terminal States



uOttawa

Ziqiao Wang<sup>1</sup> Yongyi Mao<sup>1</sup>

<sup>1</sup>School of Electrical Engineering and Computer Science, University of Ottawa

## Motivation

- Prevalent method of analyzing the generalization error of SGD via information-theoretic (IT) generalization bounds [Neu et al., 2021, Wang and Mao, 2022]:

$$\begin{aligned} \text{Gen. Err.}(\text{SGD}) &= \text{Gen. Err.}(\text{SGD}) + \text{Gen. Err.}(\text{NGD}) - \text{Gen. Err.}(\text{NGD}) \\ &\leq \text{ITBound}(\text{NGD}) + |\text{Gen. Err.}(\text{SGD}) - \text{Gen. Err.}(\text{NGD})|, \end{aligned}$$

where NGD is some noisy (stochastic) gradient descent.

- Empirical evidences [Wu et al., 2020, Li et al., 2021] show that  $|\text{Gen. Err.}(\text{SGD}) - \text{Gen. Err.}(\text{SDE})|$  is small: let NGD=SDE!
- Steady-state estimation of SDE enable us to analyze its terminal state.

## Background

- Learning algorithm  $\mathcal{A} : \mathcal{S} \rightarrow \mathcal{W}$  i.e. mapping a training sample (with size  $n$ ) to a hypothesis; Gen. Err.( $\mathcal{A}$ ) =  $\mathbb{E}[\text{Test Err.} - \text{Train Err.}]$
- SGD:  $\mathbf{w}_t = \mathbf{w}_{t-1} - \eta \tilde{\mathbf{G}}_t$ , where  $\eta$  is step size and  $\tilde{\mathbf{G}}_t$  is the mini-batch gradient with batch size  $b$ .
- SDE:  $\mathbf{w}_t = \mathbf{w}_{t-1} - \eta \mathbf{G}_t + \eta \mathbf{C}_t^{1/2} \mathbf{N}_t$ , where  $\mathbf{G}_t$  is the full-batch gradient,  $\mathbf{N}_t \sim \mathcal{N}(0, \mathbf{I}_d)$  and  $\mathbf{C}_t$  is gradient noise covariance (GNC):

$$\mathbf{C}_t \triangleq \frac{n-b}{b(n-1)} \left( \frac{1}{n} \sum_{i=1}^n \nabla \ell_i \nabla \ell_i^\top - \mathbf{G}_t \mathbf{G}_t^\top \right)$$

- Information-theoretic generalization bounds:

**Lemma 1.** For a subGaussian loss, Gen. Err.  $\leq \mathcal{O}\left(\sqrt{\frac{I(\mathcal{W}; \mathcal{S})}{n}}\right)$ .

**Lemma 2.** For a bounded loss, Gen. Err.  $\leq \mathcal{O}\left(\sqrt{D_{\text{KL}}(\mathbf{Q}_{\mathcal{W}|\mathcal{S}} || \mathbf{P}_{\mathcal{W}|\mathcal{S}})}\right)$ , where  $\mathcal{S}_J$  is a random subset of  $\mathcal{S}$ ,  $\mathbf{Q}_{\mathcal{W}|\mathcal{S}}$  is the posterior induced by  $\mathcal{A}$  and  $\mathbf{P}_{\mathcal{W}|\mathcal{S}}$  is a data-dependent prior.

## Generalization Bounds Via Full Trajectories

Recall  $I(\mathcal{X}; \mathcal{Y}) \leq \mathbb{E}_{\mathcal{X}}[D_{\text{KL}}(\mathbf{Q}_{\mathcal{Y}|\mathcal{X}} || \mathbf{P}_{\mathcal{Y}})]$ ,  $\mathbf{P}_{\mathcal{Y}}$  is some arbitrary prior.

- Using an isotropic Gaussian as prior, we have

**Theorem 1.** Let  $\Sigma_t^\mu \triangleq \mathbb{E}[\nabla \ell \nabla \ell^\top] - \mathbb{E}[\nabla \ell] \mathbb{E}[\nabla \ell]^\top$  be the population GNC. Assume  $\Sigma_t^\mu, \mathbf{C}_t \succ 0$ ,

$$\text{Gen. Err.} \lesssim \sqrt{\frac{1}{n} \sum_{t=1}^T \mathbb{E} \left[ d \log \frac{\text{tr}\{\Sigma_t^\mu\}}{bd} - \mathbb{E}[\text{tr} \log \mathbf{C}_t] \right]}.$$

**Remark.**  $\text{tr}\{\Sigma_t^\mu\} = \mathbb{E}[||\mathbf{G}_t - \mathbb{E}[\nabla \ell]||^2 + \text{tr}\{\mathbf{C}_t\}] \Rightarrow$

- First term: the sensitivity of  $\mathbf{G}_t$  to some variation of the training set  $\mathcal{S}$ .
- Second term: the gradient noise magnitude induced by mini-batch.
- By-product: recovering a bound for Gradient Langevin dynamics

**Corollary 1.** If  $\mathbf{C}_t = \mathbf{I}_d$ , then

$$\text{Gen. Err.} \lesssim \sqrt{\frac{d}{n} \sum_{t=1}^T \mathbb{E} \log \left( \mathbb{E}[||\mathbf{G}_t - \mathbb{E}[\nabla \ell]||^2]/d + 1 \right)}.$$

**Remark.** Not necessarily depends on  $d$  (by  $\log(x+1) \leq x$ ).

- Using an anisotropic Gaussian as prior, we have

**Theorem 2.** Under the same conditions in **Theorem 1.**,

$$\text{Gen. Err.} \lesssim \sqrt{\sum_{t=1}^T \frac{\mathbb{E}[\text{tr} \log (\Sigma_t^\mu \mathbf{C}_t^{-1}/b)]}{n}}.$$

**Remark.** **Theorem 2.** is tighter than **Theorem 1.**

Let  $\Sigma_t = b\mathbf{C}_t$ , then  $\Sigma_t^\mu \Sigma_t^{-1}$  is small  $\iff$  SGD is insensitive to the randomness of  $\mathcal{S}$ . Same intuition with  $I(\mathcal{W}; \mathcal{S})$  in **Lemma 1**.

## Take Home Messages

- Trajectories-based bounds need less assumptions but are time-dependent.
- Terminal-state-based bounds are time-independent but require additional assumptions and approximations.

## Generalization Bounds Via Terminal State

Quadratic loss:  $\mathbf{w} \rightarrow$  local minimum  $\mathbf{w}^*$ , let  $\mathbf{H}_{\mathbf{w}^*}$  be Hessian at  $\mathbf{w}^*$ ,

$$\text{Loss of } \mathbf{w} = \text{Loss of } \mathbf{w}^* + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^\top \mathbf{H}_{\mathbf{w}^*}(\mathbf{w} - \mathbf{w}^*).$$

$\xrightarrow{T \rightarrow \infty}$  given a  $\mathcal{S}$  and its local minimum  $\mathbf{w}_s^*$ ,  $\mathbf{W}_T \sim \mathcal{N}(\mathbf{w}_s^*, \Lambda_{\mathbf{w}_s^*})$ .

$\xrightarrow{W_s^* \sim Q_{W_s^*|s}}$   $Q_{W_s^*|s} = \mathbb{E}_{\mathbf{W}_s^*} [\mathcal{N}(\mathbf{W}_s^*, \Lambda_{\mathbf{w}_s^*})]$  is a mixture of Gaussian.

- Lemma 3.**  $\Lambda_{\mathbf{w}^*} \mathbf{H}_{\mathbf{w}^*} + \mathbf{H}_{\mathbf{w}^*} \Lambda_{\mathbf{w}^*} - \eta \mathbf{H}_{\mathbf{w}^*} \Lambda_{\mathbf{w}^*} \mathbf{H}_{\mathbf{w}^*} = \eta \mathbf{C}_T$ .

- Hessian-based Result

**Theorem 3.** Let  $\Lambda_{\mathbf{w}_\mu^*} \triangleq \mathbb{E}[(\mathbf{W} - \mathbb{E}[\mathbf{W}_s^*])(\mathbf{W} - \mathbb{E}[\mathbf{W}_s^*])^\top]$ . Under some mild assumptions,

$$\text{Gen. Err.} \lesssim \sqrt{\frac{1}{n} \mathbb{E} \left[ \text{tr} \log \left( [\mathbf{H}_{\mathbf{w}^*} \mathbf{C}_T^{-1}] \Lambda_{\mathbf{w}_\mu^*} \right) \right]}.$$

**Remark.** Alignment between a population and a sample stationary dist.

- Norm-based Result

**Theorem 4.** Let  $\hat{\mathbf{w}}$  be a reference vector. Under some mild assumptions,

$$\text{Gen. Err.} \lesssim \sqrt{\frac{d}{n} \log \left( \frac{b}{\eta d} \mathbb{E} ||\mathbf{W}_s^* - \hat{\mathbf{w}}||^2 + 1 \right)}.$$

**Remark.** i)  $\hat{\mathbf{w}} = \mathbb{E}[\mathbf{W}_s^*] \Rightarrow$  Optimal; ii)  $\hat{\mathbf{w}} = \mathbf{w}_0 \Rightarrow$  "Distance from initialization"; iii)  $\hat{\mathbf{w}} = 0 \Rightarrow$  Weight Decay.

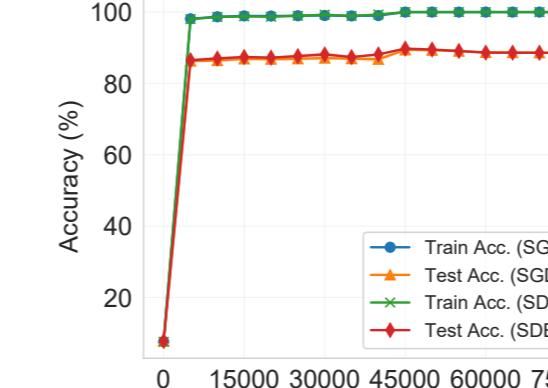
- Stability-based Result

**Theorem 5.** Recall **Lemma 2.** and under some mild assumptions,

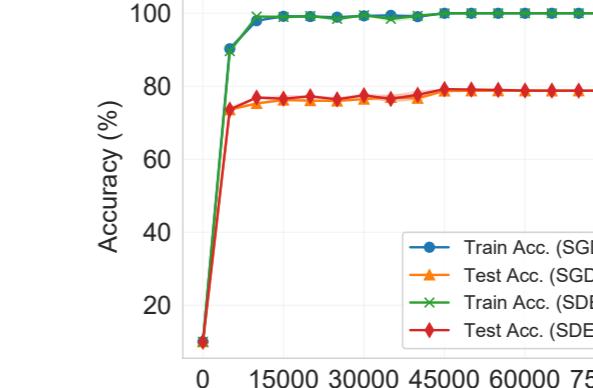
$$\text{Gen. Err.} \lesssim \sqrt{\frac{b}{\eta} \mathbb{E} ||\mathbf{W}_s^* - \mathbf{W}_{s_J}^*||^2}.$$

**Remark.** No Lipschitz constant contained; Fast-rate in some cases.

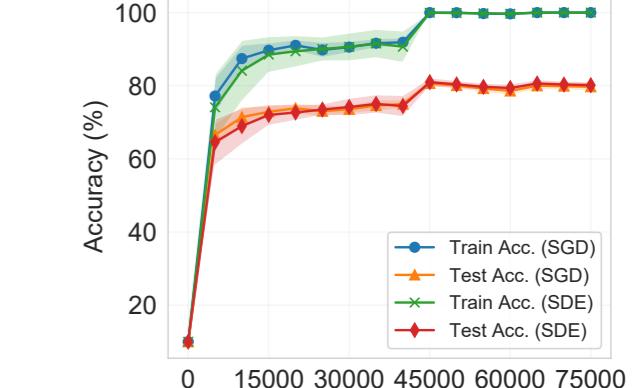
## Empirical Results



(a) VGG on (small) SVHN

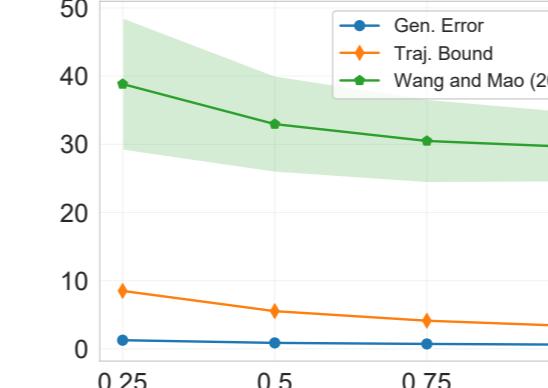


(b) VGG on CIFAR10

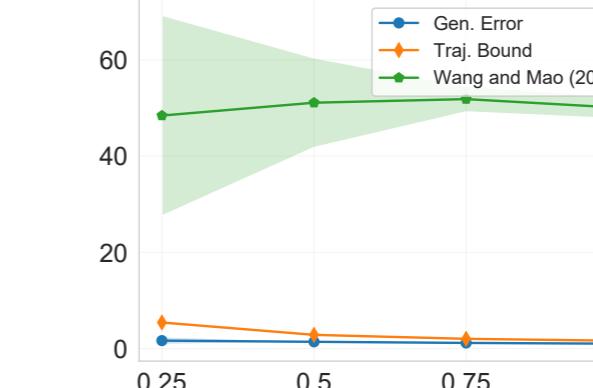


(c) ResNet on CIFAR10

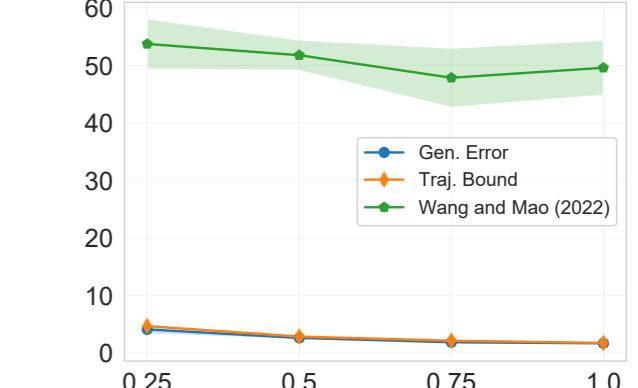
Figura 1: Performance of VGG-11 and ResNet-18 trained with SGD and SDE.



(a) VGG on (small) SVHN

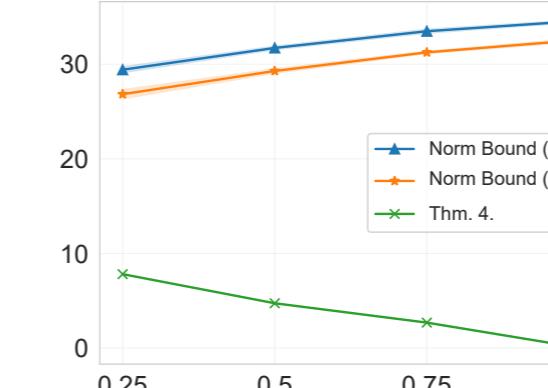


(b) VGG on CIFAR10

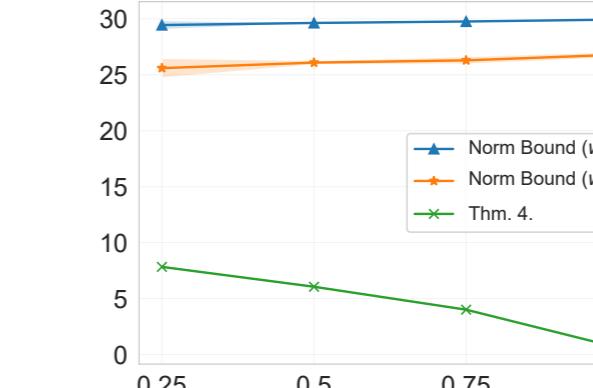


(c) ResNet on CIFAR10

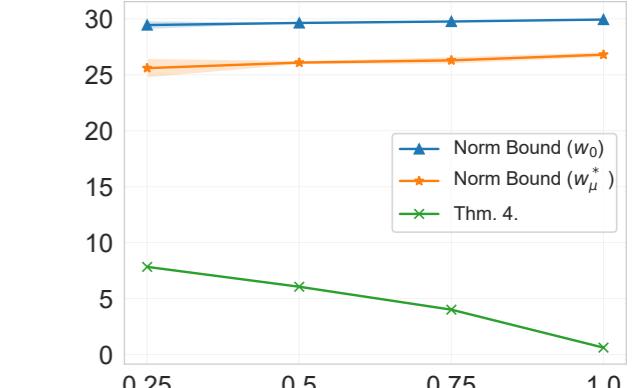
Figura 2: Scaled trajectories-based bound. Compared with Wang and Mao [2022].



(a) VGG on (small) SVHN



(b) VGG on CIFAR10



(c) ResNet on CIFAR10

## Reference

- Zhiyuan Li et al. On the validity of modeling sgd with stochastic differential equations (sdes). *NeurIPS*, 2021.
- Gergely Neu et al. Information-theoretic generalization bounds for stochastic gradient descent. In *COLT*, 2021.
- Ziqiao Wang and Yongyi Mao. On the generalization of models trained with SGD: Information-theoretic bounds and implications. In *ICLR*, 2022.
- Jingfeng Wu et al. On the noisy gradient descent that generalizes as sgd. In *ICML*, 2020.