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

Ziqiao Wang [1]    Yongyi Mao [1]

[1]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]:

  Gen. Err.(SGD) = Gen. Err.(SGD) + Gen. Err.(NGD) − Gen. Err.(NGD)
  ≤ ITBound(NGD) + |Gen. Err.(SGD) − Gen. Err.(NGD)|,

  where NGD is some noisy (stochastic) gradient descent.
- Empirical evidences [Wu et al., 2020, Li et al., 2021] show that |Gen. Err.(SGD) − Gen. Err.(SDE)| is small: let NGD=SDE!
- Steady-state estimation of SDE enable us to analyze its terminal state.

## Background

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

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

- Information-theoretic generalization bounds:

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

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

## Generalization Bounds Via Full Trajectories

Recall $\boldsymbol{I}(\boldsymbol{X}; \boldsymbol{Y}) \leq \mathbb{E}_{\boldsymbol{X}}\left[\mathbf{D}_{\mathrm{KL}}\left(\boldsymbol{Q}_{Y|X} || \boldsymbol{P}_Y\right)\right]$, $\boldsymbol{P}_Y$ is some arbitrary prior.

- Using an isotropic Gaussian as prior, we have

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

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

  **Remark.** $\boldsymbol{tr}\{\Sigma_t^\mu\} = \mathbb{E}\left[ ||\boldsymbol{G}_t - \mathbb{E}[\nabla \ell]||^2 + \boldsymbol{tr}\{\boldsymbol{C}_t\}\right] \Longrightarrow$

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

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

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

  **Remark.** Not necessarily depends on $\boldsymbol{d}$ (by $\log(x+1) \leq \boldsymbol{x}$).
- Using an anisotropic Gaussian as prior, we have

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

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

  **Remark. Theorem 2.** is tighter than **Theorem 1.**
  Let $\Sigma_t = \boldsymbol{b}\boldsymbol{C}_t$, then $\Sigma_t^\mu \Sigma_t^{-1}$ is small $\iff$ SGD is insensitive to the randomness of $\boldsymbol{S}$. Same intuition with $I(\boldsymbol{W}; \boldsymbol{S})$ in **Lemma 1**.

## Take-Home Messages

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

## Generalization Bounds Via Terminal State

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

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

$\overset{T \to \infty}{\Longrightarrow}$ given a $\boldsymbol{s}$ and its local minimum $\boldsymbol{w}_s^*$, $\boldsymbol{W}_T \sim \mathcal{N}(\boldsymbol{w}_s^*, \Lambda_{w_s^*})$.
$\overset{W_s^* \sim Q_{W_s^*|s}}{\Longrightarrow} \boldsymbol{Q}_{W_T|s} = \mathbb{E}_{W_s^*}^s\left[\mathcal{N}(\boldsymbol{W}_s^*, \Lambda_{W_s^*})\right]$ is a mixture of Gaussian.

- **Lemma 3.** $\Lambda_{w^*}\boldsymbol{H}_{w^*} + \boldsymbol{H}_{w^*}\Lambda_{w^*} - \eta \boldsymbol{H}_{w^*}\Lambda_{w^*}\boldsymbol{H}_{w^*} = \eta \boldsymbol{C}_T.$
- Hessian-based Result

  **Theorem 3.** Let $\Lambda_{w_\mu^*} \triangleq \mathbb{E}\left[\left(\boldsymbol{W} - \mathbb{E}\left[\boldsymbol{W}_S^*\right]\right)\left(\boldsymbol{W} - \mathbb{E}\left[\boldsymbol{W}_S^*\right]\right)^{\mathsf{T}}\right].$ Under some mild assumptions,

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

  **Remark.** Alignment between a population and a sample stationary dist.
- Norm-based Result

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

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

  **Remark.** i) $\hat{\boldsymbol{w}} = \mathbb{E}\left[\boldsymbol{W}_S^*\right] \Rightarrow$ Optimal; ii) $\hat{\boldsymbol{w}} = \boldsymbol{w}_0 \Rightarrow$ "Distance from initialization"; iii) $\hat{\boldsymbol{w}} = 0 \Rightarrow$ Weight Decay.
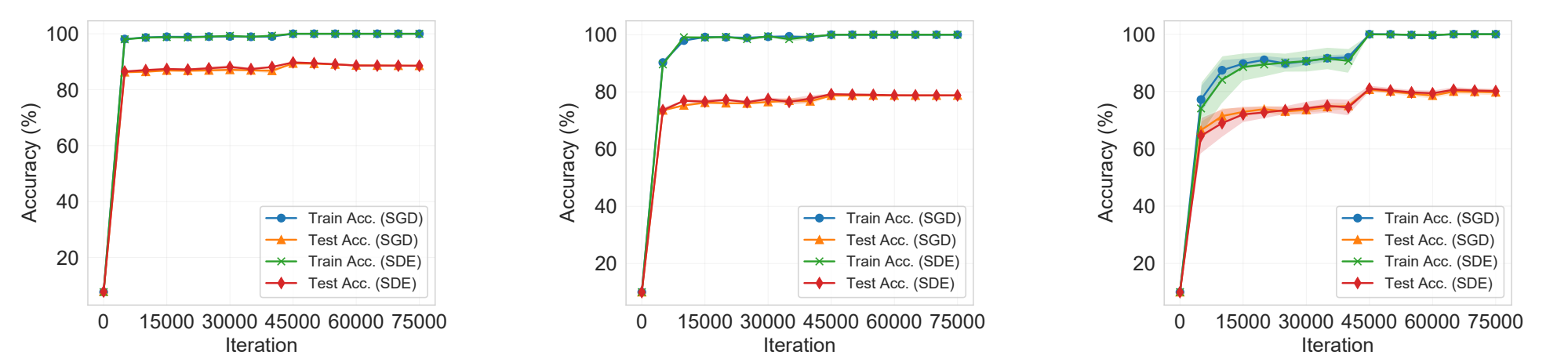- Stability-based Result

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

$$\text{Gen. Err.} \precsim \sqrt{\frac{\boldsymbol{b}}{\eta}\mathbb{E}||\boldsymbol{W}_S^* - \boldsymbol{W}_{S_J}^*||^2}.$$
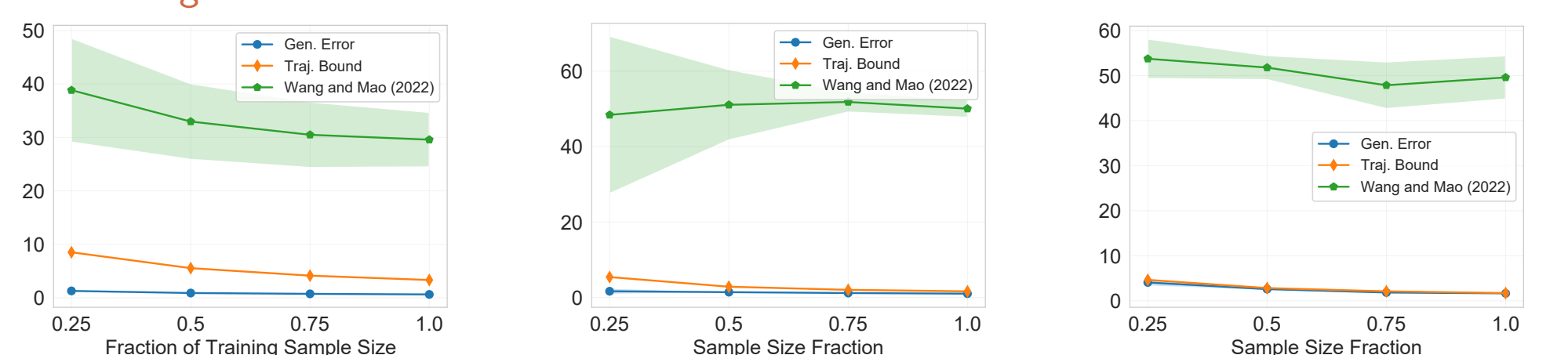
  **Remark.** No Lipschitz constant contained; Fast-rate in some cases.
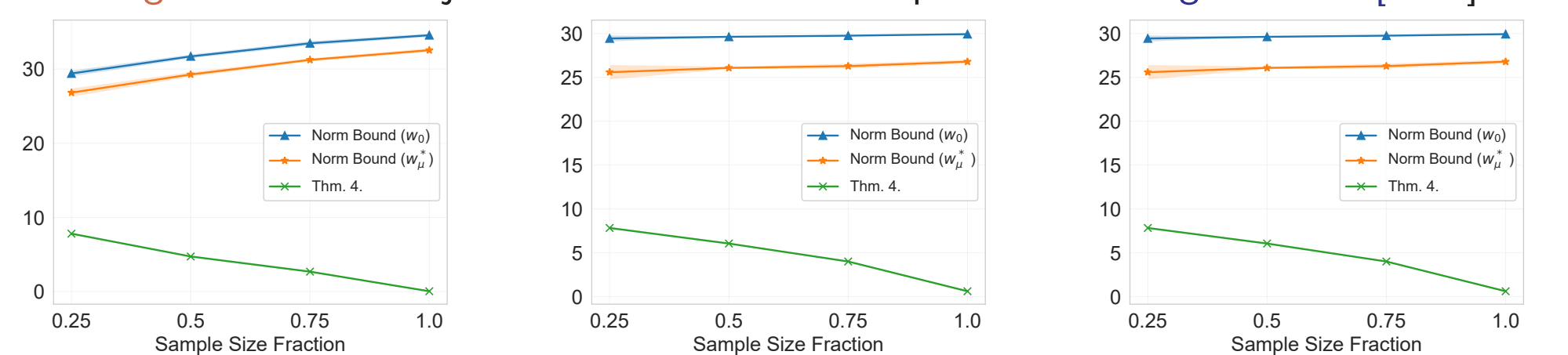
## Empirical Results



(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

Figura 3: Scaled terminal-state based bound.

## 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.