

IDENTIFYING STABILITY REGIONS OF SGD WITH CONSTANT LEARNING RATES

Anonymous authors

Paper under double-blind review

ABSTRACT

The trade-off inherent in constant learning rate stochastic gradient descent (SGD) has been well-documented empirically: larger learning rates often yield faster convergence, but risk the possibility of exploding. However, the relevant question of an appropriate choice of learning rate has rarely received systematic treatment; one often chooses learning schedules based on domain knowledge and preliminary numerical experiments without theoretical guidance. This question is intimately related to the concept of “edge of stability”, which refers to a regime where the chain neither converges nor explodes. Despite rich literature on deterministic gradient descent, the rigorous characterization of “edge of stability” for the more ubiquitous SGD chains, remains an open question. In this paper, we formalize the notion of the stability region, and develop theoretical guarantees for estimating the stability region for SGD for a wide class of strongly convex objectives. We introduce a stochastic version of Lyapunov exponent for SGD, which yields a practical, data-driven threshold for admissible learning rates. Moreover, all of our theoretical results are backed by extensive experiments. Collectively, these findings demonstrate a practically implementable as well as theoretically valid way of choosing learning rate parameters in various problems, while also paving the way to potential generalization to more complicated nonconvex landscapes.

1 INTRODUCTION

The dynamics of stochastic gradient descent (SGD) and related optimization methods have been studied extensively from the perspective of stability, generalization, and convergence. Foundational analyses such as Hardt et al. (2016) established stability guarantees for SGD and connected them to generalization, while subsequent works have investigated SGD as an approximate Bayesian inference procedure (Mandt et al., 2017) and as a stochastic process with heavy-tailed gradient noise (Simsekli et al., 2019). More recently, SGD has also been analyzed as a random dynamical system with almost sure convergence properties (Daneshmand et al., 2024) and from a nonlinear time series perspective (Li et al., 2025). However, a consistent theme with the majority of these literature is the lack of principled guidelines on how to choose the (small enough) step-size that ensures the stability of the system. On the other hand, choosing a learning rate that is too small leads to excruciatingly slow convergence. Edge of stability analysis reflects the sweet spot between stability and convergence.

However, until recently, the *edge of stability* literature has largely focused on deterministic gradient descent (GD). Conventional theoretical analyses typically focus on the inverted problem of the stability threshold—namely, convergence guarantees at the sharpness threshold (i.e., the maximum eigenvalue of the Hessian) that guarantees stability for a GD algorithm with a given step size. The practically relevant problem of determining a problem and data-dependent threshold of learning rate that ensures stability, is much less explored. Moreover, often stochastic gradient descent is used over vanilla GD in an online setting, and much less is known about the edge-of-stability threshold for the SGD algorithms. In this article, we bridge this gap between theory and practice by proposing a theoretically valid, as well as practically implementable data-driven estimate of edge-of-stability for SGD algorithms in strongly convex setting. Our main contributions are as follows.

1.1 MAIN CONTRIBUTIONS

Maximal expansion parameter As a stepping stone to the notion of edge-of-stability, we analyze the geometric moment contraction of the SGD dynamics and define the *maximal expansion parameter* $L^\ell(\gamma)$ as the maximal Lipschitz parameter for ℓ -step SGD dynamics given $\ell \in \mathbb{N}_+$ and step size $\gamma > 0$. This parameter can be understood as the value

of the weakest possible contraction of the SGD functional with step-size γ . Leveraging tools from time-series theory, we provide asymptotic theory for estimating $L^\ell(\gamma)$ uniformly over γ ;

Theorem 1.1 (Theorem 3.1, informal). *Under standard regularity conditions, it follows that $\sup_{\gamma \in \Gamma} |\hat{L}^\ell(\gamma) - L^\ell(\gamma)| = O_{\mathbb{P}}(\frac{\log n}{\sqrt{n}})$, where Γ is a compact set.*

Towards the development of this result, we also borrow insights from high-dimensional statistics literature to provide a sharp uniform moment bound on the partial sums of i.i.d. random functions. We expect this result to be of independent interest.

Conceptual development and estimation of edge-of-stability for SGD. Developing on the concept of maximal expansion parameter, we rigorously characterize the *edge-of-stability*, denoted by γ_ℓ , in terms of the smallest learning rate that pushes the ℓ -step maximal expansion parameter beyond 1, thereby making the chain explode. Our definition leads to a natural estimation strategy for this *edge-of-stability* threshold, denoted by $\hat{\gamma}_{\ell,n}$. **Theory for the estimator $\hat{\gamma}_n$** To the best of our knowledge, this work is the first one to provide finite-sample error bounds on the convergence property of $\gamma_{\ell,n}$; in particular, we prove the following theorem.

Theorem 1.2 (Theorem 4.3, informal). *Under standard regularity conditions, it follows that $|\hat{\gamma}_{\ell,n} - \gamma_\ell| = O_{\mathbb{P}}(\frac{\log n}{\sqrt{n}})$.*

Here we present two examples on linear regression and expectile regression respectively. The detailed settings are deferred to Remark 2.2 and Section 5. In particular, the exact forms of the learning-rate boundary can be provided in the linear regression model, which are $\gamma = 2/3$ and $\gamma = 10/3$ for the random samples generated from standard normal distribution and standard uniform distribution, respectively, with $p = 2$ and $d = 1$. As shown in Figure 1, by our proposed methodology, we can very accurately hit the boundary that we derived theoretically (denoted by vertical dashed lines in Figure 1(a)).

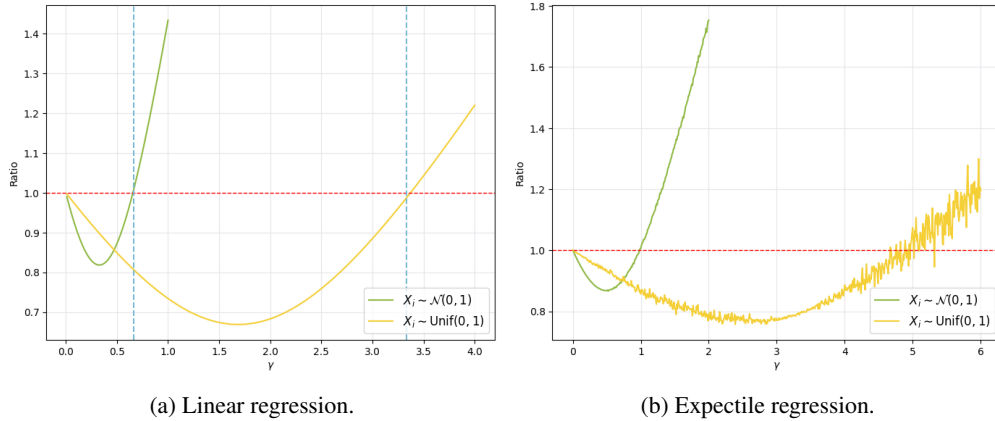


Figure 1: Examples for edge of stability. Green: $\mathcal{N}(0, 1)$; Yellow: $\text{Unif}[0, 1]$. All the experiments are repeated 30 times. The detailed setting is provided in Section 5.

Connections with Lyapunov theory. Our framework admits a natural interpretation in terms of Lyapunov theory once we adopt an asymptotic point of view on edge-of-stability. Indeed, by letting $\ell \rightarrow \infty$ in the definition of the maximal expansion parameter, submultiplicativity and Fekete’s lemma ensure that the limit:

$$\lambda_p(\gamma) := \lim_{\ell \rightarrow \infty} \frac{1}{\ell} \log L_\ell^p(\gamma)$$

exists. This quantity is precisely the maximal *Lyapunov exponent* associated with the stochastic dynamics of SGD at learning rate γ : it measures the exponential rate at which moment distances between trajectories grow (if positive) or decay (if negative). In particular,

$$\begin{aligned} \lambda_p(\gamma) < 0 &\Rightarrow \text{exponential contraction of } p\text{-th moments,} \\ \lambda_p(\gamma) > 0 &\Rightarrow \text{exponential expansion of } p\text{-th moments.} \end{aligned}$$

Accordingly, the oracle edge of stability can be equivalently characterized as the zero-crossing of this exponent:

$$\gamma_\infty(p) := \inf\{\gamma \in \Gamma \mid \lambda_p(\gamma) \geq 0\}.$$

This viewpoint places our notion of edge-of-stability squarely within the classical Lyapunov framework: SGD dynamics remain stable as long as the maximal Lyapunov exponent is negative, and instability begins exactly at the point where it reaches zero. The construction aligns with classical work on Lyapunov exponents for products of random matrices and random dynamical systems, beginning with Oseledets' multiplicative ergodic theorem (Oseledets, 1968) and subsequent developments in the monographs of Bougerol and Lacroix Bougerol & Lacroix (1985) and Arnold Arnold (1998). In those settings, the sign of the maximal Lyapunov exponent governs long-run stability of the system. Our moment-based definition $\lambda_p(\gamma)$ can be interpreted as an analogue tailored to stochastic approximation schemes such as SGD, and places the edge-of-stability phenomenon within the same analytical framework.

Notation. In this paper, we denote the set $\{1, \dots, n\}$ by $[n]$. The d -dimensional Euclidean space is \mathbb{R}^d . For a vector $a \in \mathbb{R}^d$, $|a|$ denotes its Euclidean norm. For a matrix $M \in \mathbb{R}^{d \times m}$, $|A|$ denotes its Euclidean operator norm. For a random vector $X \in \mathbb{R}^d$, we denote $\|X\| := \sqrt{\mathbb{E}[|X|^2]}$. We also denote in-probability convergence, and stochastic boundedness by $o_{\mathbb{P}}$ and $O_{\mathbb{P}}$ respectively. We write $a_n \lesssim b_n$ if $a_n \leq Cb_n$ for some constant $C > 0$, and $a_n \asymp b_n$ if $C_1b_n \leq a_n \leq C_2b_n$ for some constants $C_1, C_2 > 0$. Often we denote $a_n \lesssim b_n$ by $a_n = O(b_n)$. Additionally, if $a_n/b_n \rightarrow 0$, we write $a_n = o(b_n)$. For a compact convex set $\Gamma \subset \mathbb{R}^d$, we denote by $\text{int}(\Gamma) := \{x \in \Gamma : \exists \varepsilon > 0 \text{ such that } B_\varepsilon(x) \subset \Gamma\}$, where $B_\varepsilon(x) := \{y : |x - y| < \varepsilon\}$ is the ε -ball around $x \in \mathbb{R}^d$. In particular, we denote the closed unit ball in \mathbb{R}^d by $\mathcal{B} := \overline{B_1(0)}$.

2 EDGE OF STABILITY: PRELIMINARIES

For a function $G : \mathbb{R}^d \mapsto \mathbb{R}$, consider the following optimization problem:

$$\theta^* = \arg \min_{\theta \in \mathcal{D}} G(\theta), \quad \mathcal{D} \subset \mathbb{R}^d \text{ is compact and convex,}$$

and let $\xi_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{P}$ be the innovations. Subsequently, all the probability statements are carried out on the same measure space as \mathcal{P} . Define $F \in \mathcal{C}^1$. With an online stream of ξ_1, ξ_2, \dots , the classical SGD algorithm estimates θ^* via the recursion

$$\theta_i = F_{\xi_i}^\gamma(\theta_{i-1}), \text{ with } F_{\xi_i}^\gamma(\theta) = \theta - \gamma \nabla g(\theta, \xi_i), \quad i = 1, 2, \dots, \quad (1)$$

where g is a measurable function, and $g(\cdot, x) \in \mathcal{C}^2$ satisfies $\mathbb{E}[\nabla g(\theta, \xi)] = \nabla G(\theta)$. Here $\gamma > 0$ is the constant learning rate. Before proceeding further, we introduce two key assumptions that are ubiquitous in SGD literature, as well as heavily used throughout our article.

Assumption 2.1 (μ -strong convexity). *There exists a $\mu > 0$ such that g is μ -strongly convex; in other words, for all $\theta, \theta' \in \mathbb{R}^d$,*

$$\langle m(\theta) - m(\theta'), \theta - \theta' \rangle \geq \mu |\theta - \theta'|^2,$$

where $m(\theta) := \mathbb{E}[\nabla g(\theta, \xi)]$, $\xi \sim \mathcal{P}$.

Strong convexity is a textbook assumption in the stochastic approximation literature (Ruppert, 1988; Polyak & Juditsky, 1992; Bottou et al., 2018a). It guarantees uniqueness of the minimizer and provides a quadratic lower bound that underlies contraction arguments. This assumption is standard in convex SGD theory, and is satisfied by canonical problems such as linear or regularized logistic regression. While it does not extend to general nonconvex objectives, it is well aligned with our focus on strongly convex settings.

Assumption 2.2 (Stochastic Lipschitz continuity). *Let $p \geq 1$. There exists some constant $N_p > 0$ such that, for all $\theta, \theta' \in \mathbb{R}^d$,*

$$\|\nabla g(\theta, \xi) - \nabla g(\theta', \xi)\|_p \leq N_p |\theta - \theta'|.$$

Strong convexity guarantees uniqueness of the minimizer and provides a quadratic lower bound on the objective. This ensures that the SGD iterates are attracted toward a single point rather than drifting among multiple optima, and it underlies the contraction arguments that follow. On the other hand, stochastic Lipschitz-ness controls the variability of the stochastic gradients across different parameter values. This assumption enables us to bound deviations of the stochastic gradients uniformly, which is essential when passing from local to global statements in a concentration analysis. We remark that Assumptions 2.1 and 2.2 are standard features of statistical analysis of convex stochastic optimization, and have appeared extensively in Ruppert (1988); Polyak & Juditsky (1992); Bottou et al. (2018b); Chen et al. (2020); Zhu et al. (2023); Wei et al. (2023); Li et al. (2024).

2.1 MAXIMAL EXPANSION PARAMETER: INTRODUCTION

As discussed in §1, the learning rate $\gamma > 0$ plays a fundamental role in the performance of SGD; a larger value of γ may lead to θ_i being divergent. However, one can preclude the possibility of explosion by theoretically analyzing the maximum possible contraction after a given number of iterates from the current instance. We formalize this insight by borrowing the notion of contractive maps in dynamic systems defined by Wu (2005);

Definition 2.1 (MEP- ℓ). *The p -th Maximal Expansion Parameter of lag 1 (MEP-1) is defined as*

$$L_p(\gamma) := \sup_{\theta \neq \theta'} \frac{\mathbb{E} \left[\left| F_{\xi_i}^\gamma(\theta) - F_{\xi_i}^\gamma(\theta') \right|^p \right]}{|\theta - \theta'|^p}. \quad (2)$$

Generalizing (2), for $\ell \in \mathbb{N}_+$, the ℓ -lag maximal expansion (MEP- ℓ) can be defined as:

$$L_p^\ell(\gamma) := \sup_{\theta \neq \theta' \in \mathcal{D}} \frac{\mathbb{E} \left[\left| F_{\xi_{i+\ell-1}:\xi_i}^\gamma(\theta) - F_{\xi_{i+\ell-1}:\xi_i}^\gamma(\theta') \right|^p \right]}{|\theta - \theta'|^p},$$

where the composite map $F_{(a+b):a}^\gamma(\cdot) := F_{a+b}^\gamma \circ \dots \circ F_{a+1}^\gamma \circ F_a^\gamma(\cdot)$.

The quantity $L_p(\gamma)$ can be interpreted as the maximal possible value of the Lipschitz constant in equation (17) of Li et al. (2025); as we will discuss in §4, this interpretation readily leads to a notion of edge-of-stability through the need to ensure geometric moment contraction. However, before proceeding further, we take a pause here to make a crucial observation regarding the tractability of the maximal expansion parameter.

The maximal expansion parameter, as is defined, concerns computing a supremum over pairs of distinct points θ, θ' . This form may appear cumbersome for both analysis, as well as any direct approach to estimation. However, in Lemma 2.2, we transform the corresponding sample version into a tractable quantity through equivalent characterization through $\nabla_\theta F_{\xi_i}^\gamma(\theta)$ for all ξ_i and θ .

Lemma 2.2. *Let $\mathcal{D} \subset \mathbb{R}^d$ be compact convex set and $\gamma > 0$ be given. Suppose $F_{\xi_i}^\gamma(\theta)$ be as in Equation (1). Then, under Assumption 2.2 it follows that:*

$$\sup_{\theta \neq \theta' \in \mathcal{D}} \frac{1}{n} \sum_{i=1}^n \frac{\left| F_{\xi_i}^\gamma(\theta) - F_{\xi_i}^\gamma(\theta') \right|^p}{|\theta - \theta'|^p} = \sup_{\theta \in \mathcal{D}} \sup_{u: |u|=1} \frac{1}{n} \sum_{i=1}^n \left| \nabla_\theta F_{\xi_i}^\gamma(\theta) u \right|^p. \quad (3)$$

Additionally, it follows that

$$\sup_{\theta \neq \theta' \in \mathcal{D}} \frac{\mathbb{E} \left[\left| F_{\xi_i}^\gamma(\theta) - F_{\xi_i}^\gamma(\theta') \right|^p \right]}{|\theta - \theta'|^p} = \sup_{\theta \in \mathcal{D}} \sup_{u \in \mathbb{R}^d: |u|=1} \mathbb{E} \left[\left| \nabla_\theta F_{\xi_i}^\gamma(\theta) u \right|^p \right].$$

Remark 2.1. *Virtually the same arguments as Lemma 2.2 allow us to write*

$$\sup_{\theta \in \mathcal{D}} \sup_{u \in \mathbb{R}^d: |u|=1} \frac{1}{n} \sum_{i=1}^n \left[\left| \nabla_\theta F_{\xi_i}^\gamma(\theta) u \right|^p \right] = \sup_{\theta \in \mathcal{D}} \lim_{\delta \rightarrow 0} \sup_{v: |v|=1} \frac{1}{n} \sum_{i=1}^n \frac{|F_{\xi_i}^\gamma(\theta) - F_{\xi_i}^\gamma(\theta + \delta v)|^p}{|\delta|^p}. \quad (4)$$

Equation (4) is especially useful in situations where the computation of $\nabla_\theta F_{\xi_i}^\gamma(\theta)$ is intractable. It allows us perform numerical differentiation by considering a fine-grained mesh around θ in different directions.

Remark 2.2 (Range of γ in linear regression). *Consider the linear regression model*

$$Y_i = X_i^T \theta + \epsilon_i,$$

where $\theta \in \mathbb{R}^d$ is the population parameter vector of interest and $\epsilon_i \in \mathbb{R}$ are i.i.d. random noise independent of $\{X_i\}_{i \geq 1}$. In this setting, the Assumption 2.1 and Assumption 2.2 holds with $\mu = \lambda_{\min}\{\mathbb{E}(X_i X_i^T)\}$ and $N_2 = \sup_{\delta \in \mathbb{R}^d: |\delta|=1} \|X_i X_i^T \delta\|_2^2$, where $\lambda_{\min}\{\cdot\}$ refers to the smallest eigenvalue. Consequently, Theorem 2.2 in Li et al. (2025) ensures $L_p(\gamma) < 1$ as long as

$$0 < \gamma < \frac{2\lambda_{\min}\{\mathbb{E}(X_i X_i^T)\}}{\sup_{\delta \in \mathbb{R}^d: |\delta|=1} \|X_i X_i^T \delta\|_2^2}.$$

It can be demonstrated that this range reaches the optimum in general. By this expression, for $d = 1$, the boundary of γ is $2/3 \approx 0.67$ when X_i follows the standard normal distribution $\mathcal{N}(0, 1)$ and is $10/3 \approx 3.3$ when X_i follows the standard uniform distribution $U[0, 1]$.

Lemma 2.2 allows us to consider a supremum over a single parameter, boosting tractability by eliminating dependence on arbitrary pairs. In lieu of Lemma 2.2, one can approach estimating $L_p(\gamma)$ (and in general $L_p^\ell(\gamma)$) by way of the corresponding empirical versions:

$$\hat{L}_p^n(\gamma) := \sup_{\theta \in \mathcal{D}} \sup_{u: |u|=1} \frac{1}{n} \sum_{i=1}^n \left| \nabla_{\theta} F_{\xi_i}^{\gamma}(\theta) u \right|^p,$$

and in general

$$\hat{L}_p^{\ell,n}(\gamma) := \sup_{\theta \neq \theta' \in \mathcal{D}} \frac{\frac{1}{n} \sum_{i=1}^n |F_{\xi_{i+\ell-1}:\xi_i}^{\gamma}(\theta) - F_{\xi_{i+\ell-1}:\xi_i}^{\gamma}(\theta')|^p}{|\theta - \theta'|^p}, \ell \in \mathbb{N}. \quad (5)$$

Following from Lemma 2.2, we would also like to introduce a similar notion for $L_p^\ell(\gamma)$ and its sample version $\hat{L}_p^{\ell,n}(\gamma)$:

$$\begin{aligned} L_p^\ell(\gamma) &= \sup_{\theta \in \mathcal{D}} \sup_{u: |u|=1} \mathbb{E} \left[\left| \nabla_{\theta} (F_{(i+\ell-1):i}(\theta)) u \right|^p \right] = \sup_{\theta \in \mathcal{D}} \sup_{u: |u|=1} \mathbb{E} \left[\left| \left(\prod_{k=1}^{\ell} \nabla_{\theta} F_{\xi_{i+\ell-k}}^{\gamma}(\theta^{l-k}) \right) u \right|^p \right] \\ \hat{L}_p^{\ell,n}(\gamma) &= \sup_{\theta \in \mathcal{D}} \sup_{u: |u|=1} \frac{1}{n} \sum_{i=1}^n \left| \nabla_{\theta} (F_{(i+\ell-1):i}(\theta)) u \right|^p = \sup_{\theta \in \mathcal{D}} \sup_{u: |u|=1} \frac{1}{n} \sum_{i=1}^n \left| \left(\prod_{k=1}^{\ell} \nabla_{\theta} F_{\xi_{i+\ell-k}}^{\gamma}(\theta^{l-k}) \right) u \right|^p \end{aligned}$$

where $\theta^{(0)} = \theta$ and for $k > 0$, $\theta^{(k)} = F_{\xi_{i+k-1}}^{\gamma}(\theta^{(k-1)})$.

Subsequently, we primarily focus on $L_p(\gamma)$ and its estimator $\hat{L}_p^n(\gamma)$. A naive treatment of the general ℓ -case can be understood to be quite similar; however, we mention another interesting property of the function $\text{MEP-}\ell$ that renders the general case practically trivial after one has considered the $\ell = 1$ scenario. In particular, it follows that the sequence $\{L_p^\ell(\gamma)\}_{\ell \in \mathbb{N}_+}$ is submultiplicative.

Proposition 1. *Set $p \geq 1$ and $\gamma \in \Gamma$, and let $k, \ell \in \mathbb{N}$. Then:*

$$L_p^{\ell+k}(\gamma) \leq L_p^k(\gamma) \cdot L_p^\ell(\gamma).$$

Armed with these additional insights, in the next section we develop an asymptotic theory for $\hat{L}_p^n(\gamma)$.

3 THEORETICAL RESULTS ON MAXIMAL EXPANSION PARAMETERS

Before stating our main results, we collect a set of regularity assumptions that ensure both well-posedness of the optimization problem and tractability of the analysis. Some of these are standard in the study of SGD, but we briefly comment on their roles.

Assumption 3.1 (Compact and convex domains). *The parameters θ and γ are confined to compact convex domains $\mathcal{D} \subset \mathbb{R}^d$ and $\Gamma = [a, b]$, for $b > a > 0$.*

Compactness of the parameter and learning-rate domains is not intrinsic to SGD, but serves as a standard technical device in empirical process theory. It guarantees well-posedness when taking suprema over continuous index sets and facilitates the use of covering arguments and δ -nets. Although unconstrained optimization problems such as linear regression are typically posed on \mathbb{R}^d , in practice SGD iterates remain bounded due to regularization, explicit projection, or simply because divergence leads to algorithmic instability (see, e.g., projection-based variants of SGD in Nemirovski et al. (2009), Lan (2012)). Finally, the assumption that $a \wedge b > 0$ excludes the trivial case $L_p^{\ell,n}(\gamma) = 1$, where the SGD chain does not move at all.

Assumption 3.2 (Lipschitz property). *We assume that there exists a constant $K_p < \infty$ such that the operator norm of $\nabla_{\theta} F_{\xi}^{\gamma}(\theta)$ adheres to the following property with respect to θ and γ :*

$$\mathbb{E} \left[\sup_{(\theta, \gamma, u) \neq (\theta', \gamma', u') \in \mathcal{D} \times \Gamma \times \mathcal{B}} \frac{\left| \left| \nabla_{\theta} F_{\xi_i}^{\gamma}(\theta) u \right|^p - \left| \nabla_{\theta} F_{\xi_i}^{\gamma'}(\theta') u' \right|^p \right|}{(|\theta - \theta'| + |\gamma - \gamma'| + |u - u'|)^p} \right] \leq K_p$$

Bounding higher-order derivatives of the stochastic update map is not a universal assumption, but is a reasonable strengthening of smoothness. In the SGD chain, its contraction dynamics are characterized by its the first derivative of the iterate function. In order to control this derivative, we must bound second order derivative behavior of the function, giving rise Assumption 3.2. Although quite strong, this condition is satisfied by many smooth models of practical interest (e.g. generalized linear models), and rules out only highly irregular loss landscapes.

Assumption 3.3 ($2p$ -moment bound). *Fix $p \geq 1$. Assume*

$$A := \mathbb{E} \left[\sup_{\theta \in \mathcal{D}, \gamma \in \Gamma} \sup_{u: |u|=1} \left| \nabla_{\theta} F_{\xi_i}^{\gamma}(\theta) u \right|^{2p} \right] < \infty.$$

Finite $2p$ -th moments of the stochastic gradients strengthen Assumption 2.1 and are standard when deriving concentration inequalities for SGD. Higher-moment assumptions of this type are routinely employed in empirical process theory (see, e.g., Chernozhukov et al. (2018)) to obtain exponential tail bounds, and they also appear in modern analyses of statistical inference for SGD (Chen et al., 2020). In our setting, this condition ensures that deviation inequalities for the empirical expansion parameter hold with high probability, which is essential for establishing nonasymptotic confidence statements about the edge of stability. While stronger than bounded variance, this requirement remains reasonable in practice for smooth models where gradients have sub-Gaussian or sub-exponential tails.

Assumption 3.4 (Differentiability). *Fix $p \geq 1$. We assume that $\frac{\partial}{\partial \gamma} L_p^{\ell}(\gamma)$ is defined for all $\gamma \in \Gamma$, and that there exists some $K_p > 0$ such that $\sup_{\gamma \in \Gamma} \left| \frac{\partial}{\partial \gamma} L_p^{\ell}(\gamma) \right| \leq K_p$.*

Differentiability of $L_{\ell}^{\gamma}(p)$ with respect to γ ensures that the stability threshold behaves regularly in a neighborhood of the edge. This smoothness enables a first-order expansion around $\gamma_{\ell}(p)$, which is the key step in transferring concentration of $\hat{L}_{\ell,n}^{\gamma}(p)$ into consistency of $\hat{\gamma}_{\ell,n}(p)$.

3.1 ASYMPTOTICS OF MEP- ℓ

In this section, we control the estimation error of $\hat{L}_p^n(\gamma)$ uniformly over $\gamma \in \Gamma$, setting the stage of eventual estimation of the edge-of-stability upon its definition. To that end, we recognize that $\hat{L}_p^n(\gamma) = \sup_{\theta \in \mathcal{D}} \frac{1}{n} \sum_{i=1}^n \left| \nabla_{\theta} F_{\xi_i}^{\gamma}(\theta) \right|^p$ as the \mathcal{L}_{∞} norm of mean of random functions. Subsequently, adapting the tools of Chernozhukov et al. (2018), we provide a general result controlling the partial sums of i.i.d. random functions.

Theorem 3.1. *Let $\Phi \subset \mathbb{R}^d$ be a compact convex set and $X_1(\varphi), \dots, X_n(\varphi)$ be i.i.d. random functions with $X_i : \mathbb{R}^d \mapsto \mathbb{R}^m$ for some $d, m \geq 1$. For $p \geq 1$, denote*

$$K_{\Phi} := \mathbb{E} \left[\sup_{\varphi \neq \varphi' \in \Phi} \frac{|X_i(\varphi) - X_i(\varphi')|^p}{|\varphi - \varphi'|^p} \right] < \infty, \text{ and} \quad (6)$$

$$A_{\Phi,p} := \mathbb{E} \left[\sup_{\varphi \in \Phi} |X_i(\varphi)|^{2p} \right] < \infty. \quad (7)$$

Let the n -th partial sum be defined as $S_{n,p}(\varphi) := \sum_{i=1}^n |X_i(\varphi)|^p$. Then it holds that:

$$\mathbb{E} \left[\sup_{\varphi \in \Phi} |S_{n,p}(\varphi) - \mathbb{E}[S_{n,p}(\varphi)]| \right] = O(\sqrt{n \log n}),$$

where $O(\cdot)$ hides constants solely related to p, d, m, Φ and μ .

Theorem 3.1, to the best of our knowledge, is the *first such result* controlling the moments of sums of random function. As such, it may be of independent interest. Importantly, due to the compactness of Φ , the mean discrepancy between the \mathcal{L}_{∞} norm of empirical and oracle average, decays at the near-parametric rate $O(\sqrt{(\log n)/n})$.

This general result serves as the workhorse for bounding the estimation error of $\hat{L}_p^n(\gamma)$. As an application of Theorem 3.1, we recover the following guarantee on $\hat{L}_p^n(\gamma)$, and more generally $\hat{L}_p^{\ell,n}(\gamma)$.

Theorem 3.2. *Fix $\ell \in \mathbb{N}_+$, and recall $L_p^{\ell}(\gamma)$ and $\hat{L}_p^{\ell,n}(\gamma)$ from Definition 2.1 and (5) respectively. Then, under Assumptions 3.1-3.3, it holds that:*

$$\mathbb{E} \left[\sup_{\gamma \in \Gamma} \left| \hat{L}_p^{\ell,n}(\gamma) - L_p^\ell(\gamma) \right| \right] = O \left(\frac{\log n}{\sqrt{n}} \right).$$

Establishing such guarantees is essential: without quantitative control of the estimation error, any attempt to approximate the edge of stability would remain heuristic. Theorem 3.2 provides precisely this control, paving the way for our eventual goal: precisely estimating edge-of-stability.

4 EDGE OF STABILITY: DEFINITION AND ESTIMATION

In this section, we endeavor to precisely characterize the edge of stability through the explosions of MEP- ℓ . Subsequently, we propose a corresponding version of data-driven edge-of-stability, and provide finite sample error bounds.

Definition 4.1 (EOS- ℓ). Fix $\ell \in \mathbb{N}_+$. The oracle edge-of-stability of lag ℓ (EOS- ℓ) is defined as

$$\gamma_\ell(p) := \inf \{ \gamma > 0 \mid L_p^\ell(\gamma) \geq 1 \}.$$

Clearly, $L_p^\ell(0) = 1$. By recalling Assumption 3.1, $\gamma_\ell(p)$ can be interpreted as smallest $\gamma > 0$ such that the geometric moment contraction no longer holds for the SGD dynamics. As with $L_p^\ell(\gamma)$, we ignore the subscript ℓ whenever $\ell = 1$.

Remark 4.1. It is not yet evident why $\gamma_\ell(p)$ even exists. To ensure its existence, we proceed via the following argument.

1. Recall Theorem 2.2 in Li et al. (2025); under Assumptions 2.1-2.2, there exists a function $\kappa : \mathbb{R}_+ \mapsto \mathbb{R}_+$, such that for $0 < \gamma < \kappa(p)$, we have $L_p(\gamma) < 1$. Here we remark that Li et al. (2025) dealt with the $p > 1$ case; which however can imply the case with $p = 1$ by Hölder’s inequality as shown in Wu & Shao (2004).
2. Since $g(\cdot, \xi) \in C^2$, by the Lebesgue Dominate Convergence Theorem (DCT), $\lim_{\gamma \rightarrow \infty} L_p(\gamma)/\gamma^p = \sup_{\theta} \sup_{u \in \mathbb{R}^d: |u|=1} \mathbb{E}[\|\nabla_{\theta}^2 g(\theta, \xi)u\|^p]$.

Conditions 1 and 2 above ensure that $\gamma_\ell(p) \in \Gamma$ exists.

Definition of the empirical version of $\gamma_\ell(p)$, denoted by $\hat{\gamma}_{\ell,n}(p)$, is not straight-forward, since the guarantees in Li et al. (2025) extend only to $L_p^\ell(\gamma)$, and not to its empirical version. However, Theorem 3.2 ensures that for all $\gamma \in \Gamma$ for any compact set Γ , $L_p^\ell(\gamma)$ is closely approximated by its empirical version $\hat{L}_p^{\ell,n}(\gamma)$. Therefore, it is conceivable to leverage Theorem 3.2 to obtain a precisely-defined compact convex set Γ , such that, with high probability, $\hat{L}_p^{\ell,n}(\gamma)$ crosses 1 on $\text{int}(\Gamma)$. More formally, by Theorem 2.2 in Li et al. (2025) and continuity of $L_p^\ell(\cdot) < 1$, there exists some $\delta > 0$ and $\gamma_0 > 0$ such that $L_p^\ell(\gamma) < 1$ for all $\gamma \in B_\delta(\gamma_0)$. On the other hand, let $\gamma_\ell^\dagger(p) := \inf \{ \gamma > 0 \mid L_p^\ell(\gamma) > 2 \}$. Similar to Remark 4.1, $\gamma_\ell^\dagger(p)$ is well-defined. Then we proceed to define the edge-of-stability at lag ℓ .

Definition 4.2. Denote $\Gamma := [\gamma_0, \gamma_\ell^\dagger(p)]$. The oracle EOS- ℓ is defined as

$$\hat{\gamma}_{\ell,n}(p) := \min \left\{ \gamma \in \Gamma \mid \hat{L}_p^{\ell,n}(\gamma) \geq 1 \right\}.$$

Note that, by definition of Γ , it also follows that $\gamma_\ell(p) \in \text{int}(\Gamma)$. We establish the following conventions for the edge-cases: if $\sup_{\gamma \in \Gamma} \hat{L}_p^{\ell,n}(\gamma) < 1$, then $\hat{\gamma}_{\ell,n}(p) = 0$; on the other hand, if $\inf_{\gamma > 0} \hat{L}_p^{\ell,n}(\gamma) = 1$, then $\hat{\gamma}_{\ell,n}(p) = \infty$. In fact, in the following we prove that these edge cases have vanishing probability, and consequently, we recover the asymptotic consistency of $\gamma_{\ell,n}(p)$ as an estimator of $\gamma_\ell(p)$.

Theorem 4.3. Fix $\ell \in \mathbb{N}_+$, and recall $\gamma_\ell(p)$ and $\hat{\gamma}_{\ell,n}(p)$ from Definitions 4.1 and 4.2 respectively. Then, under Assumptions 3.1-3.3,

$$\mathbb{P}(\hat{\gamma}_{\ell,n}(p) \in \text{int}(\Gamma)) \rightarrow 1 \text{ as } n \rightarrow \infty.$$

Additionally, it holds that:

$$|\hat{\gamma}_{\ell,n}(p) - \gamma_\ell(p)| = O_{\mathbb{P}} \left(\frac{\log n}{\sqrt{n}} \right).$$

To the best of our knowledge, Theorem 4.3 provides the *only*, provably consistent estimator of $\text{EOS-}\ell$ in the context of SGD. Beyond theoretical interest, the practical relevance of estimator cannot be overstated; $\hat{\gamma}_{\ell,n}(p)$ indicates a data-driven threshold of the learning rate, beyond which the SGD dynamics explode with high probability.

5 SIMULATION

In this section, we empirically characterize the edge-of-stability region, and assess its optimality. Across a suite of synthetic settings (linear and expectile regression with varying dimension, lag, and data distributions), we estimate the contraction ratio $L_p(\gamma)^{1/p}$ as a function of γ and identify the smallest γ at which contraction fails. The resulting empirical boundary closely matches our theoretical prediction, demonstrating that the proposed “edge of stability” is tight. Taken together, these results validate the theory and provide actionable guidance for selecting constant step sizes that guarantee convergence in practice.

We first demonstrate our result focusing on the following data generating mechanism:

$$Y_i = X_i^\top \theta^* + \epsilon_i,$$

and let $\xi_i = \{X_i, y_i\}_{i \in \mathbb{N}_+}$ denote the observed sequential data and θ^* is the unknown population parameter of interest. We study two convex models: (i) linear regression with squared loss

$$G_1(\theta) = \mathbb{E}_{\xi_i = (X_i, y_i) \sim \Pi_2} (X_i^\top \theta - y_i)^2 / 2,$$

where $F_{\xi_i}^\gamma(\theta)$ takes the following form:

$$F_{\xi_i}^\gamma(\theta) = \theta - \gamma X_i (X_i^\top \theta - y_i),$$

and (ii) expectile regression with the asymmetric least-square loss

$$G_2(\theta) = \mathbb{E}_{\xi_i = (X_i, y_i) \sim \Pi_2} |w - 1_{\{X_i^\top \theta - y_i > 0\}}| (X_i^\top \theta - y_i)^2 / 2, \quad \text{with weight } w \in (0, 1),$$

and corresponding $F_{\xi_i}^\gamma(\theta)$ is given by

$$F_{\xi_i}^\gamma(\theta) = \theta - |w - 1_{\{X_i^\top \theta - y_i > 0\}}| X_i (X_i^\top \theta - y_i).$$

The feature vector $X \in \mathbb{R}^d$ is drawn either from a Gaussian design $\mathcal{N}(0, I_d)$ or a product Uniform design $\text{Unif}([0, 1]^d)$ and the noise ξ is drawn from standard Gaussian distribution, independent of $\{X_i\}_{i \in \mathbb{N}}$. We vary the ambient dimension $d \in \{1, 2, 3, 5, 10\}$, the composition lag $\ell \in \{1, 5, 10\}$, and the moment index $p \in \{2, 4\}$. For linear regression we sweep γ on a grid $\Gamma_{\text{norm}} = \{0.01, 0.02, \dots, 1.00\}$ under the Gaussian design, and for the Uniform design we use $\Gamma_{\text{unif}} = \{0.01, 0.02, \dots, 4.00\}$ to account for the different curvature scales observed in practice.

Across all configurations, the mapping $\gamma \mapsto L_p(\gamma)^{1/p}$ exhibits a pronounced elbow shape, where the estimated ratio initially declines from 1, reaches a minimum, and then reverses, crossing 1 at the stability edge; beyond the crossing it grows rapidly, ultimately diverging. The transition occurs well within the plotted range, so the edge $\hat{\gamma}$ is visually stable and can be localized to a narrow interval.

In Figure 2 and Figure 4 (in the Appendix §B), subplots(a) demonstrate, that increasing the moment index $p \in \{1, 2, 3, 5, 10\}$ shifts the crossing leftward while keeping the minimum shallow. This indicates that heavier emphasis on tail deviations tightens the admissible step-size, which aligns with the result proposed in Li et al. (2025). Varying the lag ℓ in subplots(b) of Figure 2 and Figure 4 primarily enlarge the edge of stability as lag ℓ increases: for any fixed γ , sub-multiplicative gives $L_p^\ell(\gamma) \leq L_p(\gamma)^p$, so increasing ℓ pushes ratios further below on the stable side and further above 1 on the unstable side. Thus the increase of ℓ allows larger γ to ensure the contraction. The dimensional study in subplots(c) of Figure 2 and Figure 4 shows the early contraction of the stable region as d increases. In addition, the empirical edge $\hat{\gamma}_{\ell,n}(p)$ extracted at the yellow curve in subplots(a) and red curve in subplots(b) closely matches the theoretical boundary proposed in Remark 2.2 for $d = 1$ and $\ell = 1$ case, where $\hat{\gamma}_{\ell,n}(p) = \frac{2}{3}$ for $X_i \sim \mathcal{N}(0, 1)$ and $\hat{\gamma}_{\ell,n}(p) = \frac{10}{3}$ for $X_i \sim \text{Unif}[0, 1]$. As a conclusion, the results displayed in Figure 2 and Figure 4 validate that the stability set $\gamma : L_p(\gamma) < 1$ is a single interval starting at 0, its boundary is accurately captured by the unique intersection with level 1, and its dependence on p , ℓ , and d follows the theoretical predictions.

Figure 3 shows that expectile regression mirrors the linear case: the edge of stability (the unique crossing of $L_p(\gamma)^{1/p}$ with level 1) decreases as the moment index p increases and increases as the lag ℓ grows. The first trend follows the

p -sensitivity of the contraction metric via Hölder’s inequality. The second follows from the sub-multiplicativity of the maximal expansion parameter (MEP), $L_{p,\ell+k}(\gamma) \leq L_{p,\ell}(\gamma)L_{p,k}(\gamma)$, which strengthens contraction on the stable side and steepens growth on the unstable side with right shifting the crossing in γ . The same qualitative dependencies appear for expectile regression for dimension d : the edge moves left as d grows (a smaller stable γ). Taken together, these curves confirm that the qualitative and quantitative dependence of the stability edge on p, ℓ and d persists beyond squared loss.

6 CONCLUSIONS AND DISCUSSION

In this work, we provided a principled characterization of the stability region of SGD with constant learning rates. By introducing the notion of the maximal expansion parameter and connecting it to Lyapunov exponents, we established a rigorous definition of the edge-of-stability and developed a consistent, data-driven estimator for identifying admissible learning rates. Our theoretical results, complemented by extensive simulations on linear and expectile regression, confirm that the proposed framework accurately captures the transition from stable to unstable regimes. These findings supply both a theoretical foundation and a practical tool for selecting constant step sizes in online learning algorithms.

Looking ahead, the observed dependence of stability thresholds on factors such as dimension, lag, and moment index underscores the importance of adaptive, data-driven tuning strategies, rather than relying on fixed heuristics. Moreover, by situating SGD stability with the Lyapunov exponent in dynamic systems, our work lays the groundwork for unifying deterministic and stochastic stability analyses, potentially leading to sharper guidelines for learning rate selection across a broad range of optimization problems.

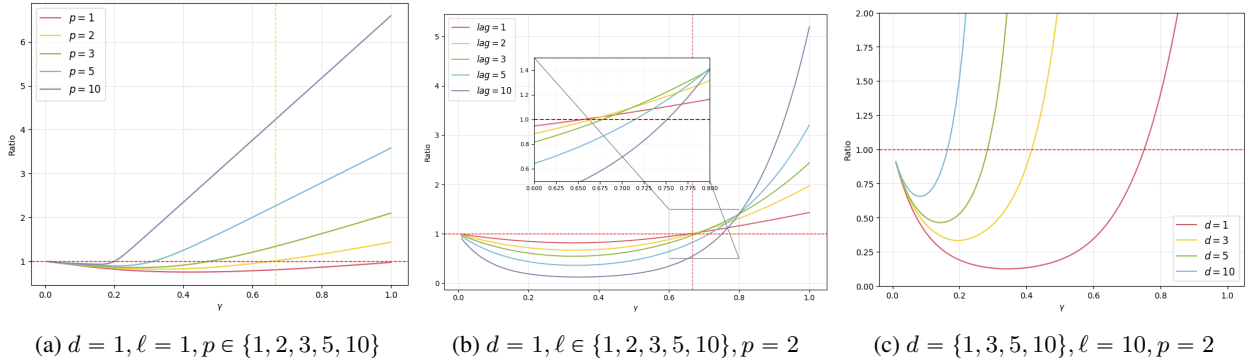


Figure 2: **Linear regression with $X_i \sim \mathcal{N}(0, I_d)$.** Each panel plots $\hat{L}_p^\ell(\gamma)^{1/p}$ versus the constant step size γ for linear regression. Experimental factors and grids follow the setup marked in subplot labels.

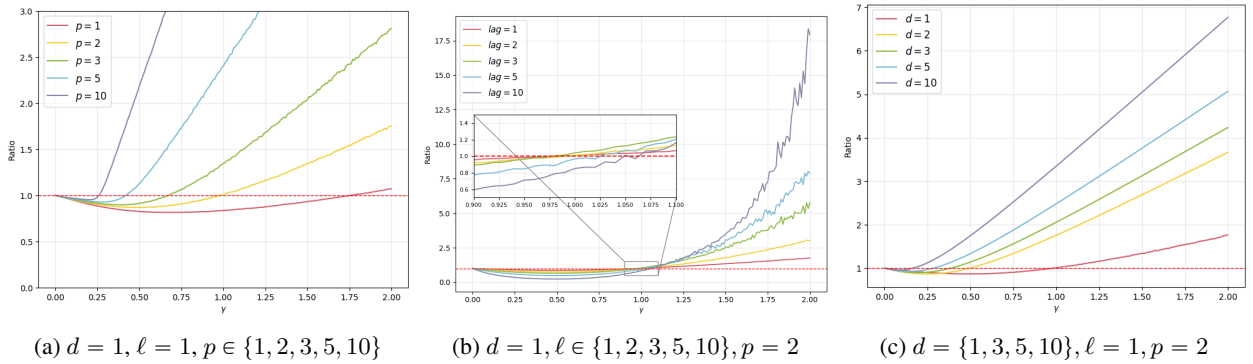


Figure 3: **Expectile regression with $X_i \sim \mathcal{N}\{0, 1\}$ and weight $\omega = 0.2$.** Each panel plots $\hat{L}_p^\ell(\gamma)^{1/p}$ versus the constant step size γ for expectile regression and is averaged over 30 experiments. Experimental factors and grids follow the setup marked in subplot labels.

ETHICS STATEMENT

The research follows all ethical guidelines. No human data or ethically sensitive content is involved. All potential limitations and justifications are adequately addressed. We do not anticipate any negative impacts, and as such the paper does not include a dedicated speculative discussion of broader societal impacts.

REPRODUCIBILITY STATEMENT

All the relevant reproducible codes can be found in the anonymous Github repository. All the theoretical results and assumptions are rigorously proved and validated in the Appendix §C-§E.

AUTHOR CONTRIBUTIONS

All the authors contributed equally to this research.

REFERENCES

- A. Agarwala and J. Pennington. High dimensional analysis reveals conservative sharpening and a stochastic edge of stability, 2025. URL <https://arxiv.org/abs/2404.19261>.
- Kwangjun Ahn, Sebastien Bubeck, Sinho Chewi, Yin Tat Lee, Felipe Suarez, and Yi Zhang. Learning threshold neurons via edge of stability. In *Advances in Neural Information Processing Systems*, volume 36, pp. 19540–19569, 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/3e592c571de69a43d7a870ea89c7e33a-Paper-Conference.pdf.
- Arseniy Andreyev and Pierfrancesco Beneventano. Edge of Stochastic Stability: Revisiting the Edge of Stability for SGD, 2025. URL <https://arxiv.org/abs/2412.20553>.
- Ludwig Arnold. *Random Dynamical Systems*. Springer, Berlin, 1998.
- S. Arora, Z. Li, and A. Panigrahi. Understanding Gradient Descent on the Edge of Stability in Deep Learning. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162, pp. 948–1024, 2022.
- F. Bach and E. Moulines. Non-strongly-convex smooth stochastic approximation with convergence rate $O(1/n)$. In *Advances in Neural Information Processing Systems*, pp. 773–781, 2013.
- D. Barrett and B. Dherin. Implicit Gradient Regularization. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=3q5IqUrkcF>.
- P. L. Bartlett, P. M. Long, G. Lugosi, and A. Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- L. Bottou, F. E. Curtis, and J. Nocedal. Optimization methods for large-scale machine learning. *SIAM Review*, 60(2): 223–311, 2018a.
- Léon Bottou, Frank E. Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. *SIAM Rev.*, 60(2):223–311, 2018b. ISSN 0036-1445,1095-7200. doi: 10.1137/16M1080173. URL <https://doi.org/10.1137/16M1080173>.
- Philippe Bougerol and Jean Lacroix. *Products of Random Matrices with Applications to Schrödinger Operators*. Birkhäuser, Boston, 1985.
- L. Chen and J. Bruna. Beyond the edge of stability via two-step gradient updates. In *Proceedings of the 40th International Conference on Machine Learning*, 2023.
- Xi Chen, Jason D. Lee, Xin T. Tong, and Yichen Zhang. Statistical inference for model parameters in stochastic gradient descent. *Ann. Statist.*, 48(1):251–273, 2020. ISSN 0090-5364,2168-8966. doi: 10.1214/18-AOS1801. URL <https://doi.org/10.1214/18-AOS1801>.

- V. Chernozhukov, D. Chetverikov, and K. Kato. Inference on Causal and Structural Parameters Using Many Moment Inequalities. *arXiv:1312.7614v6*, 2018.
- J. Cohen, S. Kaur, Y. Li, J. Z. Kolter, and A. Talwalkar. Gradient Descent on Neural Networks Typically Occurs at the Edge of Stability. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=jh-rTtvkGeM>.
- J. M. Cohen, B. Ghorbani, S. Krishnan, N. Agarwal, S. Medapati, M. Badura, D. Suo, D. Cardoze, Z. Nado, G. E. Dahl, and J. Gilmer. Adaptive Gradient Methods at the Edge of Stability, 2024. URL <https://arxiv.org/abs/2207.14484>.
- A. Damian, E. Nichani, and J. D. Lee. Self-Stabilization: The Implicit Bias of Gradient Descent at the Edge of Stability. In *ICLR*, 2023. URL <https://openreview.net/forum?id=nhKHA59gXz>.
- H. Daneshmand, N. Le Roux, and A. M. Bronstein. Stochastic Gradient Descent as a Dynamical System: Almost Sure Convergence Analysis. *Journal of Machine Learning Research*, 25(1):1–35, 2024.
- S. Dereich, R. Graeber, A. Jentzen, and A. Riekert. Asymptotic stability properties and a priori bounds for Adam and other gradient descent optimization methods, 2025. URL <https://arxiv.org/abs/2509.10476>.
- M. Even, S. Pesme, S. Gunasekar, and N. Flammarion. (S)GD over Diagonal Linear Networks: Implicit bias, Large Stepsizes and Edge of Stability. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 29406–29448. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/5da6ce80e97671b70c01a2e703b868b3-Paper-Conference.pdf.
- A. Ghosh, S. M. Kwon, R. Wang, S. Ravishankar, and Q. Qu. Learning Dynamics of Deep Matrix Factorization Beyond the Edge of Stability. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=J4Dvxv7WnG>.
- M. Hardt, B. Recht, and Y. Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International Conference on Machine Learning*, pp. 1225–1234, 2016.
- P. Jain, S. M. Kakade, R. Kidambi, P. Netrapalli, and A. Sidford. Parallelizing stochastic gradient descent for least squares regression: Mini-batching, averaging, and model misspecification. *Journal of Machine Learning Research*, 19(1):1–42, 2018.
- Guanghui Lan. An optimal method for stochastic composite optimization. *Mathematical Programming*, 133(1):365–397, 2012.
- J. Li, Z. Lou, S. Richter, and W. B. Wu. The Stochastic Gradient Descent from a Nonlinear Time Series Perspective. *Manuscript*, TBD(TBD):TBD, 2025.
- Jiaqi Li, Zhipeng Lou, Stephan Richter, and Wei-Biao Wu. The stochastic gradient descent from a nonlinear time series perspective. *preprint*, 2024.
- R. Li, J. Chen, and T. Liang. How does Batch Normalization Help Optimization in Deep Learning: A Look from Stability Analysis. *Transactions on Machine Learning Research*, 2023.
- C. Liu, L. Zhu, and M. Belkin. Loss landscapes and optimization in over-parameterized non-linear systems and neural networks. *Foundations and Trends in Machine Learning*, 15(1):1–159, 2022.
- Liming Liu, Zixuan Zhang, Simon Du, and Tuo Zhao. A Minimalist Example of Edge-of-Stability and Progressive Sharpening, 2025. URL <https://arxiv.org/abs/2503.02809>.
- Philip M. Long and Peter L. Bartlett. Sharpness-aware minimization and the edge of stability. *J. Mach. Learn. Res.*, 25(1), 2024. ISSN 1532-4435.
- A. Ly and P. Gong. Optimization on multifractal loss landscapes explains a diverse range of geometrical and dynamical properties of deep learning. *Nature Communications*, 16, 04 2025. doi: 10.1038/s41467-025-58532-9.
- S. Mandt, M. D. Hoffman, and D. M. Blei. Stochastic gradient descent as approximate Bayesian inference. *Journal of Machine Learning Research*, 18(1):4873–4907, 2017.

- Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–1609, 2009.
- V. I. Oseledets. A multiplicative ergodic theorem. characteristic Lyapunov exponents of dynamical systems. *Trans. Moscow Math. Soc.*, 19:197–231, 1968.
- Boris T. Polyak and Anatoli B. Juditsky. Acceleration of stochastic approximation by averaging. *SIAM Journal on Control and Optimization*, 30(4):838–855, 1992.
- V. Roulet, A. Agarwala, J.-B. Grill, G. Swirszcz, M. Blondel, and F. Pedregosa. Stepping on the edge: Curvature aware learning rate tuners. In *Advances in Neural Information Processing Systems*, volume 37, pp. 47708–47740, 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/555479a201da27c97aaeed842d16ca49-Paper-Conference.pdf.
- David Ruppert. Efficient estimations from a slowly convergent robbins-monro process. *Technical Report*, 1988.
- U. Simsekli, L. Sagun, and M. Gurbuzbalaban. A tail-index analysis of stochastic gradient noise in deep neural networks. In *International Conference on Machine Learning*, pp. 5834–5843, 2019.
- M. Song and C. Yun. Trajectory Alignment: Understanding the Edge of Stability Phenomenon via Bifurcation Theory. In *Advances in Neural Information Processing Systems*, volume 36, pp. 71632–71682, 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/e2a9256bd816ab9e082dfaa22f1f62a2-Paper-Conference.pdf.
- Z. Wang, Z. Li, and J. Li. Analyzing sharpness along GD trajectory: Progressive sharpening and edge of stability. In *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=thgItcQrJ4y>.
- Ziyang Wei, Wanrong Zhu, and Wei Biao Wu. Weighted averaged stochastic gradient descent: Asymptotic normality and optimality. *arXiv preprint arXiv:2307.06915*, 2023.
- J. Wu, V. Braverman, and J. D. Lee. Implicit Bias of Gradient Descent for Logistic Regression at the Edge of Stability. In *Advances in Neural Information Processing Systems*, volume 36, pp. 74229–74256, 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/eb189151ced0ff808abafdl6a51fec92-Paper-Conference.pdf.
- L. Wu, C. Ma, and E. Weinan. How SGD selects the global minima in over-parameterized learning: A dynamical stability perspective. In *Advances in Neural Information Processing Systems*, pp. 8279–8288, 2018.
- Wei Biao Wu. Nonlinear system theory: another look at dependence. *Proc. Natl. Acad. Sci. USA*, 102(40):14150–14154, 2005. ISSN 0027-8424,1091-6490. doi: 10.1073/pnas.0506715102. URL <https://doi.org/10.1073/pnas.0506715102>.
- Wei Biao Wu and Xiaofeng Shao. Limit theorems for iterated random functions. *J. Appl. Probab.*, 41(2):425–436, 2004. ISSN 0021-9002,1475-6072. doi: 10.1239/jap/1082999076. URL <https://doi.org/10.1239/jap/1082999076>.
- Y. Yang, Z. Hao, X. Zhang, and W. Wang. A Statistical Perspective on Adaptive Learning Rates in Deep Learning. In *International Conference on Learning Representations*, 2023.
- S. Zeng and Y. Lei. Stability and Generalization Analysis of Decentralized SGD: Sharper Bounds Beyond Lipschitzness and Smoothness. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=g4eTrS2U8o>.
- C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 65(6):107–115, 2022.
- J. Zhang, S. P. Karimireddy, A. Veit, S. Kim, S. J. Reddi, S. Kumar, and S. Sra. Why ADAM beats SGD for attention models. In *Advances in Neural Information Processing Systems*, volume 33, pp. 22863–22876, 2020.
- Ruiqi Zhang, Jingfeng Wu, Licong Lin, and Peter L. Bartlett. Minimax Optimal Convergence of Gradient Descent in Logistic Regression via Large and Adaptive Stepsizes, 2025. URL <https://arxiv.org/abs/2504.04105>.

- Wanrong Zhu, Xi Chen, and Wei Biao Wu. Online covariance matrix estimation in stochastic gradient descent. *J. Amer. Statist. Assoc.*, 118(541):393–404, 2023. ISSN 0162-1459,1537-274X. doi: 10.1080/01621459.2021.1933498. URL <https://doi.org/10.1080/01621459.2021.1933498>.
- X. Zhu, Z. Wang, X. Wang, M. Zhou, and R. Ge. Understanding Edge-of-Stability Training Dynamics with a Minimalist Example. *ArXiv*, abs/2210.03294, 2022. URL <https://api.semanticscholar.org/CorpusID:252762329>.
- D. Zou, Y. Cao, D. Zhou, and Q. Gu. Gradient descent optimizes over-parameterized deep ReLU networks. *Machine Learning*, 110:247–284, 2021.

Appendix

This appendix is devoted to additional discussion, collection of mathematical arguments and additional simulation results. In particular, in §A we discuss some other approaches to edge-of-stability analysis, as well as the existing gaps in the literature. In §C-§E, we provide detailed proofs to our theoretical results.

A RELATED LITERATURE

The edge of stability phenomenon was first systematically identified by Cohen et al. (2021) in the context of neural networks, showing empirically that GD trajectories typically operate at the stability threshold. Subsequent work such as Ahn et al. (2023) extended these insights to simple neuron models, establishing that edge of stability behavior arises even in minimal architectures. These contributions built on a long line of optimization analyses (Bottou et al., 2018a; Bach & Moulines, 2013) that emphasized the importance of step-size selection and convergence guarantees.

Several papers seek to isolate the mechanisms behind the edge of stability using simplified or tractable models. Arora et al. (2022) developed a theoretical framework for GD at the edge of stability, while Zhu et al. (2022) and Liu et al. (2025) employed minimalist examples to clarify the core dynamics. Variants such as diagonal linear networks (Even et al., 2023) and two-step updates (Chen & Bruna, 2023) further illuminate how the phenomenon arises across different formulations. Parallel lines of work have also explored how normalization or regularization mechanisms affect optimization stability, e.g. Li et al. (2023) on batch normalization and Barrett & Dherin (2021) on implicit gradient regularization.

Another strand of work interprets the edge of stability through the geometry of the loss landscape. Progressive sharpening along training trajectories was analyzed by Wang et al. (2022), while Song & Yun (2023) provided a bifurcation-theoretic view. More recent work has refined these ideas via high-dimensional analysis (Agarwala & Pennington, 2025), sharpness-aware methods (Long & Bartlett, 2024), and curvature-aware learning-rate tuning (Roulet et al., 2024). These developments resonate with broader optimization perspectives on adaptive learning rates (Yang et al., 2023) and comparisons of adaptive methods with SGD (Zhang et al., 2020).

Beyond stability, the edge of stability has been connected to implicit bias and generalization. For instance, Wu et al. (2023) and Damian et al. (2023) study logistic regression at the edge of stability, highlighting the implicit regularization induced by GD. Related work considers minimax optimal convergence (Zhang et al., 2025) and generalization in decentralized SGD settings (Zeng & Lei, 2025). This complements a broader literature on benign overfitting and generalization in over-parameterized models (Bartlett et al., 2020; Zhang et al., 2022; Zou et al., 2021; Liu et al., 2022), where stability considerations play a central role.

While the majority of results concern deterministic GD, several papers have begun exploring extensions. Andreyev & Beneventano (2025) revisited the notion of stability under stochastic gradient descent, whereas Cohen et al. (2024) and Dereich et al. (2025) examined adaptive and Adam-type methods, respectively. Other directions extend edge of stability analysis to deep linear networks (Ghosh et al., 2025) and multi-fractal loss landscapes (Ly & Gong, 2025). These developments connect naturally to classical work on stochastic approximation (Wu et al., 2018; Jain et al., 2018) and continue the trend of relating stochastic dynamics to stability properties.

Despite this growing body of work, the focus has remained predominantly on GD. By contrast, our work develops a systematic analysis of edge of stability in the context of *stochastic gradient descent*, providing a sharper understanding of how stochasticity modifies, stabilizes, or destabilizes the classical GD picture. In this way, we broaden the scope of the edge of stability framework to settings of practical relevance.

B ADDITIONAL SIMULATION RESULTS

In addition to the plot we have in the main context, we report here the simulation results for the linear regression with uniformly distributed random samples. The discussion can be referred to in the main text Section 5.

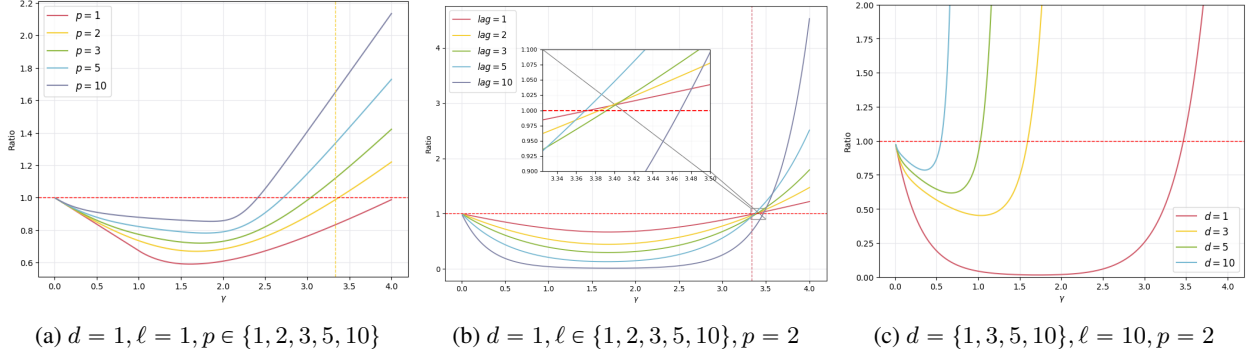


Figure 4: **Linear regression with $X_i \sim \text{Unif}([0, 1]^d)$.** Each panel plots $\hat{L}_p^\ell(\gamma)^{1/p}$ versus the constant step size γ for linear regression. Experimental factors and grids follow the setup marked in subplot labels.

C PROOFS OF §2

C.1 PROOF OF LEMMA 2.2

Proof. W.l.o.g., we consider the case $\ell = 1$; the case for general $\ell \in \mathbb{N}$ is similar. Note that since $F_\xi^\gamma(\cdot) \in C^1$, it follows that

$$\sup_{\theta \in \mathcal{D}} \sup_{u: |u|=1} \frac{1}{n} \sum_{i=1}^n |\nabla_\theta F_{\xi_i}^\gamma(\theta) u|^p \leq \sup_{\theta \in \mathcal{D}} \frac{1}{n} \sum_{i=1}^n \limsup_{\theta' \rightarrow \theta} \frac{|F_{\xi_i}^\gamma(\theta) - F_{\xi_i}^\gamma(\theta')|^p}{|\theta - \theta'|^p} \leq \sup_{\theta \neq \theta' \in \mathcal{D}} \frac{1}{n} \sum_{i=1}^n \frac{|F_{\xi_i}^\gamma(\theta) - F_{\xi_i}^\gamma(\theta')|^p}{|\theta - \theta'|^p}. \quad (8)$$

On the other hand, by Jensen's inequality and the convexity of \mathcal{D} ,

$$\begin{aligned} \sup_{\theta \neq \theta' \in \mathcal{D}} \frac{1}{n} \sum_{i=1}^n \frac{|F_{\xi_i}^\gamma(\theta) - F_{\xi_i}^\gamma(\theta')|^p}{|\theta - \theta'|^p} &= \sup_{\theta \neq \theta' \in \mathcal{D}} \frac{1}{n} \sum_{i=1}^n \frac{\left| \int_0^1 \frac{\partial}{\partial t} F_{\xi_i}^\gamma(\theta' + t(\theta - \theta')) dt \right|^p}{|\theta - \theta'|^p} \\ &\leq \sup_{\theta \neq \theta' \in \mathcal{D}} \frac{1}{n} \sum_{i=1}^n \frac{\int_0^1 \left| \frac{\partial}{\partial t} F_{\xi_i}^\gamma(\theta' + t(\theta - \theta')) \right|^p dt}{|\theta - \theta'|^p} \\ &= \sup_{\theta \neq \theta' \in \mathcal{D}} \frac{1}{n} \sum_{i=1}^n \frac{\int_0^1 \left| \nabla_\theta F_{\xi_i}^\gamma(\theta' + t(\theta - \theta')) (\theta - \theta') \right|^p dt}{|\theta - \theta'|^p} \\ &\leq \sup_{\theta \neq \theta' \in \mathcal{D}, t \in [0, 1]} \sup_{u: |u|=1} \frac{1}{n} \sum_{i=1}^n \left| \nabla_\theta F_{\xi_i}^\gamma(\theta' + t(\theta - \theta')) u \right|^p \\ &\leq \sup_{\theta \in \mathcal{D}} \sup_{u: |u|=1} \frac{1}{n} \sum_{i=1}^n \left| \nabla_\theta F_{\xi_i}^\gamma(\theta) u \right|^p. \end{aligned} \quad (9)$$

Equations (8) and (9) jointly conclude the proof of (3). In lieu of $\sup_{\theta \neq \theta' \in \mathcal{D}} \mathbb{E} \left[\frac{|F_{\xi_i}^\gamma(\theta) - F_{\xi_i}^\gamma(\theta')|^p}{|\theta - \theta'|^p} \right] < \infty$ from Assumption 3.1, Dominated Convergence Theorem entails (2.2).

□

C.2 PROOF OF PROPOSITION 1

Proof. We denote $H_\ell(\theta) := F_{i+\ell-1:i}(\theta)$ and $\mathcal{F}_{\ell-1} := \sigma(\xi_i, \dots, \xi_{i+\ell-1})$. Then:

$$\begin{aligned}
L_p^{\ell+k}(\gamma) &= \sup_{\theta \neq \theta' \in \mathcal{D}} \frac{\mathbb{E} \left[|F_{i+\ell+k-1:i}(\theta) - F_{i+\ell+k-1:i}(\theta')|^p \right]}{|\theta - \theta'|^p} = \sup_{\theta \neq \theta' \in \mathcal{D}} \mathbb{E} \left[\frac{|F_{i+\ell+k-1:i}(\theta) - F_{i+\ell+k-1:i}(\theta')|^p}{|\theta - \theta'|^p} \right] \\
&= \sup_{\theta \neq \theta' \in \mathcal{D}} \mathbb{E} \left[\frac{|F_{i+\ell+k-1:i+\ell}(H_\ell(\theta)) - F_{i+\ell+k-1:i+\ell}(H_\ell(\theta'))|^p}{|H_\ell(\theta) - H_\ell(\theta')|^p} \cdot \frac{|H_\ell(\theta) - H_\ell(\theta')|^p}{|\theta - \theta'|^p} \right] \\
&= \sup_{\theta \neq \theta' \in \mathcal{D}} \mathbb{E} \left[\mathbb{E} \left[\frac{|F_{i+\ell+k-1:i+\ell}(H_\ell(\theta)) - F_{i+\ell+k-1:i+\ell}(H_\ell(\theta'))|^p}{|H_\ell(\theta) - H_\ell(\theta')|^p} \cdot \frac{|H_\ell(\theta) - H_\ell(\theta')|^p}{|\theta - \theta'|^p} \mid \mathcal{F}_{\ell-1} \right] \right] \\
&= \sup_{\theta \neq \theta' \in \mathcal{D}} \mathbb{E} \left[\mathbb{E} \left[\frac{|F_{i+\ell+k-1:i+\ell}(H_\ell(\theta)) - F_{i+\ell+k-1:i+\ell}(H_\ell(\theta'))|^p}{|H_\ell(\theta) - H_\ell(\theta')|^p} \mid \mathcal{F}_{\ell-1} \right] \cdot \frac{|H_\ell(\theta) - H_\ell(\theta')|^p}{|\theta - \theta'|^p} \right].
\end{aligned}$$

Conditionally on $\mathcal{F}_{\ell-1}$, $F_{i+\ell+k-1:i+\ell}$ is driven by k new i. i. d. innovations which are independent of $\mathcal{F}_{\ell-1}$. Therefore we deduce that:

$$\mathbb{E} \left[\frac{|F_{i+\ell+k-1:i+\ell}(H_\ell(\theta)) - F_{i+\ell+k-1:i+\ell}(H_\ell(\theta'))|^p}{|H_\ell(\theta) - H_\ell(\theta')|^p} \mid \mathcal{F}_{\ell-1} \right] \leq L_p^k(\gamma).$$

Therefore:

$$L_p^{\ell+k}(\gamma) \leq \sup_{\theta \neq \theta' \in \mathcal{D}} \mathbb{E} \left[L_p^k(\gamma) \cdot \frac{|H_\ell(\theta) - H_\ell(\theta')|^p}{|\theta - \theta'|^p} \right] = L_p^k(\gamma) \cdot \sup_{\theta \neq \theta' \in \mathcal{D}} \frac{|H_\ell(\theta) - H_\ell(\theta')|^p}{|\theta - \theta'|^p} = L_p^k(\gamma) \cdot L_p^\ell(\gamma).$$

□

D PROOFS OF §3

Before we proceed to the key arguments behind the theoretical results of §3, it is instrumental to introduce a two key result that serves as the backbone of our arguments. This result originate from Chernozhukov et al. (2018), and serves as sharp probabilistic controls on the fluctuations of empirical sums indexed by high-dimensional parameter sets. We restate it here in a form adapted to our setting.

Lemma D.1. *Let $X_1, \dots, X_n \in \mathbb{R}^p$ be independent random vectors with $p \geq 2$. Define $M := \max_{1 \leq i \leq n, 1 \leq j \leq p} |X_{ij}|$ and $\sigma^2 := \max_{1 \leq j \leq p} \sum_{i=1}^n \mathbb{E}[X_{ij}^2]$. Then:*

$$\mathbb{E} \left[\max_{1 \leq j \leq p} \left| \sum_{i=1}^n (X_{ij} - \mathbb{E}[X_{ij}]) \right| \right] \leq K(\sigma \sqrt{\log p} + \sqrt{\mathbb{E}[M^2]} \log p),$$

where $K > 0$ is a universal constant.

This lemma complements the previous one by providing an expectation bound for the same maximal deviation and quantifies the typical size of the deviation, showing that it scales as $O(\sqrt{\log p})$ up to constants depending on variance and maximal moments. In summary, it provides the empirical process tools that underpin our general moment bound in Theorem 3.1. We note that the for the sake of brevity, the results are proved for $\ell = 1$; the general ℓ -cases follow by a simple conditional argument akin to Proposition 1.

D.1 PROOF OF THEOREM 3.1

The key idea of Theorem 3.1 is to discretize the set Φ with suitably selected grid, before applying Lemma D.1 to control the deviations of functions evaluated on those grid-points. This grid is carefully chosen to have appropriate packing radius, that allows us to move seamlessly into the compact set Φ while maintaining the rate derived on the grid-points. We formalize this idea through a novel technique leveraging ε -nets.

Proof. Let $N := n^c$ for some $c > p/2$. For a given $\varphi \in \Phi$, we denote $[\varphi]_N := \frac{1}{N} ([N\varphi^1], \dots, [N\varphi^d])$, with φ^k being the k th coordinate of φ . Then, by compactness and convexity of Φ , $\mathcal{N} := \{[\varphi]_N \mid \varphi \in \Phi\}$ is a δ_n -net for Φ , where $\delta := \delta_n \leq L_\Phi n^{-c}$ for some constant $L_\Phi > 0$ that depends only on Φ . Enumerate its elements as $\{\varphi_1, \dots, \varphi_J\}$ and observe $J \leq L_\Phi \cdot N^d$. Recall $A_{\Phi,p}$ defined in Theorem 3.1, and set $X_{ij} := |X_i(\varphi_j)|^p$. Clearly, with $\sigma^2 := \max_{1 \leq j \leq J} \sum_{i=1}^n \mathbb{E}[X_{ij}^2]$, we obtain, via (7),

$$\sigma^2 \leq n \mathbb{E} \left[\sup_{\varphi \in \Phi} |X_i(\varphi)|^{2p} \right] = n \cdot A_{\Phi,p}. \quad (10)$$

On the other hand, letting $M^2 := \max_{1 \leq i \leq n, 1 \leq j \leq J} |X_{ij}|^2$, it follows

$$\mathbb{E}[M^2] = \mathbb{E} \left[\max_{1 \leq i \leq n} \sup_{\varphi \in \Phi} |X_i(\varphi)|^2 \right] \leq \sum_{i=1}^n \mathbb{E} \left[\sup_{\varphi \in \Phi} |X_i(\varphi)|^2 \right] = n \cdot A_{\Phi,p}. \quad (11)$$

In view of (10) and (11), Lemma D.1 entails

$$\begin{aligned} \mathbb{E} \left[\max_{1 \leq j \leq J} \left| \sum_{i=1}^n (X_{ij} - \mathbb{E}[X_{ij}]) \right| \right] &\leq K \left(\sigma \sqrt{\log J} + \sqrt{\mathbb{E}[M^2] \log J} \right) \\ &= K \left(\sqrt{n \cdot A_{\Phi,p}} \sqrt{\log L_\Phi + cd \log n} + \sqrt{n \cdot A_{\Phi,p}} (\log L_\Phi + cd \log n) \right) \\ &\leq B \cdot \sqrt{n} \log n, \end{aligned} \quad (12)$$

where $K > 0$ is a universal constant and $B > 0$ depends only on $A_{\Phi,p}$, c and d . With this necessary derivations taken care of, we proceed towards the main arguments. By definition, $|\varphi - [\varphi]_N| < \delta$. Recall $S_{n,p}(\cdot)$ from the statement of Theorem 3.1. Note that

$$\begin{aligned} &\mathbb{E} \left[\sup_{\varphi \in \Phi} |S_{n,p}(\varphi) - \mathbb{E}[S_{n,p}(\varphi)]| \right] \\ &\leq \mathbb{E} \left[\max_{1 \leq j \leq J} |S_{n,p}([\varphi]_N) - \mathbb{E}[S_{n,p}([\varphi]_N)]| \right] + \mathbb{E} \left[\sup_{\varphi \in \Phi} |S_{n,p}(\varphi) - S_{n,p}([\varphi]_N)| \right] + \sup_{\phi \in \Phi} |\mathbb{E}[S_{n,p}(\phi)] - \mathbb{E}[S_{n,p}([\phi]_N)]| \\ &:= T_1 + T_2 + T_3. \end{aligned} \quad (13)$$

We tackle (13) one-by-one. Equation (12) instructs that $T_1 = O(\sqrt{n} \log n)$. Next, moving on to T_2 , we observe that

$$\begin{aligned} \mathbb{E} \left[\sup_{\varphi \in \Phi} |S_{n,p}(\varphi) - S_{n,p}([\varphi]_N)| \right] &\leq n \cdot \mathbb{E} \left[\sup_{\varphi \in \Phi} |X_i^p(\varphi) - X_i^p([\varphi]_N)| \right] \\ &\leq np \cdot \mathbb{E} \left[2 \sup_{\varphi \in \Phi} |X_i(\varphi)|^{p-1} \cdot \sup_{\varphi \in \Phi} |X_i(\varphi) - X_i([\varphi]_N)| \right] \end{aligned} \quad (14)$$

$$\leq 2np \left(\mathbb{E} \left[\sup_{\varphi \in \Phi} |X_i(\varphi)|^p \right] \right)^{\frac{p-1}{p}} \left(\mathbb{E} \left[\sup_{\varphi \in \Phi} |X_i(\varphi) - X_i([\varphi]_N)|^p \right] \right)^{\frac{1}{p}} \quad (15)$$

$$\leq 2np \sqrt[p]{A_{\Phi,p}^{\frac{p-1}{p}}} (K_\Phi \delta)^{\frac{1}{p}} = O(n \cdot \delta^{-c/p}) = O(n^{1-c/p}), \quad (16)$$

where, (14) follows due to the elementary inequality $||a|^p - |b|^p| \leq p(|a|^{p-1} + |b|^{p-1}) \cdot |a - b|$, for $p \geq 1$, $a, b \in \mathbb{R}$; (15) involves an application of Hölder's inequality, and finally, (16) invokes (6) and (7). Note that, trivially $T_3 \leq T_2$. Therefore, (13), along with $\delta = O(n^{-c})$ with $c > p/2$, begets,

$$\mathbb{E} \left[\sup_{\varphi \in \Phi} |S_{n,p}(\varphi) - \mathbb{E}[S_{n,p}(\varphi)]| \right] \lesssim \sqrt{n} \log n + n^{1-c/p} = O(\sqrt{n} \log n),$$

where \lesssim hides constants pertaining p, d and φ . This completes the proof. \square

D.2 PROOF OF THEOREM 3.2

The key idea behind Theorem 3.2 is to express the data-driven MEP's as supremum of random functions, before invoking Theorem 3.1.

Proof. For $\theta \in \mathcal{D}$, $\gamma \in \Gamma$ and $u \in \mathcal{B}$, denote

$$M_n(\theta, \gamma, u) := \frac{1}{n} \sum_{i=1}^n \left| \nabla_{\theta} F_{\xi_i}^{\gamma}(\theta) u \right|^p, \text{ and, } M(\theta, \gamma, u) := \mathbb{E} \left[\left| \nabla_{\theta} F_{\xi_i}^{\gamma}(\theta) u \right|^p \right].$$

We start off by establishing

$$\mathbb{E} \left[\sup_{\theta \in \mathcal{D}, \gamma \in \Gamma, u \in \mathcal{B}} |M_n(\theta, \gamma, u) - M(\theta, \gamma, u)| \right] = O \left(\frac{\log n}{\sqrt{n}} \right). \quad (17)$$

Observe that $\Phi := \mathcal{D} \times \Gamma \times \mathcal{B}$ is a compact set, and $F_{\xi_i}^{\gamma}(\theta)$ are i. i. d. random functions taking values in $\varphi \in \Phi$. Moreover, Assumptions 3.2 and 3.3 correspond to (6) and (7) respectively. Therefore, a direct application of Theorem 3.1 entails (17). Finally, in lieu of Lemma 2.2, (17) yields

$$\begin{aligned} \mathbb{E} \left[\sup_{\gamma \in \Gamma} \left| \hat{L}_p^{\ell, n}(\gamma) - L_p^{\ell}(\gamma) \right| \right] &= \mathbb{E} \left[\sup_{\gamma \in \Gamma} \left| \sup_{\theta \in \mathcal{D}} \sup_{u: |u|=1} \frac{1}{n} \sum_{i=1}^n \left| \nabla_{\theta} F_{\xi_i}^{\gamma}(\theta) u \right|^p - \sup_{\theta \in \mathcal{D}} \sup_{u: |u|=1} \mathbb{E} \left[\left| \nabla_{\theta} F_{\xi_i}^{\gamma}(\theta) u \right|^p \right] \right| \right] \\ &\leq \frac{1}{n} \mathbb{E} \left[\sup_{\theta \in \mathcal{D}, \gamma \in \Gamma, u \in \mathcal{B}} \left| \sum_{i=1}^n \left(\left| \nabla_{\theta} F_{\xi_i}^{\gamma}(\theta) u \right|^p - \mathbb{E} \left[\left| \nabla_{\theta} F_{\xi_i}^{\gamma}(\theta) u \right|^p \right] \right) \right| \right] \\ &= O \left(\frac{\log n}{\sqrt{n}} \right), \end{aligned}$$

which completes the proof. \square

E PROOF OF THEOREM 4.3

Proof. We provide the proof for $\ell = 1$. By definition, $\hat{\gamma}_n(p) \in \Gamma$. Fix some $M > 0$ such that $L_p(\gamma_0) + M < 1$. Therefore, invoking Theorem 3.2, it follows,

$$\mathbb{P} \left(\hat{L}_p^n(\gamma_0) < 1 \right) \geq \mathbb{P} \left(\left| \hat{L}_p^n(\gamma_0) - L_p(\gamma_0) \right| < M \right) \rightarrow 1 \text{ as } n \rightarrow \infty. \quad (18)$$

Additionally, suppose $0 < M' < 1$. By the continuity of $L_p(\cdot)$, $L_p(\gamma^+(p)) = 2$, hence, yet another application of Theorem 3.2 entails that

$$\mathbb{P} \left(\hat{L}_p^n(\gamma^+(p)) > 1 \right) \geq \mathbb{P} \left(\left| \hat{L}_p^n(\gamma^+(p)) - 2 \right| < M' \right) \rightarrow 1 \text{ as } n \rightarrow \infty. \quad (19)$$

In view of continuity of $\hat{L}_p^n(\cdot)$, equations (18) and (19) combined, yield that

$$\mathbb{P}(\hat{\gamma}_n(p) \in \text{int}(\Gamma)) \geq \mathbb{P} \left(\hat{L}_p^n(\gamma_0) < 1, \hat{L}_p^n(\gamma^+(p)) > 1 \right) \rightarrow 1 \text{ as } n \rightarrow \infty. \quad (20)$$

This completes the proof of our first assertion. We leverage (20) en route to our second assertion. To that end, observe that following from Assumption 3.4, $L_p(\cdot)$ is differentiable at $\gamma(p)$ with its derivative bounded by K_p . So there exists some $K \leq K_p$, such that we can use it to write out first order Taylor expansion of $L(\cdot)$ about $\gamma(p)$:

$$L(\gamma) - L(\gamma(p)) = K(\gamma - \gamma(p)) + o(\gamma - \gamma(p)). \quad (21)$$

From Theorem 3.2, it follows given $\varepsilon > 0$ that there exist some $G_{\varepsilon} > 0$ and $N_{\varepsilon} > 0$ such that for all $n > N_{\varepsilon}$:

$$\mathbb{P} \left(\sup_{\gamma \in \Gamma} |L_n(\gamma) - L(\gamma)| > G_{\varepsilon} \frac{\log n}{\sqrt{n}} \right) \leq \varepsilon.$$

If $\hat{\gamma}_{\ell, n} \in \Gamma$, then following from the continuity of L_n , we have $L_n(\hat{\gamma}_{\ell, n}) = 1 = L(\gamma_{\ell})$. Therefore,

$$\begin{aligned} \mathbb{P} \left(\sup_{\gamma \in \Gamma} |L_n(\gamma) - L(\gamma)| > G_{\varepsilon} \frac{\log n}{\sqrt{n}} \right) &\geq \mathbb{P} \left(\hat{\gamma}_{\ell, n} \in \Gamma, |L_n(\hat{\gamma}_{\ell, n}) - L(\hat{\gamma}_{\ell, n})| > G_{\varepsilon} \frac{\log n}{\sqrt{n}} \right) \\ &\geq \mathbb{P} \left(\hat{\gamma}_{\ell, n} \in \Gamma, |\hat{\gamma}_{\ell, n} - \gamma_{\ell}| > K_1 \cdot \frac{\log n}{\sqrt{n}} \right), \end{aligned} \quad (22)$$

where, in (22), we invoke (21). Combined with (21), (22) yields

$$\mathbb{P} \left(\hat{\gamma}_{\ell, n} \in \Gamma, |\hat{\gamma}_{\ell, n} - \gamma_{\ell}| > \frac{G_{\varepsilon}}{K_1} \cdot \frac{\log n}{\sqrt{n}} \right) \leq \varepsilon, \quad (23)$$

where, $K_1 := \frac{G_{\varepsilon}}{K}$. Equations (20), (23) jointly concludes the proof. \square