

Asymptotics of Stochastic Gradient Descent with Dropout Regularization in Linear Models

Jiaqi Li

*Department of Statistics
University of Chicago
Chicago, IL, USA*

JQLI@UCHICAGO.EDU

Johannes Schmidt-Hieber

*Department of Applied Mathematics
University of Twente
Enschede, Netherlands*

A.J.SCHMIDT-HIEBER@UTWENTE.NL

Wei Biao Wu

*Department of Statistics
University of Chicago
Chicago, IL, USA*

WBWU@UCHICAGO.EDU

Editor: Michael Mahoney

Abstract

This paper proposes an asymptotic theory for online inference of the stochastic gradient descent (SGD) iterates with dropout regularization in linear regression. Specifically, we establish the geometric-moment contraction (GMC) for constant step-size SGD dropout iterates to show the existence of a unique stationary distribution of the dropout recursive function. Based on the GMC property, we use the functional dependence measure to provide quenched central limit theorems (CLT) for the gradient descent iterates with dropout regularization. Moreover, we obtain CLTs for the Ruppert-Polyak averaged GD (AGD) and averaged SGD (ASGD) iterates with dropout. Based on these asymptotic normality results, we further introduce an online estimator for the long-run covariance matrix of ASGD dropout to facilitate inference in a recursive manner with efficiency in computational time and memory. The numerical experiments demonstrate that for large samples, the proposed confidence intervals for ASGD with dropout achieve the nominal coverage probability.

Keywords: stochastic gradient descent, dropout regularization, ℓ^2 -regularization, online inference, quenched central limit theorems

1. Introduction

Dropout regularization is a popular method in deep learning (Hinton et al., 2012; Krizhevsky et al., 2012; Srivastava et al., 2014). During each training iteration, each hidden unit is randomly masked with probability $1 - p$. This ensures that a hidden unit cannot rely on the presence of another hidden unit. Dropout therefore provides an incentive for different units to act more independently and avoids co-adaptation, which means that different units do the same.

There is a rich literature contributing to the theoretical understanding of dropout regularization. As pointed out in Srivastava et al. (2014), the core idea of dropout is to artificially

introduce stochasticity to the training process, preventing the model from learning statistical noise in the data. Starting with the connection of dropout and ℓ^2 -regularization that appeared already in the original dropout article Srivastava et al. (2014), numerous works investigated the statistical properties of dropout by marginalizing the loss functions over dropout noises and linking them with explicit regularization (Arora et al., 2021; Baldi and Sadowski, 2013; Cavazza et al., 2018; McAllester, 2013; Mianjy and Arora, 2019; Mianjy et al., 2018; Senen-Cerda and Sanders, 2022; Srivastava et al., 2014; Wager et al., 2013). The empirical study in Wei et al. (2020) concluded that adding dropout noise to gradient descent also introduces implicit effects, which cannot be characterized by connections between the gradients of marginalized loss functions and explicit regularizers. For the linear regression model and fixed learning rates, Clara et al. (2024) proved that the implicit effect of dropout adds noise to the iterates and that for a large class of design matrices, this implicit noise does not vanish in the limit.

Though the convergence theory of dropout in fixed design and full gradients has been widely investigated, an analysis of dropout with random design or sequential observations is still lacking, not to mention online statistical inference. To bridge this gap, we provide a theoretical framework for dropout applied to stochastic gradient descent (SGD). In particular, we establish the geometric-moment contraction (GMC) for the SGD dropout iterates for a range of constant learning rates α . We provide useful and sharp moment inequalities to prove the q -th moment convergence of SGD dropout for any $q > 1$.

Besides the convergence and error bounds of SGD dropout, statistical inference of SGD-based estimators is also gaining attention (Fang, 2019; Fang et al., 2019; Liang and Su, 2019; Su and Zhu, 2023; Zhong et al., 2024). Instead of focusing on point estimators using dropout regularization, we quantify the uncertainty of the estimates through their confidence intervals or confidence regions (Chen et al., 2020; Zhu et al., 2023). Nevertheless, it is challenging to derive asymptotic normality for SGD dropout or its variants, such as averaged SGD (ASGD) (Ruppert, 1988; Polyak and Juditsky, 1992). The reason is that the SGD iterates are dependent and the initialization makes the SGD iterates non-stationary. In this paper, we leverage the GMC property of SGD dropout and show quenched central limit theorems (CLT) for both SGD and ASGD dropout estimates. Additionally, we propose an online estimator for the long-run covariance matrix of ASGD dropout to facilitate the online inference.

Contributions. This study employs powerful techniques from time series analysis to derive a general asymptotic theory for the SGD iterates with dropout regularization. Specifically, the key contributions can be summarized as follows.

- (1) We establish the geometric-moment contraction (GMC) of the non-stationary SGD dropout iterates, whose recursion can be viewed as a vector auto-regressive process (VAR). The possible range of learning rates that ensures GMC can be related to the condition number of the design matrix with dropout.
- (2) The GMC property guarantees the existence of a unique stationary distribution of the SGD iterates with dropout, and leads to the L^q -convergence, the asymptotic normality, and the Gaussian approximation rate of the SGD dropout estimates and their Ruppert-Polyak averaged version.

- (3) We derive a new moment inequality in Lemma 27, proving that for any two random vectors \mathbf{x}, \mathbf{y} of the same length, the q -th moment $\mathbb{E}\|\mathbf{x} + \mathbf{y}\|_2^q$ can have a sharp bound in terms of $\mathbb{E}\|\mathbf{x}\|_2^q$, $\mathbb{E}\|\mathbf{y}\|_2^q$ and $\mathbb{E}(\mathbf{x}^\top \mathbf{y})$, without the condition $\mathbb{E}[\mathbf{y} \mid \mathbf{x}] \stackrel{\text{a.s.}}{=} 0$ required in previous results (Rio, 2009). The derived moment inequality is also applicable to many other L^q -convergence problems in machine learning.
- (4) An online statistical inference method is introduced to construct joint confidence intervals for averaged SGD dropout iterates. The coverage probability is shown to be asymptotically accurate in theory and simulation studies.

The rest of the paper is organized as follows. We introduce the dropout regularization in gradient descent in Section 2. Followed by Section 3, we establish the geometric-moment contraction for dropout in gradient descent and provide the asymptotic normality. In Section 4, we generalize the theory to stochastic gradient descent. In Section 5, we provide an online inference algorithm for the ASGD dropout with theoretical guarantees. Finally, we present simulation studies in Section 6. All the technical proofs are postponed to the Appendix.

1.1 Background

Dropout regularization. After its introduction by Hinton et al. (2012); Srivastava et al. (2014), dropout regularization was found to be closely related to ℓ^2 -regularization in linear regression and generalized linear models. See also Baldi and Sadowski (2013); McAllester (2013). Wager et al. (2013) extended this connection to more general injected forms of noise, showing that dropout induces an ℓ^2 -penalty after rescaling the data by the estimated inverse diagonal Fisher information. In neural networks with a single hidden layer, dropout noise marginalization leads to a nuclear norm regularization, as studied in matrix factorization (Cavazza et al., 2018), linear neural networks (Mianjy et al., 2018), deep linear neural networks (Mianjy and Arora, 2019) and shallow ReLU-activated networks (Arora et al., 2021). Moreover, Gal and Ghahramani (2016b) showed that dropout can be interpreted as a variational approximation to the posterior of a Bayesian neural network. Gal and Ghahramani (2016a) applied this new variational inference based dropout technique in recurrent neural networks (RNN) and long-short term memory (LSTM) models. Additional research has explored the impact of dropout on convolutional neural networks (Wu and Gu, 2015) and generalization properties via Rademacher complexity bounds (Arora et al., 2021; Gao and Zhou, 2016; Wan et al., 2013; Zhai and Wang, 2018). Dropout has been successfully applied in various domains, including image classification (Krizhevsky et al., 2012), handwriting recognition (Pham and Le, 2021) and heart sound classification (Kay and Agarwal, 2016).

Stochastic gradient descent. To learn from huge datasets, stochastic gradient descent (SGD) (Robbins and Monro, 1951; Kiefer and Wolfowitz, 1952) is a computationally attractive variant of the gradient descent method. While dropout and SGD have been studied separately, only little theory has been developed so far for SGD training with dropout regularization. Mianjy and Arora (2020) showed the necessary number of SGD iterations to achieve suboptimality in ReLU shallow neural networks for classification tasks, which is independent of the dropout probability due to a strict condition on data structures. Senen-

Cerda and Sanders (2023) extended this to more generic results without assuming any specific data structures, focusing instead on reaching stationarity in non-convex functions using dropout-like SGD. Furthermore, Senen-Cerda and Sanders (2022) analyzed the gradient flow of dropout in shallow linear networks and studied the asymptotic convergence rate of dropout by marginalizing the dropout noise in a shallow network. However, a theoretical convergence analysis or inference theory of SGD dropout iterates without marginalization has not been explored yet in the literature.

1.2 Notation

We denote column vectors in \mathbb{R}^d by lowercase bold letters, that is, $\mathbf{x} := (x_1, \dots, x_d)^\top$ and write $\|\mathbf{x}\|_2 := (\mathbf{x}^\top \mathbf{x})^{1/2}$ for the Euclidean norm. The expectation and covariance of random vectors are respectively denoted by $\mathbb{E}[\cdot]$ and $\text{Cov}(\cdot)$. For two positive number sequences (a_n) and (b_n) , we say $a_n = O(b_n)$ (resp. $a_n \asymp b_n$) if there exists $c > 0$ such that $a_n/b_n \leq c$ (resp. $1/c \leq a_n/b_n \leq c$) for all large n , and say $a_n = o(b_n)$ if $a_n/b_n \rightarrow 0$ as $n \rightarrow \infty$. Let (x_n) and (y_n) be two sequences of random variables. Write $x_n = O_{\mathbb{P}}(y_n)$ if for $\forall \epsilon > 0$, there exists $c > 0$ such that $\mathbb{P}(|x_n/y_n| \leq c) > 1 - \epsilon$ for all large n , and say $x_n = o_{\mathbb{P}}(y_n)$ if $x_n/y_n \rightarrow 0$ in probability as $n \rightarrow \infty$.

We denote matrices by uppercase letters. The $d \times d$ identity matrix is symbolized by I_d . Given matrices A and B of compatible dimension, their matrix product is denoted by juxtaposition. Write A^\top for the transpose of A and define $\mathbb{A} := A^\top A$. When A and B are of the same dimension, the Hadamard product $A \odot B$ is given by element-wise multiplication $(A \odot B)_{ij} = A_{ij}B_{ij}$. For any $A \in \mathbb{R}^{d \times d}$, let $\text{Diag}(A) := I_d \odot A$ denote the diagonal matrix with the same main diagonal as A . Given $p \in (0, 1)$, define the matrices

$$\begin{aligned} \bar{A} &:= A - \text{Diag}(A), \\ A_p &:= pA + (1 - p)\text{Diag}(A). \end{aligned}$$

In particular, $A_p = p\bar{A} + \text{Diag}(A)$, so A_p results from re-scaling the off-diagonal entries of A by p . For a matrix A , the operator norm induced by the Euclidean norm $\|\cdot\|_2$ is the spectral norm and will always be written without sub-script, that is, $\|A\| := \|A\|_{\text{op}}$.

2. Dropout Regularization

The stochasticity of dropout makes it challenging to analyze the asymptotic properties of dropout in stochastic gradient descent. To address the complex stochastic structure, we investigate in Section 2 the dropout regularization in gradient descent, and then consider stochastic gradient descent in Section 4.1.

Consider a linear regression model with fixed design matrix $X \in \mathbb{R}^{n \times d}$ and outcome $\mathbf{y} \in \mathbb{R}^n$, that is,

$$\mathbf{y} = X\boldsymbol{\beta}^* + \boldsymbol{\epsilon}, \tag{1}$$

with unknown regression vector $\boldsymbol{\beta}^* \in \mathbb{R}^d$, and random noise $\boldsymbol{\epsilon} \in \mathbb{R}^n$. The task is to recover $\boldsymbol{\beta}^*$ from the observed data (\mathbf{y}, X) . Moreover, we suppose that $\mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{0}$ and $\text{Cov}(\boldsymbol{\epsilon}) = I_n$. We highlight that the noise distribution of $\boldsymbol{\epsilon}$ is often explicitly modeled as multivariate normal, but this is not necessary for this analysis. We also assume that the design matrix X has no zero columns. Because of that we also say that model (1) is in reduced form. We

Notation	Definition	Equation	Index Range
$\tilde{\beta}_k$	GD iterates with dropout	Eq. (2)	$k \in \mathbb{N}$
$\tilde{\beta}_k^\circ$	stationary GD iterates with dropout	Eq. (14)	$k \in \mathbb{Z}$
$\tilde{\beta}$	ℓ^2 regularizer in the GD setting with dropout	Eq. (3)	/
β_k^\dagger	stationary affine GD iterates with dropout	Eq. (15)	$k \in \mathbb{Z}$
$\check{\beta}_k$	SGD iterates with dropout	Eq. (32)	$k \in \mathbb{N}$
$\check{\beta}_k^\circ$	stationary SGD iterates with dropout	Eq. (42)	$k \in \mathbb{Z}$
$\check{\beta}$	ℓ^2 regularizer in the SGD setting with dropout	Eq. (34)	/
$\bar{\beta}_k^{\text{gd}}$	AGD iterates with dropout	Eq. (24)	$k \in \mathbb{N}$
$\bar{\beta}_k^{\text{sgd}}$	ASGD iterates with dropout	Eq. (48)	$k \in \mathbb{N}$

Table 1: List of the sequences defined in the paper.

can always bring the model into reduced form, since zero columns and the corresponding regression coefficients have no effect on the outcome \mathbf{y} and can thus be eliminated from the model.

We consider the least-squares criterion $\frac{1}{2}\|\mathbf{y} - X\beta\|_2^2$ for the estimation of β^* . For the minimization, we adopt a constant learning-rate gradient descent algorithm with random dropouts in each iteration. Following the seminal work on dropout by Srivastava et al. (2014), we call a $d \times d$ random diagonal matrix D a p -dropout matrix if its diagonal entries satisfy $D_{ii} \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(p)$, with some retaining probability $p \in (0, 1)$. On average, D has pd diagonal entries equal to 1 and $(1-p)d$ diagonal entries equal to 0. For simplicity, the dependence of D on p will only be stated if unclear from the context. For a sequence of independent and identically distributed (i.i.d.) dropout matrices D_k , $k = 1, 2, \dots$, and some constant learning rate $\alpha > 0$, the k -th step gradient descent iterate with dropout takes the form

$$\begin{aligned} \tilde{\beta}_k(\alpha) &= \tilde{\beta}_{k-1}(\alpha) - \alpha \nabla_{\tilde{\beta}_{k-1}} \frac{1}{2} \|\mathbf{y} - XD_k \tilde{\beta}_{k-1}(\alpha)\|_2^2 \\ &= \tilde{\beta}_{k-1}(\alpha) + \alpha D_k X^\top (\mathbf{y} - XD_k \tilde{\beta}_{k-1}(\alpha)), \end{aligned} \quad (2)$$

with $\tilde{\beta}_0$ a given initial vector. Taking the expectation with respect to D_k , one obtains the gradient descent iterates with respect to the loss function $\mathbb{E}[\frac{1}{2}\|\mathbf{y} - XD\beta\|_2^2 | \mathbf{y}, X]$, where the expectation is taken only over the stochasticity in the dropout matrix D . The minimizer of the loss is denoted by

$$\tilde{\beta} := \arg \min_{\beta \in \mathbb{R}^d} \mathbb{E} \left[\frac{1}{2} \|\mathbf{y} - XD\beta\|_2^2 \mid \mathbf{y}, X \right]. \quad (3)$$

Throughout Section 3, we assume that the sample (\mathbf{y}, X) is given, and all the expectations are taken only over dropout matrix D , i.e., $\mathbb{E}[\cdot] = \mathbb{E}_D[\cdot]$. In fact, the vector $\tilde{\beta}$ has a closed form expression. To see this, we denote the Gram matrix by

$$\mathbb{X} = X^\top X$$

and recall

$$\bar{\mathbb{X}} = \mathbb{X} - \text{Diag}(\mathbb{X}), \quad \mathbb{X}_p = p\mathbb{X} + (1-p)\text{Diag}(\mathbb{X}). \quad (4)$$

Note that $D^2 = D$, $\text{Diag}(\mathbb{X}) = \mathbb{X}_p - p\bar{\mathbb{X}}$, and that diagonal matrices always commute. Since the fixed design matrix X is assumed to be in reduced form with $\min_i \mathbb{X}_{ii} > 0$, one can show that solving the gradient for the minimizer $\tilde{\beta}$ in (3) (Srivastava et al., 2014; Clara et al., 2024) leads to the closed form expression

$$\tilde{\beta} = p \left(p^2 \mathbb{X} + p(1-p)\text{Diag}(\mathbb{X}) \right)^{-1} X^\top \mathbf{y} = \mathbb{X}_p^{-1} X^\top \mathbf{y}. \quad (5)$$

If the columns of X are orthogonal, then \mathbb{X} is a diagonal matrix, $\mathbb{X}_p = \mathbb{X}$ and $\tilde{\beta}$ coincides with the classical least-squares estimator $\mathbb{X}^{-1} X^\top \mathbf{y}$. We refer to Section 4.1 for a counterpart of $\tilde{\beta}$ using stochastic gradient. While previous studies have linked ℓ^2 -regularization and GD with dropout by marginalizing the training loss over the dropout noise (Wager et al. (2013); McAllester (2013); Srivastava et al. (2014)), Theorem 7 of Clara et al. (2024) uncovers a persistent gap in their covariance structures under a constant learning rate. This proves that dropout moreover incurs an implicit effect beyond what an explicit ℓ^2 -penalty can capture. Consequently, it is essential to investigate the GD dropout iterates $\tilde{\beta}_k(\alpha)$.

A crucial argument in the analysis of the dropout iterate $\tilde{\beta}_k$ is to rewrite the dropout update formula as

$$\tilde{\beta}_k(\alpha) - \tilde{\beta} = \underbrace{(I_d - \alpha D_k \mathbb{X} D_k)}_{=: A_k(\alpha)} (\tilde{\beta}_{k-1}(\alpha) - \tilde{\beta}) + \underbrace{\alpha D_k \bar{\mathbb{X}} (pI_d - D_k)}_{=: b_k(\alpha)} \tilde{\beta}, \quad (6)$$

see Section 4.1 in Clara et al. (2024).

3. Asymptotic Properties of Dropout in GD

To study the asymptotic properties of gradient descent with dropout, we first establish the geometric-moment contraction for the GD dropout sequence. Subsequently, we derive the quenched central limit theorems for both iterative dropout estimates and their Ruppert-Polyak averaged variants. Furthermore, we provide the quenched invariance principle for the Ruppert-Polyak averaged dropout with the optimal Gaussian approximation rate.

3.1 Geometric-Moment Contraction (GMC)

First, we extend the geometric-moment contraction in Wu and Shao (2004) to the cases where the inputs of iterated random functions are i.i.d. random matrices.

Definition 1 (Geometric-moment contraction) *For i.i.d. $d \times d$ random matrices Ψ_i, Ψ'_j , $i, j \in \mathbb{Z}$, consider a stationary causal process*

$$\boldsymbol{\theta}_k = g(\Psi_k, \dots, \Psi_1, \Psi_0, \Psi_{-1}, \dots), \quad k \in \mathbb{Z}, \quad (7)$$

for a measurable function $g(\cdot)$ such that the d -dimensional random vector $\boldsymbol{\theta}_k$ has a finite q -th moment $\mathbb{E} \|\boldsymbol{\theta}_k\|_2^q < \infty$, for some $q \geq 1$. We say that $\boldsymbol{\theta}_k$ is geometric-moment contracting if there exists some constant $r_q \in (0, 1)$ such that

$$(\mathbb{E} \|\boldsymbol{\theta}_k - \boldsymbol{\theta}'_k\|_2^q)^{1/q} = O(r_q^k), \quad \text{for all } k = 1, 2, \dots, \quad (8)$$

where $\theta'_k = g(\Psi_k, \dots, \Psi_1, \Psi'_0, \Psi'_{-1}, \dots)$ is a coupled version of θ_k with Ψ_i , $i \leq 0$, replaced by i.i.d. copies Ψ'_i .

In general, an iterated random function satisfies the geometric-moment contraction property under regularity conditions on convexity and stochastic Lipschitz continuity, see Section B.1 in the Appendix for details. Here, we focus on the contraction property with $\Psi_k = D_k$, the k -th dropout matrix. Setting

$$f_D(\mathbf{u}) := \mathbf{u} + \alpha D X^\top (\mathbf{y} - X D \mathbf{u}),$$

we can rewrite the recursion of the dropout gradient descent iterate $\tilde{\beta}_k(\alpha)$ in (2) as

$$\tilde{\beta}_k(\alpha) = \tilde{\beta}_{k-1}(\alpha) + \alpha D_k X^\top (\mathbf{y} - X D_k \tilde{\beta}_{k-1}(\alpha)) =: f_{D_k}(\tilde{\beta}_{k-1}(\alpha)). \quad (9)$$

We shall show that, under quite general conditions on the constant learning rate $\alpha > 0$, this process satisfies the geometric-moment contraction in Definition 1, and converges weakly to a unique stationary distribution π_α on \mathbb{R}^d , that is, for any continuous function $h \in \mathcal{C}(\mathbb{R}^d)$ with $\|h\|_\infty < \infty$, $\mathbb{E}[h(\tilde{\beta}_k(\alpha))] \rightarrow \int h(\mathbf{u}) \pi_\alpha(d\mathbf{u})$ as $k \rightarrow \infty$. We then write $\tilde{\beta}_k(\alpha) \Rightarrow \pi_\alpha$. Set

$$r_{\alpha,q} := \left(\sup_{\mathbf{v} \in \mathbb{R}^d: \|\mathbf{v}\|_2=1} \mathbb{E} \left\| (I_d - \alpha D_1 \mathbb{X} D_1) \mathbf{v} \right\|_2^q \right)^{1/q}. \quad (10)$$

In particular, for $q = 2$, we can rewrite the squared norm and obtain $r_{\alpha,2}^2 = \lambda_{\max}(\mathbb{E}(I_d - \alpha D_1 \mathbb{X} D_1)^2)$ with $\lambda_{\max}(\cdot)$ the largest eigenvalue. The exponential stability of online algorithms with constant step size has also been studied by Durmus et al. (2021) in the context of linear stochastic approximation and by Samsonov et al. (2024) to establish convergence of the stationary law in Wasserstein distance for temporal-difference learning. Theorem 3 below complements these results by proving contraction in the Euclidean distance to show the existence and uniqueness of the stationary distribution for GD with dropout and constant learning rate.

Lemma 2 *If $q > 1$ and $\alpha \|\mathbb{X}\| < 2$, then, $r_{\alpha,q} < 1$.*

As we only assumed that the design matrix X has no zero column, $\mathbb{X} = X^\top X$ can be singular. The previous lemma shows that even for singular \mathbb{X} , contraction coefficient $r_{\alpha,q} < 1$ is possible. Without dropout, for any \mathbf{v} in the kernel of \mathbb{X} , one has $\|(I_d - \alpha \mathbb{X})\mathbf{v}\|_2 = \|\mathbf{v}\|_2$, implying that $\sup_{\mathbf{v} \in \mathbb{R}^d: \|\mathbf{v}\|_2=1} \mathbb{E} \|(I_d - \alpha \mathbb{X})\mathbf{v}\|_2^q \geq 1$.

Theorem 3 (Geometric-moment contraction of GD dropout) *Let $q > 1$. Choose a positive learning rate α satisfying $\alpha \|\mathbb{X}\| < 2$. For two dropout sequences $\tilde{\beta}_k(\alpha), \tilde{\beta}'_k(\alpha)$, $k = 0, 1, \dots$, generated by the recursion (6) with the same dropout matrices but possibly different initial vectors $\tilde{\beta}_0, \tilde{\beta}'_0$, we have*

$$\left(\mathbb{E} \|\tilde{\beta}_k(\alpha) - \tilde{\beta}'_k(\alpha)\|_2^q \right)^{1/q} \leq r_{\alpha,q}^k \|\tilde{\beta}_0 - \tilde{\beta}'_0\|_2. \quad (11)$$

Moreover, there exists a unique stationary distribution π_α which does not depend on the initialization $\tilde{\beta}_0$, such that $\tilde{\beta}_k(\alpha) \Rightarrow \pi_\alpha$ as $k \rightarrow \infty$.

Regarding the condition $\alpha\|\mathbb{X}\| < 2$, numerical experiments in Section 6.1 show that, for a wide range of settings, choosing the learning rate α slightly larger than $2/\|\mathbb{X}\|$, the empirical version of $r_{\alpha,q}$ will be larger than 1, and consequently, the GMC property in (11) will no longer hold. As we only require $\alpha\|\mathbb{X}\| < 2$, the dimension d influences the GMC property in the previous theorem only through the structure of the Gram matrix $\mathbb{X} = X^\top X$.

As mentioned before, $r_{\alpha,2}^2 = \lambda_{\max}(\mathbb{E}(I_d - \alpha D_1 \mathbb{X} D_1)^2) < 1$. A special case of Theorem 3 is thus $\mathbb{E}\|\tilde{\beta}_k(\alpha) - \tilde{\beta}'_k(\alpha)\|_2^2 \leq (\lambda_{\max}(\mathbb{E}(I_d - \alpha D_1 \mathbb{X} D_1)^2))^k \|\tilde{\beta}_0 - \tilde{\beta}'_0\|_2^2$. Theorem 3 indicates that although the GD dropout sequence $\{\tilde{\beta}_k(\alpha)\}_{k \in \mathbb{N}}$ is non-stationary due to the initialization, it is asymptotically stationary and approaches the unique stationary distribution π_α at an exponential rate. Such geometric-moment contraction result is fundamental to establish a central limit theorem for the iterates.

Another consequence of Theorem 3 is that if β is drawn from the stationary distribution π_α and D is an independently sampled dropout matrix, then also $f_D(\beta) \sim \pi_\alpha$. We refer to Corollary 23 in the Appendix for the detailed arguments. This result means that if the initialization $\tilde{\beta}_0^\circ$ is sampled from the stationary distribution π_α , then, the marginal distribution of any of the GD dropout iterates $\tilde{\beta}_k^\circ(\alpha)$ will follow this stationary distribution as well.

We can also define the GD dropout iterates $\tilde{\beta}_k^\circ(\alpha)$ for negative integers k by considering i.i.d. dropout matrices D_k for all integers $k \in \mathbb{Z}$ and observing that the limit

$$\tilde{\beta}_k^\circ(\alpha) := \lim_{m \rightarrow \infty} f_{D_k} \circ f_{D_{k-1}} \circ \cdots \circ f_{D_{k-m}}(\beta) =: h_\alpha(D_k, D_{k-1}, \dots), \quad (12)$$

exists almost surely and does not depend on β (see Corollary 23). Then, it follows that $\tilde{\beta}_k^\circ(\alpha) = f_{D_k}(\tilde{\beta}_{k-1}^\circ(\alpha))$ also holds for negative integers and the geometric-moment contraction in Definition 1 is satisfied for $\tilde{\beta}_k^\circ(\alpha)$, that is,

$$\left(\mathbb{E} \|h_\alpha(D_k, \dots, D_1, D_0, D_{-1}, \dots) - h_\alpha(D_k, \dots, D_1, D'_0, D'_{-1}, \dots)\|_2^q \right)^{1/q} = O(r_{\alpha,q}^k), \quad (13)$$

for some $q \geq 1$, $r_{\alpha,q} \in (0, 1)$ defined in (10), and i.i.d. dropout matrices D_k, D'_ℓ , $k, \ell \in \mathbb{Z}$.

3.2 Iterative Dropout Schemes

Equation (6) rewrites the GD dropout iterates into $\tilde{\beta}_k(\alpha) - \tilde{\beta} = A_k(\alpha)(\tilde{\beta}_{k-1}(\alpha) - \tilde{\beta}) + \mathbf{b}_k(\alpha)$. If the initial vector $\tilde{\beta}_0^\circ$ is sampled from the stationary distribution π_α , we also have

$$\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta} = A_k(\alpha)(\tilde{\beta}_{k-1}^\circ(\alpha) - \tilde{\beta}) + \mathbf{b}_k(\alpha), \quad (14)$$

and for any $k = 0, 1, \dots$, $\tilde{\beta}_k^\circ(\alpha) \sim \pi_\alpha$. We can see that $\{\tilde{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ is a stationary vector autoregressive process (VAR) with random coefficients. While $(A_k(\alpha), \mathbf{b}_k(\alpha))$ are i.i.d., $A_k(\alpha)$ and $\mathbf{b}_k(\alpha)$ are dependent. This poses challenges to prove asymptotic normality of the dropout iterates. An intermediate recursion is obtained by replacing $A_k(\alpha) = I_d - \alpha D_k \mathbb{X} D_k$ by its expectation $\mathbb{E}[A_k(\alpha)] = I_d - \alpha p \mathbb{X}_p$. This gives the recursion

$$\beta_k^\dagger(\alpha) - \tilde{\beta} = (I_d - \alpha p \mathbb{X}_p)(\beta_{k-1}^\dagger(\alpha) - \tilde{\beta}) + \mathbf{b}_k(\alpha). \quad (15)$$

The condition $\alpha \in (0, 2/\|\mathbb{X}\|)$ implies the operator norm bound $\|I_d - \alpha p \mathbb{X}_p\| < 1$ (see Step 1 in the proof of Lemma 4). Thus GMC holds and the sequence $\{\beta_k^\dagger(\alpha)\}_{k \in \mathbb{N}}$ converges to

a unique stationary distribution as $k \rightarrow \infty$. For brevity, throughout the paper, we shall assume that the initialization β_0^\dagger follows the stationary distribution, and therefore the affine sequence $\{\beta_k^\dagger(\alpha)\}_{k \in \mathbb{N}}$ is stationary. The proof then derives the asymptotic normality for $\beta_k^\dagger(\alpha)$, and shows that the difference between $\tilde{\beta}_k(\alpha)$ and $\beta_k^\dagger(\alpha)$ is negligible, in the sense that for $q \geq 2$, $(\mathbb{E}\|\tilde{\beta}_k(\alpha) - \beta_k^\dagger(\alpha)\|_2^q)^{1/q} = O(\alpha + r_{\alpha,q}^k \|\tilde{\beta}_0 - \tilde{\beta}_0^\circ\|_2)$, where the first part is due to the affine approximation in Lemma 4 and the second part results from the GMC property in Theorem 3.

Lemma 4 (Affine approximation) *If $\alpha \in (0, 2/\|\mathbb{X}\|)$, then the difference sequence $\delta_k(\alpha) = \tilde{\beta}_k^\circ(\alpha) - \beta_k^\dagger(\alpha)$ satisfies $\mathbb{E}[\delta_k(\alpha)] = 0$ and for any $q \geq 2$, $\max_k (\mathbb{E}\|\delta_k(\alpha)\|_2^q)^{1/q} = O(\alpha)$.*

Lemma 5 (Moment convergence of iterative GD dropout) *Let $q \geq 2$. For the stationary GD dropout sequence $\{\tilde{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ defined in (6), if $\alpha \in (0, 2/\|\mathbb{X}\|)$, we have*

$$\max_k (\mathbb{E}\|\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}\|_2^q)^{1/q} = O(\sqrt{\alpha}). \quad (16)$$

If the Gram matrix \mathbb{X} is ill-conditioned in high dimensions, then the allowed range for α is more restrictive. While a smaller α improves the moment convergence rate $O(\sqrt{\alpha})$ by Lemma 5, it causes the contraction constant $r_{\alpha,q}$ to approach 1, thereby requiring a larger number of iterations k for $\tilde{\beta}_k(\alpha)$ to converge to the stationary distribution π_α . If $\alpha \downarrow 0$, the stationary distribution π_α converges to the normal distribution in a sense that is formally stated in the next result. Recall that $\pi_\alpha(\mathcal{A}) = \mathbb{P}(\tilde{\beta}_1^\circ(\alpha) \in \mathcal{A})$ is the stationary distribution defined in (14).

Theorem 6 (Quenched CLT of iterative GD dropout) *Consider the iterative gradient descent dropout sequence $\{\tilde{\beta}_k(\alpha)\}_{k \in \mathbb{N}}$ in (6) and the ℓ^2 -regularized estimator $\tilde{\beta}$ in (5). Assume that the constant learning rate α satisfies $\alpha \in (0, 2/\|\mathbb{X}\|)$, $\tilde{\beta} \neq 0$, and suppose that for every $l = 1, \dots, d$, there exists $m \neq l$ such that $\mathbb{X}_{lm} \neq 0$. For any set $\mathcal{B} = (-\infty, b_1] \times \dots \times (-\infty, b_d]$ with real numbers b_1, \dots, b_d ,*

$$\pi_\alpha(\tilde{\beta} + \sqrt{\alpha}\Xi^{1/2}(\alpha)\mathcal{B}) \rightarrow \mathbb{P}(z \in \mathcal{B}), \quad \text{as } \alpha \rightarrow 0, \quad (17)$$

with $z \sim \mathcal{N}(0, I_d)$ following the standard d -variate normal distribution, and

$$\Xi(\alpha) := \text{Cov}\left(\frac{\beta_1^\dagger(\alpha) - \tilde{\beta}}{\sqrt{\alpha}}\right) = \frac{\mathbb{E}[(\beta_1^\dagger(\alpha) - \tilde{\beta})(\beta_1^\dagger(\alpha) - \tilde{\beta})^\top]}{\alpha}. \quad (18)$$

In (17), the centering term is $\tilde{\beta}$ since one can show that the affine sequence $\beta_k^\dagger(\alpha)$ satisfies $\mathbb{E}[\beta_k^\dagger(\alpha) - \tilde{\beta}] = 0$ (see the proof of Theorem 6), which along with Lemma 4 yields $\mathbb{E}[\tilde{\beta}_1^\circ(\alpha) - \tilde{\beta}] = 0$. The result also implies convergence on all hyper-rectangles $(a_1, b_1] \times \dots \times (a_d, b_d]$. One can derive more explicit expressions of $\Xi(\alpha)$ and $\Xi(0) := \lim_{\alpha \downarrow 0} \Xi(\alpha)$. Reshaping a $d \times s$ matrix $U = (\mathbf{u}_1, \dots, \mathbf{u}_s)$ with d -dimensional column vectors $\mathbf{u}_1, \dots, \mathbf{u}_s$ into a ds -dimensional column vector gives $\text{vec}(U) := (\mathbf{u}_1^\top, \dots, \mathbf{u}_s^\top)^\top$. Moreover, for any two matrices $A \in \mathbb{R}^{p \times q}$ and $B \in \mathbb{R}^{m \times n}$, the Kronecker product $A \otimes B$ is the $pm \times qn$ block matrix, with each block given by $(A \otimes B)_{ij} = A_{ij}B$. Following Theorem 1 in Pflug (1986),

and assuming that $\Xi(\alpha)$ is differentiable with respect to α , $\Xi(\alpha)$ becomes the solution of a classical Lyapunov equation

$$\Xi(\alpha)(p\mathbb{X}_p) + (p\mathbb{X}_p)\Xi(\alpha) = S,$$

that is,

$$\Xi(\alpha) = V_0 + \alpha B_p, \quad (19)$$

where the $d \times d$ matrices S, V_0 and B_p are respectively defined as

$$S = \frac{1}{\alpha^2} \text{Cov}(\mathbf{b}_1(\alpha)) = \text{Cov}(D_1 \bar{\mathbb{X}}(pI_d - D_1) \tilde{\boldsymbol{\beta}}), \quad (20)$$

$$\text{vec}(V_0) = (I_d \otimes p\mathbb{X}_p + p\mathbb{X}_p \otimes I_d)^{-1} \cdot \text{vec}(S), \quad (21)$$

$$\text{vec}(B_p) = (I_d \otimes p\mathbb{X}_p + p\mathbb{X}_p \otimes I_d)^{-1} \cdot \text{vec}(p^2 \mathbb{X}_p V_0 \mathbb{X}_p). \quad (22)$$

The matrix V_0 satisfies the equation $V_0 p \mathbb{X}_p + p \mathbb{X}_p V_0 = S$. By definition, the matrix S is independent of α and $\bar{\mathbb{X}} \neq 0$ since there exist non-zero diagonal and off-diagonal elements by assumptions in Theorem 6. Let $S_0 = \tilde{\boldsymbol{\beta}} \tilde{\boldsymbol{\beta}}^\top$. By the proof of Theorem 6, we can express S in terms of p, X and $\tilde{\boldsymbol{\beta}}$ as follows,

$$\begin{aligned} S = & p^3 (\bar{\mathbb{X}} S_0 \bar{\mathbb{X}})_p - 2p \left(p \bar{\mathbb{X}}_p (S_0 \bar{\mathbb{X}})_p + p^2 (1-p) \text{Diag}(\bar{\mathbb{X}} S_0 \bar{\mathbb{X}}) \right) \\ & + p \bar{\mathbb{X}}_p (S_0)_p \bar{\mathbb{X}}_p + p^2 (1-p) \left(\text{Diag}(\bar{\mathbb{X}} (S_0)_p \bar{\mathbb{X}}) + 2 \bar{\mathbb{X}}_p \text{Diag}(\bar{S}_0 \bar{\mathbb{X}}) + (1-p) \bar{\mathbb{X}} \odot \bar{S}_0^\top \odot \bar{\mathbb{X}} \right). \end{aligned} \quad (23)$$

One can see that, $\Xi(0) = \lim_{\alpha \downarrow 0} \Xi(\alpha) = V_0$, and in particular, for small p , $\text{vec}(V_0)$ can be approximated by $(I_d \otimes \mathbb{X}_p + \mathbb{X}_p \otimes I_d)^{-1} \cdot \text{vec}(\bar{\mathbb{X}}_p [\tilde{\boldsymbol{\beta}} \tilde{\boldsymbol{\beta}}^\top]_p \bar{\mathbb{X}}_p)$. By the expansion in (19), the CLT in (17) still holds with $\Xi(\alpha)$ therein replaced by $\Xi(0) = V_0$.

By Theorems 3 and 6, the GD dropout iterates $\tilde{\boldsymbol{\beta}}_k(\alpha)$ exhibit asymptotic normality via a two-step limit. First, for any fixed constant learning rate $\alpha > 0$, the GMC property (cf. Theorem 3) guarantees that the GD dropout iterates will first converge to the stationary distribution π_α as the iteration number k grows. Decreasing then the learning rate $\alpha \rightarrow 0$, π_α converges to the normal distribution (cf. Theorem 6). Similar two-step schemes are also explored in the literature; for example, Defazio et al. (2024) proposed a modern learning schedule which consists of an initial warm-up stage with relatively large learning rates followed by rapid learning rate annealing near the end of training.

3.3 Dropout with Ruppert-Polyak Averaging

To reduce the variance of the gradient descent iterates $\tilde{\boldsymbol{\beta}}_k(\alpha)$ introduced by the random dropout matrix D_k , we now consider the averaged GD dropout (AGD) iterate

$$\bar{\boldsymbol{\beta}}_k^{\text{gd}}(\alpha) = \frac{1}{k} \sum_{i=1}^k \tilde{\boldsymbol{\beta}}_i(\alpha), \quad (24)$$

following the averaging scheme in Ruppert (1988); Polyak and Juditsky (1992). We derive the asymptotic normality of $\bar{\boldsymbol{\beta}}_k^{\text{gd}}(\alpha)$ in the following theorem.

Theorem 7 (Quenched CLT of averaged GD dropout) For a constant learning rate $\alpha \in (0, 2/\|\mathbb{X}\|)$ and any fixed initial vector $\tilde{\beta}_0$, the averaged GD dropout sequence satisfies

$$\sqrt{k}(\bar{\beta}_k^{\text{gd}}(\alpha) - \tilde{\beta}) \Rightarrow \mathcal{N}(0, \Sigma(\alpha)), \quad \text{as } k \rightarrow \infty, \quad (25)$$

with $\Sigma(\alpha) = \sum_{i=-\infty}^{\infty} \text{Cov}(\tilde{\beta}_0^{\circ}(\alpha), \tilde{\beta}_i^{\circ}(\alpha))$ the long-run covariance matrix of the stationary process $\tilde{\beta}_i^{\circ}(\alpha) \sim \pi_{\alpha}$.

One can choose a few learning rates, say $\alpha_1, \dots, \alpha_s$, and run gradient descent for each of these learning rates in parallel by computing $\tilde{\beta}_k(\alpha_1), \dots, \tilde{\beta}_k(\alpha_s)$ for $k = 1, 2, \dots$. An example is federated learning where data are distributed across different clients (Dean et al., 2012; Karimireddy et al., 2020; Zinkevich et al., 2010).

Corollary 8 (Quenched CLT of parallel averaged GD dropout) Let $s \geq 1$. Consider constant learning rates $\alpha_1, \dots, \alpha_s \in (0, 2/\|\mathbb{X}\|)$. Then, for any initial vectors $\tilde{\beta}_0(\alpha_1), \dots, \tilde{\beta}_0(\alpha_s)$,

$$\sqrt{k} \cdot \text{vec}(\bar{\beta}_k^{\text{gd}}(\alpha_1) - \tilde{\beta}, \dots, \bar{\beta}_k^{\text{gd}}(\alpha_s) - \tilde{\beta}) \Rightarrow \mathcal{N}(0, \Sigma^{\text{vec}}), \quad \text{as } k \rightarrow \infty, \quad (26)$$

with $\text{vec}(\mathbf{u}_1, \dots, \mathbf{u}_s) = (\mathbf{u}_1^{\top}, \dots, \mathbf{u}_s^{\top})^{\top} \in \mathbb{R}^{ds}$ for d -dimensional vectors $\mathbf{u}_1, \dots, \mathbf{u}_s$, and the long-run covariance matrix

$$\Sigma^{\text{vec}} = \sum_{i=-\infty}^{\infty} \text{Cov}\left(\text{vec}(\tilde{\beta}_0^{\circ}(\alpha_1), \dots, \tilde{\beta}_0^{\circ}(\alpha_s)), \text{vec}(\tilde{\beta}_i^{\circ}(\alpha_1), \dots, \tilde{\beta}_i^{\circ}(\alpha_s))\right).$$

Next, we provide a stronger Gaussian approximation result, namely the rate for the Komlós–Major–Tusnády (KMT) approximation (Komlós et al., 1975, 1976; Berkes et al., 2014). In the quenched invariance principle below, we show that one can achieve the optimal Gaussian approximation rate for the averaged GD dropout process.

Theorem 9 (Quenched invariance principle of averaged GD dropout) Assume the constant learning rate satisfies $\alpha \in (0, 2/\|\mathbb{X}\|)$. Define a partial sum process $(S_i^{\circ}(\alpha))_{1 \leq i \leq t}$ for $t \in \mathbb{N}_+$ with

$$S_i^{\circ}(\alpha) = \sum_{k=1}^i (\tilde{\beta}_k^{\circ}(\alpha) - \tilde{\beta}). \quad (27)$$

Then, there exists a (richer) probability space $(\Omega^*, \mathcal{A}^*, \mathbb{P}^*)$ on which one can define d -dimensional random vectors $\tilde{\beta}_k^*$, the associated partial sum process $S_i^*(\alpha) = \sum_{k=1}^i (\tilde{\beta}_k^*(\alpha) - \tilde{\beta})$, and a Gaussian process $G_i^* = \sum_{k=1}^i \mathbf{z}_k^*$, with independent Gaussian random vectors $\mathbf{z}_k^* \sim \mathcal{N}(0, I_d)$, such that $(S_i^{\circ}(\alpha))_{1 \leq i \leq t} \stackrel{\mathcal{D}}{=} (S_i^*(\alpha))_{1 \leq i \leq t}$ and

$$\max_{1 \leq i \leq t} \|S_i^* - \Sigma^{1/2}(\alpha) G_i^*\|_2 = o_{\mathbb{P}}(t^{1/q}), \quad \text{in } (\Omega^*, \mathcal{A}^*, \mathbb{P}^*), \quad (28)$$

where $\Sigma(\alpha)$ is the long-run covariance matrix defined in Theorem 7. In addition, this approximation holds for all $(S_i^{\tilde{\beta}_0}(\alpha))_{1 \leq i \leq t}$ given any arbitrary initial vector $\tilde{\beta}_0 \in \mathbb{R}^d$, where

$$S_i^{\tilde{\beta}_0}(\alpha) = \sum_{k=1}^i (\tilde{\beta}_k(\alpha) - \tilde{\beta}). \quad (29)$$

Theorem 9 shows that one can approximate the averaged GD dropout sequence by Brownian motions. Specifically, for any fixed initial vector $\check{\beta}_0 \in \mathbb{R}^d$, the partial sum process converges in the Euclidean norm, uniformly over u ,

$$\{t^{-1/2}S_{[tu]}^{\check{\beta}_0}(\alpha), 0 \leq u \leq 1\} \Rightarrow \{\Sigma^{1/2}(\alpha)\mathbb{B}(u), 0 \leq u \leq 1\}, \quad (30)$$

where $[m] = \max\{i \in \mathbb{Z} : i \leq m\}$, and $\mathbb{B}(u)$ is the standard d -dimensional Brownian motion, that is, it can be represented as a d -dimensional vector of independent standard Brownian motions. According to the arguments in Karmakar and Wu (2020), the KMT approximation rate $o_{\mathbb{P}}(t^{1/q})$ is optimal for fixed-dimension time series. Since we can view the GD dropout sequence as a VAR(1) process, the approximation rate in Theorem 9 is optimal for the partial sum process $(S_i^{\check{\beta}_0}(\alpha))_{1 \leq i \leq t}$.

4. Asymptotic Properties of Dropout in SGD

In the previous section, we considered a fixed design matrix and (full) gradient descent with dropout. Computing the gradient over the entire dataset can be computationally expensive, especially with large datasets. We now investigate stochastic gradient descent with dropout regularization.

4.1 Dropout Regularization in SGD

Consider i.i.d. covariate vectors $\mathbf{x}_k \in \mathbb{R}^d$, $k = 1, 2, \dots$, from some distribution Π , and the realizations $y_k | \mathbf{x}_k$ from a linear regression model

$$y_k = \mathbf{x}_k^\top \beta^* + \epsilon_k, \quad (31)$$

with unknown regression vector $\beta^* \in \mathbb{R}^d$. We assume that the model is in reduced form, which here means that $\min_i (\mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top])_{ii} > 0$. In addition, we assume that the i.i.d. random noises ϵ_k satisfy $\mathbb{E}[\epsilon_k] = 0$ and $\text{Var}(\epsilon_k) = 1$. In this paper, we focus on the classical case where the SGD computes the gradient based on an individual observation (y_k, \mathbf{x}_k) and a dropout matrix D_k . For constant learning rate α and initialization $\check{\beta}_0$, the k -th step SGD iterate with Bernoulli dropout is

$$\begin{aligned} \check{\beta}_k(\alpha) &= \check{\beta}_{k-1}(\alpha) - \alpha \nabla_{\check{\beta}_{k-1}} \frac{1}{2} (y_k - \mathbf{x}_k^\top D_k \check{\beta}_{k-1}(\alpha))^2 \\ &= \check{\beta}_{k-1}(\alpha) + \alpha D_k \mathbf{x}_k (y_k - \mathbf{x}_k^\top D_k \check{\beta}_{k-1}(\alpha)). \end{aligned} \quad (32)$$

This is a sequential estimation, or online learning scheme, as computing $\check{\beta}_k(\alpha)$ from $\check{\beta}_{k-1}(\alpha)$ only requires the k -th sample (y_k, \mathbf{x}_k) and the dropout matrix D_k . To study the contraction property of the SGD dropout iterates $\check{\beta}_k(\alpha)$, we express the recursion (32) by an iterated random function $\check{f} : \mathbb{R}^{d \times d} \times \mathbb{R} \times \mathbb{R}^d \times \mathbb{R}^d \mapsto \mathbb{R}^d$ with

$$f_{D,(y,\mathbf{x})}(\mathbf{u}) = \mathbf{u} + \alpha D \mathbf{x} (y - \mathbf{x} D \mathbf{u}),$$

that is,

$$\begin{aligned} \check{\beta}_k(\alpha) &= \check{\beta}_{k-1}(\alpha) + \alpha D_k \mathbf{x}_k (y_k - \mathbf{x}_k^\top D_k \check{\beta}_{k-1}(\alpha)) \\ &=: \check{f}_{D_k,(y_k,\mathbf{x}_k)}(\check{\beta}_{k-1}(\alpha)). \end{aligned} \quad (33)$$

We shall show that this iterated random function \check{f} is geometrically contracting under suitable conditions, and therefore, there exists a unique stationary distribution $\check{\pi}_\alpha$ such that $\check{\beta}_k(\alpha) \Rightarrow \check{\pi}_\alpha$, where \Rightarrow denotes the convergence in distribution.

From now on, let (y, \mathbf{x}) be a sample with the same distribution as (y_k, \mathbf{x}_k) . By marginalizing over all randomness, we can view the SGD dropout in (32) as a minimizer of the ℓ^2 -regularized least-squares loss

$$\check{\beta} := \arg \min_{\beta \in \mathbb{R}^d} \mathbb{E}_{(y, \mathbf{x})} \mathbb{E}_D \left[\frac{1}{2} (y - \mathbf{x}^\top D \beta)^2 \right]. \quad (34)$$

Here, the expectation is taken over both the random sample (y, \mathbf{x}) and the dropout matrix D . Throughout the rest of the paper, we shall write $\mathbb{E}[\cdot] = \mathbb{E}_{(y, \mathbf{x})} \mathbb{E}_D[\cdot]$ when no confusion should be caused.

Denote the $d \times d$ Gram matrix by $\mathbb{X}_k = \mathbf{x}_k \mathbf{x}_k^\top$, and define

$$\bar{\mathbb{X}}_k = \mathbb{X}_k - \text{Diag}(\mathbb{X}_k), \quad \mathbb{X}_{k,p} = p\mathbb{X}_k + (1-p)\text{Diag}(\mathbb{X}_k). \quad (35)$$

By Lemma 28 in the Appendix, we have a closed form solution for $\check{\beta}$ as follows

$$\check{\beta} = p \left(p^2 \mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top] + p(1-p) \text{Diag}(\mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top]) \right)^{-1} \mathbb{E}[y_1 \mathbf{x}_1] = (\mathbb{E}[\mathbb{X}_{1,p}])^{-1} \mathbb{E}[y_1 \mathbf{x}_1],$$

and thus, we obtain the relationship $\mathbb{E}[\mathbb{X}_{1,p}] \check{\beta} = \mathbb{E}[y_1 \mathbf{x}_1]$. To study the SGD with dropout, we now focus on the difference process $\{\check{\beta}_k(\alpha) - \check{\beta}\}_{k \in \mathbb{N}}$. As in the case of gradient descent, this process can be written in autoregressive form,

$$\check{\beta}_k(\alpha) - \check{\beta} = \underbrace{(I_d - \alpha D_k \mathbb{X}_k D_k)}_{=: \check{A}_k(\alpha)} (\check{\beta}_{k-1}(\alpha) - \check{\beta}) + \underbrace{\alpha D_k \mathbf{x}_k (y_k - \mathbf{x}_k^\top D_k \check{\beta})}_{=: \check{b}_k(\alpha)}. \quad (36)$$

4.2 GMC of Dropout in SGD

Establishing the geometric-moment contraction (GMC) property to the stochastic gradient descent iterates with dropout is non-trivial as the randomness of $\check{\beta}_k(\alpha)$ not only comes from the dropout matrix D_k , but also the random sample (y_k, \mathbf{x}_k) . Recall that $\mathbb{X}_{k,p}$ is defined in (35) and by Lemma 21(ii), $\mathbb{E}[D \mathbb{X}_k D] = p \mathbb{E}[\mathbb{X}_{k,p}]$.

Lemma 10 (Learning-rate range in SGD dropout) *Assume that for some $q \geq 2$, the q -th moment $\mu_q(\mathbf{v}) = (\mathbb{E} \|D_k \mathbb{X}_k D_k \mathbf{v}\|_2^q)^{1/q} < \infty$ for all unit vectors $\mathbf{v} \in \mathbb{R}^d$. If the learning rate $\alpha > 0$ satisfies*

$$\frac{\alpha(q-1)}{2} \sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \frac{(1 + \alpha \mu_q(\mathbf{v}))^{q-2} \mu_q(\mathbf{v})^2}{p \mathbf{v}^\top \mathbb{E}[\mathbb{X}_{k,p}] \mathbf{v}} < 1, \quad (37)$$

then, for a dropout matrix D ,

$$\check{r}_{\alpha,q}^q := \sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \mathbb{E} \|(I_d - \alpha D \mathbb{X}_k D) \mathbf{v}\|_2^q < 1. \quad (38)$$

This provides a sufficient condition for the learning rate α which ensures contraction of $I_d - \alpha D \mathbb{X}_k D$ for moments $q \geq 2$. This will lead to L^q -convergence of the SGD dropout iterates and determines the convergence rate in the Gaussian approximation in Theorem 17.

For the special and important case $q = 2$, the identities $\mu_2(\mathbf{v})^2 = \mathbb{E}\|D_k \mathbb{X}_k D_k \mathbf{v}\|_2^2$ and $\mathbb{E}(D_k \mathbb{X}_k D_k) = p \mathbb{E}[\mathbb{X}_{k,p}]$ imply that condition (37) can be rewritten into

$$\frac{\alpha}{2} \sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \frac{\mu_2(\mathbf{v})^2}{p \mathbf{v}^\top \mathbb{E}[\mathbb{X}_{k,p}] \mathbf{v}} < 1, \quad (39)$$

and

$$0 < \alpha < \inf_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \frac{2 \mathbf{v}^\top \mathbb{E}(D_k \mathbb{X}_k D_k) \mathbf{v}}{\mathbb{E}\|D_k \mathbb{X}_k D_k \mathbf{v}\|_2^2}. \quad (40)$$

For $q = 2$, Lemma 29 in the Appendix states that the conclusion of the previous lemma is also implied by the condition $\mathbb{E}[2\mathbb{X}_k - \alpha \mathbb{X}_k^2] > 0$.

Remark 11 *The upper bound in condition (40) can be interpreted as a generalized eigenvalue. For a given pair of real symmetric matrices (A, B) and a given non-zero real vector \mathbf{u} , the generalized Rayleigh quotient is defined as*

$$R(A, B; \mathbf{u}) := \frac{\mathbf{u}^\top A \mathbf{u}}{\mathbf{u}^\top B \mathbf{u}}.$$

One can show that $\sup_{\mathbf{u} \neq 0} R(A, B; \mathbf{u}) = \max\{\lambda : \det(A - \lambda B) = 0\}$ is the largest generalized eigenvalue of (A, B) . Recall $\mathbb{X}_k = \mathbf{x}_k \mathbf{x}_k^\top$, let $\tilde{\mathbf{x}}_k = D_k \mathbf{x}_k$, and define the matrices

$$\Sigma_1 := \mathbb{E}(D_k \mathbb{X}_k D_k) = \mathbb{E}(\tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top), \quad \Sigma_2 := \mathbb{E}(D_k \mathbb{X}_k D_k \mathbb{X}_k D_k) = \mathbb{E}(\tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top \tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top) = \mathbb{E}(\tilde{\mathbf{x}}_k \tilde{\mathbf{x}}_k^\top)^2.$$

Letting $\mathbf{u} = \mathbf{v}$, $A = \Sigma_2$ and $B = \Sigma_1$, condition (40) is equivalent to

$$0 < \alpha < \inf_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \frac{2 \mathbf{v}^\top \Sigma_1 \mathbf{v}}{\mathbf{v}^\top \Sigma_2 \mathbf{v}} = \frac{2}{\sup_{\mathbf{v}} R(\Sigma_2, \Sigma_1; \mathbf{v})}.$$

For the geometric-moment contraction for the SGD dropout sequence, we impose the following moment conditions.

Assumption 1 (Finite moment) *Assume that for some $q \geq 2$, the random noises ϵ and the random sample \mathbf{x} in model (31) have finite $2q$ -th moment $\mathbb{E}[|\epsilon|^{2q}] + \|\mathbf{x}\|_2^{2q} < \infty$.*

Lemma 30 in the Appendix shows that this assumption ensures the finite q -th moment of the stochastic gradient in (32) evaluated at the true parameter β^* and the ℓ^2 -minimizer $\check{\beta}$ in model (31), that is,

$$\left(\mathbb{E} \left\| \nabla_{\beta^*} \frac{1}{2} (y - \mathbf{x}^\top D \beta^*)^2 \right\|_2^q \right)^{1/q} = \left(\mathbb{E} \|D \mathbf{x} (y - \mathbf{x}^\top D \beta^*)\|_2^q \right)^{1/q} < \infty,$$

and $(\mathbb{E} \|\nabla_{\check{\beta}} \frac{1}{2} (y - \mathbf{x}^\top D \check{\beta})^2\|_2^q)^{1/q} < \infty$. Assumption 1 moreover guarantees $(\mathbb{E} \|\nabla_{\beta^*} \frac{1}{2} (y - \mathbf{x}^\top D \beta^*)^2\|_2^q)^{1/q} < \infty$ for all $\beta \in \mathbb{R}^d$.

Theorem 12 (Geometric-moment contraction of SGD dropout) *Let $q \geq 2$. Suppose that Assumption 1 holds and the learning rate α satisfies (37). For two dropout sequences $\check{\beta}_k(\alpha)$ and $\check{\beta}'_k(\alpha)$, $k = 0, 1, \dots$, that are generated by the recursion (32) with the same sequence of dropout matrices $\{D_k\}_{k \in \mathbb{N}}$ but possibly different initializations $\check{\beta}_0, \check{\beta}'_0$, we have*

$$\left(\mathbb{E}\|\check{\beta}_k(\alpha) - \check{\beta}'_k(\alpha)\|_2^q\right)^{1/q} \leq \check{r}_{\alpha,q}^k \|\check{\beta}_0 - \check{\beta}'_0\|_2, \quad \text{for all } k = 1, 2, \dots, \quad (41)$$

with

$$\check{r}_{\alpha,q} = \left(\sup_{\mathbf{v} \in \mathbb{R}^d: \|\mathbf{v}\|_2=1} \mathbb{E}\|\check{A}_1(\alpha)\mathbf{v}\|_2^q\right)^{1/q} < 1.$$

Moreover, for any initial vector $\check{\beta}_0 \in \mathbb{R}^d$, there exists a unique stationary distribution $\check{\pi}_\alpha$ which does not depend on $\check{\beta}_0$, such that $\check{\beta}_k(\alpha) \Rightarrow \check{\pi}_\alpha$ as $k \rightarrow \infty$.

By Theorem 12, initializing $\check{\beta}_0^\circ \sim \check{\pi}_\alpha$ leads to the stationary SGD dropout sequence $\{\check{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ by following the recursion

$$\check{\beta}_k^\circ(\alpha) - \check{\beta} = \check{A}_k(\alpha)(\check{\beta}_{k-1}^\circ(\alpha) - \check{\beta}) + \check{\mathbf{b}}_k(\alpha), \quad k = 1, 2, \dots, \quad (42)$$

where the ℓ^2 -regularized minimizer $\check{\beta}$ is defined in (34), and the random coefficients $\check{A}_k(\alpha) = I_d - \alpha D_k \mathbb{X}_k D_k$ and $\check{\mathbf{b}}_k(\alpha) = \alpha D_k \mathbf{x}_k (y_k - \mathbf{x}_k^\top D_k \check{\beta})$ are defined in (36). Furthermore, recall the iterated random function $\check{f}_{D,(y,\mathbf{x})}(\beta)$ defined in (33). As a direct consequence of Theorem 12, we have

$$\check{\beta}_k^\circ(\alpha) = \check{f}_{D_k,(y_k,\mathbf{x}_k)}(\check{\beta}_{k-1}^\circ(\alpha)), \quad (43)$$

which holds for all $k \in \mathbb{Z}$. To see the case with $k \leq 0$, we only need to notice that, for any $\beta \in \mathbb{R}^d$, we have the limit

$$\check{\beta}_k^\circ(\alpha) := \lim_{m \rightarrow \infty} \check{f}_{\xi_k} \circ \dots \circ \check{f}_{\xi_{k-m}}(\beta) =: \check{h}_\alpha(\xi_k, \xi_{k-1}, \dots), \quad (44)$$

where \check{h}_α is a measurable function that depends on α , and we use ξ_k to denote all the new-coming random parts in the k -th iteration, that is,

$$\xi_k = (D_k, (y_k, \mathbf{x}_k)), \quad k \in \mathbb{Z}. \quad (45)$$

For $k \leq 0$, ξ_k can be viewed as an i.i.d. copy of ξ_j for some $j \geq 1$. The limit $\check{h}_\alpha(\xi_k, \xi_{k-1}, \dots)$ exists almost surely and does not depend on β . Therefore, the SGD dropout iteration $\check{\beta}_k^\circ(\alpha) = \check{f}_{D_k,(y_k,\mathbf{x}_k)}(\check{\beta}_{k-1}^\circ(\alpha))$ in (43) holds for all $k \in \mathbb{Z}$.

As for GD with dropout, the parameter dimension d influences the contraction constant $\check{r}_{\alpha,q}$ only indirectly through the matrix $\mathbb{E}[\mathbf{x}_k \mathbf{x}_k^\top]$. For i.i.d. Gaussian entries a more explicit expression is derived in the next lemma.

Lemma 13 (Dependence on d in SGD with dropout) *If the covariate vectors \mathbf{x}_k are independently drawn from a standard multivariate Gaussian, then,*

$$\check{r}_{\alpha,2}^2 = 1 - 2\alpha p + \alpha^2[(d-1)p^2 + 3p]. \quad (46)$$

Hence, $\alpha < 2/((d-1)p + 3)$ ensures $\check{r}_{\alpha,2} < 1$.

4.3 Asymptotics of Dropout in SGD

In this section, we provide the asymptotics for the k -th iterate of SGD dropout and the Ruppert-Polyak averaged version.

Lemma 14 (Moment convergence of iterative SGD dropout) *Let $q \geq 2$ and suppose that Assumption 1 holds. For the stationary SGD dropout sequence $\{\check{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ defined in (42) with learning rate α satisfying (37), we have*

$$\max_k (\mathbb{E} \|\check{\beta}_k^\circ(\alpha) - \check{\beta}\|_2^q)^{1/q} = O(\sqrt{\alpha}). \quad (47)$$

Besides the stochastic order of the last iterate of SGD dropout $\check{\beta}_k(\alpha)$, we are also interested in the limiting distribution of the Ruppert-Polyak averaged SGD dropout, which can effectively reduce the variance and keep the online computing scheme. In particular, we define

$$\bar{\beta}_k^{\text{sgd}}(\alpha) = \frac{1}{k} \sum_{i=1}^k \check{\beta}_i(\alpha). \quad (48)$$

Theorem 15 (Quenched CLT of averaged SGD dropout) *Under Assumption 1, if the learning rate α satisfies (37), then,*

$$\sqrt{k}(\bar{\beta}_k^{\text{sgd}}(\alpha) - \check{\beta}) \Rightarrow \mathcal{N}(0, \check{\Sigma}(\alpha)), \quad \text{as } k \rightarrow \infty, \quad (49)$$

with $\check{\Sigma}(\alpha) := \sum_{i=-\infty}^{\infty} \mathbb{E}[(\check{\beta}_0^\circ(\alpha) - \check{\beta})(\check{\beta}_i^\circ(\alpha) - \check{\beta})^\top]$ the long-run covariance matrix of the stationary process $\check{\beta}_i^\circ(\alpha) \sim \tilde{\pi}_\alpha$.

The CLTs for fixed design (Theorem 7) and random design (Theorem 15) have different covariance matrices in the limit $k \rightarrow \infty$. In particular, the limiting covariance matrix of the averaged SGD iterates with dropout depends on the random design.

As discussed above Corollary 8, one can also choose different learning rates $\alpha_1, \dots, \alpha_s$ and then run the SGD dropout sequences $\check{\beta}_k(\alpha_1), \dots, \check{\beta}_k(\alpha_s)$ in parallel. For d -dimensional vectors $\mathbf{u}_1, \dots, \mathbf{u}_s$, recall that $\text{vec}(\mathbf{u}_1, \dots, \mathbf{u}_s) := (\mathbf{u}_1^\top, \dots, \mathbf{u}_s^\top)^\top$ is the ds -dimensional concatenation.

Corollary 16 (Quenched CLT of parallel averaged SGD dropout) *Under Assumption 1, consider a sequence of constant learning rates $\alpha_1, \dots, \alpha_s$, for $s \geq 1$, satisfying the condition in (37). Then, for any initial vectors $\check{\beta}_0(\alpha_1), \dots, \check{\beta}_0(\alpha_s)$,*

$$\sqrt{k} \cdot \text{vec}(\bar{\beta}_k^{\text{sgd}}(\alpha_1) - \check{\beta}, \dots, \bar{\beta}_k^{\text{sgd}}(\alpha_s) - \check{\beta}) \Rightarrow \mathcal{N}(0, \check{\Sigma}^{\text{vec}}), \quad \text{as } k \rightarrow \infty, \quad (50)$$

with the long-run covariance matrix

$$\check{\Sigma}^{\text{vec}} = \sum_{i=-\infty}^{\infty} \text{Cov}(\text{vec}(\check{\beta}_0^\circ(\alpha_1), \dots, \check{\beta}_0^\circ(\alpha_s)), \text{vec}(\check{\beta}_i^\circ(\alpha_1), \dots, \check{\beta}_i^\circ(\alpha_s))).$$

Theorem 17 (Quenched invariance principle of averaged SGD dropout) *Let Assumption 1 hold for some $q > 2$ and consider the learning rate α satisfying (37). Define a partial sum process $(\check{S}_i^\circ(\alpha))_{1 \leq i \leq t}$ for $t \in \mathbb{N}_+$ with*

$$\check{S}_i^\circ(\alpha) = \sum_{k=1}^i (\check{\beta}_k^\circ(\alpha) - \check{\beta}). \quad (51)$$

Then, there exists a (richer) probability space $(\check{\Omega}^, \check{\mathcal{A}}^*, \check{\mathbb{P}}^*)$ on which one can define d -dimensional random vectors $\check{\beta}_k^*$, the associated partial sum process $\check{S}_i^*(\alpha) = \sum_{k=1}^i (\check{\beta}_k^*(\alpha) - \check{\beta})$, and a Gaussian process $\check{G}_i^* = \sum_{k=1}^i \check{z}_k^*$, with independent Gaussian random vectors $\check{z}_k \sim \mathcal{N}(0, I_d)$, such that $(\check{S}_i^\circ)_{1 \leq i \leq t} \stackrel{D}{=} (\check{S}_i^*)_{1 \leq i \leq t}$ and*

$$\max_{1 \leq i \leq t} \|\check{S}_i^* - \check{\Sigma}^{1/2}(\alpha) \check{G}_i^*\|_2 = o_{\mathbb{P}}(t^{1/q}), \quad \text{in } (\check{\Omega}^*, \check{\mathcal{A}}^*, \check{\mathbb{P}}^*), \quad (52)$$

where $\check{\Sigma}(\alpha)$ is the long-run covariance matrix defined in Theorem 15. In addition, this approximation holds for all $(\check{S}_i^{\check{\beta}_0}(\alpha))_{1 \leq i \leq t}$ given any arbitrary initialization $\check{\beta}_0 \in \mathbb{R}^d$, where

$$\check{S}_i^{\check{\beta}_0}(\alpha) = \sum_{k=1}^i (\check{\beta}_k(\alpha) - \check{\beta}). \quad (53)$$

5. Online Inference for SGD with Dropout

The long-run covariance matrix $\check{\Sigma}(\alpha)$ of the averaged SGD dropouts is usually unknown and needs to be estimated. We now propose an online estimation method for $\check{\Sigma}(\alpha)$, and establish theoretical guarantees.

The key idea is to adopt the non-overlapping batched means (NBM) method (Lahiri, 1999, 2003; Xiao and Wu, 2011), which resamples blocks of observations to estimate the long-run covariance of dependent data. Essentially, a sequences of non-overlapping blocks are pre-specified. When the block sizes are large enough, usually increasing as the the sample size grows, the block sums shall behave similar to independent observations and therefore can be used to estimate the long-run covariance. In this paper, to facilitate the online inference of the dependent SGD dropout iterates $\{\check{\beta}_k(\alpha)\}_{k \in \mathbb{N}}$, we shall extend the offline NBM estimators to online versions by only including the past SGD dropout iterates in each batch. The overlapped batch-means (OBM) methods are also investigated in literature; see for example Xiao and Wu (2011); Zhu et al. (2023). We shall only focus on the NBM estimates in this study given its simpler structure.

Let η_1, η_2, \dots be a strictly increasing integer-valued sequence satisfying $\eta_1 = 1$ and $\eta_{m+1} - \eta_m \rightarrow \infty$ as $m \rightarrow \infty$. For each m , we let B_m denote the block

$$B_m = \{\eta_m, \eta_m + 1, \dots, \eta_{m+1} - 1\}. \quad (54)$$

For the k -th SGD dropout iteration, denote by $\psi(k)$ the largest index m such that $\eta_m \leq k$. For any d -dimensional vector $\mathbf{v} = (v_1, \dots, v_d)^\top$, the Kronecker product is the $d \times d$ matrix $\mathbf{v}^{\otimes 2} = (v_i v_j)_{i,j=1}^d$. Based on the non-overlapping blocks $\{B_m\}_{m \in \mathbb{N}}$, for the k -th iteration, we

can estimate the long-run covariance matrix $\check{\Sigma}(\alpha)$ in Theorem 15 by

$$\hat{\Sigma}_k(\alpha) = \frac{1}{k} \sum_{m=1}^{\psi(k)-1} \left(\sum_{i \in B_m} [\check{\beta}_i(\alpha) - \bar{\beta}_k^{\text{sgd}}(\alpha)] \right)^{\otimes 2} + \frac{1}{k} \left(\sum_{i=\eta_{\psi(k)}}^k [\check{\beta}_i(\alpha) - \bar{\beta}_k^{\text{sgd}}(\alpha)] \right)^{\otimes 2}. \quad (55)$$

The estimator $\hat{\Sigma}_k(\alpha)$ is composed of two parts. The first part takes the sum within each block and then estimates the sample covariances of these centered block sums. The second part accounts for the remaining observations, which can be viewed as the estimated covariance of the tail block.

For the recursive computation of $\hat{\Sigma}_k(\alpha)$, we need to rewrite (55) such that, in the k -th iteration, we can update $\hat{\Sigma}_k(\alpha)$ based on the information from the $(k-1)$ -th step and the latest iterate $\check{\beta}_k(\alpha)$. To this end, we denote the number of iterates included in the tail part (i.e., the second part) in (55) by

$$\delta_\eta(k) = k - \eta_{\psi(k)} + 1, \quad (56)$$

and define two partial sums

$$\mathcal{S}_m(\alpha) = \sum_{i \in B_m} \check{\beta}_i(\alpha) \quad \text{and} \quad \mathcal{R}_k(\alpha) = \sum_{i=\eta_{\psi(k)}}^k \check{\beta}_i(\alpha). \quad (57)$$

Then, we notice that the estimator $\hat{\Sigma}_k(\alpha)$ in (55) can be rewritten as follows,

$$\begin{aligned} \hat{\Sigma}_k(\alpha) &= \frac{1}{k} \left[\left(\sum_{m=1}^{\psi(k)-1} \mathcal{S}_m(\alpha)^{\otimes 2} + \mathcal{R}_k(\alpha)^{\otimes 2} \right) + \left(\sum_{m=1}^{\psi(k)-1} |B_m|^2 + |\delta_\eta(k)|^2 \right) \bar{\beta}_k^{\text{sgd}}(\alpha)^{\otimes 2} \right. \\ &\quad - \left(\sum_{m=1}^{\psi(k)-1} |B_m| \mathcal{S}_m(\alpha) + \delta_\eta(k) \mathcal{R}_k(\alpha) \right) \bar{\beta}_k^{\text{sgd}}(\alpha)^\top \\ &\quad \left. - \bar{\beta}_k^{\text{sgd}}(\alpha) \left(\sum_{m=1}^{\psi(k)-1} |B_m| \mathcal{S}_m(\alpha) + \delta_\eta(k) \mathcal{R}_k(\alpha) \right)^\top \right] \\ &=: \frac{1}{k} \left[\mathcal{V}_k(\alpha) + K_k \bar{\beta}_k^{\text{sgd}}(\alpha)^{\otimes 2} - H_k(\alpha) \bar{\beta}_k^{\text{sgd}}(\alpha)^\top - \bar{\beta}_k^{\text{sgd}}(\alpha) H_k(\alpha)^\top \right]. \quad (58) \end{aligned}$$

As such, the estimation of $\check{\Sigma}(\alpha)$ reduces to computing $\{\mathcal{V}_k(\alpha), K_k, H_k(\alpha), \bar{\beta}_k^{\text{sgd}}(\alpha)\}$ recursively with respect to k . We provided the pseudo codes of the recursion in Algorithm 1. We shall further establish the convergence rate of the proposed online estimator $\hat{\Sigma}_k(\alpha)$ in Theorem 18.

The rationale behind Algorithm 1 is as follows: if $k+1 < \eta_{\psi(k)+1}$, then the index $k+1$ still belongs to the block $B_{\psi(k)}$ and $\psi(k+1) = \psi(k)$. Also we have $\mathcal{R}_{k+1}(\alpha) = \mathcal{R}_k(\alpha) + \check{\beta}_{k+1}(\alpha)$ and $\delta_\eta(k+1) = \delta_\eta(k) + 1$. Consequently, $\{K_{k+1}, \mathcal{V}_{k+1}(\alpha), H_{k+1}(\alpha)\}$ can be recursively updated via

$$\begin{aligned} K_{k+1} &= K_k - |\delta_\eta(k)|^2 + |\delta_\eta(k+1)|^2, \\ \mathcal{V}_{k+1}(\alpha) &= \mathcal{V}_k(\alpha) - \mathcal{R}_k(\alpha)^{\otimes 2} + \mathcal{R}_{k+1}(\alpha)^{\otimes 2}, \\ H_{k+1}(\alpha) &= H_k(\alpha) - \delta_\eta(k) \mathcal{R}_k(\alpha) + \delta_\eta(k+1) \mathcal{R}_{k+1}(\alpha). \end{aligned}$$

Algorithm 1: Online estimation of long-run covariance matrices of ASGD dropout

Data: Sequential random samples $(y_1, \mathbf{x}_1), \dots, (y_k, \mathbf{x}_k)$; sequential dropout matrices D_1, \dots, D_k ; constant learning rate α ; predefined sequences $\{\eta_m\}_{m \in \mathbb{N}}$

Result: ASGD dropout $\bar{\beta}_{k+1}^{\text{sgd}}(\alpha)$; estimated long-run covariance matrix $\hat{\Sigma}_{k+1}(\alpha)$

Initialize $\check{\beta}_0(\alpha) = \bar{\beta}_0^{\text{sgd}}(\alpha) = \mathcal{R}_0(\alpha) \leftarrow 0$,
 $\psi(0) \leftarrow 1, \delta_\eta(0) \leftarrow 1, K_0 = H_0(\alpha) \leftarrow 1, \mathcal{V}_0(\alpha) \leftarrow 0$

for $k = 0, 1, 2, 3, \dots$ **do**

$\check{\beta}_{k+1}(\alpha) \leftarrow \check{\beta}_k(\alpha) + \alpha D_k \mathbf{x}_k (y_k - \mathbf{x}_k^\top D_k \check{\beta}_{k-1}(\alpha))$; /* SGD dropout */

$\bar{\beta}_{k+1}^{\text{sgd}}(\alpha) \leftarrow \{k \bar{\beta}_k^{\text{sgd}}(\alpha) + \check{\beta}_{k+1}(\alpha)\} / (k+1)$; /* ASGD dropout */

if $k+1 < \eta_{\psi(k)+1}$ **then**

$\mathcal{R}_{k+1}(\alpha) \leftarrow \mathcal{R}_k(\alpha) + \check{\beta}_{k+1}(\alpha), \delta_\eta(k+1) \leftarrow \delta_\eta(k) + 1$;

$K_{k+1} \leftarrow K_k - \delta_\eta^2(k) + \delta_\eta^2(k+1), \psi(k+1) \leftarrow \psi(k)$;

$H_{k+1}(\alpha) \leftarrow H_k(\alpha) - \delta_\eta(k) \mathcal{R}_k(\alpha) + \delta_\eta(k+1) \mathcal{R}_{k+1}(\alpha)$;

$\mathcal{V}_{k+1}(\alpha) \leftarrow \mathcal{V}_k(\alpha) - \mathcal{R}_k(\alpha)^{\otimes 2} + \mathcal{R}_{k+1}(\alpha)^{\otimes 2}$;

else

$\mathcal{R}_{k+1}(\alpha) \leftarrow \check{\beta}_{k+1}(\alpha), \delta_\eta(k+1) \leftarrow 1$;

$\psi(k+1) \leftarrow \psi(k) + 1$;

$K_{k+1} \leftarrow K_k + 1, H_{k+1}(\alpha) \leftarrow H_k(\alpha) + \mathcal{R}_{k+1}(\alpha)$;

$\mathcal{V}_{k+1}(\alpha) \leftarrow \mathcal{V}_k(\alpha) + \mathcal{R}_{k+1}(\alpha)^{\otimes 2}$;

end

$\hat{\Sigma}_{k+1}(\alpha) \leftarrow$
 $[\mathcal{V}_{k+1}(\alpha) + K_{k+1} \bar{\beta}_{k+1}^{\text{sgd}}(\alpha)^{\otimes 2} - H_{k+1}(\alpha) \bar{\beta}_{k+1}^{\text{sgd}}(\alpha)^\top - \bar{\beta}_{k+1}^{\text{sgd}}(\alpha) H_{k+1}(\alpha)^\top] / (k+1)$;

/* Estimated long-run covariance matrix */

end

Otherwise, if $k+1 = \eta_{\psi(k)}$, we have $\psi(k+1) = \psi(k) + 1$. Hence $\mathcal{R}_{k+1}(\alpha) = \check{\beta}_{k+1}(\alpha)$ and $\delta_\eta(k+1) = 1$. In this case, $\{K_{k+1}, \mathcal{V}_{k+1}(\alpha), H_{k+1}(\alpha)\}$ can be recursively updated as follows,

$$\begin{aligned} K_{k+1} &= K_k + 1, \\ \mathcal{V}_{k+1}(\alpha) &= \mathcal{V}_k(\alpha) + \mathcal{R}_{k+1}(\alpha)^{\otimes 2}, \\ H_{k+1}(\alpha) &= H_k(\alpha) + \mathcal{R}_{k+1}(\alpha). \end{aligned}$$

As such, given $\check{\beta}_1(\alpha), \dots, \check{\beta}_k(\alpha)$, the estimator $\hat{\Sigma}_k(\alpha)$ for the long-run covariance matrix $\check{\Sigma}(\alpha)$ can be updated online, requiring only $O(1)$ memory storage.

Theorem 18 (Precision of $\hat{\Sigma}_k(\alpha)$) *Let $\eta_m = \lfloor cm^\zeta \rfloor$ for some $c > 0$ and $\zeta > 1$. If the conditions of Theorem 12 are satisfied with $q \geq 4$, we have*

$$\mathbb{E} \|\hat{\Sigma}_k(\alpha) - \check{\Sigma}(\alpha)\|_F \lesssim dk^{(1/\zeta - 1)\vee(-1/(2\zeta))},$$

where $\|\cdot\|_F$ denotes the Frobenius norm, and the constants in \lesssim are independent of k and d . Consequently, for any dimension d , if the number of iterations k satisfies

$$k \gtrsim d^{\frac{\zeta}{\zeta-1} \vee \frac{1}{2\zeta}},$$

then $\mathbb{E}\|\hat{\Sigma}_k(\alpha) - \check{\Sigma}(\alpha)\|_F = o(1)$.

In particular, for $\zeta = 3/2$,

$$\mathbb{E}\|\hat{\Sigma}_k(\alpha) - \check{\Sigma}(\alpha)\|_F \lesssim dk^{-1/3}. \quad (59)$$

This rate is optimal among long-run covariance estimators in the fixed-dimensional settings, even when comparing to offline estimation; see Xiao and Wu (2011) for details. This result indicates that for the dimension d , at least $n \gtrsim d^3$ training budget is needed to ensure consistency. The slow scaling of the algorithm in the dimension makes it, however, challenging to run the procedure for $d \gtrsim 1000$ requiring $\gtrsim 10^9$ training samples.

By the estimation procedure summarized in Algorithm 1, we can asymptotically estimate the long-run covariance matrix of the SGD dropout iterates $\bar{\beta}_k^{\text{sgd}}(\alpha)$ for any arbitrarily fixed initial vector. For some given confidence level $\omega \in (0, 1)$, in the k -th iteration the online confidence interval for each coordinate $\check{\beta}_j$, $j = 1, \dots, d$, of the vector $\check{\beta}$ in (34) is

$$\text{CI}_{\omega,k,j} := \left[\bar{\beta}_{k,j}^{\text{sgd}}(\alpha) - z_{1-\omega/2} \sqrt{\hat{\sigma}_{k,jj}(\alpha)/k}, \bar{\beta}_{k,j}^{\text{sgd}}(\alpha) + z_{1-\omega/2} \sqrt{\hat{\sigma}_{k,jj}(\alpha)/k} \right], \quad (60)$$

with $z_{1-\omega/2}$ denoting the $(1 - \omega/2)$ -percentile of the standard normal distribution. Here, $\hat{\sigma}_{k,jj}(\alpha)$ is the j -th diagonal of the proposed online long-run covariance estimator $\hat{\Sigma}_k(\alpha)$ in (55), and $\bar{\beta}_{k,j}^{\text{sgd}}(\alpha)$ is the j -th coordinate of the averaged SGD dropout estimate $\bar{\beta}_k^{\text{sgd}}(\alpha)$. Furthermore, the online joint confidence regions for the vector $\check{\beta}$ is

$$\text{CI}_{\omega,k} := \left\{ \beta \in \mathbb{R}^d : k(\bar{\beta}_k^{\text{sgd}}(\alpha) - \beta)^\top \hat{\Sigma}_k^{-1}(\alpha) (\bar{\beta}_k^{\text{sgd}}(\alpha) - \beta) \leq \chi_{d,1-\omega/2}^2 \right\}, \quad (61)$$

where $\chi_{d,1-\omega/2}^2$ is the $(1 - \omega/2)$ -percentile of the χ_d^2 distribution with d degrees of freedom.

Corollary 19 (Asymptotic coverage probability) *Suppose that Assumption 1 holds and the learning rate α satisfies (37). Given $\omega \in (0, 1)$ and $\eta_m = \lfloor cm^\zeta \rfloor$ for some $c > 0$ and $\zeta > 1$, $\text{CI}_{\omega,k,j}$ defined in (60), and $\text{CI}_{\omega,k}$ defined in (61) are asymptotic $100(1 - \omega)\%$ confidence intervals, that is, $\mathbb{P}(\check{\beta}_j \in \text{CI}_{\omega,k,j}) \rightarrow 1 - \omega$ for all $j = 1, \dots, d$, and $\mathbb{P}(\check{\beta} \in \text{CI}_{\omega,k}) \rightarrow 1 - \omega$, as $k \rightarrow \infty$. More generally, for any d -dimensional unit-length vector \mathbf{v} with $\|\mathbf{v}\|_2 = 1$, and $z_{1-\omega/2}$ the $(1 - \omega/2)$ -quantile of the standard normal distribution,*

$$\text{CI}_{\omega,k}^{\text{Proj}} := \left[\mathbf{v}^\top \bar{\beta}_k^{\text{sgd}}(\alpha) - z_{1-\omega/2} \sqrt{k^{-1}(\mathbf{v}^\top \hat{\Sigma}_k(\alpha) \mathbf{v})}, \mathbf{v}^\top \bar{\beta}_k^{\text{sgd}}(\alpha) + z_{1-\omega/2} \sqrt{k^{-1}(\mathbf{v}^\top \hat{\Sigma}_k(\alpha) \mathbf{v})} \right] \quad (62)$$

is an asymptotic $100(1 - \omega)\%$ confidence interval for the one-dimensional projection $\mathbf{v}^\top \check{\beta}$, that is, $\mathbb{P}(\mathbf{v}^\top \check{\beta} \in \text{CI}_{\omega,k}^{\text{Proj}}) \rightarrow 1 - \omega$, as $k \rightarrow \infty$.

By the quenched CLT of the averaged SGD dropout sequence $\{\bar{\beta}_k^{\text{sgd}}(\alpha)\}_{k \in \mathbb{N}}$ in Theorem 17 and the consistency of $\hat{\Sigma}_k(\alpha)$ in Theorem 18, we can apply Slutsky's theorem and obtain the results in Corollary 19. In Section 6, we shall validate the proposed online inference method by examining the estimation accuracy of the proposed online estimator $\hat{\Sigma}_k(\alpha)$ and the coverage probability of $\text{CI}_{\omega,k}^{\text{Proj}}$ under different settings.

6. Numerical Experiments

In this section, we present the results of the numerical experiments to demonstrate the validity of the proposed online inference methodology. For a real-data application see Appendix H. The codes for reproducing all results and figures can be found online¹.

6.1 On the Range of the Learning Rate

The GD dropout iterates can be defined via the recursion (6), $\tilde{\beta}_k(\alpha) - \tilde{\beta} = A_k(\alpha)(\tilde{\beta}_{k-1}(\alpha) - \tilde{\beta}) + \mathbf{b}_k(\alpha)$, and the derived theory requires the learning rate α to satisfy $\alpha\|\mathbb{X}\| < 2$. Via a simulation study we show that this range is close to sharp to guarantee that the contraction constant

$$r_{\alpha,2}^2 = \sup_{\mathbf{v} \in \mathbb{R}^d: \|\mathbf{v}\|_2=1} \mathbb{E}\|A_1(\alpha)\mathbf{v}\|_2^2 = \lambda_{\max}(\mathbb{E}[A_1^\top(\alpha)A_1(\alpha)]) < 1. \quad (63)$$

This then indicates that the condition $\alpha\|\mathbb{X}\| < 2$ in Lemma 2 is fairly sharp.

For the $n \times d$ full design matrix X , we independently generate each entry of X from the standard normal distribution. Since $\mathbb{X} = X^\top X$, the upper bound $2/\lambda_{\max}(\mathbb{X})$ of the learning rate α can be computed. Then, we independently generate $N = 500$ dropout matrices D_i , $i = 1, \dots, N$ with retaining probability p . The simulation study evaluates the empirical contraction constant

$$\hat{r}_{\alpha,2}^2 := \lambda_{\max}\left(N^{-1} \sum_{i=1}^N A_i^\top(\alpha)A_i(\alpha)\right), \quad (64)$$

for different training budget n , dimension d , retaining probability p , and learning rate α . Table 2 shows that even if the learning rate α exceeds the upper bound $2/\lambda_{\max}(\mathbb{X})$ by a

n, d	$2/\lambda_{\max}(\mathbb{X})$	p	α	$\hat{r}_{\alpha,2}^2$
100, 5	0.0151	0.9	0.0150	0.97
			0.0154	1.02
100, 50	0.0068	0.9	0.0067	0.93
			0.0072	1.01
		0.8	0.0068	0.90
100, 100	0.0052	0.9	0.0075	1.02
			0.0050	0.93
		0.5	0.0057	1.15
			0.0050	0.94
			0.0075	1.06

Table 2: Effects of the learning rate α on the geometric moment contraction of the GD iterates with dropout.

small margin, the contraction will not hold any more since $\hat{r}_{\alpha,2}^2 > 1$. This indicates that the condition $\alpha\|\mathbb{X}\| < 2$ is close to sharp.

1. https://github.com/jiaqili97/Dropout_SGD

6.2 Estimation of Long-Run Covariance Matrix

In this section, we provide the simulation results of the proposed long-run covariance estimator $\hat{\Sigma}_k(\alpha)$ defined in (55), and its online version (58).

Figure 1 shows the convergence of the GD and SGD iterates with dropout. The coordinates of the true regression vector β^* are equidistantly spaced between 0 and 1. One can see that the initialization is quickly forgotten in both GD and SGD algorithms.

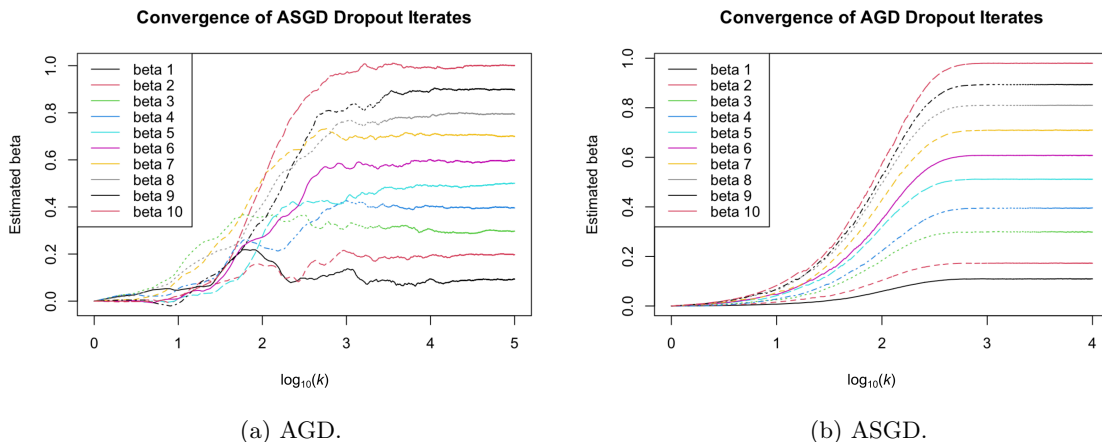


Figure 1: Convergence traces of AGD and ASGD iterates with dropout regularization based on a single run, with dimension $d = 10$ and initialization at zero. The coordinates of the true parameter β^* are equidistantly spaced between 0 and 1, the learning rate $\alpha = 0.01$, and the retaining probability $p = 0.9$. Each curve represents the convergence trace of one coordinate.

Figures 2–3 evaluate the performance of the online long-run covariance matrix estimator $\hat{\Sigma}_k(\alpha)$. Suppose that the training budget is n . The parameters in the blocks $B_m = \{\eta_m, \eta_{m+1}, \dots, \eta_{m+1} - 1\}$, $m = 1, \dots, M$, defined in (54), are chosen as $M = \lfloor \sqrt{n} \rfloor$ and $\eta_m = m^2$. In Figure 2, we can see that the long-run variances of each coordinate of the ASGD dropout iterates $\bar{\beta}_k^{\text{sgd}}(\alpha)$ converges as the number of iterations k grows. Since there is no closed-form expression for the true long-run covariance matrix $\check{\Sigma}(\alpha)$ defined in Theorem 15, in Figure 3, we set the true long-run covariance matrix as the average of the 10^8 -th iteration over 20 replications of the experiment, while we plot the first 10^7 iterations only. Figure 3(a) shows that the empirical convergence of the Frobenius norm $\mathbb{E}\|\hat{\Sigma}_k(\alpha) - \check{\Sigma}(\alpha)\|_F$ for dimension d between 100 and 1000 agrees with the theoretical convergence rate $\propto k^{-1/3}$ in (59). For $v^\top = d^{-1/2}(1, \dots, 1)^\top$, the length of the confidence interval for the one-dimensional projection $v^\top \beta^*$ is displayed in Figure 3(b). The same v is used for the simulation results summarized in Table 3. In the next section, we shall show that by using these estimated long-run variances, the online confidence intervals achieve asymptotically the nominal coverage probability.

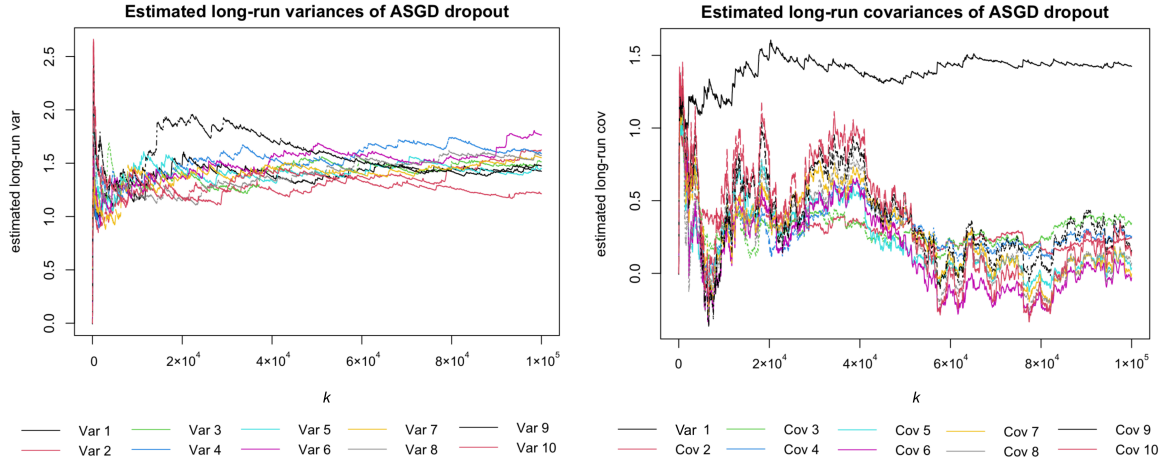


Figure 2: Under the same setting as in Figure 1, the left panel displays the estimated long-run variances of ASGD dropout iterates, that is, the diagonal entries of the estimated long-run covariance matrix $\hat{\Sigma}_k(\alpha)$; the right panel displays the first row of $\hat{\Sigma}_k(\alpha)$.

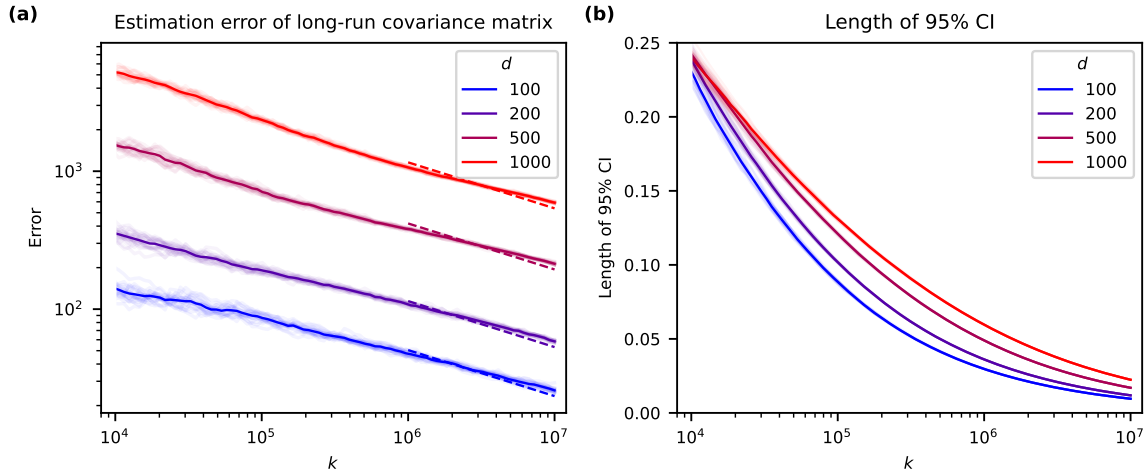


Figure 3: Performance of the long-run covariance matrix estimator $\hat{\Sigma}_k(\alpha)$ for learning rate $\alpha = 1/d$ and $p = 0.9$. (a) Empirical estimation errors $\mathbb{E}\|\hat{\Sigma}_k(\alpha) - \hat{\Sigma}(\alpha)\|_F$ averaged over 20 replications. The dashed lines show the theoretical decay rate $\propto k^{-1/3}$. The coordinates of β^* are equidistantly spaced between 0 and 1. (b) Length of the confidence intervals $\text{CI}_{\omega,k}^{\text{Proj}}$ in (62).

6.3 Online Confidence Intervals of ASGD Dropout Iterates

Recall the $100(1 - \omega)\%$ online confidence interval $\text{CI}_{\omega,k}^{\text{Proj}}$ in (62) for the one-dimensional projection of the true parameter β^* , i.e., $v^\top \beta^*$. For dimension $d \in \{50, 100, 200, 500\}$, a similar performance in convergence of the coverage probabilities is observed in Figure 4. In Table 3, we report the coverage probabilities of the confidence intervals $\text{CI}_{\omega,k}^{\text{Proj}}$ for dimension $d \in \{50, 100, 200, 500\}$, dropout probabilities $p \in \{0.5, 0.9\}$, and constant learning rates α ranging from 0.002 to 0.02. The results demonstrate that the online confidence intervals can achieve the target coverage probability as the number of iterations increases.

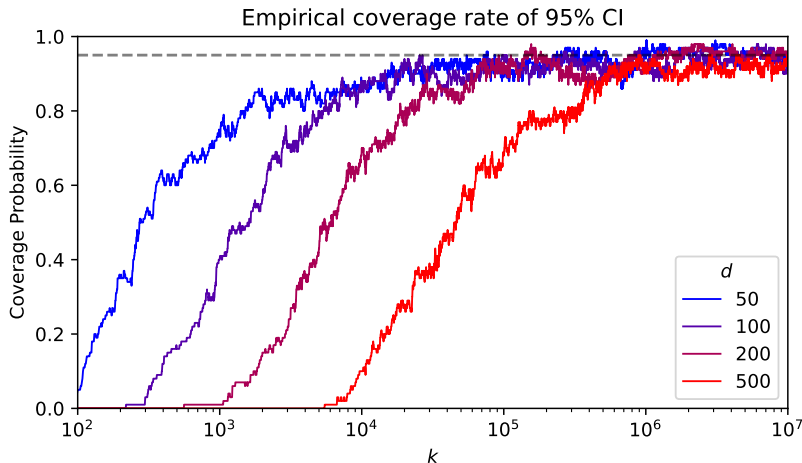


Figure 4: Coverage probabilities of 95% confidence intervals for one-dimensional projection of ASGD dropout iterates based on 100 independent repetitions, dropout probability $p = 0.9$, learning rate $\alpha = 1/d$, and coordinates of β^* equidistantly spaced between 0 and 1 and initialized at zero. The gray dashed line shows the nominal 0.95 coverage.

$p = 0.9$					
d	$k = 10^5$	$k = 5 * 10^5$	$k = 10^6$	$k = 5 * 10^6$	$k = 10^7$
50	0.91 (0.0286)	0.95 (0.0218)	0.95 (0.0218)	0.95 (0.0218)	0.92 (0.0271)
100	0.94 (0.0237)	0.92 (0.0271)	0.91 (0.0286)	0.94 (0.0237)	0.92 (0.0271)
200	0.91 (0.0286)	0.89 (0.0313)	0.92 (0.0271)	0.93 (0.0255)	0.96 (0.0196)
500	0.68 (0.0466)	0.89 (0.0313)	0.92 (0.0271)	0.93 (0.0255)	0.95 (0.0218)
$p = 0.5$					
d	$k = 10^5$	$k = 5 * 10^5$	$k = 10^6$	$k = 5 * 10^6$	$k = 10^7$
50	0.94 (0.0237)	0.92 (0.0271)	0.92 (0.0271)	0.95 (0.0218)	0.96 (0.0196)
100	0.94 (0.0237)	0.96 (0.0196)	0.93 (0.0255)	0.94 (0.0237)	0.96 (0.0196)
200	0.81 (0.0392)	0.91 (0.0286)	0.93 (0.0255)	0.95 (0.0218)	0.93 (0.0255)
500	0.55 (0.0497)	0.86 (0.0347)	0.90 (0.0300)	0.96 (0.0196)	0.96 (0.0196)

Table 3: Empirical coverage probabilities of 95% confidence intervals based on 100 independent repetitions (with standard errors in the brackets) and learning rate $\alpha = 1/d$.

Acknowledgments

We are grateful to the referees and the associate editor for constructive feedback that improved the work considerably. Moreover, the authors want to thank Haoxiong Yan for the support of large-scale computation, and thank Gabriel Clara for several helpful suggestions. Jiaqi Li’s research is partially supported by the NSF (Grant NSF/DMS-2515926). Johannes Schmidt-Hieber has received funding from the Dutch Research Council (NWO) via the Vidi grant VI.Vidi.192.021. Wei Biao Wu’s research is partially supported by the NSF (Grants NSF/DMS-2311249, NSF/DMS-2027723).

References

- Raman Arora, Peter Bartlett, Poorya Mianjy, and Nathan Srebro. Dropout: Explicit Forms and Capacity Control. In *Proceedings of the 38th International Conference on Machine Learning*, pages 351–361. PMLR, 2021.
- Pierre Baldi and Peter Sadowski. Understanding Dropout. In *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013.
- Roja Bandari, Sitaram Asur, and Bernardo Huberman. The pulse of news in social media: forecasting popularity. *Proceedings of the International AAAI Conference on Web and Social Media*, 6(1):26–33, August 2021.
- István Berkes, Weidong Liu, and Wei Biao Wu. Komlós–Major–Tusnády approximation under dependence. *The Annals of Probability*, 42(2):794–817, 2014.
- Andreas Brandt. The Stochastic Equation $Y_{n+1} = A_n Y_n + B_n$ with Stationary Coefficients. *Advances in Applied Probability*, 18(1):211–220, 1986.
- Peter J Brockwell and Richard A Davis. *Time Series: Theory and Methods*. Springer, second edition edition, 1991.
- Donald L. Burkholder. Sharp inequalities for martingales and stochastic integrals. In *Colloque Paul Lévy sur les processus stochastiques*, number 157-158 in Astérisque, pages 75–94. Société mathématique de France, 1988.
- Jacopo Cavazza, Pietro Morerio, Benjamin Haeffele, Connor Lane, Vittorio Murino, and Rene Vidal. Dropout as a Low-Rank Regularizer for Matrix Factorization. In *Proceedings of the 21st International Conference on Artificial Intelligence and Statistics*, pages 435–444. PMLR, 2018.
- Likai Chen and Wei Biao Wu. Stability and asymptotics for autoregressive processes. *Electronic Journal of Statistics*, 10(2):3723–3751, 2016.
- Xi Chen, Jason D. Lee, Xin T. Tong, and Yichen Zhang. Statistical inference for model parameters in stochastic gradient descent. *The Annals of Statistics*, 48(1):251–273, 2020.
- Giovanni Luca Ciampaglia, Azadeh Nematzadeh, Filippo Menczer, and Alessandro Flammini. How algorithmic popularity bias hinders or promotes quality. *Scientific Reports*, 8(1):15951, October 2018.

- Gabriel Clara, Sophie Langer, and Johannes Schmidt-Hieber. Dropout regularization versus l_2 -penalization in the linear model. *Journal of Machine Learning Research*, 25(204):1–48, 2024.
- Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Marc’ aurelio Ranzato, Andrew Senior, Paul Tucker, Ke Yang, Quoc Le, and Andrew Ng. Large Scale Distributed Deep Networks. In *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012.
- Aaron Defazio, Ashok Cutkosky, Harsh Mehta, and Konstantin Mishchenko. Optimal Linear Decay Learning Rate Schedules and Further Refinements. *arXiv preprint*, October 2024. arXiv:2310.07831.
- Persi Diaconis and David Freedman. Iterated Random Functions. *SIAM Review*, 41(1): 45–76, 1999.
- Alain Durmus, Eric Moulines, Alexey Naumov, Sergey Samsonov, Kevin Scaman, and Hoi-To Wai. Tight High Probability Bounds for Linear Stochastic Approximation with Fixed Stepsize. In *Advances in Neural Information Processing Systems*, volume 34, pages 30063–30074. Curran Associates, Inc., 2021.
- Yixin Fang. Scalable statistical inference for averaged implicit stochastic gradient descent. *Scandinavian Journal of Statistics*, 46(4):987–1002, 2019.
- Yixin Fang, Jinfeng Xu, and Lei Yang. Online Bootstrap Confidence Intervals for the Stochastic Gradient Descent Estimator. *Journal of Machine Learning Research*, 19:1–21, 2019.
- Kelwin Fernandes, Pedro Vinagre, Paulo Cortez, and Pedro Sernadela. Online News Popularity. UCI Machine Learning Repository, 2015. DOI: <https://doi.org/10.24432/C5NS3V>.
- Yarin Gal and Zoubin Ghahramani. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. In *Proceedings of the 33rd International Conference on Machine Learning*, pages 1050–1059. PMLR, 2016a.
- Yarin Gal and Zoubin Ghahramani. A Theoretically Grounded Application of Dropout in Recurrent Neural Networks. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016b.
- Wei Gao and Zhi-Hua Zhou. Dropout Rademacher complexity of deep neural networks. *Science China Information Sciences*, 59(7):072104, 2016.
- Geoffrey E. Hinton, Nitish Srivastava, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint*, 2012. arXiv:1207.0580.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. SCAFFOLD: Stochastic Controlled Averaging for Federated Learning. In *Proceedings of the 37th International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020.

- Sayar Karmakar and Wei Biao Wu. Optimal Gaussian Approximation for Multiple Time Series. *Statistica Sinica*, 30(3):1399–1417, 2020.
- Edmund Kay and Anurag Agarwal. DropConnected neural network trained with diverse features for classifying heart sounds. In *2016 Computing in Cardiology Conference (CinC)*, pages 617–620, 2016.
- J. Kiefer and J. Wolfowitz. Stochastic Estimation of the Maximum of a Regression Function. *The Annals of Mathematical Statistics*, 23(3):462–466, 1952.
- J. Komlós, P. Major, and G. Tusnády. An approximation of partial sums of independent RV’s, and the sample DF. I. *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 32(1):111–131, 1975.
- J. Komlós, P. Major, and G. Tusnády. An approximation of partial sums of independent RV’s, and the sample DF. II. *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 34(1):33–58, 1976.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. In *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012.
- Soumendra N. Lahiri. Theoretical comparisons of block bootstrap methods. *The Annals of Statistics*, 27(1):386–404, 1999.
- Soumendra N. Lahiri. *Resampling Methods for Dependent Data*. Springer Series in Statistics. Springer, New York, NY, 2003.
- Tengyuan Liang and Weijie J. Su. Statistical Inference for the Population Landscape via Moment-Adjusted Stochastic Gradients. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 81(2):431–456, 2019.
- David McAllester. A PAC-Bayesian tutorial with a dropout bound. *arXiv preprint*, 2013. arXiv:1307.2118.
- Poorya Mianjy and Raman Arora. On Dropout and Nuclear Norm Regularization. In *Proceedings of the 36th International Conference on Machine Learning*, pages 4575–4584. PMLR, 2019.
- Poorya Mianjy and Raman Arora. On Convergence and Generalization of Dropout Training. In *Advances in Neural Information Processing Systems*, volume 33, pages 21151–21161. Curran Associates, Inc., 2020.
- Poorya Mianjy, Raman Arora, and Rene Vidal. On the Implicit Bias of Dropout. In *Proceedings of the 35th International Conference on Machine Learning*, pages 3540–3548. PMLR, 2018.
- Georg Ch. Pflug. Stochastic Minimization with Constant Step-Size: Asymptotic Laws. *SIAM Journal on Control and Optimization*, 24(4):655–666, 1986.

- Hieu Pham and Quoc Le. AutoDropout: Learning Dropout Patterns to Regularize Deep Networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(11):9351–9359, 2021.
- B. T. Polyak and A. B. Juditsky. Acceleration of Stochastic Approximation by Averaging. *SIAM Journal on Control and Optimization*, 30(4):838–855, 1992.
- Emmanuel Rio. Moment inequalities for sums of dependent random variables under projective conditions. *Journal of Theoretical Probability*, 22(1):146–163, 2009.
- Herbert Robbins and Sutton Monro. A Stochastic Approximation Method. *The Annals of Mathematical Statistics*, 22(3):400–407, 1951.
- David Ruppert. Efficient Estimations from a Slowly Convergent Robbins-Monro Process, 1988. Technical report, Cornell University Operations Research and Industrial Engineering.
- Sergey Samsonov, Daniil Tiapkin, Alexey Naumov, and Eric Moulines. Improved High-Probability Bounds for the Temporal Difference Learning Algorithm via Exponential Stability. In *Proceedings of Thirty Seventh Conference on Learning Theory*, pages 4511–4547. PMLR, June 2024.
- Albert Senen-Cerda and Jaron Sanders. Asymptotic Convergence Rate of Dropout on Shallow Linear Neural Networks. *SIGMETRICS Perform. Eval. Rev.*, 50(1):105–106, 2022.
- Albert Senen-Cerda and Jaron Sanders. Almost Sure Convergence of Dropout Algorithms for Neural Networks. *arXiv preprint*, 2023. arXiv:2002.02247.
- Nitish Srivastava, Geoffrey E. Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research*, 15(56):1929–1958, 2014.
- Weijie Su and Yuancheng Zhu. HiGrad: Uncertainty Quantification for Online Learning and Stochastic Approximation. *Journal of Machine Learning Research*, 24:1–53, 2023.
- Gabor Szabo and Bernardo A. Huberman. Predicting the popularity of online content. *Commun. ACM*, 53(8):80–88, August 2010.
- Stefan Wager, Sida Wang, and Percy S Liang. Dropout Training as Adaptive Regularization. In *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013.
- Li Wan, Matthew Zeiler, Sixin Zhang, Yann Le Cun, and Rob Fergus. Regularization of Neural Networks using DropConnect. In *Proceedings of the 30th International Conference on Machine Learning*, pages 1058–1066. PMLR, 2013.
- Colin Wei, Sham Kakade, and Tengyu Ma. The Implicit and Explicit Regularization Effects of Dropout. In *Proceedings of the 37th International Conference on Machine Learning*, pages 10181–10192. PMLR, 2020.

- Haibing Wu and Xiaodong Gu. Towards dropout training for convolutional neural networks. *Neural Networks*, 71:1–10, 2015.
- Wei Biao Wu. Nonlinear system theory: Another look at dependence. *PNAS*, 102(40):14150–14154, 2005.
- Wei Biao Wu. Asymptotic theory for stationary processes. *Statistics and Its Interface*, 4(2):207–226, 2011.
- Wei Biao Wu and Xiaofeng Shao. Limit Theorems for Iterated Random Functions. *Journal of Applied Probability*, 41(2):425–436, 2004.
- Han Xiao and Wei Biao Wu. A Single-Pass Algorithm for Spectrum Estimation With Fast Convergence. *IEEE Transactions on Information Theory*, 57(7):4720–4731, 2011.
- Ke Zhai and Huan Wang. Adaptive dropout with rademacher complexity regularization. In *International Conference on Learning Representations*, 2018.
- Yanjie Zhong, Jiaqi Li, and Soumendra N. Lahiri. Probabilistic guarantees of stochastic recursive gradient in non-convex finite sum problems. In *Advances in Knowledge Discovery and Data Mining*, pages 142–154. Springer Nature Singapore, 2024.
- Wanrong Zhu, Xi Chen, and Wei Biao Wu. Online Covariance Matrix Estimation in Stochastic Gradient Descent. *Journal of the American Statistical Association*, 118(541):393–404, 2023.
- Martin Zinkevich, Markus Weimer, Lihong Li, and Alex Smola. Parallelized Stochastic Gradient Descent. In *Advances in Neural Information Processing Systems*, volume 23. Curran Associates, Inc., 2010.

Appendix A. Technical Lemmas

Lemma 20 (Burkholder, 1988; Rio, 2009) *Let $q > 1$, $q' = \min\{q, 2\}$, and define $M_T = \sum_{t=1}^T \xi_t$, where ξ_t are martingale differences with a finite q -th moment. Then*

$$(\mathbb{E}\|M_T\|_2^q)^{q'/q} \leq K_q^{q'} \sum_{t=1}^T (\mathbb{E}\|\xi_t\|_2^q)^{q'/q}, \quad \text{where } K_q = \max\{(q-1)^{-1}, \sqrt{q-1}\}.$$

Lemma 21 (Clara et al., 2024) *For any matrices A and B in $\mathbb{R}^{d \times d}$, $p \in (0, 1)$, and a diagonal matrix $D \in \mathbb{R}^{d \times d}$, the following results hold:*

(i) $\overline{AD} = \overline{AD}$, $\overline{DA} = D\overline{A}$, and $\overline{A_p} = p\overline{A} = \overline{A_p}$;

If in addition, the diagonal matrix D is random and independent of A and B , with the diagonal entries satisfying $D_{ii} \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(p)$, $1 \leq i \leq d$, then,

(ii) $\mathbb{E}[DAD] = pA_p$, where $A_p = pA + (1-p)\text{Diag}(A)$;

(iii) $\mathbb{E}[DADBD] = pA_p B_p + p^2(1-p)\text{Diag}(\overline{AB})$, where $\overline{A} = A - \text{Diag}(A)$;

(iv) $\mathbb{E}[DADBD C D] = pA_p B_p C_p + p^2(1-p)[\text{Diag}(\overline{AB_p C}) + A_p \text{Diag}(\overline{BC}) + \text{Diag}(A\overline{B})C_p + (1-p)A \odot \overline{B}^\top \odot C]$, where \odot denotes the Hadamard product.

Lemma 22 (Properties of operator norm) *Let $A = (a_{ij})_{1 \leq i, j \leq d}$ be a real $d \times d$ matrix. View A as a linear map $\mathbb{R}^d \mapsto \mathbb{R}^d$ and denote its operator norm by $\|A\|$.*

(i) *(Inequalities for variants of A). $\|\text{Diag}(A)\| \leq \|A\|$, $\|A_p\| \leq \|A\|$, and if in addition, A is positive semi-definite, then also $\|\overline{A}\| \leq \|A\|$;*

(ii) *(Frobenius norm). $\|A\| = \sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \|A\mathbf{v}\|_2 \leq \|A\|_F$, where $\|A\|_F$ denotes the Frobenius norm, i.e., $\|A\|_F = (\sum_{i,j=1}^d |a_{ij}|^2)^{1/2}$;*

(iii) *(Largest magnitude of eigenvalues). For a symmetric $d \times d$ matrix A , we have $\max_{1 \leq i \leq d} |\lambda_i(A)| = \|A\|$, where $\lambda_i(A)$ denotes the i -th largest eigenvalue of A . If in addition, A is positive semi-definite, then also $\lambda_{\max}(A) = \|A\|$.*

Proof The inequalities in (i) follow directly from Lemma 19 in Clara et al. (2024). For (ii), we notice that for any unit vector $\mathbf{v} \in \mathbb{R}^d$, one can find a basis $\{\mathbf{e}_1, \dots, \mathbf{e}_d\}$ and write \mathbf{v} into $\mathbf{v} = \sum_{j=1}^d c_j \mathbf{e}_j$, with $\mathbf{e}_j \in \mathbb{R}^d$, and the real coefficients c_j satisfying $\sum_{j=1}^d c_j^2 = 1$. Then, it follows from the orthogonality of \mathbf{e}_j and the Cauchy-Schwarz inequality that

$$\begin{aligned} \|A\mathbf{v}\|_2^2 &= \left\| \sum_{j=1}^d c_j A\mathbf{e}_j \right\|_2^2 \leq \left(\sum_{j=1}^d |c_j| \|A\mathbf{e}_j\|_2 \right)^2 \leq \left(\sum_{j=1}^d |c_j|^2 \right) \left(\sum_{j=1}^d \|A\mathbf{e}_j\|_2^2 \right) \\ &= \sum_{j=1}^d \|A\mathbf{e}_j\|_2^2 = \|A\|_F^2. \end{aligned} \tag{65}$$

Since this result holds for any unit vector $\mathbf{v} \in \mathbb{R}^d$, the desired result in (ii) is achieved.

(iii) is well-known. Nevertheless, we provide a proof here for completeness. For any eigenvalue $\lambda_i(A)$, denote its associated unit eigenvector by \mathbf{v} . Then, $\|A\mathbf{v}\|_2 = |\lambda_i(A)| \|\mathbf{v}\|_2 = |\lambda_i(A)|$, which further yields $|\lambda_i(A)| \leq \sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \|A\mathbf{v}\|_2 = \|A\|$, uniformly over i . Hence, the inequality can be obtained. If in addition, A is symmetric, then A can be diagonalized by an orthogonal matrix Q and a diagonal matrix Λ such that $A = Q^\top \Lambda Q$.

Therefore, $\sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \|A\mathbf{v}\|_2 = \max_{1 \leq i \leq d} |\lambda_i(A)|$, which completes the proof. \blacksquare

Appendix B. Proofs in Section 3.1

This section is devoted to the proofs of the geometric-moment contraction (GMC) for the dropout iterates with gradient descent (GD), i.e., $\{\tilde{\beta}_k(\alpha) - \beta\}_{k \in \mathbb{N}}$. We first extend the results in Wu and Shao (2004) to the cases where the inputs of iterated random functions are i.i.d. random matrices. Then, we present the proof for the sufficient condition of the GMC in terms of the constant learning rate α in Lemma 2, and showcase the GMC of $\{\tilde{\beta}_k(\alpha) - \beta\}_{k \in \mathbb{N}}$ in Theorem 3.

B.1 GMC – Random Matrix Version

Let (\mathcal{Y}, ρ) be a complete and separable metric space, endowed with its Borel sets \mathbb{Y} . Consider an iterated random function on the state space $\mathcal{Y} \subset \mathbb{R}^d$, for some fixed $d \geq 1$, with the form

$$\mathbf{y}_k = f(\mathbf{y}_{k-1}, X_k) = f_{X_k}(\mathbf{y}_{k-1}), \quad k \in \mathbb{N}, \quad (66)$$

where $f(\mathbf{y}, \cdot)$ is the \mathbf{y} -section of a jointly measurable function $f : \mathcal{Y} \times \mathcal{X} \mapsto \mathcal{Y}$; the random matrices X_k , $k \in \mathbb{N}$, take values in a second measurable space $\mathcal{X} \subset \mathbb{R}^{d \times d}$, and are independently distributed with identical marginal distribution H . The initial point $\mathbf{y}_0 \in \mathcal{Y}$ is independent of all X_k .

We are interested in the sufficient conditions on $f_X(\mathbf{y})$ such that there is a unique stationary probability π on \mathcal{Y} with $\mathbf{y}_k \Rightarrow \pi$ as $k \rightarrow \infty$. To this end, define a composite function

$$\mathbf{y}_k(\mathbf{y}) = f_{X_k} \circ f_{X_{k-1}} \circ \cdots \circ f_{X_1}(\mathbf{y}), \quad \text{for } \mathbf{y} \in \mathcal{Y}. \quad (67)$$

We say that \mathbf{y}_k is *geometric-moment contracting* if for any two independent random vectors $\mathbf{y} \sim \pi$ and $\mathbf{y}' \sim \pi$ in \mathcal{Y} , there exist some $q > 0$, $C_q > 0$ and $r_q \in (0, 1)$, such that for all $k \in \mathbb{N}$,

$$\mathbb{E}[\rho^q(\mathbf{y}_k(\mathbf{y}), \mathbf{y}_k(\mathbf{y}'))] \leq C_q r_q^k. \quad (68)$$

Wu and Shao (2004) provided the sufficient conditions for (68) when X_k are random variables and \mathbf{y}_k and \mathbf{y} are one-dimensional. Their results can be directly extended to the random matrix version and we state them here for the completeness of this paper.

Assumption 2 (Finite moment) *Assume that there exists a fixed vector $\mathbf{y}^* \in \mathcal{Y}$ and some $q > 0$ such that*

$$I(q, \mathbf{y}^*) := \mathbb{E}_{X \sim H}[\rho^q(\mathbf{y}^*, f_X(\mathbf{y}^*))] = \int_{\mathcal{X}} \rho^q(\mathbf{y}^*, f_X(\mathbf{y}^*)) H(dX) < \infty.$$

Assumption 3 (Stochastic Lipschitz continuity) *Assume that there exists some $q > 0$ and some $\mathbf{y}_0 \in \mathcal{Y}$ such that*

$$L_q := \sup_{\mathbf{y}_0 \in \mathcal{Y}, \mathbf{y}_0 \neq \mathbf{y}'_0} \frac{\mathbb{E}_{X \sim H}[\rho^q(f_X(\mathbf{y}_0), f_X(\mathbf{y}'_0))]}{\rho^q(\mathbf{y}_0, \mathbf{y}'_0)} < 1,$$

where $L_q = L_q(\mathbf{y}_0)$ is a local Lipschitz constant.

Because of the affine form of the recursion in (36), the ergodic theory in Diaconis and Freedman (1999) also applies here for the SGD sequence with dropout.

Corollary 23 (GMC – random matrix version) *Suppose that Assumptions 2 and 3 hold. Define a backward iteration process*

$$\mathbf{z}_k(\mathbf{y}) = \mathbf{z}_{k-1}(f_{X_k}(\mathbf{y})) = f_{X_1} \circ f_{X_2} \circ \cdots \circ f_{X_k}(\mathbf{y}), \quad \text{for } \mathbf{y} \in \mathcal{Y}. \quad (69)$$

Then, $\mathbf{z}_k(\mathbf{y}) \stackrel{\mathcal{D}}{=} \mathbf{y}_k(\mathbf{y})$, and there exists a random vector $\mathbf{z}_\infty \in \mathcal{Y}$ such that for any $\mathbf{y} \in \mathcal{Y}$,

$$\mathbf{z}_k(\mathbf{y}) \xrightarrow{\text{a.s.}} \mathbf{z}_\infty.$$

The limit \mathbf{z}_∞ is measurable with respect to the σ -algebra $\sigma(X_1, X_2, \dots)$ and does not depend on \mathbf{y} . In addition,

$$\mathbb{E}[\rho^q(\mathbf{z}_k(\mathbf{y}), \mathbf{z}_\infty)] \leq Cr_q^k, \quad \text{for all } k \in \mathbb{N}, \quad (70)$$

where the constant $C > 0$ only depends on q and \mathbf{y}^* in Assumption 2, L_q in Assumption 3 and \mathbf{y}_0 . Additionally, (68) holds.

Before showing the detailed proof of Corollary 23, we first introduce a key technique for this proof in the following remark.

Remark 24 (Backward iteration) *We shall comment on the intuition for defining the backward iteration \mathbf{z}_k in (69). Recall the i.i.d. random samples X_1, \dots, X_n . Clearly, for any fixed initial point $\mathbf{y}_0 \in \mathcal{Y}$, for all $k \in \mathbb{N}$, we have the relations*

$$\begin{aligned} \mathbf{y}_{k+1}(\mathbf{y}_0) &= f_{X_{k+1}}(\mathbf{y}_k(\mathbf{y}_0)), \\ \mathbf{z}_{k+1}(\mathbf{y}_0) &= \mathbf{z}_k(f_{X_{k+1}}(\mathbf{y}_0)). \end{aligned}$$

To prove the existence of the limit for $\mathbf{y}_k = f_{X_k} \circ f_{X_{k-1}} \circ \cdots \circ f_{X_1}(\mathbf{y}_0)$, we need to make use of the contracting property of the function $f_X(\cdot)$ stated in Assumption 3. However, we cannot directly apply it to the forward iteration, because by the Markov property, given the present position of the chain, the conditional distribution of the future does not depend on the past. This indicates

$$\mathbb{E}[\rho^q(\mathbf{y}_{k+1}(\mathbf{y}^*), \mathbf{y}_k(\mathbf{y}^*))] = \mathbb{E}[\mathbb{E}[\rho^q(f_{X_{k+1}}(\mathbf{y}_k(\mathbf{y}^*)), \mathbf{y}_k(\mathbf{y}^*)) \mid X_{k+1}]], \quad (71)$$

where the two parts inside of $\rho(\cdot, \cdot)$ are operated by two different functions, which are $f_{X_{k+1}}(\cdot)$ and $f_{X_k} \circ f_{X_{k-1}} \circ \cdots \circ f_{X_1}(\cdot)$ respectively. In fact, as pointed out by Diaconis and Freedman (1999), the forward iteration \mathbf{y}_k moves ergodically through \mathcal{Y} , which behaves quite differently from the backward iteration $\mathbf{z}_k(\cdot) = f_{X_1} \circ f_{X_2} \circ \cdots \circ f_{X_k}(\cdot)$ in (69), which does converge to a limit. To see this, we note that by Assumptions 2 and 3, there exists some $\mathbf{y}^* \in \mathcal{Y}$ such that

$$\begin{aligned} \mathbb{E}[\rho^q(\mathbf{z}_{k+1}(\mathbf{y}^*), \mathbf{z}_k(\mathbf{y}^*))] &= \mathbb{E}[\mathbb{E}[\rho^q(\mathbf{z}_k(f_{X_{k+1}}(\mathbf{y}^*)), \mathbf{z}_k(\mathbf{y}^*)) \mid X_{k+1}]] \\ &\leq L_q^k \mathbb{E}[\rho^q(f_{X_{k+1}}(\mathbf{y}^*), \mathbf{y}^*)] \\ &= L_q^k I(q, \mathbf{y}^*), \end{aligned} \quad (72)$$

which is summable over k by Assumption 3. Since \mathcal{Y} is a complete space, we can mimic the idea of a Cauchy sequence to prove the existence of the limit \mathbf{z}_∞ and further show $\mathbf{z}_k \xrightarrow{a.s.} \mathbf{z}_\infty$, by applying the Borel-Cantelli lemma. Since X_1, \dots, X_n are i.i.d. and thus exchangeable, we have $\mathbf{z}_k(\mathbf{y}_0) \stackrel{D}{=} \mathbf{y}_k(\mathbf{y}_0)$. Hence, we can show that \mathbf{y}_k also converges to \mathbf{z}_∞ in distribution.

Proof Let $q \in (0, 1]$ such that both Assumptions 2 and 3 hold. We will only show the desired results for this choice of q , since if Assumptions 2 and 3 are satisfied for some $q > 1$, then they are also valid for all $q \leq 1$ by Hölder's inequality (Wu and Shao, 2004). Recall the definition of integral $I(q, \mathbf{y}^*)$ in Assumption 2. Let $\mathbf{y}_0 \in \mathcal{Y}$ satisfy Assumption 3. Then,

$$\begin{aligned} I(q, \mathbf{y}_0) &= \mathbb{E}[\rho^q(\mathbf{y}_0, f_X(\mathbf{y}_0))] \\ &\leq \mathbb{E}[\rho(\mathbf{y}_0, \mathbf{y}^*) + \rho(\mathbf{y}^*, f_X(\mathbf{y}^*)) + \rho(f_X(\mathbf{y}^*), f_X(\mathbf{y}_0))]^q \\ &\leq \rho^q(\mathbf{y}_0, \mathbf{y}^*) + I(q, \mathbf{y}^*) + \mathbb{E}[\rho^q(f_X(\mathbf{y}^*), f_X(\mathbf{y}_0))] \\ &\leq \rho^q(\mathbf{y}_0, \mathbf{y}^*) + I(q, \mathbf{y}^*) + L_q \rho^q(\mathbf{y}^*, \mathbf{y}_0) < \infty, \end{aligned}$$

where the first inequality follows from the triangle inequality, the second one is by Assumption 3 and Jensen's inequality, and the last one is due to Assumption 2. A similar argument as in (72) yields

$$\mathbb{E}[\rho^q(\mathbf{z}_{k+1}(\mathbf{y}_0), \mathbf{z}_k(\mathbf{y}_0))] \leq L_q^k I(q, \mathbf{y}_0) =: \delta_k, \quad (73)$$

where $\delta_k = \delta_k(q, \mathbf{y}_0)$ solely depends on k , q , L_q and \mathbf{y}_0 . By Markov's inequality, we have

$$\mathbb{P}\left(\rho(\mathbf{z}_{k+1}(\mathbf{y}_0), \mathbf{z}_k(\mathbf{y}_0)) \geq \delta_k^{1/(2q)}\right) \leq \delta_k^{1/2}. \quad (74)$$

Since $\sum_{k=1}^{\infty} \delta_k^{1/2} < \infty$, it follows from the first Borel-Cantelli lemma that

$$\mathbb{P}\left(\rho(\mathbf{z}_{k+1}(\mathbf{y}_0), \mathbf{z}_k(\mathbf{y}_0)) \geq \delta_k^{1/(2q)} \text{ for infinitely many } k\right) = 0. \quad (75)$$

Again, since $\delta_k^{1/2}$ is summable, \mathbf{z}_k is a Cauchy sequence in space \mathcal{Y} , which together with the completeness of \mathcal{Y} gives that almost surely, there exists a random vector $\mathbf{z}_\infty \in \mathcal{Y}$ such that

$$\mathbf{z}_k(\mathbf{y}_0) \xrightarrow{a.s.} \mathbf{z}_\infty, \quad \text{as } k \rightarrow \infty,$$

where \mathbf{z}_∞ is $\sigma(X_1, X_2, \dots)$ -measurable. Let π be the probability distribution of \mathbf{z}_∞ .

Furthermore, it follows from the triangle inequality and Jensen's inequality that for any fixed $\mathbf{y}_0 \in \mathcal{Y}$,

$$\begin{aligned} \mathbb{E}[\rho^q(\mathbf{z}_k(\mathbf{y}_0), \mathbf{z}_\infty)] &\leq \mathbb{E}\left[\sum_{l=0}^{\infty} \rho(\mathbf{z}_{k+1+l}(\mathbf{y}_0), \mathbf{z}_{k+l}(\mathbf{y}_0))\right]^q \\ &\leq \sum_{l=0}^{\infty} \mathbb{E}[\rho^q(\mathbf{z}_{k+1+l}(\mathbf{y}_0), \mathbf{z}_{k+l}(\mathbf{y}_0))] \\ &\leq \delta_k / (1 - L_q), \end{aligned} \quad (76)$$

For any $\mathbf{y} \in \mathcal{Y}$, by Assumption 3 and triangle inequality,

$$\begin{aligned} \mathbb{E}[\rho^q(\mathbf{z}_k(\mathbf{y}), \mathbf{z}_\infty)] &\leq \mathbb{E}[\rho^q(\mathbf{z}_k(\mathbf{y}), \mathbf{z}_k(\mathbf{y}_0))] + \mathbb{E}[\rho^q(\mathbf{z}_k(\mathbf{y}_0), \mathbf{z}_\infty)] \\ &\leq L_q^k \rho^q(\mathbf{y}_0, \mathbf{y}) + \delta_k / (1 - L_q). \end{aligned} \quad (77)$$

Recall that $\delta_k = L_q^k I(q, \mathbf{y}_0)$. Let $C = I(q, \mathbf{y}_0) / (1 - L_q) + \rho^q(\mathbf{y}_0, \mathbf{y})$ and we have shown result (70) with $r_q = L_q$. Since $C r_q^k$ in (70) is summable over k , it again follows from Borel-Cantelli lemma that for any $\mathbf{y} \in \mathcal{Y}$,

$$\mathbf{z}_k(\mathbf{y}) \xrightarrow{a.s.} \mathbf{z}_\infty, \quad \text{as } k \rightarrow \infty,$$

and therefore, the limit

$$\mathbf{v}_k(\mathbf{y}) = \lim_{m \rightarrow \infty} f_{X_{k+1}} \circ f_{X_{k+2}} \circ \cdots \circ f_{X_{k+m}}(\mathbf{y}) \quad (78)$$

exists almost surely.

Finally, we notice that for any two independent random vectors $\mathbf{y} \sim \pi$ and $\mathbf{y}' \sim \pi$,

$$\begin{aligned} \mathbb{E}[\rho^q(\mathbf{y}_k(\mathbf{y}), \mathbf{y}_k(\mathbf{y}'))] &\leq \mathbb{E}[\rho^q(\mathbf{y}_k(\mathbf{y}), \mathbf{y}_k(\mathbf{y}_0))] + \mathbb{E}[\rho^q(\mathbf{y}_k(\mathbf{y}_0), \mathbf{y}_k(\mathbf{y}'))] \\ &= 2\mathbb{E}[\rho^q(\mathbf{z}_k(\mathbf{v}_k), \mathbf{z}_k(\mathbf{y}_0))] \\ &= 2\mathbb{E}[\rho^q(\mathbf{z}_\infty, \mathbf{z}_k(\mathbf{y}_0))] \leq 2\delta_k / (1 - L_q), \end{aligned} \quad (79)$$

where the first equation follows from the observation that \mathbf{v}_k has the identical distribution as $\mathbf{z}_\infty = \mathbf{z}_k(\mathbf{v}_k(\mathbf{y})) \sim \pi$ and is independent of i.i.d. random matrices X_1, \dots, X_k because \mathbf{v}_k as defined in (78) only depends on X_i for large $i \geq k + 1$. The desired result in (68) has been achieved. \blacksquare

The recursion $\mathbf{y}_k = f(\mathbf{y}_{k-1}, X_k)$ is only defined for positive integers k . Nevertheless, Corollary 23 guarantees that for $k = 0, -1, \dots$ the relation $\mathbf{y}_k = f(\mathbf{y}_{k-1}, X_k)$ also holds. See Remark 2 in Wu and Shao (2004) for a simple way to define \mathbf{y}_k when $k = 0, -1, \dots$ in the one-dimensional case. The vector versions can be similarly constructed.

B.2 Proof of Lemma 2

Proof Let D be a dropout matrix with the same distribution as D_1 . Since $\mathbb{X} = X^\top X$ is positive semi-definite and by assumption $\alpha \|\mathbb{X}\| < 2$, we have $-I_d < I_d - \alpha D \mathbb{X} D \leq I_d$ and consequently $\|I_d - \alpha D \mathbb{X} D\| \leq 1$. Thus for a unit vector \mathbf{v} , $\|(I_d - \alpha D \mathbb{X} D) \mathbf{v}\|_2 \leq 1$. This means that for $q \geq 2$, we can use $\|\cdot\|_2^q = \|\cdot\|_2^2 \cdot \|\cdot\|_2^{q-2}$ to bound

$$\begin{aligned} r_{\alpha, q}^q &\leq \sup_{\mathbf{v} \in \mathbb{R}^d: \|\mathbf{v}\|_2=1} \mathbb{E} \left\| (I_d - \alpha D \mathbb{X} D) \mathbf{v} \right\|_2^2 \\ &= \sup_{\mathbf{v} \in \mathbb{R}^d: \|\mathbf{v}\|_2=1} \mathbf{v}^\top \mathbb{E} \left[(I_d - \alpha D \mathbb{X} D)^2 \right] \mathbf{v} \\ &= \left\| \mathbb{E} \left[(I_d - \alpha D \mathbb{X} D)^2 \right] \right\|. \end{aligned} \quad (80)$$

For a $d \times d$ and positive semi-definite matrix A , we have $A^2 \leq \|A\|A$. To see this, let \mathbf{v}_j be the eigenvectors of A with corresponding eigenvalues λ_j . Any vector \mathbf{w} can be written as $\mathbf{w} = \gamma_1 \mathbf{v}_1 + \dots + \gamma_d \mathbf{v}_d$ with coefficients $\gamma_1, \dots, \gamma_d$. Now $\mathbf{w}^\top A^2 \mathbf{w} = \gamma_1^2 \lambda_1^2 + \dots + \gamma_d^2 \lambda_d^2 \leq (\max_j \lambda_j)(\gamma_1^2 \lambda_1 + \dots + \gamma_d^2 \lambda_d) = \mathbf{w}^\top \|A\|A \mathbf{w}$. Since \mathbf{w} was arbitrary, this proves $A^2 \leq \|A\|A$. Moreover, recall that D_k is a diagonal matrix with diagonal entries 0 and 1. Thus $D_k^2 = D_k \leq I_d$. Because \mathbb{X} is positive semi-definite and by assumption $\Delta := 2 - \alpha \|\mathbb{X}\| > 0$, we have $\alpha^2 D_1 \mathbb{X} D_1^2 \mathbb{X} D_1 \leq \alpha^2 D_1 \mathbb{X}^2 D_1 \leq \alpha^2 D_1 \|\mathbb{X}\| \mathbb{X} D_1 \leq (2 - \Delta) \alpha D_1 \mathbb{X} D_1$. Thus,

$$(I_d - \alpha D_1 \mathbb{X} D_1)^2 = I_d - 2\alpha D_1 \mathbb{X} D_1 + \alpha^2 D_1 \mathbb{X} D_1^2 \mathbb{X} D_1 \leq I_d - \Delta \alpha D_1 \mathbb{X} D_1.$$

Taking expectation and using Lemma 21 (ii) yields $\mathbb{E}[(I_d - \alpha D_1 \mathbb{X} D_1)^2] \leq I_d - \Delta \alpha p \mathbb{X}_p$. The fact that $\mathbb{E}[(I_d - \alpha D_1 \mathbb{X} D_1)^2]$ is positive semi-definite implies that $\|\mathbb{E}[(I_d - \alpha D_1 \mathbb{X} D_1)^2]\|$ is bounded by the largest eigenvalue of $I_d - \Delta \alpha p \mathbb{X}_p$. By definition, $\mathbb{X}_p = p \mathbb{X} + (1 - p) \text{Diag}(\mathbb{X}) \geq (1 - p) \min_j \mathbb{X}_{jj} I_d$. By assumption the design is in reduced form which implies that $\min_j \mathbb{X}_{jj} > 0$. This shows that \mathbb{X}_p is positive definite and the largest eigenvalue of $I_d - \Delta \alpha p \mathbb{X}_p$ must be strictly smaller than 1. This implies $\|\mathbb{E}[(I_d - \alpha D_1 \mathbb{X} D_1)^2]\| < 1$. Combined with (80) this proves $r_{\alpha, q} < 1$.

If Lemma 2 holds for some $q \geq 2$, then by Hölder's inequality, it also holds for all $1 < q < 2$. To see this, consider a unit vector $\mathbf{v} \in \mathbb{R}^d$ and set $r(q) := (\mathbb{E} \|(I_d - \alpha D_1 \mathbb{X} D_1) \mathbf{v}\|_2^q)^{1/q}$. Then for any $1 < q' < q$, it follows from Hölder's inequality that

$$r(q')^{q'} = \mathbb{E} \left\| (I_d - \alpha D_1 \mathbb{X} D_1) \mathbf{v} \right\|_2^{q'} \leq \left(\mathbb{E} \left\| (I_d - \alpha D_1 \mathbb{X} D_1) \mathbf{v} \right\|_2^q \right)^{q'/q} = r(q)^{q'} < 1. \quad (81)$$

The desired result is obtained. \blacksquare

Throughout the rest of this paper, when there is no ambiguity, we omit the dependence on α , writing $\tilde{\beta}_k$ instead of $\tilde{\beta}_k(\alpha)$, $A_k = A_k(\alpha)$, and $\mathbf{b}_k = \mathbf{b}_k(\alpha)$.

B.3 Proof of Theorem 3

Proof Recall the recursive estimator $\tilde{\beta}_k$ defined in (6). Write $A_k = A_k(\alpha) = I_d - \alpha D_k \mathbb{X} D_k$. We consider arbitrary d -dimensional initialization vectors $\tilde{\beta}_0, \tilde{\beta}'_0 \in \mathbb{R}^d$ and write $\tilde{\beta}_k$ and $\tilde{\beta}'_k$ for the respective iterates (sharing the same dropout matrices). Now $\tilde{\beta}_k - \tilde{\beta}'_k = A_k(\tilde{\beta}_{k-1} - \tilde{\beta}'_{k-1}) =: A_k \Delta_{k-1}$, with independent A_k and Δ_{k-1} . By Lemma 2, $r := \sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} (\mathbb{E} \|A_k \mathbf{v}\|_2^q)^{1/q} < 1$, and thus, for any fixed vector \mathbf{v} , $\mathbb{E} \|A_k \mathbf{v}\|_2^q \leq r^q \|\mathbf{v}\|_2^q$. Due to the independence between A_k and Δ_{k-1} , it follows from the tower rule and the condition above that

$$\mathbb{E} \|\tilde{\beta}_k - \tilde{\beta}'_k\|_2^q = \mathbb{E} \|A_k \Delta_{k-1}\|_2^q = \mathbb{E} [\mathbb{E} [\|A_k \Delta_{k-1}\|_2^q \mid \Delta_{k-1}]] \leq \mathbb{E} [r^q \|\Delta_{k-1}\|_2^q] = r^q \mathbb{E} \|\Delta_{k-1}\|_2^q.$$

Since $A_k(\alpha)$ are i.i.d. random matrices induction on k yields the claimed geometric-moment contraction $(\mathbb{E} \|\tilde{\beta}_k(\alpha) - \tilde{\beta}'_k(\alpha)\|_2^q)^{1/q} \leq r_{\alpha, q}^k \|\tilde{\beta}_0 - \tilde{\beta}'_0\|_2$.

Finally, recall the quantity $I(q, \mathbf{y}^*)$ defined in Assumption 2. For the recursion of the sequence $\{\tilde{\beta}_k(\alpha) - \tilde{\beta}'_k\}_{k \in \mathbb{N}}$ in (6), we have $I(q, 0) = \mathbb{E} \|\mathbf{b}_k(\alpha)\|_2^q < \infty$. To see this, note that

$$\|\mathbf{b}_k(\alpha)\|_2 = \alpha \|D_k \bar{\mathbb{X}} (p I_d - D_k) \tilde{\beta}\|_2 \leq \alpha \|\bar{\mathbb{X}}\| \|\tilde{\beta}\|_2,$$

where $\|\cdot\|$ denotes the operator norm, and the last inequality holds since $\|D_k\| < 1$, $\|pI_d - D_k\| \leq \max\{p, 1-p\} \leq 1$. By Corollary 23, the geometric-moment contraction proved above implies the existence of a unique stationary distribution $\tilde{\pi}_\alpha$ of the GD dropout sequence $\tilde{\beta}_k(\alpha)$. This completes the proof. \blacksquare

Appendix C. Proofs in Section 3.2

C.1 Proof of Lemma 4

Proof Since $\tilde{\beta}_k^\circ(\alpha)$ is stationary and $\mathbb{E}[\mathbf{b}_k] = 0$ by (100), it follows that $\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}$ and $\beta_k^\dagger(\alpha) - \tilde{\beta}$ both have zero mean, and thus $\mathbb{E}[\delta_k(\alpha)] = 0$.

To prove the second claim, we first note that

$$\begin{aligned} \delta_k(\alpha) &= (I_d - \alpha D_k \mathbb{X} D_k)(\tilde{\beta}_{k-1}^\circ(\alpha) - \tilde{\beta}) + \mathbf{b}_k(\alpha) - [(I_d - \alpha p \mathbb{X}_p)(\beta_{k-1}^\dagger(\alpha) - \tilde{\beta}) + \mathbf{b}_k(\alpha)] \\ &= (I_d - \alpha p \mathbb{X}_p)\delta_{k-1} + \alpha(p \mathbb{X}_p - D_k \mathbb{X} D_k)(\tilde{\beta}_{k-1}^\circ - \tilde{\beta}) \end{aligned} \quad (82)$$

is a stationary sequence. By induction on k , we can write $\delta_k(\alpha)$ into

$$\begin{aligned} \delta_k(\alpha) &= \alpha \left[(p \mathbb{X}_p - D_k \mathbb{X} D_k)(\tilde{\beta}_{k-1}^\circ - \tilde{\beta}) + \dots \right. \\ &\quad \left. + (p \mathbb{X}_p - D_1 \mathbb{X} D_1)(I_d - \alpha p \mathbb{X}_p)^{k-1}(\tilde{\beta}_0^\circ - \tilde{\beta}) + \dots \right] \\ &= \alpha \sum_{i=1}^{\infty} (p \mathbb{X}_p - D_{k-i+1} \mathbb{X} D_{k-i+1})(I_d - \alpha p \mathbb{X}_p)^{i-1}(\tilde{\beta}_{k-i}^\circ - \tilde{\beta}) \\ &=: \alpha \sum_{i=1}^{\infty} \mathcal{M}_{k-i}(\alpha). \end{aligned} \quad (83)$$

For any $k \in \mathbb{N}$, $\{\mathcal{M}_{k-i}(\alpha)\}_{i \geq 1}$ is a sequence of martingale differences with respect to the filtration $\mathcal{F}_{k-i} = \sigma(\dots, D_{k-i-1}, D_{k-i})$, since the dropout matrix D_k is independent of $\tilde{\beta}_{k-1}^\circ$ and $\tilde{\beta}$. Therefore, we can apply Burkholder's inequality in Lemma 20 to $\sum_{i=1}^{\infty} \mathcal{M}_{k-i}(\alpha)$, and obtain, for $q \geq 2$,

$$\begin{aligned} &\left(\mathbb{E} \left\| \sum_{i=1}^{\infty} \mathcal{M}_{k-i}(\alpha) \right\|_2^q \right)^{1/q} \\ &= \left(\mathbb{E} \left\| \sum_{i=1}^{\infty} (I_d - \alpha p \mathbb{X}_p)^{i-1} (p \mathbb{X}_p - D_{k-i+1} \mathbb{X} D_{k-i+1})(\tilde{\beta}_{k-i}^\circ - \tilde{\beta}) \right\|_2^q \right)^{1/q} \\ &\lesssim \left[\sum_{i=1}^{\infty} \left(\mathbb{E} \left\| (I_d - \alpha p \mathbb{X}_p)^{i-1} (p \mathbb{X}_p - D_{k-i+1} \mathbb{X} D_{k-i+1})(\tilde{\beta}_{k-i}^\circ - \tilde{\beta}) \right\|_2^q \right)^{2/q} \right]^{1/2} \\ &\leq \left[\sum_{i=1}^{\infty} \|I_d - \alpha p \mathbb{X}_p\|^{2(i-1)} \left(\mathbb{E} \left\| (p \mathbb{X}_p - D_{k-i+1} \mathbb{X} D_{k-i+1})(\tilde{\beta}_{k-i}^\circ - \tilde{\beta}) \right\|_2^q \right)^{2/q} \right]^{1/2}, \end{aligned}$$

where the constant in \lesssim here and the rest of the proof only depends on q unless it is additionally specified.

We shall proceed the proof with two main steps. First, we show the bound $\|I_d - \alpha p \mathbb{X}_p\| < 1$ for the operator norm and thus $\sum_{i=1}^{\infty} \|I_d - \alpha p \mathbb{X}_p\|^{2(i-1)} < \infty$. Second, we provide a bound for $\mathbb{E} \|(p \mathbb{X}_p - D_k \mathbb{X} D_k)(\tilde{\beta}_{k-1}^\circ - \tilde{\beta})\|_2^q$ uniformly over k .

Step 1. Since $\alpha \|\mathbb{X}\| < 2$, it follows from Lemma 22 (i) that $\alpha \|\mathbb{X}_p\| \leq \alpha \|\mathbb{X}\| < 2$. Moreover, the assumption that the design matrix X has no zero columns guarantees that all diagonal entries of \mathbb{X} are positive and thus $\text{Diag}(\mathbb{X}) > 0$. Together with $p < 1$, this lead to $\mathbb{X}_p = p \mathbb{X} + (1-p) \text{Diag}(\mathbb{X}) \geq (1-p) \text{Diag}(\mathbb{X}) > 0$. We thus have $-I_d < I_d - \alpha p \mathbb{X}_p < I_d$. Consequently, $\|I_d - \alpha p \mathbb{X}_p\| < 1$ and

$$\sum_{i=1}^{\infty} \|I_d - \alpha p \mathbb{X}_p\|^{2(i-1)} = \frac{1}{1 - \|I_d - \alpha p \mathbb{X}_p\|^2} = O(\alpha^{-1}). \quad (84)$$

Step 2. Next, we shall bound the term $\mathbb{E} \|(p \mathbb{X}_p - D_i \mathbb{X} D_i)(\tilde{\beta}_{i-1}^\circ - \tilde{\beta})\|_2^q$. We first consider the case $q = 2$. Denote $\mathbb{M}_i = p \mathbb{X}_p - D_i \mathbb{X} D_i$. Using that $\mathbb{E}[D_i \mathbb{X} D_i] = p \mathbb{X}_p$, we find $\mathbb{E}[\mathbb{M}_i] = 0$ and by the tower rule,

$$\begin{aligned} \mathbb{E} \|(p \mathbb{X}_p - D_i \mathbb{X} D_i)(\tilde{\beta}_{i-1}^\circ - \tilde{\beta})\|_2^2 &= \mathbb{E}[(\tilde{\beta}_{i-1}^\circ - \tilde{\beta})^\top \mathbb{M}_i^\top \mathbb{M}_i (\tilde{\beta}_{i-1}^\circ - \tilde{\beta})] \\ &= \mathbb{E}[\mathbb{E}[(\tilde{\beta}_{i-1}^\circ - \tilde{\beta})^\top \mathbb{M}_i^\top \mathbb{M}_i (\tilde{\beta}_{i-1}^\circ - \tilde{\beta}) \mid \mathcal{F}_{i-1}]] \\ &\leq \|\mathbb{E}[\mathbb{M}_i^\top \mathbb{M}_i]\| \cdot \mathbb{E} \|\tilde{\beta}_{i-1}^\circ - \tilde{\beta}\|_2^2. \end{aligned} \quad (85)$$

By Lemma 5, we have $\mathbb{E} \|\tilde{\beta}_{i-1}^\circ - \tilde{\beta}\|_2^2 = O(\alpha)$. We only need to bound the operator norm $\|\mathbb{E}[\mathbb{M}_i^\top \mathbb{M}_i]\|$. To this end, we use again $\mathbb{E}[D_i \mathbb{X} D_i] = p \mathbb{X}_p$ and moreover $\mathbb{E}[D_i \mathbb{X} D_i \mathbb{X} D_i] = p \mathbb{X}_p^2 + p^2(1-p) \text{Diag}(\overline{\mathbb{X}\mathbb{X}})$, which yields,

$$\begin{aligned} \mathbb{E}[\mathbb{M}_i^\top \mathbb{M}_i] &= \mathbb{E}[(p \mathbb{X}_p - D_i \mathbb{X} D_i)^\top (p \mathbb{X}_p - D_i \mathbb{X} D_i)] \\ &= p^2 \mathbb{X}_p^2 - 2p^2 \mathbb{X}_p^2 + p \mathbb{X}_p^2 + p^2(1-p) \text{Diag}(\overline{\mathbb{X}\mathbb{X}}) \\ &= p^2(1-p) \text{Diag}(\overline{\mathbb{X}\mathbb{X}}). \end{aligned} \quad (86)$$

Recall that $\mathbb{X} = X^\top X$, where X is the fixed design matrix. Then, by Lemma 22 (i) and the sub-multiplicativity of the operator norm, we have $\|\text{Diag}(\overline{\mathbb{X}\mathbb{X}})\| \leq \|\overline{\mathbb{X}\mathbb{X}}\| \leq \|\overline{\mathbb{X}}\| \|\mathbb{X}\| \leq \|\mathbb{X}\|^2$. As a direct consequence, $\|\mathbb{E}[\mathbb{M}_i^\top \mathbb{M}_i]\| \leq p^2(1-p) \|\mathbb{X}\|^2 < \infty$, which together with Lemma 5 and (84) gives

$$\mathbb{E} \|\delta_k(\alpha)\|_2 = \left(\mathbb{E} \left\| \alpha \sum_{i=1}^{\infty} \mathcal{M}_{k-i}(\alpha) \right\|_2^2 \right)^{1/2} \lesssim \alpha \left(\sum_{i=1}^{\infty} \|I_d - \alpha p \mathbb{X}_p\|^{2(i-1)} \alpha \right)^{1/2} = O(\alpha),$$

uniformly over k . For the case with $q > 2$, we can similarly apply the tower rule and obtain

$$\begin{aligned} &\mathbb{E} \|(p \mathbb{X}_p - D_i \mathbb{X} D_i)(\tilde{\beta}_{i-1}^\circ - \tilde{\beta})\|_2^q \\ &= \mathbb{E} \left[\mathbb{E} \left[\|(p \mathbb{X}_p - D_i \mathbb{X} D_i)(\tilde{\beta}_{i-1}^\circ - \tilde{\beta})\|_2^q \mid \mathcal{F}_{i-1} \right] \right] \\ &\leq \sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \mathbb{E} \|(p \mathbb{X}_p - D_i \mathbb{X} D_i) \mathbf{v}\|_2^q \cdot \mathbb{E} \|\tilde{\beta}_{i-1}^\circ - \tilde{\beta}\|_2^q, \end{aligned} \quad (87)$$

where the last inequality can be achieved by writing $\tilde{\beta}_{i-1}^\circ - \tilde{\beta} = \|\tilde{\beta}_{i-1}^\circ - \tilde{\beta}\|_2 \mathbf{v}$. Here \mathbf{v} is the unit vector $(\tilde{\beta}_{i-1}^\circ - \tilde{\beta})/\|\tilde{\beta}_{i-1}^\circ - \tilde{\beta}\|_2$ with $\|\mathbf{v}\|_2 = 1$. In addition, recall the Frobenius norm denoted by $\|\cdot\|_F$. It follows from Lemma 22 (i) and (ii) that

$$\begin{aligned} \sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \mathbb{E} \|(p\mathbb{X}_p - D_i \mathbb{X} D_i) \mathbf{v}\|_2^q &\leq \mathbb{E} \|p\mathbb{X}_p - D_i \mathbb{X} D_i\|_F^q \\ &\lesssim \mathbb{E} (\|p\mathbb{X}_p\|_F^q + \|D_i \mathbb{X} D_i\|_F^q) \\ &\leq (p^q + 1) \|\mathbb{X}\|^q < \infty, \end{aligned}$$

where the constant in \lesssim only depends on q . Combining this with the inequality (84), we obtain $(\mathbb{E} \|\delta_k(\alpha)\|_2^q)^{1/q} = O(\alpha)$, completing the proof. \blacksquare

C.2 Proof of Lemma 5

Proof Recall that by applying induction on k to Equation (6), we can rewrite the GD dropout iterates $\tilde{\beta}_k(\alpha)$ into

$$\begin{aligned} \tilde{\beta}_k(\alpha) - \tilde{\beta} &= A_k(\alpha)(\tilde{\beta}_{k-1}(\alpha) - \tilde{\beta}) + \mathbf{b}_k(\alpha) \\ &= \sum_{i=0}^{k-1} \left(\prod_{j=k-i+1}^k A_j(\alpha) \right) \mathbf{b}_{k-i}(\alpha) + \left(\prod_{j=1}^k A_j(\alpha) \right) (\tilde{\beta}_0(\alpha) - \tilde{\beta}), \end{aligned}$$

where we set $\prod_{j=k+1}^k A_j(\alpha) = I_d$. Following Brandt (1986), since both A_k and \mathbf{b}_k are i.i.d. random coefficients, the stationary solution $\{\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}\}_{k \in \mathbb{N}}$ of this recursion can be written into

$$\begin{aligned} \tilde{\beta}_k^\circ(\alpha) - \tilde{\beta} &= A_k(\alpha)(\tilde{\beta}_{k-1}^\circ(\alpha) - \tilde{\beta}) + \mathbf{b}_k(\alpha) \\ &= \sum_{i=0}^{\infty} \left(\prod_{j=k-i+1}^k A_j(\alpha) \right) \mathbf{b}_{k-i}(\alpha) \\ &= \alpha \sum_{i=0}^{\infty} \left[\prod_{j=k-i+1}^k (I_d - \alpha D_j \mathbb{X} D_j) \right] D_{k-i} \bar{\mathbb{X}} (pI_d - D_{k-i}) \tilde{\beta} \\ &=: \alpha \sum_{i=0}^{\infty} \tilde{\mathcal{M}}_{i,k}(\alpha). \end{aligned} \tag{88}$$

We observe that, for any $k \in \mathbb{N}$, $\{\tilde{\mathcal{M}}_{i,k}(\alpha)\}_{i \in \mathbb{N}}$ is a sequence of martingale differences with respect to the filtration $\mathcal{F}_{k-i} = \sigma(D_{k-i}, D_{k-i-1}, \dots)$. Hence, it follows from Burkholder's inequality in Lemma 20 that, for $q \geq 2$,

$$\begin{aligned} (\mathbb{E} \|\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}\|_2^q)^{1/q} &= \alpha \left(\mathbb{E} \left\| \sum_{i=0}^{\infty} \tilde{\mathcal{M}}_{i,k}(\alpha) \right\|_2^q \right)^{1/q} \\ &\lesssim \alpha \left(\sum_{i=0}^{\infty} (\mathbb{E} \|\tilde{\mathcal{M}}_{i,k}(\alpha)\|_2^q)^{2/q} \right)^{1/2}, \end{aligned} \tag{89}$$

where the constant in \lesssim only depends on q . Recall H_k defined in (96), and we define a $d \times d$ matrix $B_{i,k}$ by

$$B_{i,k} = \left[\prod_{j=k-i+1}^k (I_d - \alpha D_j \mathbb{X} D_j) \right] D_{k-i} \overline{\mathbb{X}} (pI_d - D_{k-i}) = \left(\prod_{j=k-i+1}^k A_j \right) H_{k-i}. \quad (90)$$

This random matrix is independent of $\tilde{\beta}$. For $q = 2$, by the tower rule, we have

$$\begin{aligned} \mathbb{E} \|\tilde{\mathcal{M}}_{i,k}(\alpha)\|_2^2 &= \mathbb{E} \left[\mathbb{E} [\tilde{\beta}^\top B_{i,k}^\top B_{i,k} \tilde{\beta} \mid \mathcal{F}_k] \right] \\ &= \mathbb{E} \left[\mathbb{E} [\text{tr}(\tilde{\beta} \tilde{\beta}^\top B_{i,k}^\top B_{i,k}) \mid \mathcal{F}_k] \right] \\ &= \mathbb{E} [\text{tr}(\mathbb{E} [\tilde{\beta} \tilde{\beta}^\top B_{i,k}^\top B_{i,k} \mid \mathcal{F}_k])] \\ &= \mathbb{E} [\text{tr}(\tilde{\beta} \tilde{\beta}^\top B_{i,k}^\top B_{i,k})] \\ &= \text{tr}(\tilde{\beta} \tilde{\beta}^\top \mathbb{E} [B_{i,k}^\top B_{i,k}]) \\ &\leq \|\mathbb{E} [B_{i,k}^\top B_{i,k}]\| \cdot \|\tilde{\beta}\|_2^2. \end{aligned} \quad (91)$$

Following the similar arguments, we obtain for $q \geq 2$,

$$\mathbb{E} \|\tilde{\mathcal{M}}_{i,k}(\alpha)\|_2^q \leq \sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \mathbb{E} \|B_{i,k} \mathbf{v}\|_2^q \cdot \|\tilde{\beta}\|_2^q. \quad (92)$$

Moreover, we notice that by the tower rule

$$\begin{aligned} \|\mathbb{E} [B_{i,k}^\top B_{i,k}]\| &= \left\| \mathbb{E} \left[H_{k-i}^\top \left(\prod_{j=k-i+1}^k A_j \right)^\top \left(\prod_{j=k-i+1}^k A_j H_{k-i} \right) \right] \right\| \\ &\leq \|\mathbb{E} [H_{k-i}^\top A_k^\top A_k H_{k-i}]\| \cdot \left\| \mathbb{E} \left[\left(\prod_{j=k-i+1}^{k-1} A_j \right)^\top \left(\prod_{j=k-i+1}^{k-1} A_j \right) \right] \right\|. \end{aligned} \quad (93)$$

By a similar argument as Step 2 in the proof of Lemma 4, we obtain

$$\|\mathbb{E} [H_{k-i}^\top A_k^\top A_k H_{k-i}]\| \lesssim p^2 \|\mathbb{X}\|^2 < \infty, \quad (94)$$

where the constant in \lesssim is independent of α . Further, recall that A_j are i.i.d. random matrices and $\|\mathbb{E} [A_1^\top A_1]\| \leq 1 - \alpha p \lambda_{\min} [X^\top (2I_d - \alpha \mathbb{X}) X]$ by the proof of Lemma 2. When $\alpha \|\mathbb{X}\| < 2$, it follows from the sub-multiplicativity of operator norm and the similar lines as the Step 1 in the proof of Lemma 4 that

$$\begin{aligned} &\sum_{i=0}^{\infty} \left\| \mathbb{E} \left[\left(\prod_{j=k-i+1}^{k-1} A_j \right)^\top \left(\prod_{j=k-i+1}^{k-1} A_j \right) \right] \right\| \\ &= \sum_{i=0}^{\infty} \left\| \prod_{j=k-i+1}^{k-1} \mathbb{E} [A_j^\top A_j] \right\| \\ &\leq \sum_{i=2}^{\infty} \|\mathbb{E} [A_1^\top A_1]\|^{i-2} = O(1/\alpha). \end{aligned} \quad (95)$$

Therefore, $\sum_{i=0}^{\infty} \mathbb{E} \|\tilde{\mathcal{M}}_{i,k}(\alpha)\|_2^2 = O(1/\alpha)$, which yields $(\mathbb{E} \|\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}\|_2^2)^{1/2} = O(\sqrt{\alpha})$. By leveraging the inequality in (92) and the similar techniques adopted in the proof of Lemma 2 for the case with $q > 2$, we obtain that for any $q \geq 2$, $\sum_{i=0}^{\infty} (\mathbb{E} \|\tilde{\mathcal{M}}_{i,k}(\alpha)\|_2^q)^{2/q} = O(1/\alpha)$. As a direct consequence, we obtain $(\mathbb{E} \|\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}\|_2^q)^{1/q} = O(\sqrt{\alpha})$, which completes the proof. ■

C.3 Proof of Theorem 6

Proof If we can establish the asymptotic normality for the affine sequence $\{\beta_k^\dagger(\alpha) - \tilde{\beta}\}_{k \in \mathbb{N}}$, then by applying Lemma 4 and Markov's inequality, we can prove the CLT for the stationary sequence $\{\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}\}_{k \in \mathbb{N}}$. Therefore, in this proof, we shall show the CLT for $\{\beta_k^\dagger(\alpha) - \tilde{\beta}\}_{k \in \mathbb{N}}$, that is,

$$\frac{\beta_k^\dagger(\alpha) - \tilde{\beta}}{\sqrt{\alpha}} \Rightarrow \mathcal{N}(0, \Xi(\alpha)), \quad \text{as } \alpha \rightarrow 0.$$

First, we recall the random vectors $\mathbf{b}_k(\alpha)$ in (15) and let

$$\mathbf{b}_k(\alpha) =: \alpha H_k \tilde{\beta}, \quad \text{with } H_k := D_k \bar{\mathbb{X}}(pI_d - D_k). \quad (96)$$

Then, since $\{\beta_k^\dagger(\alpha) - \tilde{\beta}\}_{k \in \mathbb{N}}$ is a stationary sequence and using induction on k , we can rewrite $\beta_k^\dagger(\alpha) - \tilde{\beta}$ into

$$\begin{aligned} \beta_k^\dagger(\alpha) - \tilde{\beta} &= \alpha \left(I_d H_k + (I_d - \alpha p \bar{\mathbb{X}}_p) H_{k-1} + \cdots + (I_d - \alpha p \bar{\mathbb{X}}_p)^{k-1} H_1 + \cdots \right) \tilde{\beta} \\ &= \alpha \sum_{i=0}^{\infty} (I_d - \alpha p \bar{\mathbb{X}}_p)^i H_{k-i} \tilde{\beta}. \end{aligned} \quad (97)$$

Here, the random matrices H_k are i.i.d. and independent of β_{k-1}^\dagger . For negative index k , H_k is independently drawn from the same distribution as H_1 . In the second row of (97), we can write the recursion until $H_{-\infty}$ since in the definition of $\beta_k^\dagger(\alpha) - \tilde{\beta}$, we chose the initialization as the stationary random vector $\beta_0^\dagger = \tilde{\beta}_0^\circ \sim \pi_\alpha$. Thus, $\{\beta_k^\dagger(\alpha) - \tilde{\beta}\}_{k \in \mathbb{N}}$ is a stationary sequence. This together with the discussion at the end of Section B.1 in Appendix gives Eq. (97).

Again, since H_k are i.i.d. we shall apply the Lindeberg-Feller central limit theorem to the partial sum in (97). To this end, we first take the expectation on both sides of (97). Since the random matrices H_i are independent for all $i \in \mathbb{N}$, we obtain

$$\begin{aligned} \mathbb{E}[\beta_k^\dagger(\alpha) - \tilde{\beta}] &= \alpha \sum_{i=0}^{\infty} (I_d - \alpha p \bar{\mathbb{X}}_p)^i \mathbb{E}[H_{k-i} \tilde{\beta}] \\ &= \alpha \sum_{i=0}^{\infty} (I_d - \alpha p \bar{\mathbb{X}}_p)^i \mathbb{E}[H_{k-i}] \tilde{\beta} = 0. \end{aligned} \quad (98)$$

To see the last equality, we apply Lemma 21 (i) and (ii) and obtain $\mathbb{E}[D_k \bar{\mathbb{X}} D_k] = p \bar{\mathbb{X}}_p = p^2 \bar{\mathbb{X}}$, which gives

$$\mathbb{E}[H_k] = \mathbb{E}[D_k \bar{\mathbb{X}}(pI_d - D_k)] = p^2 \bar{\mathbb{X}} - p^2 \bar{\mathbb{X}} = 0, \quad (99)$$

As a direct consequence, by (96), we have

$$\mathbb{E}[\mathbf{b}_k(\alpha)] = \alpha \mathbb{E}[H_k] \tilde{\boldsymbol{\beta}} = 0. \quad (100)$$

Next, we shall provide a closed form of the covariance matrix $\text{Cov}(\boldsymbol{\beta}_k^\dagger(\alpha) - \tilde{\boldsymbol{\beta}})$. Notice that the random vectors $H_i \tilde{\boldsymbol{\beta}}$ are uncorrelated over different i , and $\mathbb{E}[H_i \tilde{\boldsymbol{\beta}} \tilde{\boldsymbol{\beta}}^\top H_i] = \mathbb{E}[H_1 \tilde{\boldsymbol{\beta}} \tilde{\boldsymbol{\beta}}^\top H_1]$ due to the stationarity of the sequence $\{H_i \tilde{\boldsymbol{\beta}}\}_{i \in \mathbb{N}}$. Hence, by (97), we have

$$\begin{aligned} V_\alpha &:= \text{Cov}\left(\alpha^{-1/2}(\boldsymbol{\beta}_k^\dagger(\alpha) - \tilde{\boldsymbol{\beta}})\right) \\ &= \alpha^{-1} \mathbb{E}[(\boldsymbol{\beta}_k^\dagger(\alpha) - \tilde{\boldsymbol{\beta}})(\boldsymbol{\beta}_k^\dagger(\alpha) - \tilde{\boldsymbol{\beta}})^\top] \\ &= \alpha \sum_{i=0}^{\infty} (I_d - \alpha p \mathbb{X}_p)^i \mathbb{E}[H_{k-i} \tilde{\boldsymbol{\beta}} \tilde{\boldsymbol{\beta}}^\top H_{k-i}] (I_d - \alpha p \mathbb{X}_p)^i \\ &=: \alpha \sum_{i=0}^{\infty} (I_d - \alpha p \mathbb{X}_p)^i S (I_d - \alpha p \mathbb{X}_p)^i, \end{aligned} \quad (101)$$

with $d \times d$ matrix

$$S := \mathbb{E}[H_1 \tilde{\boldsymbol{\beta}} \tilde{\boldsymbol{\beta}}^\top H_1] =: \mathbb{E}[H_1 S_0 H_1], \quad \text{where } S_0 = \tilde{\boldsymbol{\beta}} \tilde{\boldsymbol{\beta}}^\top. \quad (102)$$

Furthermore, by Lemma 21 (i), one can show that $\overline{(\mathbb{X})} = \overline{\mathbb{X}}$ and $\text{Diag}(A \overline{\mathbb{X}}) = \text{Diag}(\overline{A \mathbb{X}})$ for any matrix A . Then, by the definition of H_k in (96) and Lemma 21 (ii)–(iv), we can simplify $\mathbb{E}[H_1 S_0 H_1]$ as follows:

$$\begin{aligned} &\mathbb{E}[H_1 S_0 H_1] \\ &= \mathbb{E}[D_k \overline{\mathbb{X}} (p I_d - D_k) S_0 (p I_d - D_k) \overline{\mathbb{X}} D_k] \\ &= p^2 \mathbb{E}[D_1 \overline{\mathbb{X}} S_0 \overline{\mathbb{X}} D_1] - p \mathbb{E}[D_1 \overline{\mathbb{X}} D_1 S_0 \overline{\mathbb{X}} D_1] - p \mathbb{E}[D_1 \overline{\mathbb{X}} S_0 D_1 \overline{\mathbb{X}} D_1] + \mathbb{E}[D_1 \overline{\mathbb{X}} D_k S_0 D_k \overline{\mathbb{X}} D_1] \\ &= p^3 (\overline{\mathbb{X}} S_0 \overline{\mathbb{X}})_p - 2p \left(p \overline{\mathbb{X}}_p (S_0 \overline{\mathbb{X}})_p + p^2 (1-p) \text{Diag}(\overline{\mathbb{X}} S_0 \overline{\mathbb{X}}) \right) \\ &\quad + p \overline{\mathbb{X}}_p (S_0)_p \overline{\mathbb{X}}_p + p^2 (1-p) \left(\text{Diag}(\overline{\mathbb{X}} (S_0)_p \overline{\mathbb{X}}) + 2 \overline{\mathbb{X}}_p \text{Diag}(\overline{S_0 \mathbb{X}}) + (1-p) \overline{\mathbb{X}} \odot \overline{S_0}^\top \odot \overline{\mathbb{X}} \right). \end{aligned} \quad (103)$$

By (103), we obtain a closed form solution of S which is independent of α .

Now we are ready to solve the covariance matrix V_α in (101). We multiply the matrix $I_d - \alpha p \mathbb{X}_p$ to the left and right sides of (101) and obtain

$$(I_d - \alpha p \mathbb{X}_p) V_\alpha (I_d - \alpha p \mathbb{X}_p) = \alpha \sum_{i=0}^{\infty} (I_d - \alpha p \mathbb{X}_p)^i S (I_d - \alpha p \mathbb{X}_p)^i. \quad (104)$$

Taking the difference between V_α and $(I_d - \alpha p \mathbb{X}_p) V_\alpha (I_d - \alpha p \mathbb{X}_p)$ yields

$$V_\alpha - (I_d - \alpha p \mathbb{X}_p) V_\alpha (I_d - \alpha p \mathbb{X}_p) = \alpha S. \quad (105)$$

Denote the symmetric matrix $A_p = p \mathbb{X}_p$. By simplifying the equation above, for $\alpha > 0$, we have

$$V_\alpha A_p - A_p V_\alpha + \alpha A_p V_\alpha A_p = S. \quad (106)$$

Let $V_0 = \lim_{\alpha \rightarrow 0} V_\alpha$. As $\alpha \rightarrow 0$, the quadratic term $\alpha A_p V_\alpha A_p$ vanishes. Thus, we only need to solve the equation

$$S - V_0 A_p - A_p V_0 = 0, \quad (107)$$

to get the solution for

$$V_0 = \lim_{\alpha \rightarrow 0} V_\alpha.$$

Following Theorem 1 in Pflug (1986) and the subsequent Remark therein, we can get the closed form solution of V_0 , that is,

$$\text{vec}(V_0) = (I_d \otimes A_p + A_p \otimes I_d)^{-1} \cdot \text{vec}(S), \quad (108)$$

where the $d^2 \times d^2$ matrix $I_d \otimes A_p + A_p \otimes I_d$ is invertible since the fixed design matrix X is assumed to be in a reduced form with no zero columns. For a small $\alpha > 0$, we shall provide a similar closed form solution for $V_\alpha = V_0 + \alpha B_p$. Specifically, we need to get the closed form of the matrix B_p by solving a similar equation:

$$A_p V_0 A_p - B_p A_p - A_p B_p = 0, \quad (109)$$

which gives

$$\text{vec}(B_p) = (I_d \otimes A_p + A_p \otimes I_d)^{-1} \times \text{vec}(A_p V_0 A_p). \quad (110)$$

The deterministic matrices V_0 , A_p and B_p are all independent of α . By inserting the results of V_0 and B_p into $V_\alpha = V_0 + \alpha B_p$, we obtain

$$\Xi(\alpha) = V_\alpha = V_0 + \alpha B_p,$$

which holds uniformly over k due to the stationarity of $\{\beta_k^\dagger(\alpha) - \tilde{\beta}\}_{k \in \mathbb{N}}$.

Finally, by applying the Lindeberg-Feller central limit theorem to the partial sum in (97), we establish the asymptotic normality of $\{\beta_k^\dagger(\alpha) - \tilde{\beta}\}_{k \in \mathbb{N}}$ and complete the proof. ■

Appendix D. Proofs in Section 3.3

We first outline the main techniques for establishing the asymptotic normality of the averaged GD dropout sequence $\{\tilde{\beta}_k^{\text{gd}}(\alpha)\}_{k \in \mathbb{N}}$ defined in (24).

Recall the observation \mathbf{y} in model (1) and the dropout matrix D . For the GD dropout $\{\tilde{\beta}_k(\alpha)\}_{k \in \mathbb{N}}$ in (2), by Theorem 3, we can define a centering term as follows,

$$\tilde{\beta}_{\text{mean}}(\alpha) = \lim_{k \rightarrow \infty} \mathbb{E}_D[\tilde{\beta}_k(\alpha)] = \mathbb{E}_D[\tilde{\beta}_1^\circ(\alpha)], \quad (111)$$

where $\tilde{\beta}_1^\circ(\alpha)$ follows the stationary distribution π_α as stated in (14). According to Lemma 1 in Clara et al. (2024), we note that $\mathbb{E}_D[\tilde{\beta}_k(\alpha) - \tilde{\beta}] \neq 0$ but $\|\mathbb{E}_D[\tilde{\beta}_k(\alpha) - \tilde{\beta}]\|_2 \rightarrow 0$ if $\alpha p \|\mathbb{X}\| < 1$ with a geometric rate as $k \rightarrow \infty$. Therefore, we shall first show the central limit theorems for the partial sum of $\{\tilde{\beta}_k(\alpha) - \tilde{\beta}_{\text{mean}}(\alpha)\}_{k \in \mathbb{N}}$ and then for the one of $\{\tilde{\beta}_k(\alpha) - \tilde{\beta}\}_{k \in \mathbb{N}}$.

Next, we take a closer look at the partial sum of $\{\tilde{\beta}_k(\alpha) - \tilde{\beta}_{\text{mean}}(\alpha)\}_{k \in \mathbb{N}}$. The iterative function f defined in (9) allows us to write $\tilde{\beta}_k(\alpha) = f_{D_k}(\tilde{\beta}_{k-1}(\alpha))$ for all $k \in \mathbb{N}$. Similarly,

for the initialization $\tilde{\beta}_0^\circ(\alpha)$ that follows the unique stationary distribution π_α in Theorem 3, we can write the stationary GD dropout sequence $\{\tilde{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ into

$$\tilde{\beta}_k^\circ(\alpha) = f_{D_k}(\tilde{\beta}_{k-1}^\circ(\alpha)), \quad k \in \mathbb{N}. \quad (112)$$

Recall $\tilde{\beta}_{\text{mean}}(\alpha)$ defined in (111). Then, we can recursively rewrite $\tilde{\beta}_k^\circ(\alpha)$ using the iterative function f and obtain the partial sum

$$\begin{aligned} \tilde{S}_t^\circ(\alpha) &:= \sum_{k=1}^t [\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}_{\text{mean}}(\alpha)] \\ &= \{f_{D_1}(\tilde{\beta}_0^\circ(\alpha)) - \mathbb{E}[f_{D_1}(\tilde{\beta}_0^\circ(\alpha))]\} + \{f_{D_2} \circ f_{D_1}(\tilde{\beta}_0^\circ(\alpha)) - \mathbb{E}[f_{D_2} \circ f_{D_1}(\tilde{\beta}_0^\circ(\alpha))]\} \\ &\quad + \cdots + \{f_{D_t} \circ \cdots \circ f_{D_1}(\tilde{\beta}_0^\circ(\alpha)) - \mathbb{E}[f_{D_t} \circ \cdots \circ f_{D_1}(\tilde{\beta}_0^\circ(\alpha))]\}. \end{aligned} \quad (113)$$

Primarily, we aim to (i) prove the central limit theorem for the partial sum $t^{-1/2}\tilde{S}_t^\circ(\alpha)$, and (ii) prove the invariance principle for the partial sum process $(\tilde{S}_i^\circ(\alpha))_{1 \leq i \leq t}$. To this end, we borrow the idea of *functional dependence measure* in Wu (2005), which was further investigated in Wu (2011) to establish the asymptotic normality for sequences with *short-range dependence* (see (120) for the definition). We shall show that the GD dropout sequence $\{\tilde{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ that satisfies the geometric-moment contraction (as proved in Theorem 3) satisfies such short-range dependence condition.

Finally, we shall complete the proofs of the quenched central limit theorems by showing that, for any given constant learning rate $\alpha > 0$ satisfying the conditions in Theorem 7, and any initialization $\tilde{\beta}_0 \in \mathbb{R}^d$, the partial sum

$$\tilde{S}_t^{\tilde{\beta}_0}(\alpha) := \sum_{k=1}^t [\tilde{\beta}_k(\alpha) - \tilde{\beta}_{\text{mean}}(\alpha)] \quad (114)$$

converges to the stationary partial sum process $\tilde{S}_t^\circ(\alpha)$, in the sense that $t^{-1/2}(\mathbb{E}\|\tilde{S}_t^{\tilde{\beta}_0}(\alpha) - \tilde{S}_t^\circ(\alpha)\|_2^q)^{1/q} = o(1)$ as $t \rightarrow \infty$.

D.1 Functional Dependence Measure

Before proceeding to the proofs of Theorems 7 and 9, we first provide the detailed form of the functional dependence measure in Wu (2005) for the iterated random functions with i.i.d. random matrices as inputs. This will serve as the foundational pillar to build the asymptotic normality of averaged GD dropout iterates.

First, for any random vector $\zeta \in \mathbb{R}^d$ satisfying $\mathbb{E}\|\zeta\|_2 < \infty$, define projection operators

$$\mathcal{P}_k[\zeta] = \mathbb{E}[\zeta \mid \mathcal{F}_k] - \mathbb{E}[\zeta \mid \mathcal{F}_{k-1}], \quad k \in \mathbb{Z}, \quad (115)$$

where we recall the filtration $\mathcal{F}_i = \sigma(D_i, D_{i-1}, \dots)$ with i.i.d. dropout matrices D_i , $i \in \mathbb{Z}$. By Theorem 3 and (12), there exists a measurable function $h_\alpha(\cdot)$ such that the stationary GD dropout sequence $\{\tilde{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ can be written as the following causal process

$$\tilde{\beta}_k^\circ(\alpha) = h_\alpha(D_k, D_{k-1}, \dots) = h_\alpha(\mathcal{F}_k). \quad (116)$$

Define a coupled version of filtration \mathcal{F}_i as $\mathcal{F}_{i,\{j\}} = \sigma(D_i, \dots, D_{j+1}, D'_j, D_{j-1}, \dots)$. In addition, $\mathcal{F}_{i,\{j\}} = \mathcal{F}_i$ if $j > i$. For $q > 1$, define the *functional dependence measure* of $\tilde{\beta}_k^\circ(\alpha)$ as

$$\theta_{k,q}(\alpha) = (\mathbb{E}\|\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}_{k,\{0\}}^\circ(\alpha)\|_2^q)^{1/q}, \quad \text{where } \tilde{\beta}_{k,\{0\}}^\circ(\alpha) = h_\alpha(\mathcal{F}_{k,\{0\}}). \quad (117)$$

The above quantity can be interpreted as the dependence of $\tilde{\beta}_k^\circ(\alpha)$ on D_0 (see the discussion below Theorem 3 for the meaning of $\tilde{\beta}_k^\circ$ with $k \leq 0$), and $\tilde{\beta}_{k,\{0\}}^\circ(\alpha)$ is a coupled version of $\tilde{\beta}_k^\circ(\alpha)$ with D_0 in the latter replaced by its i.i.d. copy D'_0 . If $\tilde{\beta}_k^\circ(\alpha)$ does not functionally depend on D_0 , then $\theta_{k,q}(\alpha) = 0$.

Furthermore, if $\sum_{k=0}^{\infty} \theta_{k,q}(\alpha) < \infty$, we define the tail of the *cumulative dependence measure* as

$$\Theta_{m,q}(\alpha) = \sum_{k=m}^{\infty} \theta_{k,q}(\alpha), \quad m \in \mathbb{N}. \quad (118)$$

This can be interpreted as the cumulative dependence of $\{\tilde{\beta}_k^\circ(\alpha)\}_{k \geq m}$ on D_0 , or equivalently, the cumulative dependence of $\tilde{\beta}_0^\circ(\alpha)$ on D_j , $j \geq m$. The functional dependence measure in (117) and its cumulative variant in (118) are easy to work with and they can directly reflect the underlying data-generating mechanism of the iterative function $\tilde{\beta}_k^\circ(\alpha) = f_{D_k}(\tilde{\beta}_{k-1}^\circ(\alpha))$.

Specifically, for all $q \geq 2$, Theorem 1 in Wu (2005) pointed out a useful inequality for the functional dependence measure as follows,

$$\sum_{k=0}^{\infty} (\mathbb{E}\|\mathcal{P}_0[\tilde{\beta}_k^\circ(\alpha)]\|_2^q)^{1/q} \leq \sum_{k=0}^{\infty} \theta_{k,q}(\alpha) = \Theta_{0,q}(\alpha). \quad (119)$$

In particular, for some given learning rate $\alpha > 0$, we say the sequence $\{\tilde{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ satisfies the short-range dependence condition if

$$\Theta_{0,q}(\alpha) < \infty, \quad \text{for some } q \geq 2. \quad (120)$$

This dependence assumption has been widely adopted in the literature; see for example the invariance principle in Wu (2011); Berkes et al. (2014); Karmakar and Wu (2020). If condition (120) fails, then $\tilde{\beta}_k^\circ(\alpha)$ can be long-range dependent, and the partial sum (resp. partial sum processes) behave no longer like Gaussian random vectors (resp. Brownian motions).

Here, we introduce Theorem 3 in Wu (2011) and Theorem 2 in Karmakar and Wu (2020), which are the fundamental tools for the proofs of Theorems 7 and 9, respectively.

Lemma 25 (Asymptotic normality (Wu, 2011)) *Consider a sequence of stationary mean-zero random variables $x_k = g(\epsilon_k, \epsilon_{k-1}, \dots) \in \mathbb{R}$, for $k = 1, \dots, n$, where ϵ_k 's are i.i.d. random variables, and $g(\cdot)$ is a measurable function such that each x_k is a proper random variable. Recall the projection operator $\mathcal{P}_0[\cdot]$ defined in (115). Let $\theta_q(i) = (\mathbb{E}\|\mathcal{P}_0[x_i]\|_2^q)^{1/q}$, $q > 1$. Assume $\mathbb{E}[x_i] = 0$ and*

$$\Theta_q := \sum_{i=0}^{\infty} \theta_q(i) < \infty. \quad (121)$$

Let $S_n = \sum_{k=1}^n x_k$ and define the process $S_t = S_{[t]} + (t - [t])x_{[t]+1}$, for $t \geq 0$ and the floor function $[t] = \max\{k \in \mathbb{Z} : k \leq t\}$. Then,

(i) we have the moment inequality

$$(\mathbb{E}\|S_n\|_2^q)^{1/q} \leq \begin{cases} (q-1)^{1/2}n^{1/2}\Theta_q, & q > 2, \\ (q-1)^{-1}n^{1/q}\Theta_q, & 1 < q \leq 2. \end{cases} \quad (122)$$

(ii) If moreover (121) holds with $q = 2$, the invariance principle

$$\left\{ \frac{1}{\sqrt{n}}S_{nu}, 0 \leq u \leq 1 \right\} \Rightarrow \{\sigma\mathbb{B}(u), 0 \leq u \leq 1\}, \quad (123)$$

holds with long-run variance $\sigma^2 = \mathbb{E}\|\sum_{i=0}^{\infty} \mathcal{P}_0[x_i]\|_2^2$.

Lemma 26 (Gaussian approximation (Karmakar and Wu, 2020)) Suppose that we have a sequence of nonstationary mean-zero random vectors $\mathbf{x}_k = g_k(\epsilon_k, \epsilon_{k-1}, \dots) \in \mathbb{R}^d$, for $k = 1, \dots, n$, where the ϵ_k 's are i.i.d. random variables, and $g_k(\cdot)$ is a measurable function such that each \mathbf{x}_k is a proper random vector. Let $S_j = \sum_{k=1}^j \mathbf{x}_k$. Assume the following conditions hold for some $q > 2$:

- (i) The series $(\|\mathbf{x}_k\|_2^q)_{k \geq 1}$ is uniformly integrable: $\sup_{k \geq 1} \mathbb{E}[\|\mathbf{x}_k\|_2^q \mathbf{1}_{\|\mathbf{x}_k\|_2 \geq u}] \rightarrow 0$ as $u \rightarrow \infty$,
- (ii) The eigenvalues of covariance matrices of increment processes are lower-bounded, that is, there exists $\lambda_* > 0$ and $l_* \in \mathbb{N}$, such that for all $t \geq 1, l \geq l_*$,

$$\lambda_{\min}(\text{Cov}(S_{t+l} - S_t)) \geq \lambda_* l;$$

(iii) There exist constants $\chi > \chi_0$ and $\kappa > 0$, where

$$\chi_0 = \frac{q^2 - 4 + (q-2)\sqrt{q^2 + 20q + 4}}{8q},$$

such that the tail cumulative dependence measure

$$\Theta_{m,q}(\alpha) = \sum_{k=m}^{\infty} \theta_{k,q}(\alpha) = O\{m^{-\chi}(\log(m))^{-\kappa}\}. \quad (124)$$

Then, for all $q > 2$, there exists a probability space $(\Omega^*, \mathcal{A}^*, \mathbb{P}^*)$ on which we can define random vectors \mathbf{x}_k^* , with the partial sum process $S_i^* = \sum_{k=1}^i \mathbf{x}_k^*$ and a Gaussian process $G_i^* = \sum_{k=1}^i \mathbf{z}_k^*$. Here \mathbf{z}_k^* is a mean-zero independent Gaussian vector, such that $(S_i^*)_{1 \leq i \leq n} \stackrel{\mathcal{D}}{=} (S_i)_{1 \leq i \leq n}$ and

$$\max_{i \leq n} |S_i^* - G_i^*| = o_{\mathbb{P}}(n^{1/q}) \quad \text{in } (\Omega^*, \mathcal{A}^*, \mathbb{P}^*).$$

We notice that condition (ii) in Lemma 26 on the non-singularity is required when the sequence $\{\mathbf{x}_k\}_{k \in \mathbb{N}}$ is non-stationary. However, if the function $g_k(\cdot) \equiv g(\cdot)$, that is, the sequence $\{\mathbf{x}_k\}_{k \in \mathbb{N}}$ is stationary, then the covariance matrix of the increments is allowed to be singular. To see this, consider a stationary partial sum $S_l = (S_{l,1}, \dots, S_{l,d})^\top$ with a singular covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$ and assume $\text{rank}(\Sigma) = d - 1$. Then, there exists

a unit vector $\mathbf{v} \in \mathbb{R}^d$ such that $\Sigma \mathbf{v} = 0$, which indicates that $S_{l,1}$ can be written into a linear combination of $S_{l,2}, \dots, S_{l,d}$, and the covariance matrix of this linear combination is non-singular. Hence, condition (ii) in Lemma 26 is not required for stationary processes.

In addition, the original Theorem 3 in Wu (2011) and Theorem 2 in Karmakar and Wu (2020) considered a simple case where the i.i.d. inputs ϵ_k are one-dimensional. These two theorems still hold even if the inputs are i.i.d. random matrices such as the dropout matrices D_k in our case. In fact, as long as the inputs are i.i.d. elements, the functional dependence measure can be similarly computed as the one in one-dimensional case. The essence is that the short-range dependence condition (120) is satisfied using an appropriate norm (e.g., L^2 -norm for vectors, operator norm for matrices) by the output x_k . For example, Wu and Shao (2004) considered iterated random functions on a general metric space, and Chen and Wu (2016) assumed the ϵ_i 's to be i.i.d. random elements to derive asymptotics for x_k . We will verify this short-range dependence condition on the GD dropout vector estimates $\{\tilde{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ in the proof of Theorem 7.

D.2 Proof of Theorem 7

Proof We verify the short-range dependence condition for the stationary GD dropout sequence $\{\tilde{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$.

First, consider two different initial vectors $\tilde{\beta}_0^\circ, \tilde{\beta}_0^{\circ'} \in \mathbb{R}^d$ following the unique stationary distribution π_α in Theorem 3. Denote the two GD dropout sequences by $\{\tilde{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ and $\{\tilde{\beta}_k^{\circ'}(\alpha)\}_{k \in \mathbb{N}}$ accordingly. By the geometric-moment contraction in Theorem 3, for all $q \geq 2$, we have

$$\sup_{\tilde{\beta}_0^\circ, \tilde{\beta}_0^{\circ'} \in \mathbb{R}^d, \tilde{\beta}_0^\circ \neq \tilde{\beta}_0^{\circ'}} \frac{(\mathbb{E} \|\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}_k^{\circ'}(\alpha)\|_2^q)^{1/q}}{\|\tilde{\beta}_0^\circ - \tilde{\beta}_0^{\circ'}\|_2} \leq r_{\alpha,q}^k, \quad k \in \mathbb{N}, \quad (125)$$

for some constant $r_{\alpha,q} \in (0, 1)$. Equivalently, it can be rewritten in terms of the iterative function f defined in (9) and $h_\alpha(\cdot)$ defined in (12). That is, for all $\tilde{\beta}_0^\circ, \tilde{\beta}_0^{\circ'} \in \mathbb{R}^d$, such that $\tilde{\beta}_0^\circ \neq \tilde{\beta}_0^{\circ'}$, we have

$$\begin{aligned} & (\mathbb{E} \|f_{D_k} \circ \dots \circ f_{D_1}(\tilde{\beta}_0^\circ) - f_{D_k} \circ \dots \circ f_{D_1}(\tilde{\beta}_0^{\circ'})\|_2^q)^{1/q} \\ &= (\mathbb{E} \|h_\alpha(D_k, \dots, D_1, D_0, D_{-1}, \dots) - h_\alpha(D_k, \dots, D_1, D_0', D_{-1}', \dots)\|_2^q)^{1/q} \\ &= (\mathbb{E} \|h_\alpha(\mathcal{F}_k) - h_\alpha(\mathcal{F}_{k, \{0, -1, \dots\}})\|_2^q)^{1/q} \\ &\leq c_q r_{\alpha,q}^k, \end{aligned} \quad (126)$$

where we recall the filtration $\mathcal{F}_{i, \{j\}} = \sigma(D_i, \dots, D_{j+1}, D_j', D_{j-1}, \dots)$, and $c_q > 0$ is some constant independent of k . Moreover, since $h_\alpha(\mathcal{F}_k)$ is stationary over k and D_i and D_j' are i.i.d. random matrices, for all $i, j \in \mathbb{Z}$, it follows that

$$\begin{aligned} & \mathbb{E} \|h_\alpha(\mathcal{F}_{k, \{0\}}) - h_\alpha(\mathcal{F}_{k, \{0, -1, \dots\}})\|_2^q \\ &= \mathbb{E} \|h_\alpha(\mathcal{F}_k) - h_\alpha(\mathcal{F}_{k, \{-1, -2, \dots\}})\|_2^q \\ &= \mathbb{E} \|h_\alpha(\mathcal{F}_{k+1}) - h_\alpha(\mathcal{F}_{k+1, \{0, -1, \dots\}})\|_2^q \leq c_q' r_{\alpha,q}^k, \end{aligned} \quad (127)$$

where the constant $c'_q > 0$ is also independent of k . Hence, by (126) and (127), we can bound the functional dependence measure defined in (117) as follows

$$\begin{aligned}\theta_{k,q}(\alpha) &= (\mathbb{E}\|h_\alpha(\mathcal{F}_k) - h_\alpha(\mathcal{F}_{k,\{0\}})\|_2^q)^{1/q} \\ &\leq (\mathbb{E}\|h_\alpha(\mathcal{F}_k) - h_\alpha(\mathcal{F}_{k,\{0,-1,\dots\}})\|_2^q)^{1/q} + (\mathbb{E}\|h_\alpha(\mathcal{F}_{k,\{0,-1,\dots\}}) - h_\alpha(\mathcal{F}_{k,\{0\}})\|_2^q)^{1/q} \\ &\leq (c_q + c'_q)r_{\alpha,q}^k.\end{aligned}\tag{128}$$

As a direct result, we have finite cumulative dependence measure defined in (118), i.e.,

$$\Theta_{m,q}(\alpha) = \sum_{k=m}^{\infty} \theta_{k,q}(\alpha) = O(r_{\alpha,q}^m) < \infty.\tag{129}$$

Therefore, for the constant learning rate $\alpha > 0$ satisfying the assumptions in Theorem 3, the stationary GD dropout sequence $\{\beta_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ meets the short-range dependence requirement in (120). Consequently, the condition (121) in Lemma 25 is satisfied, which along with the Cramér-Wold device yields the central limit theorem for $\tilde{S}_t^\circ(\alpha)$ defined in (113), that is,

$$t^{-1/2}\tilde{S}_t^\circ(\alpha) \Rightarrow \mathcal{N}(0, \Sigma(\alpha)),\tag{130}$$

where the long-run covariance matrix $\Sigma(\alpha)$ is defined in Theorem 7.

Next, we bound the difference between $\tilde{S}_t^\circ(\alpha)$ and $\tilde{S}_t^{\tilde{\beta}_0}(\alpha)$ for any arbitrarily fixed $\tilde{\beta}_0 \in \mathbb{R}^d$ in the q -th moment, for all $q \geq 2$. For the constant learning rate $\alpha > 0$ satisfying $\alpha\|\mathbb{X}\| < 2$, applying Theorem 3 yields

$$\begin{aligned}& (\mathbb{E}\|\tilde{S}_t^\circ(\alpha) - \tilde{S}_t^{\tilde{\beta}_0}(\alpha)\|_2^q)^{1/q} \\ &= (\mathbb{E}\|[f_{D_1}(\tilde{\beta}_0^\circ(\alpha)) + f_{D_2} \circ f_{D_1}(\tilde{\beta}_0^\circ(\alpha)) + \dots + f_{D_t} \circ \dots \circ f_{D_1}(\tilde{\beta}_0^\circ(\alpha))] \\ &\quad - [f_{D_1}(\tilde{\beta}_0(\alpha)) + f_{D_2} \circ f_{D_1}(\tilde{\beta}_0(\alpha)) + \dots + f_{D_t} \circ \dots \circ f_{D_1}(\tilde{\beta}_0(\alpha))]\|_2^q)^{1/q} \\ &\leq \left(\sum_{k=1}^t r_{\alpha,q}^k\right) \|\tilde{\beta}_0^\circ - \tilde{\beta}_0\|_2.\end{aligned}\tag{131}$$

Since the contraction constant $r_{\alpha,q} \in (0, 1)$, we can derive the limit for the sum of the geometric series $\{r_{\alpha,q}^k\}_{k=1}^t$ as follows

$$\lim_{t \rightarrow \infty} \sum_{k=1}^t r_{\alpha,q}^k = \lim_{t \rightarrow \infty} \frac{r_{\alpha,q}(1 - r_{\alpha,q}^t)}{1 - r_{\alpha,q}} = \frac{r_{\alpha,q}}{1 - r_{\alpha,q}}.\tag{132}$$

This, together with (131) gives

$$(\mathbb{E}\|\tilde{S}_t^\circ(\alpha) - \tilde{S}_t^{\tilde{\beta}_0}(\alpha)\|_2^q)^{1/q} = O(1) = o(\sqrt{t}),\tag{133}$$

which yields the quenched central limit theorem for the partial sum $\tilde{S}_t^{\tilde{\beta}_0}(\alpha)$ defined in (114), that is, for any fixed initial point $\tilde{\beta}_0 \in \mathbb{R}^d$,

$$t^{-1/2}\tilde{S}_t^{\tilde{\beta}_0}(\alpha) \Rightarrow \mathcal{N}(0, \Sigma(\alpha)).\tag{134}$$

Finally, we shall show that $\|\sum_{k=1}^t \mathbb{E}[\tilde{\beta}_k(\alpha) - \tilde{\beta}]\|_2 = o(\sqrt{t})$. To see this, we note that given two independently chosen initial vectors $\tilde{\beta}_0$ and $\tilde{\beta}_0^\circ$, where $\tilde{\beta}_0^\circ$ follows the stationary distribution π_α while $\tilde{\beta}_0$ is an arbitrary initial point in \mathbb{R}^d , it follows from the triangle inequality that

$$\begin{aligned} \left\| \sum_{k=1}^t \mathbb{E}[\tilde{\beta}_k(\alpha) - \tilde{\beta}] \right\|_2 &= \left\| \mathbb{E} \left[\sum_{k=1}^t (\tilde{\beta}_k(\alpha) - \tilde{\beta}_k^\circ(\alpha) + \tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}) \right] \right\|_2 \\ &\leq \left\| \mathbb{E} \left[\sum_{k=1}^t (\tilde{\beta}_k(\alpha) - \tilde{\beta}_k^\circ(\alpha)) \right] \right\|_2 + \left\| \mathbb{E} \left[\sum_{k=1}^t (\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}) \right] \right\|_2 \\ &=: \text{I}_1 + \text{I}_2. \end{aligned} \quad (135)$$

We first show $\text{I}_2 = 0$. Recall the representation of $\{\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}\}_{k \in \mathbb{N}}$ in (14). Since $\mathbb{E}[A_k(\alpha)] = \mathbb{E}[I_d - \alpha D_k \mathbb{X} D_k] = I_d - \alpha p \mathbb{X}_p$ and $\mathbb{E}[\mathbf{b}_k(\alpha)] = 0$ by (100), it follows that

$$\mathbb{E}[\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}] = (I_d - \alpha p \mathbb{X}_p) \mathbb{E}[\tilde{\beta}_{k-1}^\circ(\alpha) - \tilde{\beta}]. \quad (136)$$

Thus, due to the stationarity of $\{\tilde{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ and the non-singularity of \mathbb{X}_p , we obtain that uniformly over $k \in \mathbb{N}$,

$$\mathbb{E}[\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}] = 0. \quad (137)$$

As a direct consequence,

$$\text{I}_2 = \left\| \mathbb{E} \left[\sum_{k=1}^t (\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}) \right] \right\|_2 = \left\| \sum_{k=1}^t \mathbb{E}[\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}] \right\|_2 = 0. \quad (138)$$

In addition, for the part I_1 , it follows from Jensen's inequality and (131) that

$$\begin{aligned} \text{I}_1 &= \left\| \mathbb{E} \left[\sum_{k=1}^t (\tilde{\beta}_k(\alpha) - \tilde{\beta}_k^\circ(\alpha)) \right] \right\|_2 \\ &\leq \left(\mathbb{E} \left\| \sum_{k=1}^t (\tilde{\beta}_k(\alpha) - \tilde{\beta}_k^\circ(\alpha)) \right\|_2^2 \right)^{1/2} \\ &\leq \left(\sum_{k=1}^t r_{\alpha,2}^k \right) \|\tilde{\beta}_0 - \tilde{\beta}_0^\circ\|_2. \end{aligned} \quad (139)$$

By inserting the results of parts I_1 and I_2 back to (135), we obtain

$$\left\| \sum_{k=1}^t \mathbb{E}[\tilde{\beta}_k(\alpha) - \tilde{\beta}] \right\|_2 \leq \left(\sum_{k=1}^t r_{\alpha,2}^k \right) \|\tilde{\beta}_0 - \tilde{\beta}_0^\circ\|_2. \quad (140)$$

which remains bounded as $t \rightarrow \infty$ when $\alpha \|\mathbb{X}\| < 2$ by Theorem 3. This completes the proof. \blacksquare

D.3 Proof of Corollary 8

Proof Recall the stationary GD dropout sequence $\{\tilde{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ which follows the unique stationary distribution π_α . Since this sequence satisfies the short-range dependence condition as stated in (120), it follows from Lemma 25 and the Cramér-Wold device that any fixed linear combination of the coordinates of $\tilde{S}_t^\circ(\alpha)$ in (113) converges to the corresponding linear combination of normal vectors in distribution. Then, the CLT for the averaged GD dropout with multiple learning rates holds by applying the Cramér-Wold device again, that is,

$$t^{-1/2} \text{vec}(\tilde{S}_t^\circ(\alpha_1), \dots, \tilde{S}_t^\circ(\alpha_s)) \Rightarrow \mathcal{N}(0, \Sigma^{\text{vec}}). \quad (141)$$

Then, following the similar arguments as in the proof of Theorem 7, we obtain the quenched CLT for $\text{vec}(\tilde{\beta}_k^{\text{gd}}(\alpha_1) - \tilde{\beta}, \dots, \tilde{\beta}_k^{\text{gd}}(\alpha_s) - \tilde{\beta})$. We omit the details here. \blacksquare

D.4 Proof of Theorem 9

Proof Recall the stationary GD dropout sequence $\{\tilde{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ in (14), where $\tilde{\beta}_k^\circ(\alpha)$ follows the stationary distribution π_α for all $k \in \mathbb{N}$. Also, recall the centering term $\tilde{\beta}_{\text{mean}}(\alpha) = \mathbb{E}[\tilde{\beta}_1^\circ(\alpha)]$ as defined in (111). By (137), we have $\mathbb{E}[\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}] = 0$ uniformly over $k \in \mathbb{N}$. Hence,

$$(\mathbb{E}\|\tilde{\beta}_{\text{mean}}(\alpha) - \tilde{\beta}\|_2^q)^{1/q} = (\mathbb{E}\|\mathbb{E}[\tilde{\beta}_1^\circ(\alpha) - \tilde{\beta}]\|_2^q)^{1/q} = 0. \quad (142)$$

This, along with Lemma 5 gives, for $q > 2$,

$$(\mathbb{E}\|\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}_{\text{mean}}(\alpha)\|_2^q)^{1/q} \leq (\mathbb{E}\|\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}\|_2^q)^{1/q} + (\mathbb{E}\|\tilde{\beta}_{\text{mean}}(\alpha) - \tilde{\beta}\|_2^q)^{1/q} = O(\sqrt{\alpha}). \quad (143)$$

Moreover, we notice that by Markov's inequality and (143), for any $u \in \mathbb{R}$ and $\delta > 0$, we have

$$\begin{aligned} & \sup_{k \geq 1} \mathbb{E} \left[\|\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}_{\text{mean}}(\alpha)\|_2^q \mathbf{1}_{\|\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}_{\text{mean}}(\alpha)\|_2 \geq u} \right] \\ & \leq \sup_{k \geq 1} \mathbb{E} \left[\|\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}_{\text{mean}}(\alpha)\|_2^q \cdot \|\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}_{\text{mean}}(\alpha)\|_2^\delta / u^\delta \right] \\ & = \sup_{k \geq 1} \mathbb{E} \left[\|\tilde{\beta}_k^\circ(\alpha) - \tilde{\beta}_{\text{mean}}(\alpha)\|_2^{q+\delta} \right] / u^\delta \\ & = O\{\alpha^{(q+\delta)/2} / u^\delta\}, \end{aligned} \quad (144)$$

which converges to 0 as $u \rightarrow \infty$. Therefore, condition (i) in Lemma 26 is satisfied. Since $\{\tilde{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ is stationary, following the arguments below Lemma 26, condition (ii) is not required. Regarding condition (iii), for the constant learning rate $\alpha > 0$ satisfying $\alpha \|\mathbb{X}\| < 2$, it follows from (128) that the functional dependence measure $\theta_{k,q}(\alpha) \leq c \cdot r_{\alpha,q}^k$, for all $q > 2$ and $k \in \mathbb{N}$, where the constant $c > 0$ is independent of k . Consequently, there exists a

constant $\kappa > 0$ such that the tail cumulative dependence measure of $\{\tilde{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ can be bounded by

$$\Theta_{m,q}(\alpha) = \sum_{k=m}^{\infty} \theta_{k,q}(\alpha) = O\{m^{-\chi}(\log(m))^{-\kappa}\}, \quad (145)$$

where $\chi > 0$ is some constant that can be taken to be arbitrarily large. Then, the condition (iii) in Lemma 26 is satisfied.

Thus, we can obtain the invariance principle for the stationary partial sum process $(\tilde{S}_i^\circ(\alpha))_{1 \leq i \leq t}$ defined in (113). That is, there exists a (richer) probability space $(\tilde{\Omega}^*, \tilde{\mathcal{A}}^*, \tilde{\mathbb{P}}^*)$ on which we can define random vectors $\tilde{\beta}_k^*$'s with the partial sum process $\tilde{S}_i^* = \sum_{k=1}^i (\tilde{\beta}_k^* - \tilde{\beta}_{\text{mean}})$, and a Gaussian process $\tilde{G}_i^* = \sum_{k=1}^i \tilde{z}_k^*$, where \tilde{z}_k^* 's are independent Gaussian random vectors in \mathbb{R}^d following $\mathcal{N}(0, I_d)$, such that

$$(\tilde{S}_i^*)_{1 \leq i \leq t} \stackrel{\mathcal{D}}{=} (\tilde{S}_i^\circ)_{1 \leq i \leq t}, \quad (146)$$

and

$$\max_{1 \leq i \leq t} \|\tilde{S}_i^* - \Sigma^{1/2}(\alpha) \tilde{G}_i^*\|_2 = o_{\mathbb{P}}(t^{1/q}), \quad \text{in } (\tilde{\Omega}^*, \tilde{\mathcal{A}}^*, \tilde{\mathbb{P}}^*), \quad (147)$$

where the long-run covariance matrix $\Sigma(\alpha)$ is defined in Theorem 7.

Next, recall the partial sum $\tilde{S}_i^{\tilde{\beta}_0}(\alpha) = \sum_{k=1}^i [\tilde{\beta}_k(\alpha) - \tilde{\beta}_{\text{mean}}(\alpha)]$ as defined in (114), given an arbitrarily fixed initial point $\tilde{\beta}_0 \in \mathbb{R}^d$. It follows from the triangle inequality that

$$\begin{aligned} & \left(\mathbb{E} \left[\max_{1 \leq i \leq t} \|\tilde{S}_i^{\tilde{\beta}_0}(\alpha) - \Sigma^{1/2}(\alpha) \tilde{G}_i^*\|_2 \right]^q \right)^{1/q} \\ &= \left(\mathbb{E} \left[\max_{1 \leq i \leq t} \|\tilde{S}_i^{\tilde{\beta}_0}(\alpha) - \tilde{S}_i^\circ(\alpha) + \tilde{S}_i^\circ(\alpha) - \Sigma^{1/2}(\alpha) \tilde{G}_i^*\|_2 \right]^q \right)^{1/q} \\ &\leq \left(\mathbb{E} \left[\max_{1 \leq i \leq t} \|\tilde{S}_i^{\tilde{\beta}_0}(\alpha) - \tilde{S}_i^\circ(\alpha)\|_2 + \max_{1 \leq i \leq t} \|\tilde{S}_i^\circ(\alpha) - \Sigma^{1/2}(\alpha) \tilde{G}_i^*\|_2 \right]^q \right)^{1/q} \\ &\leq \left(\mathbb{E} \left[\max_{1 \leq i \leq t} \|\tilde{S}_i^{\tilde{\beta}_0}(\alpha) - \tilde{S}_i^\circ(\alpha)\|_2 \right]^q \right)^{1/q} + \left(\mathbb{E} \left[\max_{1 \leq i \leq t} \|\tilde{S}_i^\circ(\alpha) - \Sigma^{1/2}(\alpha) \tilde{G}_i^*\|_2 \right]^q \right)^{1/q}. \end{aligned} \quad (148)$$

Therefore, to show the invariance principle for $(\tilde{S}_i^{\tilde{\beta}_0}(\alpha))_{1 \leq i \leq t}$, it suffices to bound the difference part $\max_{1 \leq i \leq t} \|\tilde{S}_i^{\tilde{\beta}_0}(\alpha) - \tilde{S}_i^\circ(\alpha)\|_2$ in terms of the q -th moment. To this end, recall the iterative function $f_D(\beta) = \beta + \alpha D X^\top (\mathbf{y} - X D \beta)$ in (9) that rewrites the GD dropout recursion (2). We note that

$$\begin{aligned} & \max_{1 \leq i \leq t} \|\tilde{S}_i^{\tilde{\beta}_0}(\alpha) - \tilde{S}_i^\circ(\alpha)\|_2 \\ &= \max_{1 \leq i \leq t} \left\| [f_{D_1}(\tilde{\beta}_0) - f_{D_1}(\tilde{\beta}_0^\circ)] + \cdots + [f_{D_i} \circ \cdots \circ f_{D_1}(\tilde{\beta}_0) - f_{D_i} \circ \cdots \circ f_{D_1}(\tilde{\beta}_0^\circ)] \right\|_2 \\ &\leq \max_{1 \leq i \leq t} \left(\|f_{D_1}(\tilde{\beta}_0) - f_{D_1}(\tilde{\beta}_0^\circ)\|_2 + \cdots + \|f_{D_i} \circ \cdots \circ f_{D_1}(\tilde{\beta}_0) - f_{D_i} \circ \cdots \circ f_{D_1}(\tilde{\beta}_0^\circ)\|_2 \right) \\ &= \|f_{D_1}(\tilde{\beta}_0) - f_{D_1}(\tilde{\beta}_0^\circ)\|_2 + \cdots + \|f_{D_t} \circ \cdots \circ f_{D_1}(\tilde{\beta}_0) - f_{D_t} \circ \cdots \circ f_{D_1}(\tilde{\beta}_0^\circ)\|_2. \end{aligned} \quad (149)$$

This, along with the triangle inequality and Theorem 3 yields

$$\begin{aligned}
& \left(\mathbb{E} \left[\max_{1 \leq i \leq t} \|\tilde{S}_i^{\tilde{\beta}_0}(\alpha) - \tilde{S}_i^\circ(\alpha)\|_2 \right]^q \right)^{1/q} \\
& \leq \left(\mathbb{E} \|f_{D_1}(\tilde{\beta}_0) - f_{D_1}(\tilde{\beta}_0^\circ)\|_2^q \right)^{1/q} + \dots + \left(\mathbb{E} \|f_{D_t} \circ \dots \circ f_{D_1}(\tilde{\beta}_0) - f_{D_t} \circ \dots \circ f_{D_1}(\tilde{\beta}_0^\circ)\|_2^q \right)^{1/q} \\
& \leq \frac{r_{\alpha,q}(1 - r_{\alpha,q}^t)}{1 - r_{\alpha,q}} \|\tilde{\beta}_0 - \tilde{\beta}_0^\circ\|_2 = o(t^{1/q}). \tag{150}
\end{aligned}$$

We insert this result back into (149), which together with (148) gives the invariance principle for the partial sum $\tilde{S}_i^{\tilde{\beta}_0}(\alpha) = \sum_{k=1}^i [\tilde{\beta}_k(\alpha) - \tilde{\beta}_{\text{mean}}(\alpha)]$.

Finally, let the partial sum $S_i^{\tilde{\beta}_0}(\alpha) = \sum_{k=1}^i [\tilde{\beta}_k(\alpha) - \tilde{\beta}]$ be as defined in Theorem 9. We shall bound the difference between $S_i^{\tilde{\beta}_0}(\alpha)$ and $\tilde{S}_i^{\tilde{\beta}_0}(\alpha)$. Since $\tilde{\beta}_{\text{mean}}(\alpha) = \mathbb{E}_D[\tilde{\beta}_1^\circ(\alpha)]$ as defined in (111) and $\tilde{\beta}_{\text{mean}}(\alpha) - \tilde{\beta} = \mathbb{E}_D[\tilde{\beta}_1^\circ(\alpha) - \tilde{\beta}] = 0$ by (137), it follows that

$$\left(\mathbb{E} \left[\max_{1 \leq i \leq t} \|\tilde{S}_i^{\tilde{\beta}_0}(\alpha) - S_i^{\tilde{\beta}_0}(\alpha)\|_2 \right]^q \right)^{1/q} = \left(\mathbb{E} \left[\max_{1 \leq i \leq t} \left\| \sum_{k=1}^i \tilde{\beta}_{\text{mean}}(\alpha) - \tilde{\beta} \right\|_2 \right]^q \right)^{1/q} = 0. \tag{151}$$

Combining this with the invariance principle for $(\tilde{S}_i^{\tilde{\beta}_0}(\alpha))_{1 \leq i \leq t}$, we obtain the same approximation rate $o_{\mathbb{P}}(t^{1/q})$ for the partial sum process $(S_i^{\tilde{\beta}_0}(\alpha))_{1 \leq i \leq t}$. This completes the proof. ■

Appendix E. Proofs in Section 4.2

Recall the SGD dropout sequence $\{\check{\beta}_k(\alpha)\}_{k \in \mathbb{N}}$ and the random coefficient $\check{A}_k(\alpha)$ in (36). To prove that $\sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \mathbb{E} \|\check{A}_k(\alpha) \mathbf{v}\|_2^q < 1$ is a sufficient condition for the geometric-moment contraction (GMC) of the SGD dropout sequence, we first introduce two useful moment inequalities in Lemma 27.

E.1 Proof of Lemma 10

Lemma 27 (Moment inequality) *Let $q \geq 2$. For any two random vectors \mathbf{x} and \mathbf{y} in \mathbb{R}^d with fixed $d \geq 1$, the following inequalities holds:*

- (i) $\mathbb{E} \|\mathbf{x} + \mathbf{y}\|_2^q - \|\mathbf{x}\|_2^q - q \|\mathbf{x}\|_2^{q-2} \mathbf{x}^\top \mathbf{y} \leq \mathbb{E} (\|\mathbf{x}\|_2 + \|\mathbf{y}\|_2)^q - \mathbb{E} \|\mathbf{x}\|_2^q - q \mathbb{E} (\|\mathbf{x}\|_2^{q-1} \|\mathbf{y}\|_2)$.
- (ii) $\mathbb{E} \|\mathbf{x} + \mathbf{y}\|_2^q - \|\mathbf{x}\|_2^q - q \|\mathbf{x}\|_2^{q-2} \mathbf{x}^\top \mathbf{y} \leq [(\mathbb{E} \|\mathbf{x}\|_2^q)^{1/q} + (\mathbb{E} \|\mathbf{y}\|_2^q)^{1/q}]^q - \mathbb{E} \|\mathbf{x}\|_2^q - q (\mathbb{E} \|\mathbf{x}\|_2^q)^{(q-1)/q} (\mathbb{E} \|\mathbf{y}\|_2^q)^{1/q}$.

Lemma 27(i) immediately follows if we can prove the inequality

$$\left| \|\mathbf{x} + \mathbf{y}\|_2^q - \|\mathbf{x}\|_2^q - q \|\mathbf{x}\|_2^{q-2} \mathbf{x}^\top \mathbf{y} \right| \leq (\|\mathbf{x}\|_2 + \|\mathbf{y}\|_2)^q - \|\mathbf{x}\|_2^q - q (\|\mathbf{x}\|_2^{q-1} \|\mathbf{y}\|_2), \tag{152}$$

which is of independent interest. The right hand side of the inequality in Lemma 27(ii) only depends on expectations of either one of the random vectors \mathbf{x} or \mathbf{y} . This makes the inequality particularly useful if \mathbf{x} and \mathbf{y} are dependent. Lemma 27(i) is more favorable in cases where \mathbf{x} and \mathbf{y} are independent, or if one vector is deterministic.

Proof [Lemma 27(i)] We can assume that $\|\mathbf{x}\|_2 > 0$ and $\|\mathbf{y}\|_2 > 0$ as otherwise, the inequality holds trivially. It is moreover sufficient to assume $\mathbf{x} = w\mathbf{e}_1$, for a positive number w and $\mathbf{e}_1 \in \mathbb{R}^d$ a unit vector. Then, we can find two numbers u, v such that

$$\mathbf{y} = u \cdot (w\mathbf{e}_1) + v\mathbf{e}_2, \quad (153)$$

where $\mathbf{e}_2 \in \mathbb{R}^d$ is a unit vector orthogonal to \mathbf{e}_1 . Let $r = \|\mathbf{y}\|_2 = \sqrt{(uw)^2 + v^2} > 0$. We note that $\mathbf{x}^\top \mathbf{y} = uw^2$ and $\|\mathbf{x} + \mathbf{y}\|_2^2 = (1+u)^2 w^2 + v^2$, which gives

$$\begin{aligned} & \|\mathbf{x} + \mathbf{y}\|_2^q - \|\mathbf{x}\|_2^q - q\|\mathbf{x}\|_2^{q-2} \mathbf{x}^\top \mathbf{y} \\ &= [(1+u)^2 w^2 + v^2]^{q/2} - w^q - qw^{q-2} uw^2 \\ &= (w^2 + 2uw^2 + r^2)^{q/2} - w^q - qw^q. \end{aligned} \quad (154)$$

Since $r = \sqrt{(uw)^2 + v^2}$, we can rewrite $uw = r\delta$ for some scalar δ with $\delta \in [-\delta^*, 1]$, where $\delta^* = (w^2 + r^2)/(2wr)$, i.e., $w^2 - 2wr\delta^* + r^2 = 0$. Here, $|\delta|$ can be viewed as the projection length of \mathbf{y} on the direction of \mathbf{x} , and the end point δ^* falls in $[0, 1]$. Then, (154) can be rewritten into

$$\varphi(\delta) := (w^2 + 2wr\delta + r^2)^{q/2} - w^q - qw^{q-1}r\delta. \quad (155)$$

Recall that $w = \|\mathbf{x}\|_2 > 0$ and $r = \|\mathbf{y}\|_2 > 0$. The first order derivative of $\varphi(\delta)$ is

$$\begin{aligned} \varphi'(\delta) &= \frac{q}{2} 2wr(w^2 + 2wr\delta + r^2)^{(q/2)-1} - qw^{q-1}r \\ &= qwr[(w^2 + 2wr\delta + r^2)^{(q/2)-1} - w^{q-2}] \\ &= qw^{q-1}r \left[\left(1 + \frac{2r\delta}{w} + \frac{r^2}{w^2} \right)^{(q/2)-1} - 1 \right]. \end{aligned} \quad (156)$$

This indicates that, for $q \geq 2$, $\varphi'(\delta) \leq 0$ when $\delta \in [-\delta^*, r/(2w)]$, and $\varphi'(\delta) > 0$ when $\delta \in (-r/(2w), 1]$. In particular, by Bernoulli's inequality, we can observe that

$$\begin{aligned} \varphi(-\delta^*) &= -w^q + qw^{q-1}r\delta^* = (q/2 - 1)w^q + qw^{q-2}r^2 > 0, \\ \varphi(-r/(2w)) &= -qw^{q-2}r^2/2 < 0, \\ \varphi(1) &= (w+r)^q - w^q - qw^{q-1}r > 0. \end{aligned} \quad (157)$$

Moreover, regarding $\varphi(-\delta^*)$ on $\delta^* \in [0, 1]$, we consider a new function $\tilde{\varphi}(s) = -w^q - qw^{q-1}rs$, which is decreasing on $s \in [-1, 0]$. Note that $\tilde{\varphi}(-1) = -w^q + qw^{q-1}r$. Thus, by comparison, we have $-\varphi(-r/(2w)) < \varphi(-\delta^*) \leq \tilde{\varphi}(-1) < \varphi(1)$. As a direct result, we obtain

$$\sup_{|\delta| \leq 1} |\varphi(\delta)| = \max\{\varphi(-\delta^*), -\varphi(-r/(2w)), \varphi(1)\} = \varphi(1). \quad (158)$$

By inserting $w = \|\mathbf{x}\|_2$ and $r = \|\mathbf{y}\|_2$ back to $\varphi(\delta)$, we obtain, for any \mathbf{x} and \mathbf{y} in \mathbb{R}^d ,

$$\left| \|\mathbf{x} + \mathbf{y}\|_2^q - \|\mathbf{x}\|_2^q - q\|\mathbf{x}\|_2^{q-2} \mathbf{x}^\top \mathbf{y} \right| \leq (\|\mathbf{x}\|_2 + \|\mathbf{y}\|_2)^q - \|\mathbf{x}\|_2^q - q\|\mathbf{x}\|_2^{q-1} \|\mathbf{y}\|_2. \quad (159)$$

The desired result holds by taking the expectation on the both sides. \blacksquare

Proof [Lemma 27(ii)] First, we define a function $\phi(t) = \|\mathbf{x} + t\mathbf{y}\|_2^q$ on $t \in [0, \infty)$. The first and second order derivatives of $\phi(t)$ are as follows

$$\begin{aligned}\phi'(t) &= \frac{d}{dt}\phi(t) = q\|\mathbf{x} + t\mathbf{y}\|_2^{q-2}(t\|\mathbf{y}\|_2^2 + \mathbf{x}^\top \mathbf{y}), \\ \phi''(t) &= \frac{d^2}{dt^2}\phi(t) = q(q-2)\|\mathbf{x} + t\mathbf{y}\|_2^{q-4}(t\|\mathbf{y}\|_2^2 + \mathbf{x}^\top \mathbf{y})^2 + q\|\mathbf{x} + t\mathbf{y}\|_2^{q-2}\|\mathbf{y}\|_2^2 \\ &= q(q-2)\|\mathbf{x} + t\mathbf{y}\|_2^{q-4}[(\mathbf{x} + t\mathbf{y})^\top \mathbf{y}]^2 + q\|\mathbf{x} + t\mathbf{y}\|_2^{q-2}\|\mathbf{y}\|_2^2 \\ &\leq q(q-1)\|\mathbf{x} + t\mathbf{y}\|_2^{q-2}\|\mathbf{y}\|_2^2,\end{aligned}\tag{160}$$

where the last inequality follows from the Cauchy-Schwarz inequality. In the previous inequality, equality holds when both random vectors \mathbf{x} and \mathbf{y} are scalars. Note that $\phi(1) = \|\mathbf{x} + \mathbf{y}\|_2^q$, $\phi(0) = \|\mathbf{x}\|_2^q$, and $\phi'(0) = q\|\mathbf{x}\|_2^{q-2}\mathbf{x}^\top \mathbf{y}$. Since

$$\begin{aligned}\phi(1) - \phi(0) - \phi'(0) &= \int_{t=0}^1 \phi'(t)dt - \phi'(0) \\ &= \int_{t=0}^1 \left(\int_{s=0}^t \phi''(s)ds + \phi'(0) \right) dt - \phi'(0) \\ &= \int_{t=0}^1 \int_{s=0}^t \phi''(s)dsdt,\end{aligned}\tag{161}$$

it follows from the upper bound of $\phi''(s)$ in (160) that

$$\|\mathbf{x} + \mathbf{y}\|_2^q - \|\mathbf{x}\|_2^q - q\|\mathbf{x}\|_2^{q-2}\mathbf{x}^\top \mathbf{y} \leq q(q-1) \int_{t=0}^1 \int_{s=0}^t \|\mathbf{x} + s\mathbf{y}\|_2^{q-2}\|\mathbf{y}\|_2^2 dsdt.\tag{162}$$

Taking the expectation yields

$$\mathbb{E}\left(\|\mathbf{x} + \mathbf{y}\|_2^q - \|\mathbf{x}\|_2^q - q\|\mathbf{x}\|_2^{q-2}\mathbf{x}^\top \mathbf{y}\right) \leq q(q-1) \int_{t=0}^1 \int_{s=0}^t \mathbb{E}(\|\mathbf{x} + s\mathbf{y}\|_2^{q-2}\|\mathbf{y}\|_2^2) dsdt.\tag{163}$$

Furthermore, by Hölder's inequality and the triangle inequality, we obtain

$$\begin{aligned}\mathbb{E}(\|\mathbf{x} + s\mathbf{y}\|_2^{q-2}\|\mathbf{y}\|_2^2) &\leq (\mathbb{E}\|\mathbf{x} + s\mathbf{y}\|_2^q)^{(q-2)/q} (\mathbb{E}\|\mathbf{y}\|_2^q)^{2/q} \\ &\leq \left[(\mathbb{E}\|\mathbf{x}\|_2^q)^{1/q} + s(\mathbb{E}\|\mathbf{y}\|_2^q)^{1/q} \right]^{q-2} (\mathbb{E}\|\mathbf{y}\|_2^q)^{2/q},\end{aligned}\tag{164}$$

which together with (163) gives

$$\begin{aligned}&\mathbb{E}\left(\|\mathbf{x} + \mathbf{y}\|_2^q - \|\mathbf{x}\|_2^q - q\|\mathbf{x}\|_2^{q-2}\mathbf{x}^\top \mathbf{y}\right) \\ &\leq q(q-1) \int_{t=0}^1 \int_{s=0}^t \left[(\mathbb{E}\|\mathbf{x}\|_2^q)^{1/q} + s(\mathbb{E}\|\mathbf{y}\|_2^q)^{1/q} \right]^{q-2} (\mathbb{E}\|\mathbf{y}\|_2^q)^{2/q} dsdt \\ &= q \int_{t=0}^1 \left\{ \left[(\mathbb{E}\|\mathbf{x}\|_2^q)^{1/q} + t(\mathbb{E}\|\mathbf{y}\|_2^q)^{1/q} \right]^{q-1} - (\mathbb{E}\|\mathbf{x}\|_2^q)^{(q-1)/q} \right\} (\mathbb{E}\|\mathbf{y}\|_2^q)^{1/q} dt \\ &= \left[(\mathbb{E}\|\mathbf{x}\|_2^q)^{1/q} + (\mathbb{E}\|\mathbf{y}\|_2^q)^{1/q} \right]^q - \mathbb{E}\|\mathbf{x}\|_2^q - q(\mathbb{E}\|\mathbf{x}\|_2^q)^{(q-1)/q} (\mathbb{E}\|\mathbf{y}\|_2^q)^{1/q}.\end{aligned}\tag{165}$$

In addition, recall that by the proof of Lemma 27(i), we have

$$\begin{aligned} \left| \|\mathbf{x} + \mathbf{y}\|_2^q - \|\mathbf{x}\|_2^q - q\|\mathbf{x}\|_2^{q-2}\mathbf{x}^\top\mathbf{y} \right| &\leq (\|\mathbf{x}\|_2 + \|\mathbf{y}\|_2)^q - \|\mathbf{x}\|_2^q - q\|\mathbf{x}\|_2^{q-1}\|\mathbf{y}\|_2 \\ &= q(q-1) \int_{t=0}^1 \int_{s=0}^t (\|\mathbf{x}\|_2 + s\|\mathbf{y}\|_2)^{q-2} \|\mathbf{y}\|_2^2 ds dt. \end{aligned}$$

Taking expectation on both sides, we obtain

$$\mathbb{E} \left| \|\mathbf{x} + \mathbf{y}\|_2^q - \|\mathbf{x}\|_2^q - q\|\mathbf{x}\|_2^{q-2}\mathbf{x}^\top\mathbf{y} \right| \leq q(q-1) \int_0^1 \int_0^t \mathbb{E} [(\|\mathbf{x}\|_2 + s\|\mathbf{y}\|_2)^{q-2} \|\mathbf{y}\|_2^2] ds dt. \quad (166)$$

It follows from Hölder's inequality and the triangle inequality that

$$\begin{aligned} \mathbb{E} [(\|\mathbf{x}\|_2 + s\|\mathbf{y}\|_2)^{q-2} \|\mathbf{y}\|_2^2] &\leq [\mathbb{E}(\|\mathbf{x}\|_2 + s\|\mathbf{y}\|_2)^q]^{(q-2)/q} (\mathbb{E}\|\mathbf{y}\|_2^q)^{2/q} \\ &\leq \left[(\mathbb{E}\|\mathbf{x}\|_2^q)^{1/q} + s(\mathbb{E}\|\mathbf{y}\|_2^q)^{1/q} \right]^{q-2} (\mathbb{E}\|\mathbf{y}\|_2^q)^{2/q}. \end{aligned} \quad (167)$$

Evaluating the double integral as in (165) yields

$$\begin{aligned} &\mathbb{E} \left| \|\mathbf{x} + \mathbf{y}\|_2^q - \|\mathbf{x}\|_2^q - q\|\mathbf{x}\|_2^{q-2}\mathbf{x}^\top\mathbf{y} \right| \\ &\leq \left[(\mathbb{E}\|\mathbf{x}\|_2^q)^{1/q} + (\mathbb{E}\|\mathbf{y}\|_2^q)^{1/q} \right]^q - \mathbb{E}\|\mathbf{x}\|_2^q - q(\mathbb{E}\|\mathbf{x}\|_2^q)^{(q-1)/q} (\mathbb{E}\|\mathbf{y}\|_2^q)^{1/q}. \end{aligned}$$

This completes the proof. ■

Proof [Lemma 10] Since a dropout matrix D_k is a diagonal matrix with values 0 and 1 on the diagonal, $D_k^2 = D_k$ and

$$\begin{aligned} \|(I_d - \alpha D_k \mathbb{X}_k D_k) \mathbf{v}\|_2^2 &= \mathbf{v}^\top (I_d - 2\alpha D_k \mathbb{X}_k D_k + \alpha^2 D_k \mathbb{X}_k D_k \mathbb{X}_k D_k) \mathbf{v} \\ &= 1 - 2\alpha \mathbf{v}^\top D_k \mathbb{X}_k D_k \mathbf{v} + \alpha^2 \mathbf{v}^\top D_k \mathbb{X}_k D_k^2 \mathbb{X}_k D_k \mathbf{v} \\ &= 1 - \alpha \mathbf{v}^\top [2D_k \mathbb{X}_k D_k - \alpha D_k \mathbb{X}_k D_k \mathbb{X}_k D_k] \mathbf{v} \\ &=: 1 - \alpha \mathbf{v}^\top M_k \mathbf{v}, \end{aligned} \quad (168)$$

with

$$M_k(\alpha) = 2D_k \mathbb{X}_k D_k - \alpha D_k \mathbb{X}_k D_k \mathbb{X}_k D_k. \quad (169)$$

Recall the condition on the learning rate α in (40), it follows that $\mathbb{E}[M_k]$ is positive definite (p.d.), which further implies the lower bound $\mathbb{E}(\alpha \mathbf{v}^\top M_k \mathbf{v}) \geq \alpha \lambda_{\min}(\mathbb{E}[M_k]) > 0$, that holds uniformly over all unit vectors \mathbf{v} . As a direct consequence,

$$\sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \mathbb{E} \|(I_d - \alpha D_k \mathbb{X}_k D_k) \mathbf{v}\|_2^2 \leq \sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} (1 - \mathbb{E}(\alpha \mathbf{v}^\top M_k \mathbf{v})) < 1, \quad (170)$$

proving the result in the case $q = 2$.

Next, we shall show that for all $q > 2$, we also have

$$\sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \mathbb{E} \|(I_d - \alpha D_k \mathbb{X}_k D_k) \mathbf{v}\|_2^q < 1.$$

In this case, the techniques in the proof of Lemma 2 cannot be directly applied due to the randomness of \mathbb{X}_k , and we need to leverage the moment inequalities in Lemma 27 instead. Specifically, let \mathbf{x} and \mathbf{y} in Lemma 27 be

$$\mathbf{x} = \mathbf{x}(\mathbf{v}) = \mathbf{v}, \quad \mathbf{y} = \mathbf{y}(\mathbf{v}) = D_k \mathbb{X}_k D_k \mathbf{v}, \quad (171)$$

respectively, for \mathbf{v} a deterministic d -dimensional unit vector. It remains to show that for any $q > 2$, $\mathbb{E}\|\mathbf{x} - \alpha\mathbf{y}\|_2^q < 1$ holds for any arbitrary unit vector \mathbf{v} . By Lemma 27,

$$\mathbb{E}\|\mathbf{x} - \alpha\mathbf{y}\|_2^q - \|\mathbf{x}\|_2^q - q\|\mathbf{x}\|_2^{q-2}\mathbb{E}(-\mathbf{x}^\top \alpha\mathbf{y}) \leq \mathbb{E}(\|\mathbf{x}\|_2 + \|\alpha\mathbf{y}\|_2)^q - \|\mathbf{x}\|_2^q - q\|\mathbf{x}\|_2^{q-1}\mathbb{E}\|\alpha\mathbf{y}\|_2, \quad (172)$$

which along with $\|\mathbf{x}\|_2 = \|\mathbf{v}\|_2 = 1$ further yields,

$$\mathbb{E}\|\mathbf{x} - \alpha\mathbf{y}\|_2^q - 1 + q\alpha\mathbb{E}(\mathbf{x}^\top \mathbf{y}) \leq \mathbb{E}(1 + \alpha\|\mathbf{y}\|_2)^q - 1 - q\alpha\mathbb{E}\|\mathbf{y}\|_2. \quad (173)$$

Therefore, to prove $\sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \mathbb{E}\|\mathbf{x} - \alpha\mathbf{y}\|_2^q < 1$, it suffices to show

$$\mathbb{E}(1 + \alpha\|\mathbf{y}\|_2)^q - 1 - q\alpha\mathbb{E}\|\mathbf{y}\|_2 < q\alpha\mathbb{E}(\mathbf{x}^\top \mathbf{y}). \quad (174)$$

By applying Lemma 27 again, we have

$$\mathbb{E}(1 + \alpha\|\mathbf{y}\|_2)^q - 1 - q\alpha\mathbb{E}\|\mathbf{y}\|_2 \leq (1 + \alpha(\mathbb{E}\|\mathbf{y}\|_2^q)^{1/q})^q - 1 - q\alpha(\mathbb{E}\|\mathbf{y}\|_2^q)^{1/q}. \quad (175)$$

Thus, we only need to show that for any d -dimensional vector \mathbf{v} ,

$$(1 + \alpha(\mathbb{E}\|\mathbf{y}\|_2^q)^{1/q})^q - 1 - q\alpha(\mathbb{E}\|\mathbf{y}\|_2^q)^{1/q} < q\alpha\mathbb{E}(\mathbf{x}^\top \mathbf{y}). \quad (176)$$

Recall the definitions of \mathbf{x} and \mathbf{y} in (171). By Lemma 21 (i), it follows that $\mathbb{E}(\mathbf{x}^\top \mathbf{y}) = p\mathbf{v}^\top \mathbb{E}[\mathbb{X}_{k,p}]\mathbf{v}$, where $\mathbb{X}_{k,p} = p\mathbb{X}_k + (1-p)\text{Diag}(\mathbb{X}_k)$. With $\mu_q = \mu_q(\mathbf{v}) = (\mathbb{E}\|\mathbf{y}(\mathbf{v})\|_2^q)^{1/q} = (\mathbb{E}\|D_k \mathbb{X}_k D_k \mathbf{v}\|_2^q)^{1/q}$, it suffices to show that

$$\sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \frac{(1 + \alpha\mu_q)^q - 1 - q\alpha\mu_q}{p\mathbf{v}^\top \mathbb{E}[\mathbb{X}_{k,p}]\mathbf{v}} < q\alpha. \quad (177)$$

Let $q > 2$. Consider a function $f : \mathbb{R}_+ \mapsto \mathbb{R}$ with $f(t) = (1+t)^q - 1 - qt$, which is strictly increasing on \mathbb{R}_+ . Note that $f''(t) = q(q-1)(1+t)^{q-2}$. Since $q > 2$, it follows that $f(t) = \int_{s=0}^t \int_{u=0}^s q(q-1)(1+u)^{q-2} du ds$. Therefore,

$$f(t) = (1+t)^q - 1 - qt \leq \frac{q(q-1)}{2}(1+t)^{q-2}t^2. \quad (178)$$

Thus, the condition (177) is satisfied if

$$\frac{q(q-1)}{2}\alpha^2 \sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \frac{(1 + \alpha\mu_q)^{q-2}\mu_q^2}{p\mathbf{v}^\top \mathbb{E}[\mathbb{X}_{k,p}]\mathbf{v}} < q\alpha. \quad (179)$$

As this is true by assumption the proof is complete. \blacksquare

E.2 Proof of Lemma 13

Proof For any fixed unit vector $\mathbf{u} \in \mathbb{R}^d$, since $\check{A}_k(\alpha) = I_d - \alpha D_k \mathbb{X}_k D_k$ as defined in (36) with $\mathbb{X}_k = \mathbf{x}_k \mathbf{x}_k^\top$, it follows that

$$\begin{aligned} \|\check{A}_k(\alpha) \mathbf{v}\|_2^2 &= \|(I_d - \alpha D_k \mathbf{x}_k \mathbf{x}_k^\top D_k) \mathbf{v}\|_2^2 \\ &= \|\mathbf{v}\|_2^2 - 2\alpha \mathbf{v}^\top D_k \mathbf{x}_k \mathbf{x}_k^\top D_k \mathbf{v} + \alpha^2 \mathbf{v}^\top D_k \mathbf{x}_k \mathbf{x}_k^\top D_k^2 \mathbf{x}_k \mathbf{x}_k^\top D_k \mathbf{v} \\ &= 1 - 2\alpha (\mathbf{x}_k^\top D_k \mathbf{v})^2 + \alpha^2 \|D_k \mathbf{x}_k\|_2^2 (\mathbf{x}_k^\top D_k \mathbf{v})^2. \end{aligned} \quad (180)$$

Write $\mathbf{v} = (v_1, \dots, v_d)^\top$ and $\mathbf{x}_k = (x_{k1}, \dots, x_{kd})^\top$. Let $D_{k,jj}$ be the j -th diagonal element in D_k . Since $\mathbf{x}_k \sim \mathcal{N}(0, I_d)$ and D_k is independent of \mathbf{x}_k , it follows that

$$\mathbb{E}(\mathbf{x}_k^\top D_k \mathbf{v})^2 = \sum_{j=1}^d v_j^2 \mathbb{E}[D_{k,jj} x_{kj}^2] = \sum_{j=1}^d v_j^2 p = p. \quad (181)$$

Moreover,

$$\begin{aligned} \|D_k \mathbf{x}_k\|_2^2 (\mathbf{x}_k^\top D_k \mathbf{v})^2 &= \left(\sum_{i=1}^d D_{k,ii} x_{ki}^2 \right) \left(\sum_{j,l=1}^d D_{k,jj} D_{k,ll} v_j v_l x_{kj} x_{kl} \right) \\ &= \sum_{i,j,l=1}^d D_{k,ii} D_{k,jj} D_{k,ll} v_j v_l x_{ki}^2 x_{kj} x_{kl}. \end{aligned} \quad (182)$$

For $j \neq l$, $\mathbb{E}[x_{ki}^2 x_{kj} x_{kl}] = 0$. For $j = l$ and $i \neq j$, $\mathbb{E}[x_{ki}^2 x_{kj}^2] = 1$. For $j = l$ and $i = j$, $\mathbb{E}[x_{ki}^4] = 3$. Therefore, by the independence of the diagonal elements in D_k , we have

$$\begin{aligned} \mathbb{E}[\|D_k \mathbf{x}_k\|_2^2 (\mathbf{x}_k^\top D_k \mathbf{v})^2] &= \sum_{j=1}^d v_j^2 \mathbb{E} \left(3D_{k,jj} + \sum_{i \neq j} D_{k,ii} D_{k,jj} \right) \\ &= \sum_{j=1}^d v_j^2 (3p + p^2(d-1)) \\ &= 3p + p^2(d-1). \end{aligned} \quad (183)$$

Finally, we obtain

$$\mathbb{E} \|\check{A}_k(\alpha) \mathbf{v}\|_2^2 = 1 - 2\alpha p + \alpha^2 [3p + (d-1)p^2],$$

for all unit vectors $\mathbf{v} \in \mathbb{R}^d$. This completes the proof. \blacksquare

E.3 Proof of Theorem 12

Proof Let the random coefficient matrix $\check{A}_k(\alpha) = I_d - \alpha D_k \mathbb{X}_k D_k$ be as defined in (36). We write $\check{A}_k(\alpha) = \check{A}_k$ exchangeably in this proof.

First, we study the case with $q = 2$. Consider two SGD dropout sequences $\{\check{\beta}_k(\alpha)\}_{k \in \mathbb{N}}$ and $\{\check{\beta}'_k(\alpha)\}_{k \in \mathbb{N}}$, given two arbitrarily fixed initial vectors $\check{\beta}_0, \check{\beta}'_0$. Let $\check{\delta} = \check{\beta}_0 - \check{\beta}'_0$. Then, it follows from the tower rule that

$$\begin{aligned}
\mathbb{E}\|\check{\beta}_k(\alpha) - \check{\beta}'_k(\alpha)\|_2^2 &= \mathbb{E}[\check{\delta}^\top \check{A}_1^\top \cdots \check{A}_k^\top \check{A}_k \cdots \check{A}_1 \check{\delta}] \\
&= \mathbb{E}[\mathbb{E}[\check{\delta}^\top \check{A}_1^\top \cdots \check{A}_k^\top \check{A}_k \cdots \check{A}_1 \check{\delta} \mid \check{A}_1, \dots, \check{A}_{k-1}]] \\
&= \mathbb{E}[\check{\delta}^\top \check{A}_1^\top \cdots \check{A}_{k-1}^\top \mathbb{E}(\check{A}_k^\top \check{A}_k) \check{A}_{k-1} \cdots \check{A}_1 \check{\delta}] \\
&\leq \|\mathbb{E}(\check{A}_k^\top \check{A}_k)\| \cdot \mathbb{E}[\check{\delta}^\top \check{A}_1^\top \cdots \check{A}_{k-1}^\top \check{A}_{k-1} \cdots \check{A}_1 \check{\delta}] \\
&\leq \prod_{i=1}^k \|\mathbb{E}(\check{A}_i^\top \check{A}_i)\| \cdot \|\check{\delta}\|_2^2.
\end{aligned} \tag{184}$$

Recall that for the constant learning rate $\alpha > 0$ satisfying the conditions in Lemma 10, we have $\|\mathbb{E}(\check{A}_i^\top \check{A}_i)\| < 1$ uniformly over $i \in \mathbb{N}$. Thus, $\prod_{i=1}^k \|\mathbb{E}(\check{A}_i^\top \check{A}_i)\| < 1$. Since the dropout matrices D_k 's are i.i.d. and are independent of the i.i.d. observations \mathbf{x}_k 's, it follows that $\prod_{i=1}^k \|\mathbb{E}(\check{A}_i^\top \check{A}_i)\| = \|\mathbb{E}(\check{A}_1^\top \check{A}_1)\|^k < 1$. This gives the desired result for the case with $q = 2$.

For $q > 2$, we note that

$$\begin{aligned}
\mathbb{E}\|\check{\beta}_k(\alpha) - \check{\beta}'_k(\alpha)\|_2^q &= \mathbb{E}\|\check{A}_k \cdots \check{A}_1 \check{\delta}\|_2^q = \mathbb{E}(\|\check{A}_k \cdots \check{A}_1 \check{\delta}\|_2^2)^{q/2} \\
&= \mathbb{E}(\check{\delta}^\top \check{A}_1^\top \cdots \check{A}_k^\top \check{A}_k \cdots \check{A}_1 \check{\delta})^{q/2}.
\end{aligned} \tag{185}$$

Similarly, it follows from the tower rule that

$$\begin{aligned}
\mathbb{E}(\check{\delta}^\top \check{A}_1^\top \cdots \check{A}_k^\top \check{A}_k \cdots \check{A}_1 \check{\delta})^{q/2} &= \mathbb{E}[\mathbb{E}[(\check{\delta}^\top \check{A}_1^\top \cdots \check{A}_k^\top \check{A}_k \cdots \check{A}_1 \check{\delta})^{q/2} \mid \check{A}_1, \dots, \check{A}_{k-1}]] \\
&\leq \sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \mathbb{E}\|\check{A}_k \mathbf{v}\|_2^q \cdot \mathbb{E}(\check{\delta}^\top \check{A}_1^\top \cdots \check{A}_{k-1}^\top \check{A}_{k-1} \cdots \check{A}_1 \check{\delta})^{q/2} \\
&\leq \prod_{i=1}^k \sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \mathbb{E}\|\check{A}_i \mathbf{v}\|_2^q \cdot \|\check{\delta}\|_2^q.
\end{aligned} \tag{186}$$

Since $\sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \mathbb{E}\|\check{A}_i \mathbf{v}\|_2^q < 1$ holds uniformly over $i \in \mathbb{N}$, we obtain the geometric-moment contraction in (41) for $q > 2$.

Finally, by Corollary 23, the geometric-moment contraction in (41) implies the existence of a unique stationary distribution $\check{\pi}_\alpha$ of the SGD dropout $\{\check{\beta}_k(\alpha)\}_{k \in \mathbb{N}}$. This completes the proof. \blacksquare

E.4 Proofs of Lemmas 28–30

Lemma 28 (Closed-form solution of the ℓ^2 minimizer) *Assume that model (31) is in reduced form, i.e., $\min_i \mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top]_{ii} > 0$. Then, for the minimizer of the ℓ^2 -regularized least-squares loss $\check{\beta} := \arg \min_{\beta \in \mathbb{R}^d} \mathbb{E}[(y - \mathbf{x}^\top D \beta)^2 / 2]$ as defined in (34), we have the closed form solution*

$$\check{\beta} = (\mathbb{E}[\mathbb{X}_{1,p}])^{-1} \mathbb{E}[y_1 \mathbf{x}_1].$$

Proof Recall the $d \times d$ Gram matrix $\mathbb{X}_k = \mathbf{x}_k \mathbf{x}_k^\top$ and

$$\bar{\mathbb{X}}_k = \mathbb{X}_k - \text{Diag}(\mathbb{X}_k), \quad \mathbb{X}_{k,p} = p\mathbb{X}_k + (1-p)\text{Diag}(\mathbb{X}_k).$$

To obtain the closed form solution, we first compute the gradient of the ℓ^2 -regularized least-squares loss as follows,

$$\begin{aligned} (y - \mathbf{x}^\top \boldsymbol{\beta})^2 &= y^2 - 2y\mathbf{x}^\top D\boldsymbol{\beta} + \boldsymbol{\beta}^\top D(\mathbf{x}\mathbf{x}^\top)D\boldsymbol{\beta}, \\ \mathbb{E}_D[(y - \mathbf{x}^\top \boldsymbol{\beta})^2] &= y^2 - 2py\mathbf{x}^\top \boldsymbol{\beta} + p^2\boldsymbol{\beta}^\top (\mathbf{x}\mathbf{x}^\top)\boldsymbol{\beta} + p(1-p)\boldsymbol{\beta}^\top \text{Diag}(\mathbf{x}\mathbf{x}^\top)\boldsymbol{\beta}, \\ \mathbb{E}_{(y,\mathbf{x})}\mathbb{E}_D[(y - \mathbf{x}^\top \boldsymbol{\beta})^2] &= \mathbb{E}[y^2] - 2p\mathbb{E}[y\mathbf{x}^\top]\boldsymbol{\beta} + p^2\boldsymbol{\beta}^\top \mathbb{E}[\mathbf{x}\mathbf{x}^\top]\boldsymbol{\beta} + p(1-p)\boldsymbol{\beta}^\top \text{Diag}(\mathbb{E}[\mathbf{x}\mathbf{x}^\top])\boldsymbol{\beta}, \\ \nabla_{\boldsymbol{\beta}}\mathbb{E}_{(y,\mathbf{x})}\mathbb{E}_D[(y - \mathbf{x}^\top \boldsymbol{\beta})^2] &= -2p\mathbb{E}[y\mathbf{x}] + 2\left(p^2\mathbb{E}[\mathbf{x}\mathbf{x}^\top] + p(1-p)\text{Diag}(\mathbb{E}[\mathbf{x}\mathbf{x}^\top])\right)\boldsymbol{\beta}. \end{aligned}$$

Recall that the i.i.d. random noise ϵ_k is independent of the i.i.d. random covariates \mathbf{x}_k . Since model (31) is assumed to be in a reduced form, i.e., $\min_i \mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top]_{ii} > 0$, the closed form solution of $\check{\boldsymbol{\beta}}$ is

$$\check{\boldsymbol{\beta}} = p\left(p^2\mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top] + p(1-p)\text{Diag}(\mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top])\right)^{-1} \mathbb{E}[y_1 \mathbf{x}_1] = (\mathbb{E}[\mathbb{X}_{1,p}])^{-1} \mathbb{E}[y_1 \mathbf{x}_1].$$

This completes the proof. \blacksquare

Recall that for any $d \times d$ matrix A , $A_p := pA + (1-p)\text{Diag}(A)$.

Lemma 29 *If the $d \times d$ matrix $\mathbb{E}[2\mathbb{X}_k - \alpha\mathbb{X}_k^2]_p$ is positive definite, then the condition on the learning rate α in (37) holds for $q = 2$.*

Proof By rewriting the condition (37) with $q = 2$, for all the unit vector $\mathbf{v} \in \mathbb{R}^d$, $\|\mathbf{v}\|_2 = 1$, we aim to show

$$0 < \alpha < \frac{2\mathbf{v}^\top \mathbb{E}(D_k \mathbb{X}_k D_k) \mathbf{v}}{\mathbb{E}\|D_k \mathbb{X}_k D_k \mathbf{v}\|_2^2}. \quad (187)$$

Since $D_k^2 = D_k \leq I_d$, it follows from Lemma 21(ii) that

$$\begin{aligned} &2\mathbf{v}^\top \mathbb{E}(D_k \mathbb{X}_k D_k) \mathbf{v} - \alpha \mathbf{v}^\top \mathbb{E}[D_k \mathbb{X}_k D_k \mathbb{X}_k D_k] \mathbf{v} \\ &\geq 2\mathbf{v}^\top \mathbb{E}(D_k \mathbb{X}_k D_k) \mathbf{v} - \alpha \mathbf{v}^\top \mathbb{E}[D_k \mathbb{X}_k^2 D_k] \mathbf{v} \\ &= \mathbf{v}^\top \mathbb{E}[D_k (2\mathbb{X}_k - \alpha\mathbb{X}_k^2) D_k] \mathbf{v} \\ &= p\mathbf{v}^\top \mathbb{E}[2\mathbb{X}_k - \alpha\mathbb{X}_k^2]_p \mathbf{v} > 0. \end{aligned}$$

As the unit vector $\mathbf{v} \in \mathbb{R}^d$ was arbitrary, condition (37) holds for $q = 2$. \blacksquare

Lemma 30 (ℓ^2 -minimizer $\check{\boldsymbol{\beta}}$ and true parameter $\boldsymbol{\beta}^*$) *Assume that $\mathbb{E}[|\epsilon|^{2q}] + \|\mathbf{x}\|_2^{2q} < \infty$. Then, the q -th moment of the gradient in (32) exists at the true parameter $\boldsymbol{\beta}^*$ in model (31), for some $q \geq 2$, that is,*

$$\left(\mathbb{E}\left\|\nabla_{\boldsymbol{\beta}^*} \frac{1}{2}(y - \mathbf{x}^\top D\boldsymbol{\beta}^*)^2\right\|_2^q\right)^{1/q} = \left(\mathbb{E}\|D\mathbf{x}(y - \mathbf{x}^\top D\boldsymbol{\beta}^*)\|_2^q\right)^{1/q} < \infty,$$

which further implies the finite q -th moment of the stochastic gradient at the ℓ^2 -minimizer $\check{\beta}$ defined in (34), that is

$$\left(\mathbb{E}\left\|\nabla_{\check{\beta}}\frac{1}{2}(y-\mathbf{x}^\top D\check{\beta})^2\right\|_2^q\right)^{1/q}=\left(\mathbb{E}\|D\mathbf{x}(y-\mathbf{x}^\top D\check{\beta})\|_2^q\right)^{1/q}<\infty,$$

Proof First, it follows from the triangle inequality that

$$\begin{aligned} & \left(\mathbb{E}\|D\mathbf{x}(y-\mathbf{x}^\top D\beta^*)\|_2^q\right)^{1/q} \\ &= \left(\mathbb{E}\|D\mathbf{x}(\mathbf{x}^\top \beta^* + \epsilon - \mathbf{x}^\top D\beta^*)\|_2^q\right)^{1/q} \\ &\leq \left(\mathbb{E}\|D\mathbf{x}\mathbf{x}^\top \beta^*\|_2^q\right)^{1/q} + \left(\mathbb{E}\|D\mathbf{x}\epsilon\|_2^q\right)^{1/q} + \left(\mathbb{E}\|D\mathbf{x}\mathbf{x}^\top D\beta^*\|_2^q\right)^{1/q}. \end{aligned} \quad (188)$$

By Assumption 1, since the dimension of β^* is fixed, we have

$$\left(\mathbb{E}\|D\mathbf{x}\mathbf{x}^\top \beta^*\|_2^q\right)^{1/q}\leq\left(\mathbb{E}\|\mathbf{x}\|_2^{2q}\right)^{1/q}\|\beta^*\|_2<\infty.$$

Due the independence between \mathbf{x} and ϵ , Assumption 1 gives

$$\left(\mathbb{E}\|D\mathbf{x}\epsilon\|_2^q\right)^{1/q}\leq\left(\mathbb{E}\|\mathbf{x}\|_2^q\right)^{1/q}\left(\mathbb{E}\|\epsilon\|_2^q\right)^{1/q}<\infty.$$

Moreover, we obtain

$$\left(\mathbb{E}\|D\mathbf{x}\mathbf{x}^\top D\beta^*\|_2^q\right)^{1/q}=\left(\mathbb{E}\|D\mathbf{x}\|_2^{2q}\right)^{1/q}\|\beta^*\|_2\leq\left(\mathbb{E}\|\mathbf{x}\|_2^{2q}\right)^{1/q}\|\beta^*\|_2<\infty. \quad (189)$$

Inserting the inequalities into (188), we obtain the finite q -th moment at the true parameter β^* . Replacing β^* in (189) by any vector $\beta \in \mathbb{R}^d$ satisfying $\|\beta\|_2 < \infty$, the result still holds.

Next, we show that the finite q -th moment of the stochastic gradient at β^* can also imply the finite q -th moment at $\check{\beta}$. Note that

$$\begin{aligned} & \left(\mathbb{E}\|D_1\mathbf{x}_1(y_1-\mathbf{x}_1^\top D_1\check{\beta})\|_2^q\right)^{1/q} \\ &\leq\left(\mathbb{E}\|D_1\mathbf{x}_1(y_1-\mathbf{x}_1^\top D_1\beta^*)\|_2^q\right)^{1/q}+\left(\mathbb{E}\|D_1\mathbb{X}_1 D_1(\check{\beta}-\beta^*)\|_2^q\right)^{1/q}. \end{aligned} \quad (190)$$

We only need to show that the second term is bounded. Since $\mathbb{X}_{1,p} = p\mathbb{X}_1 + (1-p)\text{Diag}(\mathbb{X}_1)$, $\bar{\mathbb{X}}_1 = \mathbb{X}_1 - \text{Diag}(\mathbb{X}_1)$, and $\check{\beta} = (\mathbb{E}[\mathbb{X}_{1,p}])^{-1}\mathbb{E}[y_1\mathbf{x}_1]$, it follows that

$$\begin{aligned} \mathbb{E}\|D_1\mathbb{X}_1 D_1(\check{\beta}-\beta^*)\|_2^q &= \mathbb{E}\|D_1\mathbb{X}_1 D_1((\mathbb{E}[\mathbb{X}_{1,p}])^{-1}\mathbb{E}[y_1\mathbf{x}_1]-\beta^*)\|_2^q \\ &= \mathbb{E}\|D_1\mathbb{X}_1 D_1((\mathbb{E}[\mathbb{X}_{1,p}])^{-1}\mathbb{E}[(\mathbf{x}_1^\top \beta^* + \epsilon_1)\mathbf{x}_1]-\beta^*)\|_2^q \\ &= \mathbb{E}\|D_1\mathbb{X}_1 D_1((\mathbb{E}[\mathbb{X}_{1,p}])^{-1}\mathbb{E}[\mathbb{X}_1]\beta^*-\beta^*)\|_2^q \\ &= (1-p)^q \mathbb{E}\|D_1\mathbb{X}_1 D_1(\mathbb{E}[\mathbb{X}_{1,p}])^{-1}\mathbb{E}[\bar{\mathbb{X}}_1]\beta^*\|_2^q \\ &\leq (1-p)^q \left\|(\mathbb{E}[\mathbb{X}_{1,p}])^{-1}\mathbb{E}[\bar{\mathbb{X}}_1]\right\|_2^q \sup_{\mathbf{v}\in\mathbb{R}^d,\|\mathbf{v}\|_2=1} \mathbb{E}\|D_1\mathbb{X}_1 D_1\mathbf{v}\|_2^q \cdot \|\beta^*\|_2^q. \end{aligned} \quad (191)$$

The sub-multiplicativity of the operator norm yields

$$\|(\mathbb{E}[\mathbb{X}_{1,p}])^{-1}\mathbb{E}[\overline{\mathbb{X}}_1]\| \leq \frac{\lambda_{\max}(\mathbb{E}[\overline{\mathbb{X}}_1])}{\lambda_{\min}(\mathbb{E}[\mathbb{X}_{1,p}])} < \infty, \quad (192)$$

since $\lambda_{\min}(\mathbb{E}[\mathbb{X}_{1,p}]) \geq (1-p)\lambda_{\min}(\mathbb{E}[\text{Diag}(\mathbb{X}_1)]) = (1-p) \min_i \mathbb{E}[\mathbb{X}_1]_{ii} > 0$, and $\lambda_{\max}(\mathbb{E}[\overline{\mathbb{X}}_1]) \leq \lambda_{\max}(\mathbb{E}[\mathbb{X}_1]) < \infty$. Moreover, $\|\beta^*\|_2 < \infty$. As we assume $\sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \mathbb{E}\|D_1 \mathbb{X}_1 D_1 \mathbf{v}\|_2^q < \infty$ in Lemma 10, also (191) is bounded. \blacksquare

Appendix F. Proofs in Section 4.3

F.1 Proof of Lemma 14

Proof The recursion in (36) is $\check{\beta}_k(\alpha) - \check{\beta} = \check{A}_k(\alpha)(\check{\beta}_{k-1}(\alpha) - \check{\beta}) + \check{\mathbf{b}}_k(\alpha)$, with random matrix $\check{A}_k(\alpha) = I_d - \alpha D_k \mathbb{X}_k D_k$, and random vector $\check{\mathbf{b}}_k(\alpha) = \alpha D_k \mathbf{x}_k (y_k - \mathbf{x}_k^\top D_k \check{\beta})$. Recall $\check{\beta} = \arg \min_{\beta \in \mathbb{R}^d} \mathbb{E}[(y - \mathbf{x}^\top D \beta)^2 / 2]$ in (34), where the expectation is taken over both (y, \mathbf{x}) and D . By Lemma 21 (ii),

$$\begin{aligned} \mathbb{E}_{(y,\mathbf{x})} \mathbb{E}_D [\check{\mathbf{b}}_k(\alpha)] &= \mathbb{E}_{(y,\mathbf{x})} \mathbb{E}_D [\alpha D_k \mathbf{x}_k (y_k - \mathbf{x}_k^\top D_k \check{\beta})] \\ &= \mathbb{E}_{(y,\mathbf{x})} [\alpha p I_d y_k \mathbf{x}_k - \alpha p \mathbb{X}_{k,p} \check{\beta}] \\ &= \alpha p \mathbb{E}[\mathbb{X}_{1,p}] \check{\beta} - \alpha p \mathbb{E}[\mathbb{X}_{1,p}] \check{\beta} \\ &= 0. \end{aligned} \quad (193)$$

Similar to (88), we can rewrite the stationary SGD dropout sequence $\check{\beta}_k^\circ(\alpha)$ into

$$\begin{aligned} \check{\beta}_k^\circ(\alpha) - \check{\beta} &= \sum_{i=0}^{\infty} \left(\prod_{j=k-i+1}^k \check{A}_j(\alpha) \right) \check{\mathbf{b}}_{k-i}(\alpha) \\ &= \alpha \sum_{i=0}^{\infty} \left(\prod_{j=k-i+1}^k (I_d - \alpha D_j \mathbb{X}_j D_j) \right) D_{k-i} \mathbf{x}_{k-i} (y_{k-i} - \mathbf{x}_{k-i}^\top D_{k-i} \check{\beta}) \\ &=: \alpha \sum_{i=0}^{\infty} \check{\mathcal{M}}_{i,k}(\alpha). \end{aligned} \quad (194)$$

Recall the filtration $\check{\mathcal{F}}_i = \sigma(\boldsymbol{\xi}_i, \boldsymbol{\xi}_{i-1}, \dots)$ in (205) for $i \in \mathbb{Z}$, where $\boldsymbol{\xi}_i = (y_i, \mathbf{x}_i, D_i)$. Notice that $\mathbb{E}[\check{\mathbf{b}}_k(\alpha)] = 0$ by (193), and therefore we have

$$\begin{aligned} &\mathbb{E}[\check{\mathcal{M}}_{i,k}(\alpha) \mid \check{\mathcal{F}}_{k-i+1}] \\ &= \mathbb{E} \left[\prod_{j=k-i+1}^k (I_d - \alpha D_j \mathbb{X}_j D_j) \mid \check{\mathcal{F}}_{k-i+1} \right] \mathbb{E} \left[D_{k-i} \mathbf{x}_{k-i} (y_{k-i} - \mathbf{x}_{k-i}^\top D_{k-i} \check{\beta}) \right] \\ &= \mathbb{E} \left[\prod_{j=k-i+1}^k (I_d - \alpha D_j \mathbb{X}_j D_j) \mid \check{\mathcal{F}}_{k-i+1} \right] \cdot \mathbb{E}[\check{\mathbf{b}}_{k-i}(\alpha)] \\ &= 0. \end{aligned} \quad (195)$$

Hence, for any $k \in \mathbb{N}$, $\{\check{\mathcal{M}}_{i,k}(\alpha)\}_{i \in \mathbb{N}}$ is a sequence of martingale differences with respect to the filtration $\check{\mathcal{F}}_{k-i}$. Let $t = k - i$. By applying Burkholder's inequality in Lemma 20, we have,

$$\begin{aligned} (\mathbb{E}\|\check{\beta}_k^o(\alpha) - \check{\beta}\|_2^q)^{1/q} &= \alpha \left(\mathbb{E} \left\| \sum_{t=-\infty}^k \check{\mathcal{M}}_{k-t,k}(\alpha) \right\|_2^q \right)^{1/q} \\ &\lesssim \alpha \left(\sum_{t=-\infty}^k (\mathbb{E}\|\check{\mathcal{M}}_{k-t,k}(\alpha)\|_2^q)^{2/q} \right)^{1/2}, \end{aligned} \quad (196)$$

where the constant in \lesssim only depends on q . Denote the vector $\check{\mathbf{s}}_t = D_t \mathbf{x}_t (y_t - \mathbf{x}_t^\top D_t \check{\beta})$ and the matrix product $\check{A}_{(t+1):k} = \check{A}_{(t+1):k}(\alpha) = \prod_{j=t+1}^k \check{A}_j(\alpha)$ for simplicity. Then, we can write

$$\check{\mathbf{b}}_t(\alpha) = \alpha \check{\mathbf{s}}_t \quad \text{and} \quad \check{\mathcal{M}}_{k-t,k}(\alpha) = \check{A}_{(t+1):k}(\alpha) \check{\mathbf{s}}_t. \quad (197)$$

For the case with $q = 2$, notice that $\check{A}_{(t+1):k}$ is independent of $\check{\mathbf{s}}_t$, and by the tower rule, we have

$$\begin{aligned} \mathbb{E}\|\check{\mathcal{M}}_{k-t,k}(\alpha)\|_2^2 &= \mathbb{E}\|\check{A}_{(t+1):k} \check{\mathbf{s}}_t\|_2^2 \\ &= \mathbb{E}[\mathbb{E}[\check{\mathbf{s}}_t^\top \check{A}_{(t+1):k}^\top \check{A}_{(t+1):k} \check{\mathbf{s}}_t \mid \check{\mathcal{F}}_t]] \\ &= \mathbb{E}[\mathbb{E}[\text{tr}(\check{\mathbf{s}}_t \check{\mathbf{s}}_t^\top \check{A}_{(t+1):k}^\top \check{A}_{(t+1):k}) \mid \check{\mathcal{F}}_t]] \\ &= \mathbb{E}[\text{tr}(\mathbb{E}[\check{\mathbf{s}}_t \check{\mathbf{s}}_t^\top \check{A}_{(t+1):k}^\top \check{A}_{(t+1):k} \mid \check{\mathcal{F}}_t])] \\ &= \mathbb{E}[\text{tr}(\mathbb{E}[\check{\mathbf{s}}_t \check{\mathbf{s}}_t^\top] \check{A}_{(t+1):k}^\top \check{A}_{(t+1):k})] \\ &= \text{tr}(\mathbb{E}[\check{\mathbf{s}}_t \check{\mathbf{s}}_t^\top] \mathbb{E}[\check{A}_{(t+1):k}^\top \check{A}_{(t+1):k}]) \\ &\leq \|\mathbb{E}[\check{A}_{(t+1):k}^\top \check{A}_{(t+1):k}]\| \cdot \mathbb{E}\|\check{\mathbf{s}}_t\|_2^2. \end{aligned} \quad (198)$$

Next, we shall bound the parts $\|\mathbb{E}[\check{A}_{(t+1):k}^\top \check{A}_{(t+1):k}]\|$ and $\mathbb{E}\|\check{\mathbf{s}}_t\|_2^2$ separately. First, recall $\check{A}_j(\alpha) = I_d - \alpha D_j \mathbb{X}_j D_j$ in (36), which are i.i.d. over j . By the tower rule with the induction over $j = t+1, t+2, \dots, k$, we have

$$\begin{aligned} \|\mathbb{E}[\check{A}_{(t+1):k}^\top \check{A}_{(t+1):k}]\| &= \|\mathbb{E}[\mathbb{E}[\check{A}_{(t+1):k}^\top \check{A}_{(t+1):k} \mid \check{A}_{t+1}, \check{A}_{t+2}, \dots, \check{A}_{k-1}]]\| \\ &\leq \|\mathbb{E}[\check{A}_k^\top(\alpha) \check{A}_k(\alpha)]\| \cdot \|\mathbb{E}[\check{A}_{(t+1):(k-1)}^\top \check{A}_{(t+1):(k-1)}]\| \\ &\leq \prod_{j=t+1}^k \|\mathbb{E}[\check{A}_j^\top(\alpha) \check{A}_j(\alpha)]\| \\ &= \|\mathbb{E}[\check{A}_1^\top(\alpha) \check{A}_1(\alpha)]\|^{k-t}. \end{aligned} \quad (199)$$

Moreover, recall the random matrix $M_k(\alpha) = 2D_k \mathbb{X}_k D_k - \alpha D_k \mathbb{X}_k D_k \mathbb{X}_k D_k$ as defined in (169). For any unit vector $\mathbf{v} \in \mathbb{R}^d$, by (170) in the proof of Lemma 10, we have

$$\begin{aligned} \mathbf{v}^\top \mathbb{E}[\check{A}_1^\top(\alpha) \check{A}_1(\alpha)] \mathbf{v} &\leq 1 - \alpha \mathbf{v}^\top \mathbb{E}[M_1(\alpha)] \mathbf{v} \\ &\leq 1 - \alpha \lambda_{\min}(\mathbb{E}[M_1(\alpha)]) < 1, \end{aligned} \quad (200)$$

as the constant learning rate α satisfies condition (37). In fact, condition (37) also implies that $\mathbb{E}[M_k(\alpha)]$ is positive definite for each $k \in \mathbb{N}$, which can be seen by (40).

We shall show that the term $\mathbb{E}\|\check{\mathbf{s}}_k\|_2^2 = \mathbb{E}\|D_k \mathbf{x}_k (y_k - \mathbf{x}_k^\top D_k \check{\boldsymbol{\beta}})\|_2^2$ remains bounded as $k \rightarrow \infty$. In Assumption 1, we have assumed that the stochastic gradient in the SGD dropout recursion $\nabla_{\boldsymbol{\beta}}(y - \mathbf{x}^\top D \boldsymbol{\beta})^2/2 = D \mathbf{x} (y - \mathbf{x}^\top D \boldsymbol{\beta})$, has finite q -th moment when $\boldsymbol{\beta} = \boldsymbol{\beta}^*$ for some $q \geq 2$. By Lemma 30, Assumption 1 also implies the bounded q -th moment when $\boldsymbol{\beta} = \check{\boldsymbol{\beta}}$. As a direct consequence, $\mathbb{E}\|\check{\mathbf{s}}_k\|_2^2$ is bounded as $k \rightarrow \infty$. This, along with (196), (198) and (200), yields

$$\begin{aligned}
 (\mathbb{E}\|\check{\boldsymbol{\beta}}_k^\circ(\alpha) - \check{\boldsymbol{\beta}}\|_2^2)^{1/2} &\lesssim \alpha \left(\sum_{t=-\infty}^k \mathbb{E}\|\check{\mathcal{M}}_{k-t,k}(\alpha)\|_2^2 \right)^{1/2} \\
 &\leq \alpha \left(\sum_{t=-\infty}^k \|\mathbb{E}[\check{A}_{(t+1):k}^\top \check{A}_{(t+1):k}]\| \cdot \mathbb{E}\|\check{\mathbf{s}}_t\|_2^2 \right)^{1/2} \\
 &\lesssim \alpha \left(\sum_{t=-\infty}^k (1 - \alpha \lambda_*)^{k-t} \right)^{1/2} \\
 &= \alpha \left(\sum_{i=0}^{\infty} (1 - \alpha \lambda_*)^i \right)^{1/2} \\
 &= O(\sqrt{\alpha}), \tag{201}
 \end{aligned}$$

where the last equation holds since $\sum_{i=0}^{\infty} (1 - \alpha \lambda_*)^i = \frac{1}{1 - (1 - \alpha \lambda_*)} = O(1/\alpha)$. Here, the constants in \lesssim are independent of k and α , and $\lambda_* = \lambda_{\min}(\mathbb{E}[M_1(\alpha)])$ is bounded away from zero since $\mathbb{E}[M_1(\alpha)] = \mathbb{E}[2D_1 \mathbb{X}_1 D_1 - \alpha D_1 \mathbb{X}_1 D_1 \mathbb{X}_1 D_1]$ is positive definite by condition (40).

For the case $q > 2$, following similar arguments as in (198), we obtain

$$\begin{aligned}
 \mathbb{E}\|\check{\mathcal{M}}_{k-t,k}(\alpha)\|_2^q &= \mathbb{E}(\|\check{A}_{(t+1):k} \check{\mathbf{s}}_t\|_2^2)^{q/2} \\
 &= \mathbb{E}(\check{\mathbf{s}}_t^\top \check{A}_{(t+1):k}^\top \check{A}_{(t+1):k} \check{\mathbf{s}}_t)^{q/2} \\
 &\leq \sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \mathbb{E}\|\check{A}_{(t+1):k} \mathbf{v}\|_2^q \cdot \mathbb{E}\|\check{\mathbf{s}}_t\|_2^q \\
 &\leq \left(\sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \mathbb{E}\|\check{A}_1 \mathbf{v}\|_2^q \right)^{k-t} \cdot \mathbb{E}\|\check{\mathbf{s}}_t\|_2^q. \tag{202}
 \end{aligned}$$

With $\mu_q(\mathbf{v}) = (\mathbb{E}\|D_1 \mathbb{X}_1 D_1 \mathbf{v}\|_2^q)^{1/q} < \infty$ as defined in Lemma 10 and the Equations (173) and (175) in the proof of Lemma 10, we have

$$\mathbb{E}\|\check{A}_1 \mathbf{v}\|_2^q \leq (1 + \alpha \mu_q(\mathbf{v}))^q - q \alpha \mu_q(\mathbf{v}) - q \alpha p \mathbf{v}^\top \mathbb{E}[\mathbb{X}_{1,p}] \mathbf{v}. \tag{203}$$

This, together with Taylor expansion around $\alpha = 0$ and the inequalities (196) and (202) gives finally $\max_k (\mathbb{E}\|\check{\boldsymbol{\beta}}_k^\circ(\alpha) - \check{\boldsymbol{\beta}}\|_2^q)^{1/q} = O(\sqrt{\alpha})$. \blacksquare

F.2 Proof of Theorem 15

The proofs of the quenched CLT and the invariance principle for averaged SGD dropout follow similar arguments as the ones for averaged GD dropout. The key differences lie in the functional dependence measures, because for SGD settings, both the dropout matrix D_k and the sequential observation \mathbf{x}_k are random. We shall first introduce some necessary definitions and then proceed with the rigorous proofs.

Recall the generic dropout matrix $D \in \mathbb{R}^{d \times d}$ and random sample $(y, \mathbf{x}) \in \mathbb{R} \times \mathbb{R}^d$. For the SGD dropout sequence $\{\check{\beta}_k(\alpha)\}_{k \in \mathbb{N}}$, we define the centering term

$$\check{\beta}_{\text{mean}}(\alpha) = \lim_{k \rightarrow \infty} \mathbb{E}[\check{\beta}_k(\alpha)] = \mathbb{E}[\check{\beta}_1^\circ(\alpha)], \quad (204)$$

where the expectation is taken over both (y, \mathbf{x}) and D , and $\check{\beta}_1^\circ(\alpha)$ defined in (44) follows the unique stationary distribution $\check{\pi}_\alpha$. We first use Lemma 25 to prove the CLT for the partial sum of the stationary sequence $\{\check{\beta}_k^\circ(\alpha) - \check{\beta}_{\text{mean}}(\alpha)\}_{k \in \mathbb{N}}$, and then apply the geometric-moment contraction in Theorem 12 to show the quenched CLT for the partial sum of the non-stationary one $\{\check{\beta}_k(\alpha) - \check{\beta}_{\text{mean}}(\alpha)\}_{k \in \mathbb{N}}$. Finally, we extend the quenched CLT to the partial sum of $\{\check{\beta}_k(\alpha) - \check{\beta}_{\text{mean}}(\alpha)\}_{k \in \mathbb{N}}$ by providing the upper bound of $\check{\beta}_{\text{mean}}(\alpha) - \check{\beta}$ in terms of the q -th moment for some $q \geq 2$.

Similar to Section D.1, we introduce the *functional dependence measure* in Wu (2005) for the stationary SGD dropout sequence $\{\check{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$. However, the randomness in $\check{\beta}_k^\circ(\alpha)$ is induced from both the dropout matrix D_k and the random sample (y_k, \mathbf{x}_k) . Therefore, we define a new filtration $\check{\mathcal{F}}_k$ by

$$\check{\mathcal{F}}_k = \sigma(\boldsymbol{\xi}_k, \boldsymbol{\xi}_{k-1}, \dots), \quad k \in \mathbb{Z}, \quad (205)$$

where the i.i.d. random elements $\boldsymbol{\xi}_k = (D_k, (y_k, \mathbf{x}_k))$, $k \in \mathbb{Z}$, are defined in (45). For any random vector $\zeta \in \mathbb{R}^d$ satisfying $\mathbb{E}\|\zeta\|_2 < \infty$, define projection operators

$$\check{\mathcal{P}}_k[\zeta] = \mathbb{E}[\zeta \mid \check{\mathcal{F}}_k] - \mathbb{E}[\zeta \mid \check{\mathcal{F}}_{k-1}], \quad k \in \mathbb{Z}. \quad (206)$$

By Theorem 12 and (44), there exists a measurable function $\check{h}_\alpha(\cdot)$ such that the stationary SGD dropout sequence $\{\check{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ can be represented by a causal process

$$\check{\beta}_k^\circ(\alpha) = \check{h}_\alpha(\boldsymbol{\xi}_k, \boldsymbol{\xi}_{k-1}, \dots) = \check{h}_\alpha(\check{\mathcal{F}}_k). \quad (207)$$

We denote the coupled version of $\check{\mathcal{F}}_i$ by

$$\check{\mathcal{F}}_{i, \{j\}} = \sigma(\boldsymbol{\xi}_i, \dots, \boldsymbol{\xi}_{j+1}, \boldsymbol{\xi}'_j, \boldsymbol{\xi}_{j-1}, \dots), \quad (208)$$

and let $\check{\mathcal{F}}_{i, \{j\}} = \check{\mathcal{F}}_i$ if $j > i$, where $\boldsymbol{\xi}'_j$ is an i.i.d. copy of $\boldsymbol{\xi}_i$. For $q > 1$, define the *functional dependence measure* of $\check{\beta}_k(\alpha)$ as

$$\check{\theta}_{k,q}(\alpha) = (\mathbb{E}\|\check{\beta}_k^\circ(\alpha) - \check{\beta}_{k,\{0\}}^\circ(\alpha)\|_2^q)^{1/q}, \quad \text{where } \check{\beta}_{k,\{0\}}^\circ(\alpha) = \check{h}_\alpha(\check{\mathcal{F}}_{k,\{0\}}). \quad (209)$$

In addition, if $\sum_{k=0}^{\infty} \check{\theta}_{k,q}(\alpha) < \infty$, we define the tail of *cumulative dependence measure* as

$$\check{\Theta}_{m,q}(\alpha) = \sum_{k=m}^{\infty} \check{\theta}_{k,q}(\alpha), \quad m \in \mathbb{N}. \quad (210)$$

Both $\check{\theta}_{k,q}(\alpha)$ and $\check{\Theta}_{k,q}(\alpha)$ are useful to study the dependence structure of the stationary SGD dropout iteration $\check{\beta}_k^\circ(\alpha) = f_{D_k, (y_k, \mathbf{x}_k)}(\check{\beta}_{k-1}^\circ(\alpha))$. To apply Lemma 25, we only need to show that the stationary SGD dropout sequence $\{\check{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ is short-range dependent in the sense that $\check{\Theta}_{0,q}(\alpha) < \infty$, for some $q \geq 2$.

Proof [Theorem 15] Consider two initial vectors $\check{\beta}_0^\circ, \check{\beta}_0^{\circ'}$ following the stationary distribution $\tilde{\pi}_\alpha$ defined in Theorem 12. By the recursion in (36), we obtain two stationary SGD dropout sequences $\{\check{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ and $\{\check{\beta}_k^{\circ'}(\alpha)\}_{k \in \mathbb{N}}$. It follows from Theorem 12 that for all $q \geq 2$,

$$\sup_{\check{\beta}_0^\circ, \check{\beta}_0^{\circ'} \in \mathbb{R}^d, \check{\beta}_0^\circ \neq \check{\beta}_0^{\circ'}} \frac{(\mathbb{E} \|\check{\beta}_k^\circ(\alpha) - \check{\beta}_k^{\circ'}(\alpha)\|_2^q)^{1/q}}{\|\check{\beta}_0^\circ - \check{\beta}_0^{\circ'}\|_2} \leq \check{r}_{\alpha,q}^k, \quad k \in \mathbb{N}, \quad (211)$$

with $\check{r}_{\alpha,q} = (\sup_{\mathbf{v} \in \mathbb{R}^d: \|\mathbf{v}\|_2=1} \mathbb{E} \|\check{A}_1(\alpha)\mathbf{v}\|_2^q)^{1/q}$ as defined in (38) and random matrix $\check{A}_1(\alpha) = I_d - \alpha D_1 \mathbb{X}_1 D_1$. When the constant learning rate $\alpha > 0$ satisfies the condition in (37), we have $\check{r}_{\alpha,q} \in (0, 1)$ as shown in Theorem 12. Recall the coupled filtration $\check{\mathcal{F}}_{i, \{j\}}$ = $\sigma(\xi_i, \dots, \xi_{j+1}, \xi_j', \xi_{j-1}, \dots)$ as defined in (208). Then, (44) and (211) show that

$$\begin{aligned} & (\mathbb{E} \|\check{f}_{\xi_k} \circ \dots \circ \check{f}_{\xi_1}(\check{\beta}_0^\circ) - \check{f}_{\xi_k} \circ \dots \circ \check{f}_{\xi_1}(\check{\beta}_0^{\circ'})\|_2^q)^{1/q} \\ &= (\mathbb{E} \|\check{h}_\alpha(\xi_k, \dots, \xi_1, \xi_0, \xi_{-1}, \dots) - \check{h}_\alpha(\xi_k, \dots, \xi_1, \xi_0', \xi_{-1}', \dots)\|_2^q)^{1/q} \\ &= (\mathbb{E} \|\check{h}_\alpha(\check{\mathcal{F}}_k) - \check{h}_\alpha(\check{\mathcal{F}}_{k, \{0, -1, \dots\}})\|_2^q)^{1/q} \\ &\leq \check{c}_q \check{r}_{\alpha,q}^k, \end{aligned} \quad (212)$$

for some constant $\check{c}_q > 0$ that is independent of k . Following a similar argument in (127), for all $q \geq 2$ and $k \in \mathbb{N}$, we can bound the functional dependence measure $\check{\theta}_{k,q}(\alpha)$ in (209) as follows

$$\begin{aligned} \check{\theta}_{k,q}(\alpha) &= (\mathbb{E} \|\check{h}_\alpha(\check{\mathcal{F}}_k) - \check{h}_\alpha(\check{\mathcal{F}}_{k, \{0\}})\|_2^q)^{1/q} \\ &\leq (\mathbb{E} \|\check{h}_\alpha(\check{\mathcal{F}}_k) - \check{h}_\alpha(\check{\mathcal{F}}_{k, \{0, -1, \dots\}})\|_2^q)^{1/q} + (\mathbb{E} \|\check{h}_\alpha(\check{\mathcal{F}}_{k, \{0, -1, \dots\}}) - \check{h}_\alpha(\check{\mathcal{F}}_{k, \{0\}})\|_2^q)^{1/q} \\ &\leq \check{c}'_q \check{r}_{\alpha,q}^k, \end{aligned} \quad (213)$$

for some constant $\check{c}'_q > 0$ that is independent of k . Consequently, the cumulative dependence measure $\check{\Theta}_{m,q}(\alpha)$ in (210) is also bounded for all $q \geq 2$ and $m \in \mathbb{N}$, that is,

$$\check{\Theta}_{m,q}(\alpha) = \sum_{k=m}^{\infty} \check{\theta}_{k,q}(\alpha) = O(\check{r}_{\alpha,q}^m) < \infty. \quad (214)$$

The inequality in (119) derived by Wu (2005) holds for a general class of functional dependence measures, as long as the inputs of the functional system (i.e., the measurable function $\check{h}(\xi_k, \xi_{k-1}, \dots)$) are i.i.d. elements. Thus, we can apply (119) to the projection operator $\check{\mathcal{P}}_k[\cdot]$ in (206) and obtain

$$\sum_{k=0}^{\infty} (\mathbb{E} \|\check{\mathcal{P}}_0[\check{\beta}_k^\circ(\alpha)]\|_2^q)^{1/q} \leq \sum_{k=0}^{\infty} \check{\theta}_{k,q}(\alpha) = \check{\Theta}_{0,q}(\alpha) < \infty, \quad (215)$$

implying the short-range dependence of the stationary SGD dropout sequence $\{\check{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$. Then, it follows from Lemma 25 that

$$t^{-1/2} \sum_{k=1}^t (\check{\beta}_k^\circ(\alpha) - \check{\beta}_{\text{mean}}(\alpha)) \Rightarrow \mathcal{N}(0, \check{\Sigma}(\alpha)), \quad (216)$$

where the long-run covariance matrix $\check{\Sigma}(\alpha)$ is defined in Theorem 15. Following similar arguments as in (131)–(133), for any initial vector $\check{\beta}_0 \in \mathbb{R}^d$, we can leverage the geometric-moment contraction in Theorem 12 and achieve the quenched CLT for the corresponding SGD dropout sequence $\{\check{\beta}_k(\alpha)\}_{k \in \mathbb{N}}$, that is,

$$t^{-1/2} \sum_{k=1}^t (\check{\beta}_k(\alpha) - \check{\beta}_{\text{mean}}(\alpha)) \Rightarrow \mathcal{N}(0, \check{\Sigma}(\alpha)). \quad (217)$$

Recall the ℓ^2 -minimizer $\check{\beta}$ in (34) and the centering term $\check{\beta}_{\text{mean}}(\alpha) = \lim_{k \rightarrow \infty} \mathbb{E}[\check{\beta}_k(\alpha)] = \mathbb{E}[\check{\beta}_1^\circ(\alpha)]$ in (204). We shall prove $\|\sum_{k=1}^t \mathbb{E}[\check{\beta}_k(\alpha) - \check{\beta}]\|_2 = o(\sqrt{t})$. For any two initial vectors $\check{\beta}_0$ and $\check{\beta}_0^\circ$, where $\check{\beta}_0^\circ$ follows the stationary distribution $\check{\pi}_\alpha$ in Theorem 12, while $\check{\beta}_0$ is an arbitrary initial vector in \mathbb{R}^d , it follows from the triangle inequality that

$$\begin{aligned} \left\| \sum_{k=1}^t \mathbb{E}[\check{\beta}_k(\alpha) - \check{\beta}] \right\|_2 &= \left\| \mathbb{E} \left[\sum_{k=1}^t (\check{\beta}_k(\alpha) - \check{\beta}_k^\circ(\alpha) + \check{\beta}_k^\circ(\alpha) - \check{\beta}) \right] \right\|_2 \\ &\leq \left\| \mathbb{E} \left[\sum_{k=1}^t (\check{\beta}_k(\alpha) - \check{\beta}_k^\circ(\alpha)) \right] \right\|_2 + \left\| \mathbb{E} \left[\sum_{k=1}^t (\check{\beta}_k^\circ(\alpha) - \check{\beta}) \right] \right\|_2 \\ &=: \check{\mathbb{I}}_1 + \check{\mathbb{I}}_2. \end{aligned} \quad (218)$$

For the part $\check{\mathbb{I}}_1$, Jensen's inequality and a similar argument as for (131) yield

$$\begin{aligned} \check{\mathbb{I}}_1 &= \left\| \mathbb{E} \left[\sum_{k=1}^t (\check{\beta}_k(\alpha) - \check{\beta}_k^\circ(\alpha)) \right] \right\|_2 \\ &\leq \left(\mathbb{E} \left\| \sum_{k=1}^t (\check{\beta}_k(\alpha) - \check{\beta}_k^\circ(\alpha)) \right\|_2^2 \right)^{1/2} \\ &\leq \left(\sum_{k=1}^t \check{r}_{\alpha,2}^k \right) \|\check{\beta}_0 - \check{\beta}_0^\circ\|_2. \end{aligned} \quad (219)$$

For the part $\check{\mathbb{I}}_2$, we recall the random matrix $\check{A}_1(\alpha) = I_d - \alpha D_1 \mathbb{X}_1 D_1$ in (36) with $\mathbb{X}_1 = \mathbf{x}_1 \mathbf{x}_1^\top$. Recall the notation $\mathbb{X}_{1,p} = p \mathbb{X}_1 + (1-p) \text{Diag}(\mathbb{X}_1)$ in (35). Notice that by Lemma 21 (ii), we have $\mathbb{E}_{(y,\mathbf{x})} \mathbb{E}_D[\check{A}_1(\alpha)] = \mathbb{E}_{(y,\mathbf{x})} \mathbb{E}_D[I_d - \alpha D_1 \mathbb{X}_1 D_1] = \mathbb{E}_{(y,\mathbf{x})} [I_d - \alpha p \mathbb{X}_{1,p}] = I_d - \alpha p \mathbb{E}[\mathbb{X}_{1,p}]$, which along with $\mathbb{E}[\check{\mathbf{b}}_k] = 0$ in (193) gives $\mathbb{E}[\check{\beta}_k^\circ(\alpha) - \check{\beta}] = (I_d - \alpha p \mathbb{E}[\mathbb{X}_{1,p}]) \mathbb{E}[\check{\beta}_{k-1}^\circ(\alpha) - \check{\beta}]$. Since $\{\check{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ is stationary and $\mathbb{E}[\mathbb{X}_{1,p}]$ is non-singular by the reduced-form condition $\min_i (\mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top])_{ii} > 0$ imposed on model (31), it follows that

$$\mathbb{E}[\check{\beta}_k^\circ(\alpha) - \check{\beta}] = 0, \quad \text{for all } k \in \mathbb{N}. \quad (220)$$

This further yields

$$\check{\mathbb{I}}_2 = \left\| \mathbb{E} \left[\sum_{k=1}^t (\check{\beta}_k^\circ(\alpha) - \check{\beta}) \right] \right\|_2 = \left\| \sum_{k=1}^t \mathbb{E}[\check{\beta}_k^\circ(\alpha) - \check{\beta}] \right\|_2 = 0. \quad (221)$$

By applying the results for $\check{\mathbb{I}}_1$ and $\check{\mathbb{I}}_2$ to (218), we obtain

$$\left\| \sum_{k=1}^t \mathbb{E}[\check{\beta}_k(\alpha) - \check{\beta}] \right\|_2 \leq \left(\sum_{k=1}^t \check{r}_{\alpha,2}^k \right) \|\check{\beta}_0 - \check{\beta}_0^\circ\|_2. \quad (222)$$

When the constant learning rate $\alpha > 0$ satisfies the condition in (37), $\check{r}_{\alpha,2} \in (0, 1)$ by Theorem 12, and hence (222) remains bounded as $t \rightarrow \infty$. By this result and (217), the desired quenched CLT for the partial sum $\sum_{k=1}^t (\check{\beta}_k^\circ(\alpha) - \check{\beta})$ follows. \blacksquare

Proof The proof of Corollary 16 applies the Cramér-Wold device to Theorem 15 and can be derived in the same way as the proof of Corollary 8. We omit the details here. \blacksquare

F.3 Proof of Theorem 17

Proof Consider a stationary SGD dropout sequence $\{\check{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ following the stationary distribution $\check{\pi}_\alpha$ in Theorem 12. Define the mean-zero stationary partial sum

$$\check{S}_i^\circ(\alpha) = \sum_{k=1}^i (\check{\beta}_k^\circ(\alpha) - \check{\beta}_{\text{mean}}), \quad i \in \mathbb{N}. \quad (223)$$

Recall that $\check{\beta}_{\text{mean}}(\alpha) = \mathbb{E}[\check{\beta}_1^\circ(\alpha)]$ as defined in (204). Due to the stationarity, we have $\mathbb{E}[\check{\beta}_k^\circ(\alpha) - \check{\beta}] = 0$ for all $k \in \mathbb{N}$ by (220). This gives

$$(\mathbb{E}\|\check{\beta}_{\text{mean}}(\alpha) - \check{\beta}\|_2^q)^{1/q} = (\mathbb{E}\|\mathbb{E}[\check{\beta}_1^\circ(\alpha) - \check{\beta}]\|_2^q)^{1/q} = 0. \quad (224)$$

Since we supposed in Theorem 17 that Assumption 1 holds for some $q > 2$, it follows from Lemma 14 that, for all $q > 2$,

$$(\mathbb{E}\|\check{\beta}_k^\circ(\alpha) - \check{\beta}_{\text{mean}}(\alpha)\|_2^q)^{1/q} \leq (\mathbb{E}\|\check{\beta}_k^\circ(\alpha) - \check{\beta}\|_2^q)^{1/q} + (\mathbb{E}\|\check{\beta}_{\text{mean}}(\alpha) - \check{\beta}\|_2^q)^{1/q} = O(\sqrt{\alpha}). \quad (225)$$

Following a similar argument as for (144), we can show that $\check{\beta}_k^\circ(\alpha) - \check{\beta}_{\text{mean}}(\alpha)$ satisfies condition (i) on the uniform integrability in Lemma 26. Due to the stationarity of $\{\check{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$, condition (ii) in Lemma 26 is not required (see the discussion below Lemma 26 for details). Condition (iii) is also satisfied, since when the constant learning rate $\alpha > 0$ satisfies condition (37), the stationary SGD dropout sequence $\{\check{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ is shown to be short-range dependent by (215), i.e., the tail of cumulative dependence measure $\check{\Theta}_{m,q}(\alpha) = \sum_{k=m}^{\infty} \check{\theta}_{k,q}(\alpha) < \infty$.

Hence, by Lemma 26, there exists a (richer) probability space $(\check{\Omega}^*, \check{\mathcal{A}}^*, \check{\mathbb{P}}^*)$ on which we can define random vectors $\check{\beta}_k^*$'s with the partial sum process $\check{S}_i^* = \sum_{k=1}^i (\check{\beta}_k^* - \check{\beta}_{\text{mean}})$, and a Gaussian process $\check{G}_i^* = \sum_{k=1}^i \check{z}_k^*$, where \check{z}_k^* 's are independent Gaussian random vectors in \mathbb{R}^d following $\mathcal{N}(0, I_d)$, such that

$$(\check{S}_i^*)_{1 \leq i \leq t} \stackrel{\mathcal{D}}{=} (\check{S}_i^\circ)_{1 \leq i \leq t}, \quad (226)$$

and

$$\max_{1 \leq i \leq t} \|\check{S}_i^* - \check{\Sigma}^{1/2}(\alpha) \check{G}_i^*\|_2 = o_{\mathbb{P}}(t^{1/q}), \quad \text{in } (\check{\Omega}^*, \check{\mathcal{A}}^*, \check{\mathbb{P}}^*), \quad (227)$$

where the long-run covariance matrix $\check{\Sigma}(\alpha)$ is defined in Theorem 15.

Following similar arguments as for (148)–(150), we can leverage the geometric-moment contraction in Theorem 12 to show the same Gaussian approximation rate, i.e., $o_{\mathbb{P}}(t^{1/q})$, for the partial sum sequence $(\sum_{k=1}^i (\check{\beta}_k(\alpha) - \check{\beta}_{\text{mean}}(\alpha)))_{1 \leq i \leq t}$, for any arbitrarily fixed initial vector $\check{\beta}_0 \in \mathbb{R}^d$. Finally, recall the partial sum process $\check{S}_i^{\check{\beta}_0}(\alpha) = \sum_{k=1}^i (\check{\beta}_k(\alpha) - \check{\beta})$. By a similar argument in (151), the desired Gaussian approximation result for the partial sum process $(\check{S}_i^{\check{\beta}_0}(\alpha))_{1 \leq i \leq t}$ follows. \blacksquare

Appendix G. Proofs in Section 5

G.1 Proof of Theorem 18

Lemma 31 *For any $d \times d$ symmetric matrix $S = (S_{ij})$, we have $\mathbb{E}\|S\|_F \leq \sqrt{\text{tr}\mathbb{E}(S^2)} \leq d \max_{i,j} (\mathbb{E}[S_{ij}^2])^{1/2}$.*

Proof For a symmetric matrix S , $\text{tr}(S^2) = \|S\|_F^2 = \sum_{i,j} S_{ij}^2$. Since $\sqrt{\cdot}$ is a concave function, by Jensen's inequality,

$$\mathbb{E}\|S\|_F = \mathbb{E}\sqrt{\text{tr}(S^2)} \leq \sqrt{\text{tr}[\mathbb{E}(S^2)]} = \sqrt{\sum_{i,j=1}^d \mathbb{E}[S_{ij}^2]} \leq d \max_{i,j} (\mathbb{E}[S_{ij}^2])^{1/2}.$$

This completes the proof. \blacksquare

Proof Recall the SGD dropout sequence $\{\check{\beta}_k(\alpha)\}_{k \in \mathbb{N}}$ in (36), the long-run covariance matrix $\check{\Sigma}(\alpha)$ of the averaged SGD dropout iterates in Theorem 15, and the online estimator $\hat{\Sigma}_k(\alpha)$ in (55). When there is no ambiguity, we omit the dependence on α , e.g., $\hat{\Sigma}_k = \hat{\Sigma}_k(\alpha)$ and $\check{\Sigma} = \check{\Sigma}(\alpha)$. We shall bound $\mathbb{E}\|\hat{\Sigma}_k - \check{\Sigma}\|_F$.

Recall the stationary SGD dropout sequence $\{\check{\beta}_k^\circ(\alpha)\}_{k \in \mathbb{N}}$ in (42), which follows the stationary distribution $\check{\pi}_\alpha$ in Theorem 12. For simplicity, we define $V_k(\alpha) = k\hat{\Sigma}_k(\alpha)$. By

Equation (58), we can write $V_k(\alpha)$ into

$$\begin{aligned}
 V_k(\alpha) &= \left(\sum_{m=1}^{\psi(k)-1} \mathcal{S}_m(\alpha)^{\otimes 2} + \mathcal{R}_k(\alpha)^{\otimes 2} \right) + \left(\sum_{m=1}^{\psi(k)-1} |B_m|^2 + |\delta_\eta(k)|^2 \right) \bar{\beta}_k^{\text{sgd}}(\alpha)^{\otimes 2} \\
 &\quad - \left(\sum_{m=1}^{\psi(k)-1} |B_m| \mathcal{S}_m(\alpha) + \delta_\eta(k) \mathcal{R}_k(\alpha) \right) \bar{\beta}_k^{\text{sgd}}(\alpha)^\top \\
 &\quad - \bar{\beta}_k^{\text{sgd}}(\alpha) \left(\sum_{m=1}^{\psi(k)-1} |B_m| \mathcal{S}_m(\alpha) + \delta_\eta(k) \mathcal{R}_k(\alpha) \right)^\top. \tag{228}
 \end{aligned}$$

Recall the partial sums $\mathcal{S}_m(\alpha) = \sum_{k \in B_m} \check{\beta}_k(\alpha)$ and $\mathcal{R}_k(\alpha) = \sum_{i=\eta_{\psi(k)}}^k \check{\beta}_i(\alpha)$ in (57). We similarly define

$$V_k^\circ(\alpha) = \sum_{m=1}^{\psi(k)-1} \left(\sum_{i \in B_m} \check{\beta}_i^\circ(\alpha) \right)^{\otimes 2} + \left(\sum_{i=\eta_{\psi(k)}}^k \check{\beta}_i^\circ(\alpha) \right)^{\otimes 2} =: \sum_{m=1}^{\psi(k)-1} \mathcal{S}_m^\circ(\alpha)^{\otimes 2} + \mathcal{R}_k^\circ(\alpha)^{\otimes 2}. \tag{229}$$

By the triangle inequality, we have

$$k\mathbb{E}\|\hat{\Sigma}_k - \check{\Sigma}\|_F = \mathbb{E}\|V_k - k\check{\Sigma}\|_F \leq \mathbb{E}\|V_k - V_k^\circ\|_F + \mathbb{E}\|V_k^\circ - k\check{\Sigma}\|_F. \tag{230}$$

We shall bound these two terms separately.

First, for the term $\mathbb{E}\|V_k^\circ - k\check{\Sigma}\|_F$, we shall use the results in Xiao and Wu (2011). To this end, we need to verify the assumptions on the weak dependence of SGD dropout iterates $\{\check{\beta}_k\}_{k \in \mathbb{N}}$ in (36) and the growing sizes of blocks $\{B_m\}_{m \in \mathbb{N}}$ in (54). Denote the elements of $d \times d$ matrices V_k° and $\check{\Sigma}$ respectively by

$$V_k^\circ =: (v_{ij,k}^\circ)_{1 \leq i, j \leq d} \quad \text{and} \quad \check{\Sigma} =: (\check{\sigma}_{ij})_{1 \leq i, j \leq d}. \tag{231}$$

Moreover, we write the stationary SGD dropout iterate $\check{\beta}_k^\circ$ in (42) and the ℓ^2 -minimizer $\check{\beta}$ in (34) respectively into

$$\check{\beta}_k^\circ = (\check{\beta}_{k1}^\circ, \dots, \check{\beta}_{kd}^\circ)^\top \quad \text{and} \quad \check{\beta} = (\check{\beta}_1, \dots, \check{\beta}_d)^\top. \tag{232}$$

By the short-range dependence of $\{\check{\beta}_k^\circ\}_{k \in \mathbb{N}}$ in (215), it can be shown that for any $1 \leq i, j \leq d$, $\text{Cov}(\check{\beta}_{ki}^\circ, \check{\beta}_{0j}^\circ) \leq c_{ij} \rho_{ij}^k$ for some constants $c_{ij} > 0$ and $0 \leq \rho_{ij} < 1$. For the non-overlapping blocks $\{B_m\}_{m \in \mathbb{N}}$, since the positive integers $\{\eta_m\}_{m \in \mathbb{N}}$ satisfy $\eta_{m+1} - \eta_m \rightarrow \infty$, it follows that $\eta_{m+1}/\eta_m \rightarrow 1$ as $m \rightarrow \infty$, and therefore,

$$\sum_{m=1}^M (\eta_{m+1} - \eta_m)^2 \asymp \eta_{M+1} (\eta_{M+1} - \eta_M). \tag{233}$$

Then, by Theorems 1(i) and 2(iii) in Xiao and Wu (2011), we obtain, for each $1 \leq j \leq d$,

$$\begin{aligned}
 \mathbb{E}(v_{jj,k}^\circ - k\check{\sigma}_{jj})^2 &\leq \left\{ [\mathbb{E}(v_{jj,k}^\circ - \mathbb{E}v_{jj,k}^\circ)^2]^{1/2} + [(\mathbb{E}v_{jj,k}^\circ - k\check{\sigma}_{jj})^2]^{1/2} \right\}^2 \\
 &\lesssim k^{(2/\zeta) \vee (2-1/\zeta)}. \tag{234}
 \end{aligned}$$

For $\mathbb{E}(v_{ij,k}^\circ - k\check{\sigma}_{ij})^2$ with $i \neq j$, the same rate as in (234) holds. To see this, define two new sequences $\{\beta_k^+\}_{k \in \mathbb{N}}$ and $\{\beta_k^-\}_{k \in \mathbb{N}}$ with $\beta_k^+ = (\check{\beta}_{ki}^\circ + \check{\beta}_{0j}^\circ) - (\check{\beta}_i + \check{\beta}_j)$, and $\beta_k^- = (\check{\beta}_{ki}^\circ - \check{\beta}_{0j}^\circ) - (\check{\beta}_i - \check{\beta}_j)$. Notice that

$$\begin{aligned} \check{\sigma}_{ij} &= \sum_{k=-\infty}^{\infty} \mathbb{E}(\check{\beta}_{ki}^\circ - \check{\beta}_i)(\check{\beta}_{0j}^\circ - \check{\beta}_j) \\ &= \sum_{k=-\infty}^{\infty} \mathbb{E} \frac{[(\check{\beta}_{ki}^\circ - \check{\beta}_i) + (\check{\beta}_{0j}^\circ - \check{\beta}_j)]^2 - [(\check{\beta}_{ki}^\circ - \check{\beta}_i) - (\check{\beta}_{0j}^\circ - \check{\beta}_j)]^2}{4} \\ &= \frac{1}{4} \sum_{k=-\infty}^{\infty} \mathbb{E}(\beta_k^+)^2 - \frac{1}{4} \sum_{k=-\infty}^{\infty} \mathbb{E}(\beta_k^-)^2, \end{aligned} \quad (235)$$

which can be viewed as the long-run variances of the sequences $\{\beta_k^+\}_{k \in \mathbb{N}}$ and $\{\beta_k^-\}_{k \in \mathbb{N}}$ as indicated by the last line. A similar decomposition can be applied to $v_{ij,k}^\circ$. Since the results in Xiao and Wu (2011) hold for any linear combination of weak-dependent sequences, again by the short-range dependence of $\{\check{\beta}_k^\circ\}_{k \in \mathbb{N}}$ in (215), we have

$$\mathbb{E}(v_{ij,k}^\circ - k\check{\sigma}_{ij})^2 \lesssim k^{(2/\zeta) \vee (2-1/\zeta)}, \quad \text{for all } 1 \leq i, j \leq d. \quad (236)$$

For dimension $d \geq 1$, it follows from Lemma 31 that

$$\mathbb{E}\|V_k^\circ - k\check{\Sigma}\|_F \leq d \max_{1 \leq i, j \leq d} \sqrt{\mathbb{E}(v_{ij,k}^\circ - k\check{\sigma}_{ij})^2} \lesssim dk^{(1/\zeta) \vee (1-1/(2\zeta))}, \quad (237)$$

where the constants in \lesssim are independent of k and d .

Next, we bound $\mathbb{E}\|V_k - V_k^\circ\|_F$. Similar to (237), by Lemma 31, we have

$$\mathbb{E}\|V_k - V_k^\circ\|_F \leq d \max_{1 \leq i, j \leq d} \sqrt{\mathbb{E}(v_{ij,k} - v_{ij,k}^\circ)^2}. \quad (238)$$

Thus, we only need to show the bound for the one-dimensional case. Now, consider $V_k, V_k^\circ, \mathcal{V}_k, H_k$ and $\bar{\beta}_k^{\text{sgd}}$ as scalars. Note that $\mathbb{E}\|\cdot\|_F = (\mathbb{E}[\cdot]^2)^{1/2}$ for $d = 1$. By the decomposition in (58) and applying Jensen's inequality, we have

$$(\mathbb{E}[V_k - V_k^\circ]^2)^{1/2} \leq (\mathbb{E}[\mathcal{V}_k - V_k^\circ]^2)^{1/2} + 2(\mathbb{E}[H_k \bar{\beta}_k^{\text{sgd}}]^2)^{1/2} + K_n(\mathbb{E}[\bar{\beta}_k^{\text{sgd}}]^4)^{1/2}. \quad (239)$$

We shall bound the three term respectively. First, recall the contraction constant $\check{r}_{\alpha,q}$ defined in (38). By the GMC in Theorem 12 and Hölder's inequality, it follows that

$$\begin{aligned} &(\mathbb{E}[\mathcal{V}_k - V_k^\circ]^2)^{1/2} \\ &\leq \sum_{m=1}^{\psi(k)-1} (\mathbb{E}[\mathcal{S}_m - \mathcal{S}_m^\circ]^4)^{1/4} (\mathbb{E}[\mathcal{S}_m + \mathcal{S}_m^\circ]^4)^{1/4} + (\mathbb{E}[\mathcal{R}_k - \mathcal{R}_k^\circ]^4)^{1/4} (\mathbb{E}[\mathcal{R}_k + \mathcal{R}_k^\circ]^4)^{1/4} \\ &\lesssim \sum_{m=1}^{\psi(k)-1} (\check{r}_{\alpha,q})^{\eta_m} (\eta_{m+1} - \eta_m)^{1/2} + (\check{r}_{\alpha,q})^{\eta_{\psi(k)}} (k - \eta_{\psi(k)} + 1)^{1/2} \\ &= O(1), \end{aligned} \quad (240)$$

where the constants in \lesssim and $O(\cdot)$ are independent of k . Similarly, since the dimension d is fixed and $\mathbb{E}[\check{\beta}_k^\circ] = \check{\beta}$ by (220) with the ℓ^2 -minimizer $\check{\beta}$ defined in (34), we can show that

$$(\mathbb{E}[\check{\beta}_k^{\text{sgd}}]^4)^{1/2} \lesssim \left(\mathbb{E} \left[\frac{1}{k} \sum_{i=1}^k \check{\beta}_i^\circ \right]^4 \right)^{1/2} + \left(\mathbb{E} \left[\frac{1}{k} \sum_{i=1}^k (\check{\beta}_i - \check{\beta}_i^\circ) \right]^4 \right)^{1/2} \asymp \frac{1}{k}, \quad (241)$$

and $(\mathbb{E}[H_k]^4)^{1/2} \lesssim k(k - \eta_{\psi(k)})^2$. Combining these results with the fact that $K_k \asymp k(k - \eta_{\psi(k)})$ yields

$$K_k (\mathbb{E}[\check{\beta}_k^{\text{sgd}}]^4)^{1/2} \asymp k - \eta_{\psi(k)} \quad \text{and} \quad (\mathbb{E}[H_k \check{\beta}_k^{\text{sgd}}]^2)^{1/2} \lesssim k - \eta_{\psi(k)}. \quad (242)$$

Inserting all these expressions back into (239), we obtain

$$(\mathbb{E}[V_k - V_k^\circ]^2)^{1/2} \lesssim k - \eta_{\psi(k)} \asymp k^{1-(1/\zeta)}. \quad (243)$$

Consequently, by (238), for the multi-dimensional case, it follows from (238) that

$$\mathbb{E} \|V_k - V_k^\circ\|_F \leq d \max_{1 \leq i, j \leq d} \sqrt{\mathbb{E}(v_{ij,k} - v_{ij,k}^\circ)^2} \lesssim dk^{1-1/(2\zeta)}.$$

Since $\zeta > 1$, compared to the rate of $\mathbb{E} \|V_k^\circ - k\check{\Sigma}\|_F \lesssim dk^{(1/\zeta) \vee (1-1/(2\zeta))}$ in (237), the latter dominates. This, along with (230) gives the desired result $\mathbb{E} \|\hat{\Sigma}_k(\alpha) - \check{\Sigma}(\alpha)\|_F = O(dk^{(1/\zeta) \vee (1-1/(2\zeta))})$. ■

G.2 Long-Run Covariances for Gaussian Random Samples

Consider i.i.d. covariate vectors $\mathbf{x}_k \sim \mathcal{N}(0, I_d)$, $k = 1, 2, \dots$, and the responses $y_k \mid \mathbf{x}_k$ from the linear regression model

$$y_k = \mathbf{x}_k^\top \beta^* + \epsilon_k, \quad \epsilon_k \sim \mathcal{N}(0, 1). \quad (244)$$

For the SGD iterates with dropout

$$\check{\beta}_k(\alpha) - \check{\beta} = \underbrace{(I_d - \alpha D_k \mathbb{X}_k D_k)}_{=: \check{A}_k(\alpha)} (\check{\beta}_{k-1}(\alpha) - \check{\beta}) + \underbrace{\alpha D_k \mathbf{x}_k (y_k - \mathbf{x}_k^\top D_k \check{\beta})}_{=: \check{b}_k(\alpha)},$$

defined in (36), and the minimizer of the ℓ^2 -regularized least-squares loss

$$\check{\beta} = p \left(p^2 \mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top] + p(1-p) \text{Diag}(\mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top]) \right)^{-1} \mathbb{E}[y_1 \mathbf{x}_1] = (\mathbb{E}[\mathbb{X}_{1,p}])^{-1} \mathbb{E}[y_1 \mathbf{x}_1],$$

defined in (34), we have

$$\check{\beta} = \beta^*. \quad (245)$$

This result holds since $\mathbb{E}[\mathbb{X}_{k,p}] = \mathbb{E}[D_k \mathbf{x}_k \mathbf{x}_k^\top D_k] = pI_d = \mathbb{E}[\mathbb{X}_k]$. Recall the long-run covariance matrix

$$\check{\Sigma}(\alpha) := \sum_{i=-\infty}^{\infty} \mathbb{E}[(\check{\beta}_0^\circ(\alpha) - \check{\beta})(\check{\beta}_i^\circ(\alpha) - \check{\beta})^\top],$$

defined in the quenched CLT in Theorem 15, where $\check{\beta}_i^\circ(\alpha) \sim \check{\pi}_\alpha$ denotes the stationary SGD process with dropout. We can derive a closed form expression of the limit

$$\check{\Sigma}^* := \lim_{\alpha \rightarrow 0} \check{\Sigma}(\alpha). \quad (246)$$

Lemma 32 (Long-run covariance matrices for Gaussian random samples) *For the model (244) with Gaussian random samples, suppose that conditions in Theorem 18 are satisfied. Then, $\check{\Sigma}^*$ is diagonal, with the i -th diagonal element*

$$\check{\Sigma}_{ii}(\alpha) \rightarrow \check{\Sigma}_{ii}^*, \quad \text{as } \alpha \rightarrow 0, \quad (247)$$

for all $i = 1, \dots, d$, where we write $\beta^* = (\beta_1^*, \dots, \beta_d^*)^\top$ and

$$\check{\Sigma}_{ii}^* = \frac{1}{p} + \frac{1-p}{p} \sum_{j \neq i} \beta_j^{*2}. \quad (248)$$

In particular, $\check{\Sigma}^* = I_d$ if $p = 1$.

Proof The proof consists of two steps. In **Step 1**, we construct an affine sequence with long-run covariance matrix $\check{\Sigma}$ and prove that with i.i.d. Gaussian random samples (y_k, \mathbf{x}_k) , $\check{\Sigma}^\dagger$ is diagonal, and the i -th diagonal element $\check{\Sigma}_{ii}^\dagger(\alpha) = \Sigma_{ii}^*$ for each i . Then in **Step 2**, we bound the error of this affine approximation by showing that $\|\check{\Sigma}(\alpha) - \check{\Sigma}^\dagger\| \rightarrow 0$ as $\alpha \rightarrow 0$.

Step 1. Define the affine approximation to the SGD dropout sequence $\{\check{\beta}_k(\alpha)\}_{k \in \mathbb{N}}$ as $\{\check{\beta}_k^\dagger(\alpha)\}_{k \in \mathbb{N}}$, which follows the recursion

$$\check{\beta}_k^\dagger(\alpha) - \check{\beta} = (I_d - \alpha H_p)(\check{\beta}_{k-1}^\dagger(\alpha) - \check{\beta}) + \check{\mathbf{b}}_k(\alpha), \quad (249)$$

where

$$H_p = \mathbb{E}[D_1 \mathbb{X}_1 D_1], \quad \check{\mathbf{b}}_k(\alpha) = \alpha D_k \mathbf{x}_k (y_k - \mathbf{x}_k^\top D_k \check{\beta}). \quad (250)$$

Recall the contraction constant $\check{r}_{\alpha, q}$ defined in (38). Let $q = 2$. Since $\|\cdot\|_2^2$ is convex, it follows from Jensen's inequality that

$$\mathbb{E}\|(I_d - \alpha D \mathbb{X}_1 D) \mathbf{v}\|_2^2 \geq \|\mathbb{E}[(I_d - \alpha D \mathbb{X}_1 D) \mathbf{v}]\|_2^2 = \|(I_d - \alpha H_p) \mathbf{v}\|_2^2. \quad (251)$$

Since the matrix $I - \alpha H_p$ is symmetric, we therefore obtain

$$\check{r}_{\alpha, 2}^2 = \sup_{\|\mathbf{v}\|_2=1} \mathbb{E}\|(I_d - \alpha D \mathbb{X}_1 D) \mathbf{v}\|_2^2 \geq \sup_{\|\mathbf{v}\|_2=1} \|(I_d - \alpha H_p) \mathbf{v}\|_2^2 = \|I_d - \alpha H_p\|^2. \quad (252)$$

By Theorem 12, since the constant learning rate α satisfies condition (37), we have $\|I_d - \alpha H_p\| \leq \check{r}_{\alpha, 2} < 1$. Hence, GMC also holds for the affine sequence $\{\check{\beta}_k^\dagger(\alpha)\}_{k \in \mathbb{N}}$ and $\check{\beta}_k^\dagger(\alpha)$ converges to a unique stationary distribution as $k \rightarrow \infty$. Throughout this proof, we assume that the initialization $\check{\beta}_0^\dagger$ follows this stationary distribution, and thus the sequence $\{\check{\beta}_k^\dagger(\alpha)\}_{k \in \mathbb{N}}$ is stationary.

Since $\{\check{\beta}_k^\dagger(\alpha)\}_{k \in \mathbb{N}}$ is a stationary sequence, define its long-run covariance matrix as

$$\check{\Sigma}^\dagger(\alpha) := \sum_{i=-\infty}^{\infty} \mathbb{E}[(\check{\beta}_0^\dagger(\alpha) - \check{\beta})(\check{\beta}_i^\dagger(\alpha) - \check{\beta})^\top].$$

We denote the i -th diagonal element of $\check{\Sigma}^\dagger(\alpha)$ by $\check{\Sigma}_{ii}^\dagger(\alpha)$, $i = 1, \dots, d$. Next, we show that

$$\check{\Sigma}_{ii}^\dagger(\alpha) = \check{\Sigma}_{ii}^*, \quad \text{for all } i. \quad (253)$$

Recall that $\check{\mathbf{b}}_k(\alpha)$ are i.i.d. and by (193), $\mathbb{E}[\check{\mathbf{b}}_k(\alpha)] = 0$. Thus, recursion (249) yields a VAR(1) process. By iterating (249), we obtain

$$\check{\beta}_k^\dagger(\alpha) - \check{\beta} = \sum_{j=0}^{\infty} (I_d - \alpha H_p)^j \check{\mathbf{b}}_{k-j}(\alpha), \quad (254)$$

which converges since $\|I_d - \alpha H_p\| < 1$. Define the covariance matrix

$$Q_p(\alpha) := \mathbb{E}[\check{\mathbf{b}}_1(\alpha)\check{\mathbf{b}}_1(\alpha)^\top]. \quad (255)$$

Since $\check{\mathbf{b}}_{k-j}$ are independent over j , it follows that

$$\begin{aligned} \check{\Sigma}^\dagger(\alpha) &= \sum_{j=0}^{\infty} (I_d - \alpha H_p)^j Q_p(\alpha) [(I_d - \alpha H_p)^j]^\top \\ &= [I_d - (I_d - \alpha H_p)]^{-1} Q_p(\alpha) [I_d - (I_d - \alpha H_p)]^{-1} \\ &= (\alpha H_p)^{-1} Q_p(\alpha) (\alpha H_p)^{-1}, \end{aligned} \quad (256)$$

where the sandwich form in the second equation follows from the the spectral density of $\check{\beta}_k^\dagger(\alpha) - \check{\beta}$ at frequency 0 (see Brockwell and Davis (1991), Chapter 11).

Now we compute H_p and $Q_p(\alpha)$ for i.i.d. Gaussian random samples. Since $\mathbf{x}_k \sim \mathcal{N}(0, I_d)$,

$$H_p = \mathbb{E}[D_1 \mathbb{X}_1 D_1] = p I_d, \quad (257)$$

which yields

$$\check{\Sigma}^\dagger(\alpha) = (\alpha p)^{-2} Q_p(\alpha). \quad (258)$$

Moreover, by (245), $\check{\beta} = \beta^*$, and thus $y_k = \mathbf{x}_k^\top \check{\beta} + \epsilon_k$ for all $k \in \mathbb{N}_+$, which leads to

$$\begin{aligned} Q_p(\alpha) &= \mathbb{E}[\check{\mathbf{b}}_1(\alpha)\check{\mathbf{b}}_1(\alpha)^\top] \\ &= \alpha^2 \mathbb{E}[(D_1 \mathbf{x}_1 (y_1 - \mathbf{x}_1^\top D_1 \check{\beta})) [D_1 \mathbf{x}_1 (y_1 - \mathbf{x}_1^\top D_1 \check{\beta})]^\top] \\ &= \alpha^2 \mathbb{E}[D_1 \mathbb{X}_1 D_1 (\epsilon_1 + \mathbf{x}_1^\top (I_d - D_1) \check{\beta})^2]. \end{aligned} \quad (259)$$

Expanding the square gives

$$(\epsilon_1 + \mathbf{x}_1^\top (I_d - D_1) \check{\beta})^2 = (\mathbf{x}_1^\top (I_d - D_1) \check{\beta})^2 + 2\epsilon_1 \mathbf{x}_1^\top (I_d - D_1) \check{\beta} + \epsilon_1^2.$$

We consider the three terms separately. Since ϵ_k is independent of \mathbf{x}_k and D_k ,

$$\mathbb{E}_{(y,\mathbf{x})}[D_1 \mathbb{X}_1 D_1 \epsilon_1 \mathbf{x}_1^\top (I_d - D_1) \check{\boldsymbol{\beta}} \mid D_1] = 0. \quad (260)$$

By $\mathbb{E}[\epsilon_k^2] = 1$ and $D_1^2 = D_1$, we have

$$\mathbb{E}_{(y,\mathbf{x})}[D_1 \mathbb{X}_1 D_1 \epsilon_1^2 \mid D_1] = D_1. \quad (261)$$

In addition, since $\mathbf{x}_k \sim \mathcal{N}(0, I_d)$,

$$\begin{aligned} & \mathbb{E}_{(y,\mathbf{x})}[D_1 \mathbb{X}_1 D_1 (\mathbf{x}_1^\top (I_d - D_1) \check{\boldsymbol{\beta}})^2 \mid D_1] \\ &= D_1 \mathbb{E}_{(y,\mathbf{x})}[\|(I_d - D_1) \check{\boldsymbol{\beta}}\|_2^2 I_d + 2(I_d - D_1) \check{\boldsymbol{\beta}} \check{\boldsymbol{\beta}}^\top (I_d - D_1)] D_1 \\ &= D_1 \mathbb{E}_{(y,\mathbf{x})}[\|(I_d - D_1) \check{\boldsymbol{\beta}}\|_2^2 I_d], \end{aligned} \quad (262)$$

where the last equation holds due to $D_1(I_d - D_1) = 0$. Combining all the three equations above, we obtain

$$\mathbb{E}_{(y,\mathbf{x})}[\check{\mathbf{b}}_1(\alpha) \check{\mathbf{b}}_1(\alpha)^\top \mid D_1] = \alpha^2 (1 + \|(I_d - D_1) \check{\boldsymbol{\beta}}\|_2^2) D_1. \quad (263)$$

Taking expectation over D gives the diagonal matrix $Q_p(\alpha)$ with entries

$$Q_{p,ii}(\alpha) = \alpha^2 \left(p + p(1-p) \sum_{j \neq i} \check{\beta}_j^2 \right), \quad (264)$$

and $Q_{p,ij}(\alpha) = 0$ for $i \neq j$. Inserting this back into (258) yields, for all i ,

$$\check{\Sigma}_{ii}^\dagger(\alpha) = \frac{1}{p} + \frac{1-p}{p} \sum_{j \neq i} \beta_j^{*2} = \check{\Sigma}_{ii}^*, \quad (265)$$

which is independent of α . Write $\check{\Sigma}^\dagger(\alpha) = \check{\Sigma}$.

Step 2. Finally, we complete the proof by showing that

$$\|\check{\Sigma}(\alpha) - \check{\Sigma}^\dagger\|_F \rightarrow 0, \quad \text{as } \alpha \rightarrow 0. \quad (266)$$

To this end, consider the difference

$$\check{\boldsymbol{\delta}}_k(\alpha) = \check{\boldsymbol{\beta}}_k^\circ(\alpha) - \check{\boldsymbol{\beta}}_k^\dagger(\alpha) = \check{A}_k(\alpha) \check{\boldsymbol{\delta}}_{k-1}(\alpha) - \alpha (D_k \mathbb{X}_k D_k - H_p) [\check{\boldsymbol{\beta}}_{k-1}^\dagger(\alpha) - \check{\boldsymbol{\beta}}]. \quad (267)$$

Since both initializations $\check{\boldsymbol{\beta}}_0^\circ$ and $\check{\boldsymbol{\beta}}_0^\dagger$ follow the stationary distributions, respectively, the difference sequence $\{\check{\boldsymbol{\delta}}_k(\alpha)\}_{k \in \mathbb{N}}$ is stationary. By Lemma 10 and the independence between (D_k, \mathbf{x}_k) and $\check{\boldsymbol{\beta}}_{k-1}^\dagger(\alpha)$, we obtain, for any $\eta > 0$,

$$\begin{aligned} \mathbb{E} \|\check{\boldsymbol{\delta}}_k(\alpha)\|_2^2 &\leq (1 + \eta) \mathbb{E} \|\check{A}_k(\alpha) \check{\boldsymbol{\delta}}_{k-1}(\alpha)\|_2^2 + \left(1 + \frac{1}{\eta}\right) \alpha^2 \mathbb{E} \|(D_k \mathbb{X}_k D_k - H_p) [\check{\boldsymbol{\beta}}_{k-1}^\dagger(\alpha) - \check{\boldsymbol{\beta}}]\|_2^2 \\ &\leq (1 + \eta) \check{r}_{\alpha,2}^2 \mathbb{E} \|\check{\boldsymbol{\delta}}_{k-1}(\alpha)\|_2^2 \\ &\quad + \left(1 + \frac{1}{\eta}\right) \alpha^2 \sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \mathbb{E} \|(D_k \mathbb{X}_k D_k - H_p) \mathbf{v}\|_2^2 \cdot \mathbb{E} \|\check{\boldsymbol{\beta}}_{k-1}^\dagger(\alpha) - \check{\boldsymbol{\beta}}\|_2^2. \end{aligned} \quad (268)$$

By Assumption 2, there exists some finite constant $c > 0$ such that

$$\sup_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \mathbb{E} \|(D_k \mathbb{X}_k D_k - H_p) \mathbf{v}\|_2^2 \leq c. \quad (269)$$

Moreover, by similar arguments as in the proof of Lemma 14, one can show that

$$\mathbb{E} \|\check{\boldsymbol{\beta}}_{k-1}^\dagger(\alpha) - \check{\boldsymbol{\beta}}\|_2^2 = O(\alpha). \quad (270)$$

Since $\check{r}_{\alpha,2} < 1$, by choosing η such that $(1 + \eta)\check{r}_{\alpha,2}^2 < 1$, it follows from the stationarity of $\{\check{\boldsymbol{\delta}}_k(\alpha)\}_{k \in \mathbb{N}}$ that

$$\mathbb{E} \|\check{\boldsymbol{\delta}}_k(\alpha)\|_2^2 = O(\alpha^3). \quad (271)$$

For brevity, denote

$$\mathbf{u}_k = \mathbf{u}_k(\alpha) = \check{\boldsymbol{\beta}}_k^\circ(\alpha) - \check{\boldsymbol{\beta}}, \quad \mathbf{u}_k^\dagger = \mathbf{u}_k^\dagger(\alpha) = \check{\boldsymbol{\beta}}_k^\dagger(\alpha) - \check{\boldsymbol{\beta}}. \quad (272)$$

For any lag $h \in \mathbb{Z}$, the difference of the two covariance matrices can be decomposed into

$$\mathbb{E}[\mathbf{u}_0 \mathbf{u}_h^\top] - \mathbb{E}[\mathbf{u}_0^\dagger \mathbf{u}_h^{\dagger \top}] = \mathbb{E}[\check{\boldsymbol{\delta}}_0 \mathbf{u}_h^\top] + \mathbb{E}[\mathbf{u}_0^\dagger \check{\boldsymbol{\delta}}_h^\top]. \quad (273)$$

We take the Frobenius norm and it follows from Cauchy-Schwarz inequality, Lemma 14 and expression (271) that

$$\|\mathbb{E}[\check{\boldsymbol{\delta}}_0 \mathbf{u}_h^\top]\|_F \leq (\mathbb{E} \|\check{\boldsymbol{\delta}}_0\|_2^2)^{1/2} (\mathbb{E} \|\mathbf{u}_h\|_2^2)^{1/2} = O(\alpha^{3/2}) \cdot O(\alpha^{1/2}) = O(\alpha^2). \quad (274)$$

Similarly, $\|\mathbb{E}[\mathbf{u}_0^\dagger \check{\boldsymbol{\delta}}_h^\top]\|_F = O(\alpha^2)$, which together with the expression above yields

$$\|\mathbb{E}[\mathbf{u}_0 \mathbf{u}_h^\top] - \mathbb{E}[\mathbf{u}_0^\dagger \mathbf{u}_h^{\dagger \top}]\|_F = O(\alpha^2). \quad (275)$$

For any truncation level $H \in \mathbb{N}$,

$$\begin{aligned} \|\check{\Sigma}(\alpha) - \check{\Sigma}^\dagger\|_F &= \left\| \sum_{h=-\infty}^{\infty} (\mathbb{E}[\mathbf{u}_0 \mathbf{u}_h^\top] - \mathbb{E}[\mathbf{u}_0^\dagger \mathbf{u}_h^{\dagger \top}]) \right\|_F \\ &\leq \left\| \sum_{|h| \leq H} (\mathbb{E}[\mathbf{u}_0 \mathbf{u}_h^\top] - \mathbb{E}[\mathbf{u}_0^\dagger \mathbf{u}_h^{\dagger \top}]) \right\|_F + \left\| \sum_{|h| > H} (\mathbb{E}[\mathbf{u}_0 \mathbf{u}_h^\top] - \mathbb{E}[\mathbf{u}_0^\dagger \mathbf{u}_h^{\dagger \top}]) \right\|_F \\ &\leq O(H\alpha^2) + \sum_{|h| > H} \|\mathbb{E}[\mathbf{u}_0 \mathbf{u}_h^\top]\|_F + \sum_{|h| > H} \|\mathbb{E}[\mathbf{u}_0^\dagger \mathbf{u}_h^{\dagger \top}]\|_F. \end{aligned} \quad (276)$$

Recall the functional dependence measure in (209). By GMC of \mathbf{u}_h , the coupled distance also contracts, which gives

$$\check{\theta}_{k,2}(\alpha) \leq \check{r}_{\alpha,2}^k \check{\theta}_{0,2}(\alpha) \leq 2\check{r}_{\alpha,2}^k (\mathbb{E} \|\mathbf{u}_0\|_2^2)^{1/2} = 2\check{r}_{\alpha,2}^k O(\sqrt{\alpha}). \quad (277)$$

Recall the projection operator $\check{\mathcal{P}}_k[\cdot]$ in (206). By the martingale decomposition and the orthogonality of projections, for $h > 0$,

$$\begin{aligned} \mathbb{E}[\mathbf{u}_0 \mathbf{u}_h^\top] &= \mathbb{E} \left[\left(\sum_{i=0}^{\infty} \check{\mathcal{P}}_{-i}[\mathbf{u}_0] \right) \left(\sum_{j=0}^{\infty} \check{\mathcal{P}}_{h-j}[\mathbf{u}_h] \right)^\top \right] \\ &= \sum_{i=0}^{\infty} \mathbb{E} \left[\check{\mathcal{P}}_{-i}[\mathbf{u}_0] \left(\check{\mathcal{P}}_{-i}[\mathbf{u}_h] \right)^\top \right]. \end{aligned} \quad (278)$$

By the triangle inequality, Cauchy-Schwarz inequality, and Theorem 1 in Wu (2005),

$$\begin{aligned}
\|\mathbb{E}[\mathbf{u}_0 \mathbf{u}_h^\top]\|_F &\leq \sum_{i=0}^{\infty} \left\| \mathbb{E} \left[\check{\mathcal{P}}_{-i}[\mathbf{u}_0] \left(\check{\mathcal{P}}_{-i}[\mathbf{u}_h] \right)^\top \right] \right\|_F \\
&\leq \sum_{i=0}^{\infty} \left[\|\mathbb{E}(\check{\mathcal{P}}_{-i}[\mathbf{u}_0])\|_2^2 \right]^{1/2} \left[\|\mathbb{E}(\check{\mathcal{P}}_{-i}[\mathbf{u}_h])\|_2^2 \right]^{1/2} \\
&\leq \sum_{j=0}^{\infty} \check{\theta}_{j,2}(\alpha) \check{\theta}_{h+j,2}(\alpha).
\end{aligned} \tag{279}$$

As a direct consequence,

$$\sum_{h>H} \|\mathbb{E}[\mathbf{u}_0 \mathbf{u}_h^\top]\|_F \leq \sum_{h>H} \sum_{j=0}^{\infty} \check{\theta}_{j,2}(\alpha) \check{\theta}_{h+j,2}(\alpha) = O(\alpha \check{r}_{\alpha,2}^H). \tag{280}$$

The bounds for the parts $\sum_{h<-H} \|\mathbb{E}[\mathbf{u}_0 \mathbf{u}_h^\top]\|_F$ and $\sum_{|h|>H} \|\mathbb{E}[\mathbf{u}_0^\dagger \mathbf{u}_h^{\dagger\top}]\|_F$ can follow similar arguments. Choose $H = \lceil c' \log(1/\alpha) \rceil$ with some constant $c' > 1/|\log(\check{r}_{\alpha,2})|$. Then, we have $\alpha \check{r}_{\alpha,2}^H \leq \alpha^{c' |\log(\check{r}_{\alpha,2})| + 1} = o(\alpha)$ for the tail part in (276), and $H\alpha^2 = c' \alpha^2 \log(1/\alpha)$. Both converge to 0 as $\alpha \rightarrow 0$. This completes the proof. \blacksquare

Appendix H. Real-Data Application

The *Online News Popularity* dataset from UCI ML Repository (Fernandes et al., 2015) summarizes a set of features about articles published by *Mashable* (www.mashable.com) in a period of two years. This dataset contains 39,797 samples and 58 features. We apply the efficient ASGD-dropout algorithm to predict the number of *shares* in social networks (popularity) and use our proposed online inference methodology to test the significance of features towards prediction.

To standardize the dataset, we center all the predictors and scale them to unit variance. The response is transformed as $\log(1 + \text{shares})$ and centered. This ensures comparable gradient magnitudes across features and avoids a few heavy-tailed article outliers dominating the updates. We adopt $p = 0.9$ and the dropout rate is $1 - p = 0.1$. Dropout regularization is not strictly necessary for classical linear regression, but it demonstrates how the ASGD framework can incorporate randomized feature masking, especially when scaling to large-dimensional data. The online nature of the proposed inference methodology means we never store the full design matrix or compute costly matrix inversions, useful for streaming or memory-constrained environments.

Figures 5 and 6 plot the ASGD-dropout trajectories over $n = 39,797$ iterations, showing rapid stabilization of estimated regression coefficients. By estimating the long-run variances of these coefficients and constructing marginal 95% CIs using expression (60), at the last iteration, any feature whose CI does not cover zero can be deemed statistically significant. Prominent positive predictors include *number of videos*, *number of links*, *publication at week-ends* and *global text subjectivity*. Primary negative predictors include *mid-week publication* and *average length of the words in the content*. Interestingly, for the five Latent Dirichlet

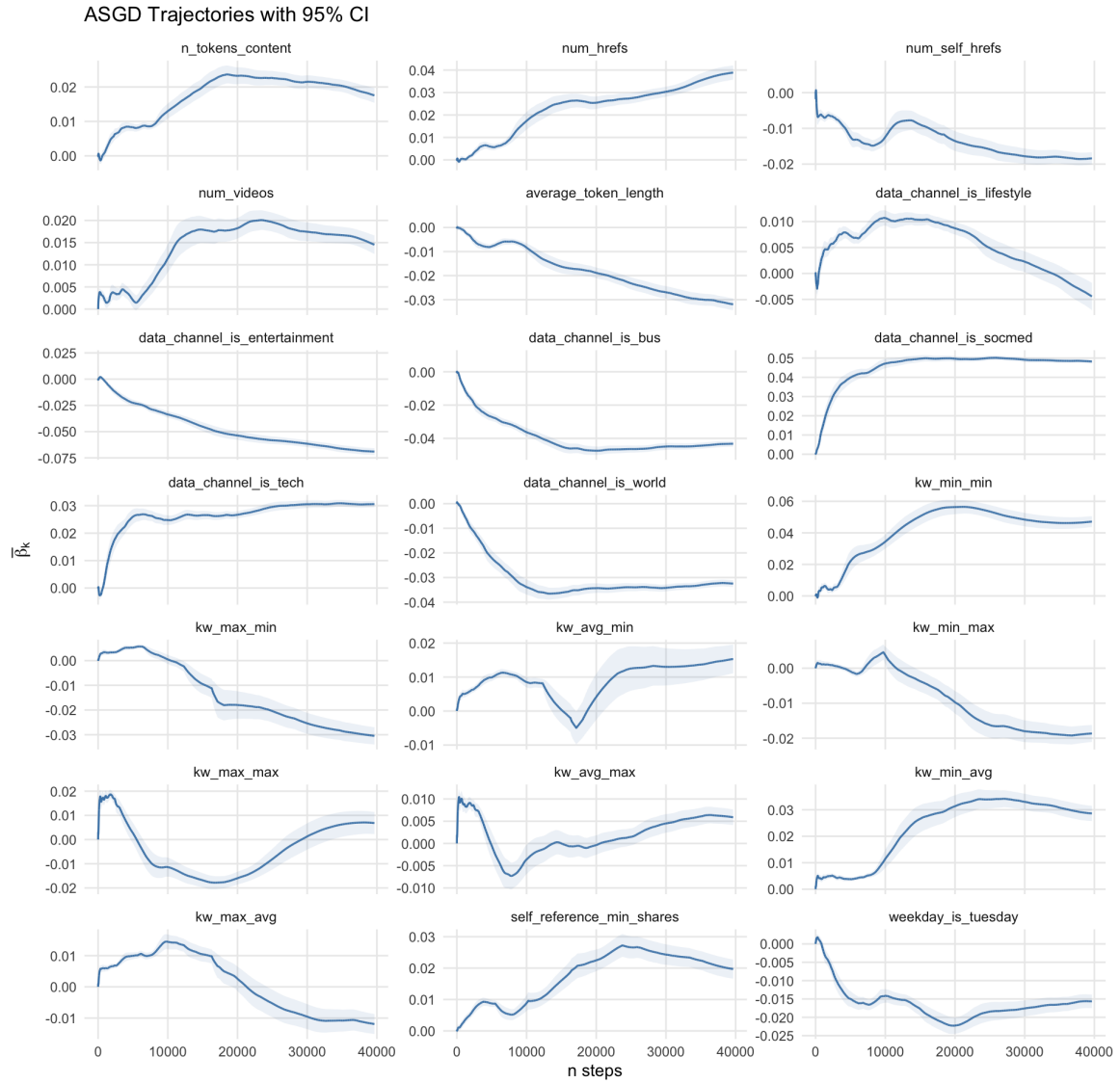


Figure 5: Convergence trace of 95% CI of ASGD iterates with dropout regularization for significant features. Each subplot represents the convergence trace of one feature.

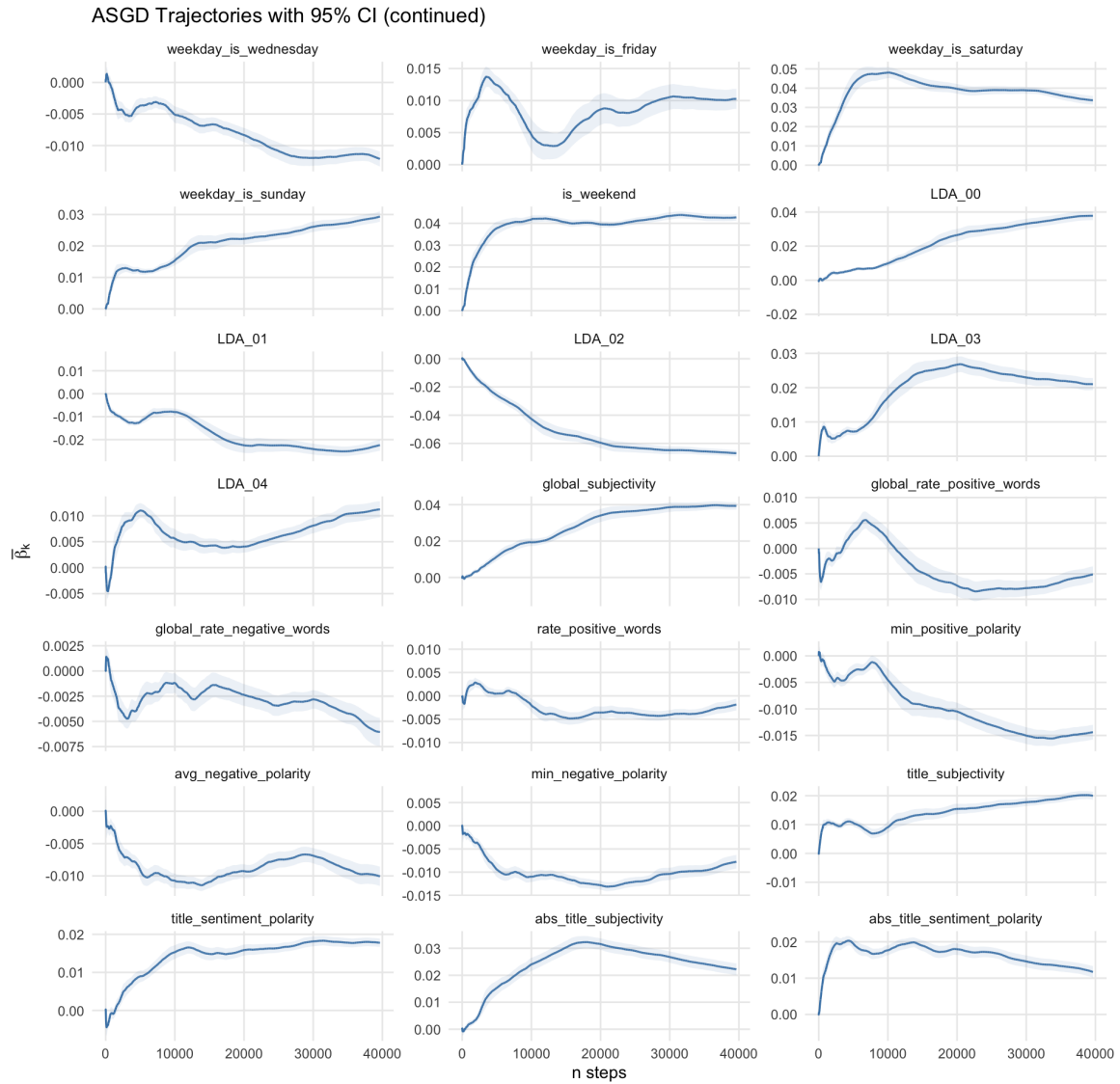


Figure 6: (Continued) Convergence trace of 95% CI of ASGD iterates with dropout regularization for significant features. Each subplot represents the convergence trace of one feature.

Allocation (LDA) topics, LDA topics 0, 3 and 4 have positive coefficients whose 95% CIs lie entirely above zero, while LDA topics 1 and 2 have coefficients whose CIs lie entirely below zero. We refer to Figures 5 and 6 for the detailed significant features, where we plot the trajectory of ASGD-dropout (solid line) with its 95% CI ribbon for the significant features.

Our findings of significance features echo the literature. For example, Ciampaglia et al. (2018) found that certain latent topics (e.g. technology, social media) boost sharing while others (e.g. politics, hard news) suppress it, matching our positive significance for certain LDA components and negative for others. Szabo and Huberman (2010) documented a “weekend effect” in which articles published on weekends tend to garner more attention, consistent with the positive coefficient on weekend publication. Moreover, Bandari et al. (2021) demonstrated that multimedia richness (links and video counts) are among the most correlated variables with news popularity.