

# Larger Datasets Can Be Repeated More: A Theoretical Analysis of Multi-Epoch Scaling in Linear Regression

Tingkai Yan<sup>\*</sup>  <sup>$\beta$</sup>  Haodong Wen<sup>\*</sup>  <sup>$\alpha$</sup>  Binghui Li<sup>\*</sup>  <sup>$\gamma$</sup>   
 Kairong Luo <sup>$\nu$</sup>  Wenguang Chen <sup>$\nu\lambda$</sup>  Kaifeng Lyu <sup>$\dagger$</sup>   <sup>$\alpha$</sup>

<sup>$\alpha$</sup> Institute for Interdisciplinary Information Sciences, Tsinghua University

<sup>$\beta$</sup> School of Mathematical Sciences, Peking University

<sup>$\gamma$</sup> Center for Machine Learning Research, Peking University

<sup>$\nu$</sup> Department of Computer Science and Technology, Tsinghua University

<sup>$\lambda$</sup> Peng Cheng Laboratory

yantingkai66@gmail.com, whd25@mails.tsinghua.edu.cn

libinghui@pku.edu.cn, luokr24@mails.tsinghua.edu.cn

{cwg, klyu}@mail.tsinghua.edu.cn

## Abstract

While data scaling laws of large language models (LLMs) have been widely examined in the one-pass regime with massive corpora, their form under limited data and repeated epochs remains largely unexplored. This paper presents a theoretical analysis of how a common workaround, training for multiple epochs on the same dataset, reshapes the data scaling laws in linear regression. Concretely, we ask: to match the performance of training on a dataset of size  $N$  for  $K$  epochs, how much larger must a dataset be if the model is trained for only one pass? We quantify this using the *effective reuse rate* of the data,  $E(K, N)$ , which we define as the multiplicative factor by which the dataset must grow under one-pass training to achieve the same test loss as  $K$ -epoch training. Our analysis precisely characterizes the scaling behavior of  $E(K, N)$  for SGD in linear regression under either strong convexity or Zipf-distributed data: (1) When  $K$  is small, we prove that  $E(K, N) \approx K$ , indicating that every new epoch yields a linear gain; (2) As  $K$  increases,  $E(K, N)$  plateaus at a problem-dependent value that grows with  $N$  ( $\Theta(\log N)$  for the strongly-convex case), implying that larger datasets can be repeated more times before the marginal benefit vanishes. These theoretical findings point out a neglected factor in a recent empirical study by Muennighoff et al. [35], which claimed that training LLMs for up to 4 epochs results in negligible loss differences compared to using fresh data at each step, *i.e.*,  $E(K, N) \approx K$  for  $K \leq 4$  in our notation. Supported by further empirical validation with LLMs, our results reveal that the maximum  $K$  value for which  $E(K, N) \approx K$  in fact depends on the data size and distribution, and underscore the need to explicitly model both factors in future studies of scaling laws with data reuse.

## 1. Introduction

Scaling laws [17, 18, 22] have emerged as a central framework for characterizing the behavior of large language model (LLM) pre-training. The Chinchilla scaling law [18] established robust empirical trends in performance as a joint function of model size and dataset size under the one-pass training paradigm, in which each data point is used at most once. This assumption, however, is

---

<sup>\*</sup> Equal contribution.

<sup>$\dagger$</sup>  Corresponding author.

becoming increasingly untenable. The quest for more capable models has driven an unprecedented escalation in data requirements: from fewer than 10 billion tokens for GPT-2, to 300 billion for GPT-3 [7], 2 trillion for Chinchilla and LLaMA-2 [18, 43], and 36 trillion for Qwen3 [53]. Projections further suggest that the pool of publicly available data may be exhausted as early as 2028 [45].

A common response to this emerging data scarcity is to train models for multiple epochs over the same dataset. Recent empirical studies have begun to examine the consequences of such repetition: for example, Muennighoff et al. [35] and Xue et al. [52] show that moderate reuse can still yield competitive pre-training performance. Yet the fundamental scaling behavior of multi-epoch training remains poorly understood—particularly from a theoretical standpoint.

In this paper, we study a fundamental question in understanding how multi-epoch training affects the data scaling laws: *To what extent does training for  $K$  epochs on  $N$  samples can be effectively seen as one-pass training with an increased number of data samples?* Formally, let  $\mathcal{L}(K, N)$  denote the expected loss of  $K$ -epoch training on  $N$  samples. We define the *effective dataset size*  $N'(K, N)$  as the minimal number of samples in one-pass training that achieves a comparable or lower loss  $\mathcal{L}(1, N') \leq \mathcal{L}(K, N)$ . In this paper, we concern about the ratio  $E(K, N) = N'(K, N)/N$ , which we term as the *effective reuse rate* of the data, a key quantity that characterizes how many times larger the dataset must grow to match the same performance as  $K$ -epoch training (see the detailed version in Definition 2.1).

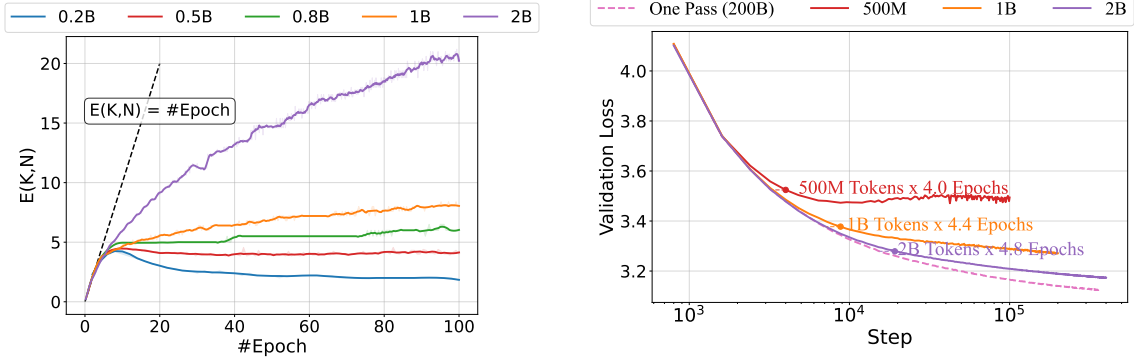
In a recent study of scaling laws for multi-epoch training, Muennighoff et al. [35] encountered this question and proposed an empirical approximation:

$$N'(K, N) = \left(1 + R^*(1 - e^{-(K-1)/R^*})\right) \cdot N, \quad (1)$$

where  $R^*$  is a fitted constant ( $R^* \approx 15.39$  in their experiments). This formula suggests that the benefit of repetition grows with  $K$  but saturates exponentially at  $(1 + R^*) \cdot N$  as  $K$  increases. While supported by some empirical evidence in their study, this approximation still leads to a noticeable gap between scaling law predictions and empirical results (see Figure 3 in their paper). Moreover, the formula implies that the ratio  $E(K, N) = N'(K, N)/N$  is independent of  $N$ , so the benefit of repeating the dataset  $K$  times is equivalent to increasing its size by a factor that depends only on  $K$ , regardless of how large  $N$  is. It remains unclear to what extent this independence holds in general.

**Our Contributions.** In this paper, we approach the above question on the effective reuse rate of data in the setting of linear regression, a setting that is simple enough to reveal the key mechanisms of data reuse, while still tractable for precise analysis under stochastic gradient descent (SGD). We provide a theoretical characterization of  $E(K, N)$  in various regimes, and point out a neglected factor in the empirical study of Muennighoff et al. [35]: the effective reuse rate depends not only on the number of epochs  $K$ , but also on the dataset size  $N$ . In fact, larger datasets can be repeated more. Our main contributions are as follows:

1. In Section 3, we study the strongly convex case of linear regression, and show that when  $K$  is small,  $E(K, N) \approx K$ , indicating that every new epoch leads to a linear gain. As  $K$  increases,  $E(K, N)$  saturates at a problem-dependent value of order  $\Theta(\log N)$ , suggesting that larger datasets can be repeated for more epochs before the marginal benefit vanishes.
2. In Section 4, we go beyond the strongly convex case and study a class of Zipf-law distributed data, and show that  $E(K, N)$  exhibits a similar scaling behavior to the strongly convex case, except that the saturation point scales as a power of  $N$  instead of  $\log N$ .



(a) The effective reuse rate  $E(K, N)$  as a function of the epoch number  $K$ . (b) Training loss as a function of training steps for different fresh data sizes.

Figure 1: The effective reuse rate  $E(K, N)$  over  $K$  and training curves in language model experiments. Figure 1(a) shows that  $E(K, N) \approx K$  when  $K$  is small, to be specific,  $K \leq 4$ . Figure 1(b) plots the points where  $E(K, N) = 0.8K$  under different configurations, and we observe that  $E(K, N)$  increases as  $N$  increases, indicating that larger datasets can be repeated more.

3. Technically, we derive the optimal learning rate (Lemma E.4) for multi-epoch SGD in linear regression and its corresponding approximation formula for the expected excess risk up to an  $o(1)$  multiplicative error (Lemma F.1). These results may be of independent interest.
4. In Section B, we conduct LLM pretraining experiments up to 200B repeated tokens, and empirically validate our theoretical predictions. The results confirm that  $E(K, N) \approx K$  for small  $K$ , and that for fixed  $K$ , the effective reuse rate increases monotonically with  $N$ . This provides direct evidence for our main conclusion: larger datasets can be repeated more.

## 2. Preliminaries

Now we introduce the problem setups in our paper; the notations are given in Section E.

**Linear Regression Problem.** We focus on a linear regression setup, where data point  $(\mathbf{x}, y) \in \mathbb{R}^d \times \mathbb{R}$  follows a joint distribution  $\mathcal{D}$  and  $\|\mathbf{x}\| \leq D$  for some constant  $D$ . W.L.O.G., we assume that the covariance matrix of data input is diagonal, i.e.,  $\mathbf{H} := \mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$ , where  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ . A direct corollary is that  $\lambda_1 \leq D^2$ . For a given data input  $\mathbf{x}$ , the label  $y$  is generated by  $y := \langle \mathbf{w}^*, \mathbf{x} \rangle + \xi$ , where  $\mathbf{w}^* \in \mathbb{R}^d$  is the ground-truth weight and  $\xi$  represents the independent random label noise with  $\mathbb{E}[\xi] = 0$  and  $\mathbb{E}[\xi^2] = \sigma^2$ . We aim to train a linear model  $f(\mathbf{x}; \mathbf{w}) := \langle \mathbf{w}, \mathbf{x} \rangle$  to predict the data label, where  $\mathbf{w} \in \mathbb{R}^d$  is the trainable parameter. We use MSE-loss  $\ell(\mathbf{w}; \mathbf{x}, y) := \frac{1}{2}(f(\mathbf{x}; \mathbf{w}) - y)^2$  to measure the fitting error. Then, the population loss is defined as  $\mathcal{L}(\mathbf{w}) := \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[\ell(\mathbf{w}; \mathbf{x}, y)]$ . Further we define the excess risk  $\mathcal{R}(\mathbf{w}) := \mathcal{L}(\mathbf{w}) - \frac{1}{2}\sigma^2$ , which is the expected population loss minus the irreducible loss  $\frac{1}{2}\sigma^2$ .

**Multi-Epoch SGD Training Algorithm.** Consider a finite training dataset with  $N$  data points  $\{(\mathbf{x}_0, y_0), (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{N-1}, y_{N-1})\}$ , where the data points  $(\mathbf{x}_i, y_i)$  are i.i.d. sampled from the distribution  $\mathcal{D}$ . We use  $K$ -epoch stochastic gradient descent (SGD) with random shuffling [1, 13,

[14, 34] to minimize the loss function. And the initial parameter  $\mathbf{w}_0$  is set to 0. Formally, we denote  $K$  independent random permutations of  $[N]$  by  $\pi_1, \dots, \pi_K$ . And we define  $j_t := \pi_{k_t}(i_t)$ , where  $i_t := t \bmod N$ ,  $k_t := \lfloor t/N \rfloor + 1$ . Then we have the update rule for  $K$ -epoch SGD with  $N$  data points:  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla_{\mathbf{w}} \ell(\mathbf{w}_t; \mathbf{x}_{j_t}, y_{j_t}) = \mathbf{w}_t - \eta \mathbf{x}_{j_t} \mathbf{x}_{j_t}^\top (\mathbf{w}_t - \mathbf{w}^*) + \eta \xi_{j_t} \mathbf{x}_{j_t}$ . Next, given a  $K$ -epoch SGD over  $N$  data points, with learning rate  $\eta$ , we define  $\mathcal{W}_{K,N,\eta}$  to be the distribution of  $\mathbf{w}_{KN}$ . The randomness within  $\mathbf{w}_{KN}$  comes from the random draw of the dataset, label noise  $\xi$ , and the shuffling in SGD. Based on this, we define the expected excess risk of a given  $K$ -epoch SGD over  $N$  data points, with learning rate  $\eta$  as  $\bar{\mathcal{R}}(K, N; \eta) := \mathbb{E}_{\mathbf{w} \sim \mathcal{W}_{K,N,\eta}}[\mathcal{R}(\mathbf{w})]$ . We assume  $\eta \leq D^{-2}$  for training stability.

**Comparing Performance under Optimal Learning Rate Regime.** To compare the performance of one-pass and multi-epoch SGD, we consider the settings where the learning rates for both methods are tuned to the optimal. Formally, we introduce the notion of the *optimal expected excess risk* of  $K$ -epoch SGD for  $N$  samples as  $\bar{\mathcal{R}}^*(K, N) := \min_{\eta \in (0, \frac{1}{D^2}]} \{\bar{\mathcal{R}}(K, N; \eta)\}$ . To calculate this value in math, we will show in the next section that we can get a learning rate choice that can approximately achieve the above optimal expected excess risk  $\bar{\mathcal{R}}^*(K, N)$  both for one-pass and multi-epoch SGD. Following our discussion in the introduction, we define the *effective reuse rate* as follows:

**Definition 2.1 (Effective Reuse Rate)** *Given  $K$ -epoch SGD trained with  $N$  fresh data samples, the effective reuse ratio is defined as  $E(K, N) := \frac{1}{N} \min\{N' \geq 0 : \bar{\mathcal{R}}^*(1, N') \leq \bar{\mathcal{R}}^*(K, N)\}$ .*

That is, the effective reuse rate measures how many times larger the dataset must grow under one-pass training to match the performance of  $K$ -epoch training, both under the optimal learning rate regime.

### 3. Multi-Epoch Scaling in Strongly Convex Linear Regression

In the study of linear regression problems, the strongly convex case is a classical and central theoretical framework, serving as the standard entry point before many relaxing to weaker conditions [12, 15]. In Section 3.1, we first give the problem setups and the main results of the effective reuse rate. In Section E.1, we give a proof sketch for our theoretical results, and the detailed proof of this section can be found in Appendix F.

#### 3.1. Main Results

We first make the following assumptions:

**Assumption 3.1 (Strong Convexity)** *We assume that  $\lambda_d \geq \mu$  for some constant  $\mu > 0$ .*

**Assumption 3.2 (Parameter Prior)** *The ground truth  $\mathbf{w}^*$  satisfies  $w_i^* \neq 0$  for all  $i \in [d]$ .*

**Assumption 3.3 (Computationally feasible number of epochs)** *We assume that the training dataset size  $N$  and number of epochs  $K$  satisfy  $K = O(N^{0.1})$ .*

To compute  $E(K, N)$ , we first precisely characterize the optimal expected excess risk. In particular, we derive asymptotic expansions for  $\bar{\mathcal{R}}^*(K, N)$  in the regimes  $K = o(\log N)$  and  $K = \omega(\log N)$ , each expressed as a leading term accompanied by an explicitly controlled higher-order remainder.

**Theorem 3.1 (Multi-Epoch Data Scaling Law)** Under [Assumptions 3.1 to 3.3](#), for multi-epoch SGD with the number of epochs  $K$ , dataset size of  $N$ , it holds that

$$\bar{\mathcal{R}}^*(K, N) = \begin{cases} \frac{\sigma^2 \text{tr}(\mathbf{H})}{8\lambda_d} (1 + o_N(1)) \cdot \frac{\log(KN)}{KN} & \text{for } K = o(\log N), \\ \frac{\sigma^2 d}{2} (1 + o_N(1)) \cdot \frac{1}{N} & \text{for } K = \omega(\log N). \end{cases}$$

[Theorem 3.1](#) describes how expected excess risk decays with number of epochs  $K$  and dataset size  $N$  when choosing the optimal learning rate. When  $K \ll \log N$ , then  $\bar{\mathcal{R}}^*(K, N) = \Theta\left(\frac{\log T}{T}\right)$  where  $T = KN$ ; by contrast, when  $K \gg \log N$ , then  $\bar{\mathcal{R}}^*(K, N) = \Theta\left(\frac{1}{N}\right)$  which does not depend on  $K$ , showing that endless data reuse turns to be useless.

Next we propose the expression of  $E(K, N)$  by applying [Theorem 3.1](#).

**Theorem 3.2** Under [Assumptions 3.1 to 3.3](#), for multi-epoch SGD with the number of epochs  $K$ , dataset size of  $N$ , it holds that

$$E(K, N) = \begin{cases} (1 + o_N(1)) \cdot K & \text{for } K = o(\log N), \\ \frac{\text{tr}(\mathbf{H})}{4\lambda_d d} (1 + o_N(1)) \cdot \log N & \text{for } K = \omega(\log N). \end{cases}$$

[Theorem 3.2](#) pinpoints two regimes for the effective reuse rate in the strongly convex case. The first one is an *effective-reuse regime*: when  $K \ll \log N$ , then  $E(K, N) = K(1 + o(1))$ . This suggests that each extra epoch is essentially as valuable as a fresh pass. The second one is a *limited-reuse regime*: when  $K \gg \log N$ , then  $E(K, N) = \frac{\text{tr}(\mathbf{H}) \log N}{4\lambda_d d} (1 + o_N(1))$ , which means additional epochs yield only logarithmic gains. This further implies that the model has effectively “seen” the dataset enough times that additional repetition is redundant.

Together, these two asymptotic descriptions expose a phase transition when the quantity  $\lim_{N \rightarrow \infty} \frac{K}{\log N}$  changes from 0 to  $\infty$ . For the former case ( $\lim_{N \rightarrow \infty} \frac{K}{\log N} = 0$ ), multi-epoch training behaves like unlimited data augmentation; for the latter ( $\lim_{N \rightarrow \infty} \frac{K}{\log N} = \infty$ ), the benefits of reusing data all but vanish, capping  $E(K, N)$  at  $\Theta(\log N)$ . This insight provides a precise guideline for practitioners: one should allocate epochs up to order  $\log N$  to maximize effective data utilization, but pushing  $K$  significantly beyond that yields rapidly diminishing returns.

#### 4. A Solvable Case with Zipf-distributed Data

Natural data distributions often exhibit power law structures. To capture this phenomenon, we go beyond the strongly convex case and analyze a stylized linear regression model with Zipf-distributed data, where the excess risk admits a closed-form expression and the effective reuse rate can be characterized explicitly.

Through this setup, we can see that the effective reuse rate exhibits a similar scaling behavior: as the number of epochs  $K$  increases,  $E(K, N)$  initially grows linearly but eventually saturates at a problem-dependent value that increases with  $N$ . In contrast to the strongly convex case, however, the saturation point does not scale as  $\sim \log N$  but instead scales as a power of  $N$ .

**Problem Setup.** We use the same notation for excess risk, one-pass and multi-epoch SGD, and *i.i.d.* training data as in [Section 2](#). We specify the data distribution as a Zipf distribution over  $d$  one-hot data points, where the  $i$ -th data point is  $\mathbf{x}^{(i)} = \mu_i \mathbf{e}_i$  for some  $\mu_i > 0$  and the probability of sampling the  $i$ -th data point is  $p_i$ . Furthermore, we define  $\Lambda_i = \mu_i^2$ ,  $\mathbf{P} = \text{diag}(p_1, p_2, \dots, p_d)$ ,

and  $\Lambda = \text{diag}(\Lambda_1, \Lambda_2, \dots, \Lambda_d)$ . The label is generated by  $y = \langle \mathbf{w}^*, \mathbf{x} \rangle$  with no label noise. The ground-truth weight  $\mathbf{w}^* \in \mathbb{R}^d$  follows an isotropic prior distribution. Here,  $\Lambda$  can have a power-law spectrum or a logarithmic power-law spectrum; the former is presented in [Section 4.1](#), and the latter is presented in [Section D.1](#).

**Assumption 4.1 (Parameter Prior)**  $\mathbf{w}^*$  is sampled from a prior distribution with  $\mathbb{E}[\mathbf{w}^* \mathbf{w}^{*\top}] = \mathbf{I}$ .

**Interpretation.** This setup can be interpreted as a simplified model of real-world data with heavy-tailed feature distributions. Each coordinate represents an atomic feature that appears with Zipf-distributed probability, mimicking the long-tailed statistics observed in domains such as text and natural language. The scaling factors  $\mu_i$  encode feature importance, which may reflect, for instance, effects introduced by feature weighting or normalization.

#### 4.1. Results on Power-Law Spectrum

**Assumption 4.2 (Power-Law Spectrum)** There exist two constants  $a, b > 0$  with  $a - b > 1$  such that the data input distribution satisfies that  $p_i = ci^{-(a-b)}$  and  $\Lambda_i = i^{-b}$ , where  $c = \left(\sum_{i=1}^d \frac{1}{i^{a-b}}\right)^{-1}$ .

Here we establish matching upper and lower bounds for  $\bar{\mathcal{R}}^*(K, N)$  in the small- $K$  and large- $K$  regimes, given the solvable model. Comparing with the strongly convex case, we observe a different scaling behavior: when  $K \ll N^{\frac{b}{a-b}}$ ,  $\bar{\mathcal{R}}^*(K, N)$  decays as a power law in  $KN$ , with exponent  $\frac{a-1}{a}$ ; whereas when  $K \gg N^{\frac{b}{a-b}}$ ,  $\bar{\mathcal{R}}^*(K, N)$  exhibits a power-law decay in  $N$  and is independent of  $K$ .

**Theorem 4.1** Consider a  $K$ -epoch SGD over  $N$  fresh data. Under Assumptions 4.1-4.2, and given the data dimension  $d = \Omega((KN)^{\frac{1}{a}})$ , it holds that

$$\bar{\mathcal{R}}^*(K, N) \asymp \begin{cases} (KN)^{-\frac{a-1}{a}} & \text{for } K = o(N^{\frac{b}{a-b}}) \\ N^{-\frac{a-1}{a-b}} & \text{for } K = \omega(N^{\frac{b}{a-b}}). \end{cases}$$

Then we derive the formula of  $E(K, N)$  by first solving the equation  $\bar{\mathcal{R}}^*(1, T') = \bar{\mathcal{R}}^*(K, N)$  based on [Theorem 4.1](#), and divide  $T'$  by  $N$ .

**Theorem 4.2 (Multi-Epoch Scaling Under Power-Law Spectrum)** Consider a  $K$ -epoch SGD over  $N$  fresh data. Under Assumptions 4.1-4.2, and given the data dimension  $d = \Omega((KN)^{\frac{1}{a}})$ , it holds that

$$E(K, N) = \begin{cases} K(1 + o(1)) & \text{for } K = o(N^{\frac{b}{a-b}}) \\ \Theta(N^{\frac{b}{a-b}}) & \text{for } K = \omega(N^{\frac{b}{a-b}}). \end{cases}$$

Under the assumption of a logarithmic power-law spectrum, the trend of the effective reuse rate as a function of  $K$  approximates the phenomena described in [Theorem 3.2](#) in the strongly convex setting and the trend described in [Theorem 3.2](#) under the power-law spectrum assumption. We still observe an effective-reuse regime ( $E(K, N) \approx K$ ) when  $K$  is relatively small ( $K \ll N^{b/(a-b)}$ ), and as  $K$  increases, the effective reuse rate undergoes a phase transition, converging to an upper bound determined by  $N$ , entering the limited-reuse regime ( $E(K, N) = \Theta(N^{b/(a-b)})$ ).

We can see that the exponent of this power of  $N$  is determined by the rate of eigenvalue decay of the Hessian and the rate of norm decay of the parameter with respect to dimension. The proofs of [Theorem 4.1](#) and [Theorem 4.2](#) are given in [Appendix H.2](#) and [Appendix H.3](#) respectively.



## References

- [1] Kwangjun Ahn, Chulhee Yun, and Suvrit Sra. Sgd with shuffling: optimal rates without component convexity and large epoch requirements. *Advances in Neural Information Processing Systems*, 33:17526–17535, 2020.
- [2] Luca Arnaboldi, Yatin Dandi, Florent Krzakala, Luca Pesce, and Ludovic Stephan. Repetita iuvant: Data repetition allows sgd to learn high-dimensional multi-index functions, 2025. URL <https://arxiv.org/abs/2405.15459>.
- [3] Alexander Atanasov, Jacob A Zavatone-Veth, and Cengiz Pehlevan. Scaling and renormalization in high-dimensional regression. *arXiv preprint arXiv:2405.00592*, 2024.
- [4] Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining neural scaling laws. *Proceedings of the National Academy of Sciences*, 121(27):e2311878121, 2024.
- [5] Blake Bordelon, Alexander Atanasov, and Cengiz Pehlevan. A dynamical model of neural scaling laws. *arXiv preprint arXiv:2402.01092*, 2024.
- [6] Blake Bordelon, Alexander Atanasov, and Cengiz Pehlevan. How feature learning can improve neural scaling laws. *arXiv preprint arXiv:2409.17858*, 2024.
- [7] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [8] François Charton and Julia Kempe. Emergent properties with repeated examples. *arXiv preprint arXiv:2410.07041*, 2024.
- [9] Yatin Dandi, Emanuele Troiani, Luca Arnaboldi, Luca Pesce, Lenka Zdeborová, and Florent Krzakala. The benefits of reusing batches for gradient descent in two-layer networks: Breaking the curse of information and leap exponents, 2024. URL <https://arxiv.org/abs/2402.03220>.
- [10] Aymeric Dieuleveut, Nicolas Flammarion, and Francis Bach. Harder, better, faster, stronger convergence rates for least-squares regression. *Journal of Machine Learning Research*, 18(101):1–51, 2017.
- [11] Gerald B Folland. *Real analysis: modern techniques and their applications*. John Wiley & Sons, 1999.
- [12] Rong Ge, Sham M Kakade, Rahul Kidambi, and Praneeth Netrapalli. The step decay schedule: A near optimal, geometrically decaying learning rate procedure for least squares. *Advances in neural information processing systems*, 32, 2019.
- [13] M. Gürbüzbalaban, A. Ozdaglar, and P. A. Parrilo. Why random reshuffling beats stochastic gradient descent. *Mathematical Programming*, 186(1–2):49–84, October 2019. ISSN 1436-4646. doi: 10.1007/s10107-019-01440-w. URL <http://dx.doi.org/10.1007/s10107-019-01440-w>.

- [14] Jeff Z. HaoChen and Suvrit Sra. Random shuffling beats sgd after finite epochs, 2019. URL <https://arxiv.org/abs/1806.10077>.
- [15] Trevor Hastie. The elements of statistical learning: data mining, inference, and prediction, 2009.
- [16] Danny Hernandez, Tom Brown, Tom Conerly, Nova DasSarma, Dawn Drain, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Tom Henighan, Tristan Hume, et al. Scaling laws and interpretability of learning from repeated data. *arXiv preprint arXiv:2205.10487*, 2022.
- [17] Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad, Md Mostofa Ali Patwary, Yang Yang, and Yanqi Zhou. Deep learning scaling is predictable, empirically. *arXiv preprint arXiv:1712.00409*, 2017.
- [18] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- [19] De Huang, Jonathan Niles-Weed, Joel A Tropp, and Rachel Ward. Matrix concentration for products. *Foundations of Computational Mathematics*, 22(6):1767–1799, 2022.
- [20] Marcus Hutter. Learning curve theory. *arXiv preprint arXiv:2102.04074*, 2021.
- [21] Ayush Jain, Andrea Montanari, and Eren Sasoglu. Scaling laws for learning with real and surrogate data, 2024. URL <https://arxiv.org/abs/2402.04376>.
- [22] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [23] Joshua Kazdan, Rylan Schaeffer, Apratim Dey, Matthias Gerstgrasser, Rafael Rafailov, David L Donoho, and Sanmi Koyejo. Collapse or thrive? perils and promises of synthetic data in a self-generating world. *arXiv preprint arXiv:2410.16713*, 2024.
- [24] Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. Deduplicating training data makes language models better. *arXiv preprint arXiv:2107.06499*, 2021.
- [25] Binghui Li, Fengling Chen, Zixun Huang, Lean Wang, and Lei Wu. Functional scaling laws in kernel regression: Loss dynamics and learning rate schedules. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025.
- [26] Houyi Li, Wenzhen Zheng, Jingcheng Hu, Qiufeng Wang, Hanshan Zhang, Zili Wang, Shijie Xuyang, Yuantao Fan, Shuigeng Zhou, Xiangyu Zhang, and Daxin Jiang. Predictable scale: Part i – optimal hyperparameter scaling law in large language model pretraining, 2025. URL <https://arxiv.org/abs/2503.04715>.
- [27] Xuheng Li and Quanquan Gu. Understanding sgd with exponential moving average: A case study in linear regression. *arXiv preprint arXiv:2502.14123*, 2025.



- [28] Junhong Lin and Lorenzo Rosasco. Optimal rates for multi-pass stochastic gradient methods, 2019. URL <https://arxiv.org/abs/1605.08882>.
- [29] Licong Lin, Jingfeng Wu, Sham M Kakade, Peter L Bartlett, and Jason D Lee. Scaling laws in linear regression: Compute, parameters, and data. *arXiv preprint arXiv:2406.08466*, 2024.
- [30] Licong Lin, Jingfeng Wu, and Peter L Bartlett. Improved scaling laws in linear regression via data reuse. *arXiv preprint arXiv:2506.08415*, 2025.
- [31] Kairong Luo, Haodong Wen, Shengding Hu, Zhenbo Sun, Zhiyuan Liu, Maosong Sun, Kaifeng Lyu, and Wenguang Chen. A multi-power law for loss curve prediction across learning rate schedules. *arXiv preprint arXiv:2503.12811*, 2025.
- [32] Alexander Maloney, Daniel A Roberts, and James Sully. A solvable model of neural scaling laws. *arXiv preprint arXiv:2210.16859*, 2022.
- [33] Alexandru Meterez, Depen Morwani, Costin-Andrei Oncescu, Jingfeng Wu, Cengiz Pehlevan, and Sham Kakade. A simplified analysis of sgd for linear regression with weight averaging. *arXiv preprint arXiv:2506.15535*, 2025.
- [34] Konstantin Mishchenko, Ahmed Khaled, and Peter Richtárik. Random reshuffling: Simple analysis with vast improvements, 2021. URL <https://arxiv.org/abs/2006.05988>.
- [35] Niklas Muennighoff, Alexander Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra Piktus, Sampo Pyysalo, Thomas Wolf, and Colin A Raffel. Scaling data-constrained language models. *Advances in Neural Information Processing Systems*, 36:50358–50376, 2023.
- [36] Preetum Nakkiran, Behnam Neyshabur, and Hanie Sedghi. The Deep Bootstrap Framework: Good Online Learners are Good Offline Generalizers. *arXiv preprint arXiv:2010.08127*, 2020.
- [37] Gergely Neu and Lorenzo Rosasco. Iterate averaging as regularization for stochastic gradient descent. In *Conference On Learning Theory*, pages 3222–3242. PMLR, 2018.
- [38] Elliot Paquette, Courtney Paquette, Lechao Xiao, and Jeffrey Pennington. 4+3 phases of compute-optimal neural scaling laws, 2025. URL <https://arxiv.org/abs/2405.15074>.
- [39] Loucas Pillaud-Vivien, Alessandro Rudi, and Francis Bach. Statistical optimality of stochastic gradient descent on hard learning problems through multiple passes, 2018. URL <https://arxiv.org/abs/1805.10074>.
- [40] Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025. URL <https://arxiv.org/abs/2412.15115>.

- [41] Utkarsh Sharma and Jared Kaplan. A neural scaling law from the dimension of the data manifold, 2020. URL <https://arxiv.org/abs/2004.10802>.
- [42] Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. Galactica: A large language model for science. *arXiv preprint arXiv:2211.09085*, 2022.
- [43] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [44] Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge University Press, 2018.
- [45] Pablo Villalobos, Anson Ho, Jaime Sevilla, Tamay Besiroglu, Lennart Heim, and Marius Hobbhahn. Will we run out of data? limits of llm scaling based on human-generated data, 2024.
- [46] Peihao Wang, Rameswar Panda, and Zhangyang Wang. Data efficient neural scaling law via model reusing. In *International Conference on Machine Learning*, pages 36193–36204. PMLR, 2023.
- [47] Alexander Wei, Wei Hu, and Jacob Steinhardt. More than a toy: Random matrix models predict how real-world neural representations generalize. In *International Conference on Machine Learning*, pages 23549–23588. PMLR, 2022.
- [48] Kaiyue Wen, Huaqing Zhang, Hongzhou Lin, and Jingzhao Zhang. From sparse dependence to sparse attention: unveiling how chain-of-thought enhances transformer sample efficiency. *arXiv preprint arXiv:2410.05459*, 2024.
- [49] Jingfeng Wu, Difan Zou, Vladimir Braverman, Quanquan Gu, and Sham M. Kakade. Last iterate risk bounds of sgd with decaying stepsize for overparameterized linear regression, 2022. URL <https://arxiv.org/abs/2110.06198>.
- [50] Jingfeng Wu, Difan Zou, Vladimir Braverman, Quanquan Gu, and Sham M. Kakade. The power and limitation of pretraining-finetuning for linear regression under covariate shift, 2022. URL <https://arxiv.org/abs/2208.01857>.
- [51] Zhangjie Xia, Chi-Hua Wang, and Guang Cheng. Data deletion for linear regression with noisy sgd. *arXiv preprint arXiv:2410.09311*, 2024.
- [52] Fuzhao Xue, Yao Fu, Wangchunshu Zhou, Zangwei Zheng, and Yang You. To repeat or not to repeat: Insights from scaling llm under token-crisis. *Advances in Neural Information Processing Systems*, 36:59304–59322, 2023.
- [53] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin

Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.

- [54] Difan Zou, Jingfeng Wu, Vladimir Braverman, Quanquan Gu, and Sham Kakade. Benign overfitting of constant-stepsizesgd for linear regression. In *Conference on Learning Theory*, pages 4633–4635. PMLR, 2021.
- [55] Difan Zou, Jingfeng Wu, Vladimir Braverman, Quanquan Gu, Dean P. Foster, and Sham M. Kakade. The benefits of implicit regularization from sgd in least squares problems, 2022. URL <https://arxiv.org/abs/2108.04552>.
- [56] Nicolas Zucchet, Francesco d’Angelo, Andrew K Lampinen, and Stephanie CY Chan. The emergence of sparse attention: impact of data distribution and benefits of repetition. *arXiv preprint arXiv:2505.17863*, 2025.

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Preliminaries</b>	<b>3</b>
<b>3</b>	<b>Multi-Epoch Scaling in Strongly Convex Linear Regression</b>	<b>4</b>
3.1	Main Results . . . . .	4
<b>4</b>	<b>A Solvable Case with Zipf-distributed Data</b>	<b>5</b>
4.1	Results on Power-Law Spectrum . . . . .	6
<b>A</b>	<b>Related Work</b>	<b>14</b>
<b>B</b>	<b>Experiments</b>	<b>15</b>
B.1	Simulations in <a href="#">Section 3</a> . . . . .	15
B.2	Simulations in <a href="#">Section 4.1</a> . . . . .	16
B.3	Empirical Verification in Large Language Models . . . . .	16
<b>C</b>	<b>Pretraining Setup</b>	<b>17</b>
<b>D</b>	<b>Results and Simulations for Logarithmic Power-Law Spectrum</b>	<b>18</b>
D.1	Results on Logarithmic Power-law Spectrum . . . . .	18
D.2	Simulations in <a href="#">Section D.1</a> . . . . .	19
<b>E</b>	<b>Notations</b>	<b>20</b>
E.1	Proof Sketch . . . . .	21
<b>F</b>	<b>Proof of Main Results in Strongly Convex Linear Regression</b>	<b>23</b>
F.1	Step I: A Concrete Version of Bias-Variance Decomposition . . . . .	23
F.2	Step II: Risk Approximation and Error Bound Analysis . . . . .	24
F.2.1	Variance Term Analysis: Proof of <a href="#">Lemma F.2</a> . . . . .	25
F.2.2	Bias Term Analysis: Proof of <a href="#">Lemma F.3</a> . . . . .	34
F.3	Step III: Narrowing the Range for Optimal Learning Rate . . . . .	38
F.3.1	A description of the Range of Optimal Learning Rate, Small- $K$ Case . . .	39
F.3.2	A description of the Range of Optimal Learning Rate, Large- $K$ Case . . .	40
F.3.3	An Approximation of the Excess Risk, Small- $K$ Case . . . . .	42
F.3.4	An Approximation of the Excess Risk, Large- $K$ Case . . . . .	43
F.4	Step IV: Deriving the Approximately Optimal Learning Rate, Proof of <a href="#">Lemma E.4</a> . . .	44
F.4.1	Proof of <a href="#">Lemma E.4</a> , small $K$ . . . . .	44
F.4.2	Proof of <a href="#">Lemma E.4</a> , large $K$ . . . . .	45
F.5	Proof of <a href="#">Theorem 3.1</a> . . . . .	45
F.6	Proof of <a href="#">Theorem 3.2</a> . . . . .	47
<b>G</b>	<b>Proof Outline for the Solvable Case with Zipf-distributed Data</b>	<b>51</b>

<b>H</b>	<b>Proof of Main Results for the Solvable Case with Zipf-distributed Data</b>	<b>52</b>
H.1	A Closed Formula for the Excess Risk: Proof of <a href="#">Lemma G.1</a>	52
H.2	Scaling Laws for Power-Law Spectrum: Proof of <a href="#">Theorem 4.1</a>	53
H.2.1	Proof of <a href="#">Theorem 4.1</a> : Small- $K$ Case	53
H.2.2	Proof of <a href="#">Theorem 4.1</a> : Large- $K$ Case	58
H.3	$E(K, N)$ for Power-Law Spectrum: Proof of <a href="#">Theorem 4.2</a>	60
H.3.1	Proof of <a href="#">Theorem 4.2</a> , small- $K$ case	60
H.3.2	Proof of <a href="#">Theorem 4.2</a> , Large- $K$ Case	60
H.4	Scaling Laws for Logarithmic Power-Law Spectrum: Proof of <a href="#">Theorem D.1</a>	61
H.4.1	Proof of <a href="#">Theorem D.1</a> : Small- $K$ Case	61
H.4.2	Proof of <a href="#">Theorem D.1</a> , Large- $K$ case	66
H.5	$E(K, N)$ for Logarithmic Power-Law Spectrum: Proof of <a href="#">Theorem D.2</a>	68
H.5.1	Proof of <a href="#">Theorem D.2</a> , Small- $K$ Case	68
H.5.2	Proof of <a href="#">Theorem D.2</a> , Large- $K$ Case	68
<b>I</b>	<b>Additional Technical lemmas</b>	<b>70</b>

## Appendix A. Related Work

**Data Reuse in LLM Pre-Training.** Empirically, there is a long debate over the effect of data reuse in LLM pre-training. Some works [16, 18, 24, 46] suggested it may be harmful, while some work [42] reported the benefit of data reusing when the number of epochs is small ( $K \leq 4$ ). Xue et al. [52] then discovered a degradation phenomenon in multi-epoch training and investigated relevant factors and regularization methods to tackle it. Muennighoff et al. [35] trained LLMs under different configurations and also found that reusing data is as good as using fresh data in the first few epochs. Yet, as the number of epochs increases, the returns for repetitions diminish. In our work, from a theoretical perspective, we rigorously analyzed the effect of data reuse using non-asymptotic techniques, and we defined and calculated the effective reuse rate under two cases, shedding light on the theoretical understanding of data reusing in LLM pre-training.

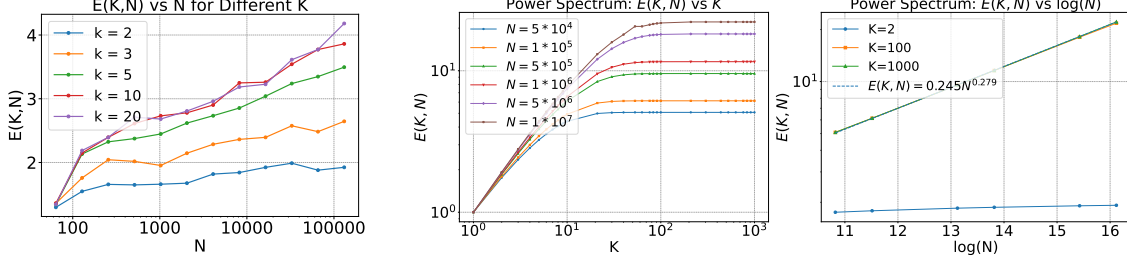
Besides the real-world LLM pre-training regime, many works also reported the improvement of data reusing under synthetic settings empirically [8, 23, 48] or theoretically [2, 9, 56].

**Empirical Findings on Scaling Laws.** Scaling laws characterize how the one-pass training loss of large-scale models scales with factors such as model size, dataset size, and compute budget. They were first identified by Hestness et al. [17] and were later popularized and systematically quantified by Kaplan et al. [22]. These scaling laws are now commonly used to guide the design and training configurations of large language models [18]. Building on this line of work, Muennighoff et al. [35] incorporated the number of training epochs into a more general scaling law, empirically modeling the effect of data reuse. In this work, we provide a multi-epoch scaling law for linear regression.

**Theoretical Explanations for Scaling Laws.** A series of studies have sought to theoretically explain the scaling behaviors for the large language model training. Early works provide theoretical insights in simplified settings, such as regression on a data manifold with fixed intrinsic dimension [41] or memorization of Zipf-distributed data [20]. Maloney et al. [32] and Bahri et al. [4] further showed how scaling laws can arise from random feature models when trained to their global optima, but without analyzing the underlying training dynamics. Another line of work [3, 5, 6, 21, 25, 29, 38, 47] have analyzed scaling laws by tracking the one-pass training dynamics of SGD using a linear model. Among them, Bordelon et al. [5] studied how the loss varies with sample size, model size and training steps. They proved that in the asymptotic regime where any two of these quantities tend to infinity following a proportional limit, the loss exhibits a power-law decay with respect to the remaining one. Lin et al. [29] recovered the power law under an infinite dimension linear regression setup, and explained the mismatch between neural scaling laws and traditional statistical learning theory. Paquette et al. [38] identified 4 phases (+3 subphases) on a random feature model with power-law data, and derived their scaling laws respectively. By comparison, our work studies how multi-epoch training reshapes these scaling laws in a linear regression setup.

**SGD Analysis in Linear Regression.** The analysis of SGD in linear regression has been extensively studied over the years, encompassing both one-pass and multi-epoch SGD. In the context of one-pass SGD, many works conducted convergence analyses under various scenarios [12, 27, 33, 49–51, 55]. For multi-epoch SGD, Lin and Rosasco [28] showed that multiple passes could act as a regularization parameter, and Pillaud-Vivien et al. [39] proved that multiple passes could achieve optimal prediction performance on a class of problems. Nonetheless, all these above works only provide the upper bound (or matching lower bound) of the loss for their convergence analysis,





(a) Strongly convex case in linear regression:  $E(K, N)$  with  $N$ . (b) The solvable case with Zipf-distributed data and power spectrum:  $E(K, N)$  versus  $K$  and  $N$ .

Figure 2: Simulation experiments for strongly-convex linear regression and the solvable case with Zipf-distributed data and power spectrum. Results show that  $E(K, N)$  is approximately proportional to some function of  $N$  when  $N$  is relatively small, and  $E(K, N) \approx K$  when  $N$  is relatively large. For the solvable case with Zipf-distributed data and power spectrum, we also fit the effective reuse rate using the formula  $E(K, N) = c_1 N^{c_2}$  suggested by Theorem 4.2, and the fitted exponent  $c_2 = 0.279 \approx \frac{b}{a-b} = \frac{2}{7}$  matches our theory.

which we cannot apply to giving an accurate characterization of  $E(K, N)$ . In this work, we discuss  $K$ -epoch SGD training in linear regression with shuffling without replacement, under an optimal constant learning rate. We derive tight risk bounds up to a  $1 + o(1)$  order, and used the results to determine  $E(K, N)$ .

**Comparison with Lin et al. [30].** A recent study on linear regression with data reusing [30] is among the most relevant to our results. They showed that when the number of epochs is relatively small (smaller than some power of the dataset size), the order of loss remains the same as one pass SGD for the same iterations, which aligns with our results. However, their results only imply that  $E(K, N) = \Theta(K)$  for small  $K$ , while our analysis directly gives the explicit loss characterization with  $o(1)$  relative error bound and a more exact description of the effective reuse rate both in the small  $K$  and large  $K$ , which reflects the data reusing scaling behavior. Our analysis further highlights the role of  $N$  in the scaling behavior of  $E(K, N)$ , a factor that was not explicitly accounted for in Muennighoff et al. [35].

## Appendix B. Experiments

### B.1. Simulations in Section 3

First, we conduct our experiments on synthetic dataset with a strongly convex linear regression to verify the characterization of effective reuse rate  $E(K, N)$  in Theorem 3.2.

**Experiments Setup.** We generate data pairs  $(x_i, y_i)$  where  $x_i \stackrel{i.i.d}{\sim} \mathcal{N}(0, \mathbf{I}_d)$  with dimension  $d = 100$ . For the label  $y_i$ , we generate it as  $y_i = \langle \mathbf{w}^*, x_i \rangle + \xi_i$ , where  $\mathbf{w}^*$  is the ground truth generate by standard Gaussian with unit variance. Also,  $\xi_i \stackrel{i.i.d}{\sim} \mathcal{N}(0, \sigma^2 \mathbf{I}_d)$ . Here in our simulation, we set  $\sigma$  to 0.1. To make our simulation aligned with the theoretical setup, we set the learning rate  $\eta \propto \frac{\log KN}{KN}$ , and we grid search the ratio  $c := \frac{\eta}{\log KN / KN}$  for the  $c^*$  which minimizes the final loss given training steps  $T = KN$ .

**Results.** As shown in Figure 2(a), we plot  $E(K, N)$  as a function of  $\log N$  for various fixed values of  $K$ . Each curve corresponds to a fixed number of epochs (e.g.,  $K = 3, 5, \dots, 20$ ) and illustrates how the effective reuse rate  $E(K, N)$  grows with dataset size. For small data size ( $\log N \ll K$ ), the effective reuse factor increases roughly linearly with  $\log N$ , indicating that adding more data substantially boosts the one-pass equivalent performance. However, as  $N$  becomes large ( $\log N \gg K$ ), each curve flattens out and approaches an asymptote at  $E(K, N) \approx K$ . In other words, once the dataset is sufficiently large relative to the number of epochs, additional passes through the same data yield no further benefit beyond a factor of  $K$ . This behavior is exactly as predicted by Theorem 3.2: when  $K$  is much smaller than  $\log N$ , we have  $E(K, N) \approx K$  (nearly full  $K$ -fold data reuse), whereas when  $K$  is large relative to  $\log N$ , the effective reuse saturates and grows only on the order of  $\log N$ .

### B.2. Simulations in Section 4.1

We now verify the predictions of Theorem 4.2 using synthetic data generated under the spectral assumptions of Section 4 with a power-law decay Hessian spectrum (Assumption 4.2). In all sub-figures of Figure 2(b), we set the data dimension  $d$  to  $10^5$  and tune all the learning rates to their optimal values. Here we set  $a = 4.5$  and 1.

**Results.** Figure 2(b) plots  $E(K, N)$  versus  $K$  and  $\log N$  for the solvable model with Zipf-distributed data. The curves depicting  $E(K, N)$  versus  $K$  show that  $E(K, N) \approx K$  when  $K$  is relatively small and saturate to some value depending on  $N$  when  $K$  is large. In the right panel, which describes the relationship between  $E(K, N)$  and  $\log N$ , we observe that when  $K$  is small (namely  $K = 2$ ),  $E(K, N)$  increases and approaches  $K$  as  $\log N$  increases, and the plots overlap when  $K$  is large. Those phenomena provide empirical confirmation of the scaling behaviors predicted by Theorem 4.2. We also fit  $E(K, N)$  in the large- $K$  regime with a power-form function as stated in Theorem 4.2. The fitted exponent is  $0.279 \approx \frac{b}{a-b} = \frac{2}{7}$ , aligning with our theory.

### B.3. Empirical Verification in Large Language Models

We conduct experiments on a large language model to empirically validate the hypothesis that larger datasets allow for more effective repetition.

**Experiments Setup.** We perform pretraining runs with fresh data sizes of 0.2B, 0.5B, 0.8B, 1.0B, and 2B tokens, each trained for 100 epochs. As a control, we also include a run with 200B fresh tokens. For each fresh dataset size  $N$  and training epoch  $K$ , we approximate the effective reuse rate  $E(K, N)$  by determining the effective fresh data size  $N_f(K, N)$  required to achieve the same validation loss after one pass through the data.

The effective reuse rate is then computed as:

$$E(K, N) = \frac{N_f(K, N)}{N}.$$

Our experiments utilize a 0.3B parameter model adapted from the Qwen2.5-0.5B architecture [40] and a subset of the DCLM dataset, totaling 200B tokens. A separate subset of the DCLM dataset is reserved for validation. Crucially, we use a constant learning rate schedule across all experiments to align with our theoretical analysis and mitigate the confounding effects of learning rate schedules, as reported in prior work [18, 31]. More details regarding the experiment setup are available in Appendix C.

**Results.** Figure 1(a) depicts the relationship between  $E(K, N)$  and  $K$ , and Figure 1(b) depicts the training curves for different data sizes, and mark the point where  $E(K, N) = 0.8K$ .

**When  $K \leq 4$ ,  $E(K, N) \approx K$ .** Our theoretical analysis indicates that  $E(K, N)$  should be close to  $K$  when  $K$  is small (e.g.,  $K \leq 4$ ). In Figure 1(a), when the epoch number is small (approximately  $\leq 5$ ), we observe that  $E(K, N)$  increases at a rate comparable to the epoch number, as indicated by the black dashed line. The fresh-data equivalent repetition for modest multi-epoch pretraining aligns with the data-constrained scaling laws [35].

**Larger Datasets Allow More Repetition.**  $E(K, N)$  increases with the number of fresh data sizes and eventually saturates for sufficiently large fresh datasets. Our results challenge the data-constrained scaling laws proposed by Muennighoff et al. [35], which assume a uniform effective number of epochs across different fresh data sizes. In Figure 1(a), similar to Figure 4 of Nakkiran et al. [36], which presents alike loss curves for experiments on CIFAR datasets, we show that at the critical points where one-pass training start to outperform multi-epoch training significantly,  $E(K, N)$  increases as  $N$  increases. This suggests the continued potential for scaling pretraining through multi-epoch training with larger datasets.

### Appendix C. Pretraining Setup

In our pretraining experiments, we employ the AdamW optimizer with a weight decay of 0.1 and a gradient clip of 1.0. We set the peak learning rate to 0.001, aligning with the approximate optimal learning rate reported by Li et al. [26]. Balancing the optimal batch size suggested by Li et al. [26] with training efficiency, we utilize a sequence batch size of 128, which corresponds to roughly 0.5M data points per batch. We adopt the vocabulary of Qwen2.5 [40] models. Our pretraining model consists of approximately 117 million non-embedding parameters, consistent with the methodology of Kaplan et al. [22], and a total of 331 million parameters following the convention of Hoffmann et al. [18]. The detailed hyperparameter configurations are presented in Table 2, and the model architecture specifications are provided in Table 1. To ensure a fair comparison by eliminating the influence of batch order variations, we fix the random seed that governs the data stream across all experiments.

Table 1: Model configurations and parameter counts.  $d_h$ : hidden dimension;  $d_f$ : feed-forward dimension;  $n_l$ : number of Transformer layers;  $n_h$ : number of attention heads;  $n_{kv}$ : number of key-value heads (for grouped-query attention); Vocab Size: size of tokenizer vocabulary; #NE params: number of non-embedding parameters (in millions); #Params: total number of model parameters (in millions).

Name	$d_h$	$d_f$	$n_l$	$n_h$	$n_{kv}$	Vocab Size	#NE params	#Params
0.5B	896	4864	24	14	2	151936	355	491
0.3B	640	3328	16	10	2	151936	117	331

Table 2: LLM Experiment Settings

Parameter	Value
<b>Data</b>	
Sequence Batch Size	128
Sequence Length	4096
<b>Learning Rate</b>	
Peak Learning Rate	0.001
Schedule	Constant
Warmup Steps	400
<b>Optimizer</b>	
Optimizer	AdamW
Weight Decay	0.1
$\beta_1$	0.9
$\beta_2$	0.95
$\epsilon$	1e-8
Gradient Clip	1.0

## Appendix D. Results and Simulations for Logarithmic Power-Law Spectrum

### D.1. Results on Logarithmic Power-law Spectrum

Further, we aim to understand under the same Hessian matrix, how the data distribution correlated with  $\mathbf{P}$  and  $\Lambda$  affects the effective reusing rate. By changing the spectrum of  $\Lambda$ , we can also obtain matching upper lower bounds for  $\bar{\mathcal{R}}^*(K, N)$  and a characterization for  $E(K, N)$ , which behave differently from the power-spectrum case.

**Assumption D.1 (Logarithmic Power-Law Spectrum)** *There exist two constants  $a > 1, b > 0$  such that the data input distribution satisfies that  $p_i = ci^{-a} \log^b(i+1)$  and  $\Lambda_i = 1/\log^b(i+1)$ , where  $c = \left(\sum_{i=1}^d i^{-a} \log^b(i+1)\right)^{-1}$ .*

**Theorem D.1** *Consider a  $K$ -epoch SGD over  $N$  fresh data. Under [Assumption 4.1](#), [Assumption D.1](#), and given the data dimension  $d = \Omega((KN)^{\frac{1}{a}})$ , it holds that*

$$\bar{\mathcal{R}}^*(K, N) \asymp \begin{cases} (KN)^{-\frac{a-1}{a}} & \text{for } K = o(\log^b N) \\ (N \log^b N)^{-\frac{a-1}{a}} & \text{for } K = \omega(\log^b N). \end{cases}$$

**Theorem D.2 (Multi-Epoch Scaling Under Logarithmic Power-Law Spectrum)** *Under [Assumption 4.1](#), [Assumption D.1](#), and given the data dimension  $d = \Omega((KN)^{\frac{1}{a}})$  for a one-pass SGD and a  $K$ -epoch SGD over  $N$  fresh data, it holds that*

$$E(K, N) = \begin{cases} K(1 + o(1)) & \text{for } K = o(\log^b N) \\ \Theta(\log^b N) & \text{for } K = \omega(\log^b N). \end{cases}$$

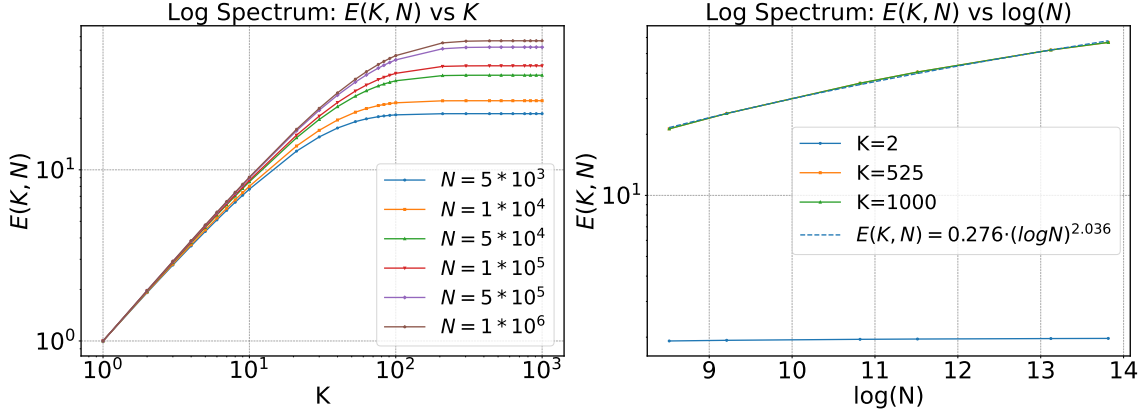


Figure 3: The solvable cases with logarithmic power-law spectrum.  $E(K, N)$  exhibits a similar behavior to that presented in Figure 2. We also fit the effective reuse rate using the formula  $E(K, N) = c_1 (\log N)^{c_2}$  suggested by Theorem 4.2, and the fitted exponent  $c_2 = 2 \approx b = 2$  matches our theory.

**The Saturation Point Varies across Different Problem Setups.** The phase transition point where the effectiveness of data reusing changes from effectively reused to limitedly reused varies across different problem setups. In strongly convex linear regression problems, this phase transition happens when the limit  $\lim_{K \rightarrow \infty} \frac{K}{\log N}$  changes from 0 to  $\infty$ . And in the above power spectrum and log-power spectrum case, the limit turns to be  $\lim_{K \rightarrow \infty} \frac{K}{N^{b/(a-b)}}$  and  $\lim_{K \rightarrow \infty} \frac{K}{\log^b N}$ .

## D.2. Simulations in Section D.1

Now we focus on validating the predictions of Theorem D.2 using synthetic data generated under the spectral assumptions of Section 4 and a log-power decay spectrum (Assumption D.1).

**Experiments Setup.** Similar to Section B.2, in all sub-figures of Figure 3, we set the data dimension  $d$  to  $10^5$  and tune all the learning rates to their optimal values. Here we set  $a = 1.5$  and  $b = 2$ .

**Simulations for the Solvable Model.** Figure 3 plots  $E(K, N)$  versus  $K$  and  $\log N$  for the solvable model. The curves depicting  $E(K, N)$  versus  $K$  and  $E(K, N)$  versus  $\log N$  show trends consistent with those in Section B.2, aligning with Theorem D.2. Furthermore, in the large- $K$  regime, we fit the exponent according to Theorem D.2 and obtain  $2.036 \approx b = 2$ , which provides strong validation of our theory.

## Appendix E. Notations

In this section, we provide the notations appeared in our paper.

**Core notations.** We use  $\|\cdot\|$  to denote the  $\ell_2$ -norm of vectors and the corresponding operator norm of matrices. For two sequences  $(A_n)_{n=0}^\infty$  and  $(B_n)_{n=0}^\infty$ , we write  $A_n = O(B_n)$ , or alternatively  $A_n \lesssim B_n$ ,  $B_n = \Omega(A_n)$ ,  $B_n \gtrsim A_n$ , if there exist constants  $C > 0$ ,  $N > 0$  such that  $|A_n| \leq C|B_n|$  for all  $n \geq N$ . We write  $A_n = \Theta(B_n)$ , or alternatively  $A_n \asymp B_n$ , if both  $A_n = O(B_n)$  and  $A_n = \Omega(B_n)$  hold. Moreover, for some variable  $n$ , we write  $A_n = o_n(B_n)$  if for every constant  $c > 0$ , there exists  $n_0 > 0$  such that  $|A_n| < c|B_n|$  for all  $n \geq n_0$ . In this paper, when  $n$  is clear from the context, we write  $A_n = o(B_n)$  for short. Furthermore, we write  $A_n = \omega(B_n)$  if  $B_n = o(A_n)$ . For matrices  $\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_n$ , we use  $\prod_{l=1}^n \mathbf{A}_l$  to denote the product  $\mathbf{A}_1 \mathbf{A}_2 \dots \mathbf{A}_n$ . Let  $\|\mathbf{u}\|_{\mathbf{S}} = \sqrt{\mathbf{u}^\top \mathbf{S} \mathbf{u}}$  for a vector  $\mathbf{u}$  and a positive semi-definite (*p.s.d*) matrix  $\mathbf{S}$ .

**Key Quantities.** We define the following key quantities to analyze the sequential updates. For each epoch  $k$ , let

$$\mathbf{A}^{(k)} := \prod_{i=N-1}^0 (\mathbf{I} - \eta \mathbf{x}_{\pi_k(i)} \mathbf{x}_{\pi_k(i)}^\top) \quad (2)$$

represent the product of sequential updates through all samples in epoch  $k$ . More generally, we define the partial product operator:

$$\mathbf{Z}_{a \rightarrow b}^{(k)} := \prod_{i=a}^b (\mathbf{I} - \eta \mathbf{x}_{\pi_k(i)} \mathbf{x}_{\pi_k(i)}^\top), \quad \text{with} \quad \mathbf{A}^{(k)} = \mathbf{Z}_{N-1 \rightarrow 0}^{(k)}.$$

We further define that  $\mathbf{Z}_{N-1 \rightarrow N}^{(k)} = \mathbf{I}$ . The cumulative effect across epochs is captured by:

$$\mathbf{T}^{(k)} := \prod_{i=K}^{k+1} \mathbf{A}^{(i)}, \quad \text{and} \quad \mathbf{T}^{(K)} = \mathbf{I}.$$

**Pseudo-expectation Notation  $\tilde{\mathbb{E}}$ .** Because matrix multiplication is non-commutative and the shuffling in training introduces statistical dependence, the expectations of the random matrices defined above cannot be written in a tractable closed form. To approximate the population excess risk, we therefore introduce the auxiliary notation  $\tilde{\mathbb{E}}$ . By construction,  $\tilde{\mathbb{E}}$  computes the expectation of each factor as if the variables were independent, deliberately neglecting the correlations. We then invoke matrix-concentration inequalities to bound the gap between this “pseudo”-expectation and the true expectation of the original dependent random variables. Specifically, for the above random matrices used in our proof, here we further define that

$$\tilde{\mathbb{E}} \mathbf{Z}_{a \rightarrow b}^{(k)} := (\mathbf{I} - \eta \mathbf{H})^{a-b+1}, \quad (3)$$

$$\tilde{\mathbb{E}} \mathbf{A}^{(k)} := (\mathbf{I} - \eta \mathbf{H})^N, \quad (4)$$

$$\tilde{\mathbb{E}} \mathbf{T}^{(k)} := (\mathbf{I} - \eta \mathbf{H})^{N(K-k)}, \quad (5)$$

$$\tilde{\mathbb{E}} \mathbf{S}_l^{(ij)} := \tilde{\mathbb{E}} \left[ \mathbf{Z}_{N-1 \rightarrow \pi_i^{-1}(l)+1}^{(i)} \right] \mathbb{E} [\mathbf{x}_l \mathbf{x}_l^\top] \tilde{\mathbb{E}} \left[ \mathbf{Z}_{\pi_j^{-1}(l)+1 \rightarrow N-1}^{(j)} \right]. \quad (6)$$



### E.1. Proof Sketch

We now provide a proof sketch of our main results. First, we need to compute the optimal expected excess risk  $\bar{\mathcal{R}}^*(K, N)$ . This requires us to compute  $\bar{\mathcal{R}}(K, N; \eta)$  and then select the optimal learning rate  $\eta^*$  that minimizes  $\bar{\mathcal{R}}(K, N; \eta)$ . However, due to the random shuffling and multi-pass processing of the training data, directly analyzing  $\bar{\mathcal{R}}(K, N; \eta)$  is intractable. To overcome this, we seek an analytic approximation of  $\bar{\mathcal{R}}(K, N; \eta)$ , which is derived through the following steps.

**Step 1: Bias-Variance Decomposition for Training Dynamics.** Following the widely-applied bias-variance decomposition approach to analyzing the dynamics of SGD training [12, 37, 49, 54], we define  $\theta_t = \mathbf{w}_t - \mathbf{w}^*$  and examine the following two processes of bias and variance:

$$\theta_{t+1}^{\text{bias}} = \theta_t^{\text{bias}} - \eta \langle \theta_t^{\text{bias}}, \mathbf{x}_{j_t} \rangle \mathbf{x}_{j_t}, \quad \theta_{t+1}^{\text{var}} = \theta_t^{\text{var}} - \eta \langle \theta_t^{\text{var}}, \mathbf{x}_{j_t} \rangle \mathbf{x}_{j_t} + \eta \xi_{j_t} \mathbf{x}_{j_t},$$

where two processes are initialized as  $\theta_0^{\text{bias}} = \mathbf{w}_0 - \mathbf{w}^*$  and  $\theta_0^{\text{var}} = \mathbf{0}$ . It follows that  $\theta_t = \theta_t^{\text{bias}} + \theta_t^{\text{var}}$ , with  $\mathbb{E}[\theta_t^{\text{var}}] = \mathbf{0}$ . We can then decompose the excess risk  $\mathcal{R}(\mathbf{w}_t)$  into two components: the *bias term* and the *variance term*, which we formalize as follows:

$$\bar{\mathcal{R}}(K, N; \eta) = \mathbb{E}_{\mathbf{w}_t \sim \mathcal{W}_{K, N, \eta}} \frac{1}{2} \|\theta_t\|_{\mathbf{H}}^2 = \mathbb{E}_{\mathbf{w}_t \sim \mathcal{W}_{K, N, \eta}} \frac{1}{2} \|\theta_t^{\text{bias}}\|_{\mathbf{H}}^2 + \mathbb{E}_{\mathbf{w}_t \sim \mathcal{W}_{K, N, \eta}} \frac{1}{2} \|\theta_t^{\text{var}}\|_{\mathbf{H}}^2.$$

**Step 2: Analytic Risk Approximation by Matrix Concentration.** A key challenge in tracking the dynamics of multi-epoch SGD training arises from the non-commutative nature of the matrices in the weight updates, which depend on randomly shuffled and multi-pass data. For example, the bias weight evolves as

$$\theta_{KN}^{\text{bias}} = \left( \prod_{k=1}^K \left( \prod_{l=1}^N \left( \mathbf{I} - \eta \mathbf{x}_{\pi_k(l)} \mathbf{x}_{\pi_k(l)}^\top \right) \right) \right) \theta_0^{\text{bias}},$$

where we can see that one data point appears more than once across different epochs. Thus, the above matrix multiplication involves massive correlated data, which makes calculating the bias term  $\mathbb{E} \left[ \|\theta_{KN}^{\text{bias}}\|_{\mathbf{H}}^2 \right]$  intractable. To resolve this issue, we borrow tools from concentration inequalities for matrix products Huang et al. [19]. Specifically, we use the following result:

**Lemma E.1 (Corollary of Theorem 7.1 in Huang et al. [19])** *Given  $n$  data points such that  $\mathbf{z}_0, \dots, \mathbf{z}_{n-1} \stackrel{i.i.d}{\sim} \mathcal{N}(0, \mathbf{H})$ , and defining  $\mathbf{A} = \prod_{j=0}^{n-1} \left( \mathbf{I} - \eta \mathbf{z}_j \mathbf{z}_j^\top \right)$ , we have  $\mathbb{E} \|\mathbf{A} - \mathbb{E} \mathbf{A}\|^l \leq \left( \sqrt{\delta_A \eta^2 n l} \right)^l$ , where  $\delta_A := \tilde{C} 8e D^4 \log d$  for some absolute constant  $\tilde{C} > 0$ .*

However, several obstacles prevent us from directly applying Lemma E.1 to our problem. For example, we actually need to control error terms like  $\mathbb{E} \left\| \prod_{i=K}^{K+1} \mathbf{A}^{(i)} - (\mathbb{E} \mathbf{A})^l \right\|$ , where  $\mathbf{A}^{(i)}$  represents the product of sequential updates through all samples in epoch  $i$  (see the formal definition in Equation (2), Appendix E). To address this, our main idea is to derive a tight upper bound for the original term, and decompose this upper bound into the sum of a series of sub-terms for which we can apply Lemma E.1. (see the detailed derivation in Appendix F.2.1 and Appendix F.2.2)

Finally, we derive an error bound on matrix deviations based on our calculations, which is a higher-order infinitesimal of the main term when  $\eta \in \left[ \Omega(T^{-1}), o(T^{-\frac{3}{4}}) \right]$  and  $K = o\left(\eta^{-1} T^{-\frac{3}{4}}\right)$ ,

with  $T := KN$  denoting the total number of training steps. This provides a theoretical guarantee for us to approximate the risk function with a tractable expression. For the bias term, we have

$$\begin{aligned}\mathbb{E} \left[ \left\| \boldsymbol{\theta}_{KN}^{\text{bias}} \right\|_{\mathbf{H}}^2 \right] &= \mathbb{E} \left[ \left\| \left( \prod_{k=1}^K \left( \prod_{l=0}^{N-1} \left( \mathbf{I} - \eta \mathbf{x}_{\pi_k(l)} \mathbf{x}_{\pi_k(l)}^\top \right) \right) \right) \boldsymbol{\theta}_0 \right\|_{\mathbf{H}}^2 \right] \\ &\approx \left\| \left( \prod_{k=1}^K \mathbb{E} \left[ \prod_{l=0}^{N-1} \left( \mathbf{I} - \eta \mathbf{x}_{\pi_k(l)} \mathbf{x}_{\pi_k(l)}^\top \right) \right] \right) \boldsymbol{\theta}_0 \right\|_{\mathbf{H}}^2 \\ &= \left\| (\mathbf{I} - \eta \mathbf{H})^{KN} \boldsymbol{\theta}_0 \right\|_{\mathbf{H}}^2,\end{aligned}$$

where the approximation step follows from [Lemma E.1](#), and the last equation follows the facts that  $\mathbb{E} \left[ \mathbf{x}_{\pi_k(l)} \mathbf{x}_{\pi_k(l)}^\top \right] = \mathbf{H}$  and  $\mathbf{x}_i$  is uncorrelated with  $\mathbf{x}_j$  for  $i \neq j$ . For the variance term, the data correlation issue is similar to what we met in the bias term case. Again, leveraging [Lemma E.1](#) and following a similar analysis, we can get an approximation formula for the variance term as shown:

$$\begin{aligned}\mathbb{E} \left[ \left\| \boldsymbol{\theta}_{KN}^{\text{var}} \right\|_{\mathbf{H}}^2 \right] &\approx \frac{2\sigma^2}{N} \text{tr} \left( \frac{(\mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^{KN}) ((\mathbf{I} - \eta \mathbf{H})^N - (\mathbf{I} - \eta \mathbf{H})^{KN})}{\mathbf{I} + (\mathbf{I} - \eta \mathbf{H})^N} \right) \\ &\quad + \eta \sigma^2 \langle \mathbf{H}, (\mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^{2KN}) (2\mathbf{I} - \eta \mathbf{H})^{-1} \rangle.\end{aligned}$$

**Step 3: Narrowing the Range for Optimal Learning Rate.** However, despite we have an analytic approximation for risk, it is important to note that this approximation holds only for a specific range of parameters. For a detailed discussion, refer to [Lemma F.1](#).

To mitigate this, we first determine a reasonable range for the optimal learning rate in two steps: First, we choose  $\tilde{\eta} = \frac{\log KN}{2\lambda_d KN}$  as a reference learning rate; Then, by comparing the losses for the reference learning rate and other candidate learning rates, we can eliminate a large range of values.

This analysis helps narrow down the potential range of learning rates ([Lemma F.5](#) for small  $K$  and [Lemma F.6](#) for large  $K$ ). Within this range, we further simplify the risk approximation to make it more tractable for optimization, as shown in the following lemmas:

**Lemma E.2 (Small  $K$ )** Let  $\mathbf{H} = \mathbf{PDP}^\top$  be the canonical form of  $\mathbf{H}$  under similarity, and let  $\tilde{\theta}_d^2 := \sum_{l=d-n_d+1}^d (\mathbf{P}\boldsymbol{\theta}_0)_l^2$ . Under [Assumption 3.1](#) and [3.3](#), for learning rate  $\eta \in \left[ \frac{\log KN}{3\lambda_d KN}, \frac{D^2 \text{tr}(\mathbf{H}) \log KN}{\lambda_d \text{tr}(\mathbf{H}^2) KN} \right]$ ,  $K = o(\log N)$ , we have  $\bar{\mathcal{R}}(K, N; \eta) = M(K, N; \eta)(1 + o(1))$  with

$$M(K, N; \eta) := \frac{1}{2} \tilde{\theta}_d^2 \lambda_d \exp(-2\lambda_d \eta KN) + \frac{\eta \text{tr}(\mathbf{H}) \sigma^2}{4}.$$

**Lemma E.3 (Large  $K$ )** We define  $\tilde{\theta}_d^2$  as the same as [Lemma E.2](#). Under [Assumption 3.1](#) and [3.3](#), for learning rate  $\eta \in \left[ \frac{\log KN}{3\lambda_d KN}, o\left(\frac{1}{N}\right) \right]$  and  $K = \omega(\log N)$ , we have  $\bar{\mathcal{R}}(K, N; \eta) = M(K, N; \eta)(1 + o(1))$  with

$$M(K, N; \eta) = \frac{1}{2} \tilde{\theta}_d^2 \lambda_d \exp(-2\lambda_d \eta KN) + \frac{\eta \text{tr}(\mathbf{H}) \sigma^2}{4} + \frac{\sigma^2 d}{2N}.$$

**Step 4: Deriving the Approximately Optimal Learning Rate.** At this point, we have narrowed down the range for the optimal learning rate and simplified the risk approximation. The next step is to approximate the optimal expected excess risk. To achieve this, we differentiate the simplified risk function  $M(K, N; \eta)$  in [Lemma E.2](#) and [Lemma E.3](#) with respect to the learning rate  $\eta$  and give the critical point  $\eta = \eta'(K, N)$ , which are presented as follows:

**Lemma E.4 (Approximately Optimal Learning Rate)** Under [Assumption 3.1](#) and [3.3](#), we consider  $K$ -epoch SGD with  $N$  fresh data and learning rate  $\eta = \eta'(K, N) = \frac{\log \rho KN}{2\lambda_d KN}$ , where  $\rho := \frac{4\hat{\theta}_d^2 \lambda_d}{\text{tr}(\mathbf{H})\sigma^2}$ . Then it holds for  $K = o(\log N)$  or  $K = \omega(\log N)$  that  $\bar{\mathcal{R}}(K, N; \eta'(K, N)) = \bar{\mathcal{R}}^*(K, N) (1 + o(1))$ .

Using Lemma E.4, we complete the proof as follows. By evaluating the risk at the approximately optimal learning rate  $\eta'(K, N) = \frac{\log \rho KN}{2\lambda_d KN}$ , we obtain an approximation of the optimal risk ([Theorem 3.1](#)), based on which we derive the effective reuse rate ([Theorem 3.2](#)).

## Appendix F. Proof of Main Results in Strongly Convex Linear Regression

In this section, we present our proof of the main results in [Section 3](#). [Section 3](#) centres on [Theorem 3.2](#), which establishes a scaling law for the effective reuse rate  $E(K, N)$  in terms of the relative magnitudes of number of epochs  $K$  and dataset size  $N$ . Its proof unfolds in four main stages, as organized in [Section E.1](#). And [Lemma F.1](#) serves as a middle product when we derive [Theorem 3.2](#).

### F.1. Step I: A Concrete Version of Bias-Variance Decomposition

Before we begin our proof, we first provide the formal version of the loss estimate for a specific range of learning rate parameters. We define a  $\hat{\mathcal{R}}(K, N, \eta)$  as the estimator of  $\bar{\mathcal{R}}(K, N; \eta)$

$$\hat{\mathcal{R}}(K, N; \eta) := \underbrace{\hat{\mathcal{R}}_1(K, N; \eta)}_{\text{bias term}} + \underbrace{\hat{\mathcal{R}}_2(K, N; \eta)}_{\text{var term across epochs}} + \underbrace{\hat{\mathcal{R}}_3(K, N; \eta)}_{\text{var term within epoch}},$$

where

$$\hat{\mathcal{R}}_1(K, N; \eta) := \frac{1}{2}(\mathbf{w}_0 - \mathbf{w}^*)^\top (\mathbf{I} - \eta \mathbf{H})^{2KN} \mathbf{H}(\mathbf{w}_0 - \mathbf{w}^*),$$

$$\hat{\mathcal{R}}_2(K, N; \eta) := \frac{\sigma^2}{N} \text{tr} \left( \frac{(\mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^{KN}) ((\mathbf{I} - \eta \mathbf{H})^N - (\mathbf{I} - \eta \mathbf{H})^{KN})}{\mathbf{I} + (\mathbf{I} - \eta \mathbf{H})^N} \right),$$

$$\hat{\mathcal{R}}_3(K, N; \eta) := \frac{\eta \sigma^2}{2} \langle \mathbf{H}, (\mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^{2KN})(2\mathbf{I} - \eta \mathbf{H})^{-1} \rangle.$$

Then, we provide the following approximation lemma.

**Lemma F.1** Under [Assumption 3.1](#) and [3.3](#), we further assume that for every  $\mathbf{x}$  in the training set,  $\|\mathbf{x}\| \leq D$  for some constant  $D > 0$ . Consider a  $K$ -epoch SGD with learning rate  $\eta \in \left[ \Omega\left(\frac{1}{T}\right), o(T^{-\frac{3}{4}}) \right]$ ,  $K = o\left(\eta^{-1}T^{-\frac{3}{4}}\right)$  and data shuffling. Then, after  $T = KN$  steps, the estimator of the expected excess risk satisfies:

$$\bar{\mathcal{R}}(K, N; \eta) = \hat{\mathcal{R}}(K, N; \eta) (1 + o(1)).$$

## F.2. Step II: Risk Approximation and Error Bound Analysis

In this section, we rigorously formulate the analytic risk approximation in [Lemma F.1](#) and provide its proof. [Lemma F.1](#) indicates that the error bound is of higher order than the main term when the parameters are restricted to a limited range of values.

Recall from [Section E.1](#) that the risk  $\bar{\mathcal{R}}(K, N; \eta)$  can be decomposed into the *bias term*  $\bar{\mathcal{R}}^{\text{bias}}(K, N; \eta) := \frac{1}{2} \|\boldsymbol{\theta}_t^{\text{bias}}\|_{\mathbf{H}}^2$  and *variance term*  $\bar{\mathcal{R}}^{\text{var}}(K, N; \eta) := \frac{1}{2} \|\boldsymbol{\theta}_t^{\text{var}}\|_{\mathbf{H}}^2$ . This implies that [Lemma F.1](#) is a direct corollary of the following two lemmas:

**Lemma F.2 (Variance Term)** *Suppose that [Assumption 3.1](#) holds. Then for a  $K$ -epoch SGD with dataset size  $N$  and learning rate  $\eta \in [\Omega(\frac{1}{T}), o(\frac{1}{T^{\frac{1}{2}}})]$  and shuffling, when  $\text{poly}(T) \gtrsim d$ , we have the estimator of the variance term  $\tilde{\mathcal{R}}^{\text{var}}(K, N; \eta) := \mathbb{E}_{\mathbf{w} \sim \mathcal{W}_{K, N, \eta}} [\mathcal{R}(\mathbf{w})^{\text{var}}]$  after  $T := KN$  steps*

$$\begin{aligned} \tilde{\mathcal{R}}^{\text{var}}(K, N; \eta) &:= \frac{\sigma^2}{N} \text{tr} \left( \frac{(\mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^{KN}) ((\mathbf{I} - \eta \mathbf{H})^N - (\mathbf{I} - \eta \mathbf{H})^{KN})}{\mathbf{I} + (\mathbf{I} - \eta \mathbf{H})^N} \right) \\ &\quad + \frac{\eta \sigma^2}{2} \langle \mathbf{H}, (\mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^{2KN}) (2\mathbf{I} - \eta \mathbf{H})^{-1} \rangle, \end{aligned}$$

where the expectation is taken on the training set and shuffle, and the estimate error is

$$\left| \tilde{\mathcal{R}}^{\text{var}}(K, N; \eta) - \bar{\mathcal{R}}^{\text{var}}(K, N; \eta) \right| = O(\eta^3 T^{\frac{3}{2}} K^2 \sqrt{\log d}).$$

when  $K \leq \frac{\log 2}{\eta \sqrt{\tilde{C} 8e D^4 T \log d}}$ .

**Lemma F.3 (Bias Term)** *Under [Assumption 3.1](#), for a  $K$ -epoch SGD with dataset size  $N$ , learning rate  $\eta$  and shuffling, when  $\text{poly}(T) \gtrsim d$ , we have the estimator of the bias term  $\tilde{\mathcal{R}}^{\text{bias}}(K, N; \eta) := \mathbb{E}_{\mathbf{w} \sim \mathcal{W}_{K, N, \eta}} [\mathcal{R}(\mathbf{w})^{\text{bias}}]$  after  $T := KN$  steps*

$$\tilde{\mathcal{R}}^{\text{bias}}(K, N; \eta) := \frac{1}{2} (\mathbf{w}_0 - \mathbf{w}^*)^\top (\mathbf{I} - \eta \mathbf{H})^{2KN} \mathbf{H} (\mathbf{w}_0 - \mathbf{w}^*).$$

Then we have the following estimate errors:

1. When  $K \geq 2$  and  $K = o\left(\frac{N^{\frac{1}{5}}}{(\log N)^{\frac{6}{5}}}\right)$ :

(a) When  $\eta \leq \frac{2 \log T}{3 \lambda_d T}$ , the estimate distance is given by

$$\left| \tilde{\mathcal{R}}^{\text{bias}}(K, N; \eta) - \bar{\mathcal{R}}^{\text{bias}}(K, N; \eta) \right| = O\left((1 - \eta \lambda_d)^{N(2K-1)} K \sqrt{\eta^2 KN}\right).$$

(b) When  $\eta \geq \frac{2 \log T}{3 \lambda_d T}$ , the estimate distance is given by

$$\left| \tilde{\mathcal{R}}^{\text{bias}}(K, N; \eta) - \bar{\mathcal{R}}^{\text{bias}}(K, N; \eta) \right| = O\left(\frac{1}{T^{\frac{4}{3}}}\right).$$

2. When  $K = 1$ :

$$\left| \tilde{\mathcal{R}}^{\text{bias}}(1, T; \eta) - \bar{\mathcal{R}}^{\text{bias}}(1, T; \eta) \right| = O\left(\eta^2 T e^{-2\lambda_d \eta T}\right).$$

F.2.1. VARIANCE TERM ANALYSIS: PROOF OF [LEMMA F.2](#)

We first recall some notations Appendix E that  $\mathbf{Z}_{a \rightarrow b}^{(k)} = \prod_{i=a}^b (\mathbf{I} - \eta \mathbf{x}_{\pi_k(i)} \mathbf{x}_{\pi_k(i)}^\top)$ ,  $\mathbf{Z}_{N-1 \rightarrow N}^{(k)} = \mathbf{I}$ ,  $\mathbf{b}^{(k)} = \sum_{l=0}^{N-1} \mathbf{Z}_{N-1 \rightarrow l+1}^{(k)} \xi_{\pi_k(l)} \mathbf{x}_{\pi_k(l)}$ ,  $\mathbf{A}^{(k)} = \mathbf{Z}_{N-1 \rightarrow 0}^{(k)}$ ,  $\mathbf{T}^{(k)} = \prod_{i=K}^{k+1} \mathbf{A}^{(i)}$ , and  $\mathbf{T}^{(K)} = \mathbf{I}$ . For simplicity, and if it does not cause confusion, we omit the superscript “var” in all the training parameters  $\boldsymbol{\theta}^{\text{var}}$  in the proof of [Lemma F.2](#). Now we derive the recursion before and after the  $k$ -th epoch.

$$\begin{aligned} \boldsymbol{\theta}_{kN} &= (\mathbf{I} - \eta \mathbf{x}_{\pi_k(N-1)} \mathbf{x}_{\pi_k(N-1)}^\top) \boldsymbol{\theta}_{(k-1)N} + \eta \xi_{\pi_k(N-1)} \mathbf{x}_{\pi_k(N-1)} \\ &= \eta \sum_{l=0}^{N-1} \mathbf{Z}_{N-1 \rightarrow l+1}^{(k)} \xi_{\pi_k(l)} \mathbf{x}_{\pi_k(l)} + \mathbf{A}^{(k)} \boldsymbol{\theta}_{(k-1)N} \\ &= \eta \mathbf{b}^{(k)} + \mathbf{A}^{(k)} \boldsymbol{\theta}_{(k-1)N}, \end{aligned}$$

where  $\pi_k(i)$  is the  $i$ -th index after the permutation  $\pi_k$  in the  $K$ -th epoch. Further writing out the above recursion gives the parameter after  $K$  epochs

$$\boldsymbol{\theta}_{KN} = \eta \sum_{k=1}^K \mathbf{A}^{(K)} \dots \mathbf{A}^{(k+1)} \mathbf{b}^{(k)}.$$

A natural move here is to replace  $\boldsymbol{\theta}_{KN}$  with the expression above in the variance term

$$\begin{aligned} \bar{\mathcal{R}}^{\text{var}}(K, N; \eta) &= \mathbb{E} \frac{1}{2} \boldsymbol{\theta}_{KN}^\top \mathbf{H} \boldsymbol{\theta}_{KN} = \mathbb{E} \frac{1}{2} \langle \mathbf{H}, \boldsymbol{\theta}_{KN} \boldsymbol{\theta}_{KN}^\top \rangle \\ &= \frac{\eta^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{\pi_1 \dots \pi_K} \sum_{i,j=1}^K \mathbf{T}^{(i)} \mathbf{b}^{(i)} (\mathbf{b}^{(j)})^\top (\mathbf{T}^{(j)})^\top \right\rangle \\ &= \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{\pi_1 \dots \pi_K} \sum_{i,j=1}^K \mathbf{T}^{(i)} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ij)} \right) (\mathbf{T}^{(j)})^\top \right\rangle \\ &= \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{\substack{i \neq j \\ i,j=1}}^K \sum_{\substack{\pi_1 \dots \pi_K \\ \text{except } \pi_i, \pi_j}} \mathbf{T}^{(i)} \sum_{\pi_i, \pi_j} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ij)} \right) (\mathbf{T}^{(j)})^\top \right\rangle \\ &\quad + \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{i=1}^K \sum_{\substack{\pi_1 \dots \pi_K \\ \text{except } \pi_i}} \mathbf{T}^{(i)} \sum_{\pi_i} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ii)} \right) (\mathbf{T}^{(i)})^\top \right\rangle. \end{aligned} \tag{7}$$

where in the third equation, we take expectations with respect to the label noise  $(\xi_l)_{l=0}^{N-1}$ , and in the last equation, we decompose the variance term into two parts, according to whether the  $\mathbf{b}^{(i)}$  and  $\mathbf{b}^{(j)}$  are from the same epoch or not.

After explicitly writing the variance term, and to get a close-form formula for it, we then take pseudo expectations of  $\mathbf{T}^{(i)}$ ,  $\mathbf{T}^{(j)}$ ,  $\mathbf{S}_l^{(ii)}$ , and  $\mathbf{S}_l^{(ij)}$  separately to get the approximation of

$\bar{\mathcal{R}}^{\text{var}}(K, N; \eta)$ , given as follows:

$$\begin{aligned} \tilde{\mathcal{R}}^{\text{var}}(K, N; \eta) &:= \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^2} \sum_{\substack{i \neq j \\ i, j=1}}^K \tilde{\mathbb{E}} \mathbf{T}^{(i)} \left( \sum_{l=0}^{N-1} \sum_{\pi_i, \pi_j} \tilde{\mathbb{E}} \mathbf{S}_l^{(ij)} \right) \tilde{\mathbb{E}} \mathbf{T}^{(i)} \right\rangle \\ &\quad + \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{N!} \sum_{i=1}^K \tilde{\mathbb{E}} \mathbf{T}^{(i)} \left( \sum_{l=0}^{N-1} \sum_{\pi_i} \tilde{\mathbb{E}} \mathbf{S}_l^{(ii)} \right) \tilde{\mathbb{E}} \mathbf{T}^{(i)} \right\rangle. \end{aligned}$$

The intuition of the ‘‘pseudo expectation’’ and the related definitions are in Appendix E. Fix  $l$ , notice that when  $i \neq j$ , by Equation (6),

$$\begin{aligned} \sum_{\pi_i, \pi_j} \tilde{\mathbb{E}} \mathbf{S}_l^{(ij)} &:= \sum_{\pi_i, \pi_j} \tilde{\mathbb{E}} \left[ \mathbf{Z}_{N-1 \rightarrow \pi_i^{-1}(l)+1}^{(i)} \mathbf{x}_l \mathbf{x}_l^\top \mathbf{Z}_{\pi_j^{-1}(l)+1 \rightarrow N-1}^{(j)} \right] \\ &:= \sum_{\pi_i, \pi_j} (\mathbf{I} - \eta \mathbf{H})^{N-1-\pi_i^{-1}(l)} \mathbf{H} (\mathbf{I} - \eta \mathbf{H})^{N-1-\pi_j^{-1}(l)}. \end{aligned}$$

For a fixed  $i$ , for all  $m \in [0, N-1]$ , there are  $(N-1)!$  permutations  $\pi_i$  that satisfies  $\pi_i(m) = l$ . So

$$\sum_{\pi_i, \pi_j} \tilde{\mathbb{E}} \mathbf{S}_l^{(ij)} = ((N-1)!)^2 \sum_{m, n=0}^{N-1} (\mathbf{I} - \eta \mathbf{H})^{N-1-m} \mathbf{H} (\mathbf{I} - \eta \mathbf{H})^{N-1-n}. \quad (8)$$

By applying a similar derivation to the  $i = j$  case, we obtain that

$$\sum_{\pi_i} \tilde{\mathbb{E}} \mathbf{S}_l^{(ii)} = (N-1)! \sum_{m=0}^{N-1} (\mathbf{I} - \eta \mathbf{H})^{N-1-m} \mathbf{H} (\mathbf{I} - \eta \mathbf{H})^{N-1-m}. \quad (9)$$

Plugging Equation (8) and Equation (9) into the expression of  $\tilde{\mathcal{R}}^{\text{var}}(K, N; \eta)$ , and we have

$$\begin{aligned} &\tilde{\mathcal{R}}^{\text{var}}(K, N; \eta) \\ &= \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{N^2} \sum_{\substack{i \neq j \\ i, j=1}}^K \tilde{\mathbb{E}} \mathbf{T}^{(i)} \left( \sum_{l=0}^{N-1} \sum_{m, n=0}^{N-1} (\mathbf{I} - \eta \mathbf{H})^{2N-2-m-n} \mathbf{H} \right) \tilde{\mathbb{E}} \mathbf{T}^{(i)} \right\rangle \\ &\quad + \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{N} \sum_{i=1}^K \tilde{\mathbb{E}} \mathbf{T}^{(i)} \left( \sum_{l=0}^{N-1} \sum_{m=0}^{N-1} (\mathbf{I} - \eta \mathbf{H})^{2N-2-2m} \mathbf{H} \right) \tilde{\mathbb{E}} \mathbf{T}^{(i)} \right\rangle \\ &= \frac{\sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{N} \sum_{\substack{i \neq j \\ i, j=1}}^K (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \underbrace{\left( \mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^N \right)^2 \mathbf{H}^{-1} (\mathbf{I} - \eta \mathbf{H})^{N(K-j)}}_{:= \Psi_1} \right\rangle \\ &\quad + \frac{\eta \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \underbrace{\sum_{i=1}^K (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \left( \mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^{2N} \right) (2\mathbf{I} - \eta \mathbf{H})^{-1} (\mathbf{I} - \eta \mathbf{H})^{N(K-i)}}_{\Psi_2} \right\rangle. \end{aligned}$$



where the second equation uses Equation (5). The quantity  $\Psi_1$  accounts for the variance term across different epochs and  $\Psi$ . Then we calculate  $\Psi_1$  and  $\Psi_2$  separately. For  $\Psi_1$ , we have

$$\begin{aligned}\Psi_1 &= \frac{\sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{N} \sum_{i,j=1}^K (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \left( \mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^N \right)^2 \mathbf{H}^{-1} (\mathbf{I} - \eta \mathbf{H})^{N(K-j)} \right\rangle \\ &\quad - \frac{\sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{N} \sum_{i=1}^K (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \left( \mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^N \right)^2 \mathbf{H}^{-1} (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \right\rangle \\ &= \frac{\sigma^2}{2N} \text{tr} \left( \left( \mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^{KN} \right)^2 \right) \\ &\quad - \frac{\sigma^2}{2N} \text{tr} \left( \left( \mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^N \right)^2 \left( \mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^{2N} \right)^{-1} \left( \mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^{2KN} \right) \right) \\ &= \frac{\sigma^2}{N} \text{tr} \left( \left( \mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^{KN} \right) \left( \mathbf{I} + (\mathbf{I} - \eta \mathbf{H})^N \right)^{-1} \left( (\mathbf{I} - \eta \mathbf{H})^N - (\mathbf{I} - \eta \mathbf{H})^{KN} \right) \right).\end{aligned}$$

The last equation is obtained by direct algebraic calculation. For  $\Psi_2$ , by direct matrix calculation, we get

$$\Psi_2 = \frac{\eta \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, (2\mathbf{I} - \eta \mathbf{H})^{-1} \left( \mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^{2KN} \right) \right\rangle.$$

Next we obtain the error bound for  $\left| \tilde{\mathcal{R}}^{\text{var}}(K, N; \eta) - \bar{\mathcal{R}}^{\text{var}}(K, N; \eta) \right|$ , which can be represented as

$$\begin{aligned}& \left| \tilde{\mathcal{R}}^{\text{var}}(K, N; \eta) - \bar{\mathcal{R}}^{\text{var}}(K, N; \eta) \right| \\ & \leq \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{\substack{i \neq j \\ i,j=1}}^K \sum_{\substack{\pi_1 \cdots \pi_K \\ \text{except } \pi_i, \pi_j}} \mathbf{T}^{(i)} \sum_{\pi_i, \pi_j} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ij)} \right) \left( \mathbf{T}^{(j)} \right)^\top \right\rangle \right. \\ & \quad \left. - \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^2} \sum_{\substack{i \neq j \\ i,j=1}}^K (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \left( \sum_{l=0}^{N-1} \sum_{\pi_i, \pi_j} \tilde{\mathbb{E}} \mathbf{S}_l^{(ij)} \right) (\mathbf{I} - \eta \mathbf{H})^{N(K-j)} \right\rangle \right| =: I_1 \\ & \quad + \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{i=1}^K \sum_{\substack{\pi_1 \cdots \pi_K \\ \text{except } \pi_i}} \mathbf{T}^{(i)} \sum_{\pi_i} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ii)} \right) \left( \mathbf{T}^{(i)} \right)^\top \right\rangle \right. \\ & \quad \left. - \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{N!} \sum_{i=1}^K (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \left( \sum_{l=0}^{N-1} \sum_{\pi_i} \tilde{\mathbb{E}} \mathbf{S}_l^{(ii)} \right) (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \right\rangle \right| =: I_2,\end{aligned}$$

where the first inequality uses the triangle inequality. The term  $I_1$  represents the error term between epochs, and  $I_2$  represents the error term within one epoch. We will bound  $I_1$  and  $I_2$  separately in the proof.

**Upper bound for  $I_1$ .** To bound  $I_1$ , a natural move here is to plug in a term that takes pseudo expectation over  $(\mathbf{T}^{(i)})_{i=1}^K$  but does not take pseudo expectation over  $(\mathbf{S}_l^{(ij)})_{l,i,j}$ , and divide  $I_1$  into two terms.

$$\begin{aligned}
 I_1 &\leq \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{\substack{i \neq j \\ i,j=1}}^K \sum_{\substack{\pi_1 \dots \pi_K \\ \text{except } \pi_i, \pi_j}} \mathbf{T}^{(i)} \sum_{\pi_i, \pi_j} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ij)} \right) (\mathbf{T}^{(j)})^\top \right\rangle \right. \\
 &\quad \left. - \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{\substack{i \neq j \\ i,j=1}}^K (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \sum_{\pi_i, \pi_j} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ij)} \right) (\mathbf{I} - \eta \mathbf{H})^{N(K-j)} \right\rangle \right| \\
 &\quad + \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^2} \sum_{\substack{i \neq j \\ i,j=1}}^K (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \sum_{\pi_i, \pi_j} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ij)} \right) (\mathbf{I} - \eta \mathbf{H})^{N(K-j)} \right\rangle \right. \\
 &\quad \left. - \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^2} \sum_{\substack{i \neq j \\ i,j=1}}^K (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \left( \sum_{l=0}^{N-1} \sum_{\pi_i, \pi_j} \tilde{\mathbb{E}} \mathbf{S}_l^{(ij)} \right) (\mathbf{I} - \eta \mathbf{H})^{N(K-j)} \right\rangle \right| \\
 &=: I_{11} + I_{12}.
 \end{aligned}$$

Next we bound the terms  $I_{11}$  and  $I_{12}$  separately. Notice that

$$\begin{aligned}
 &\sum_{\substack{i \neq j \\ i,j=1}}^K \sum_{\substack{\pi_1 \dots \pi_K \\ \text{except } \pi_i, \pi_j}} (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \sum_{\pi_i, \pi_j} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ij)} \right) (\mathbf{I} - \eta \mathbf{H})^{N(K-j)} \\
 &= (N!)^{K-2} \sum_{\substack{i \neq j \\ i,j=1}}^K (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \sum_{\pi_i, \pi_j} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ij)} \right) (\mathbf{I} - \eta \mathbf{H})^{N(K-j)} \tag{10}
 \end{aligned}$$

because the summands do not depend on the permutations except  $\pi_i, \pi_j$ , plugging [Equation \(10\)](#) into the expression of  $I_1$  we have

$$\begin{aligned}
 I_{11} &\leq \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{\substack{i \neq j \\ i,j=1}}^K \sum_{\substack{\pi_1 \dots \pi_K \\ \text{except } \pi_i, \pi_j}} \mathbf{T}^{(i)} \sum_{\pi_i, \pi_j} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ij)} \right) (\mathbf{T}^{(j)})^\top \right\rangle \right. \\
 &\quad \left. - \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{\substack{i \neq j \\ i,j=1}}^K \sum_{\substack{\pi_1 \dots \pi_K \\ \text{except } \pi_i, \pi_j}} (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \sum_{\pi_i, \pi_j} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ij)} \right) (\mathbf{I} - \eta \mathbf{H})^{N(K-j)} \right\rangle \right|.
 \end{aligned}$$

Then we use Equation (5) to split  $I_{11}$  into three terms and by triangle inequality:

$$\begin{aligned}
 I_{11} \leq & \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{\substack{i \neq j \\ i, j=1}}^K \sum_{\substack{\pi_1 \dots \pi_K \\ \text{except } \pi_i, \pi_j}} \left( \mathbf{T}^{(i)} - \tilde{\mathbb{E}} \mathbf{T}^{(i)} \right) \sum_{\pi_i, \pi_j} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ij)} \right) \tilde{\mathbb{E}} \mathbf{T}^{(i)} \right\rangle \right| \\
 & + \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{\substack{i \neq j \\ i, j=1}}^K \sum_{\substack{\pi_1 \dots \pi_K \\ \text{except } \pi_i, \pi_j}} \tilde{\mathbb{E}} \mathbf{T}^{(i)} \sum_{\pi_i, \pi_j} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ij)} \right) \left( \mathbf{T}^{(j)} - \tilde{\mathbb{E}} \mathbf{T}^{(j)} \right) \right\rangle \right| \\
 & + \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{\substack{i \neq j \\ i, j=1}}^K \sum_{\substack{\pi_1 \dots \pi_K \\ \text{except } \pi_i, \pi_j}} \left( \mathbf{T}^{(i)} - \tilde{\mathbb{E}} \mathbf{T}^{(i)} \right) \sum_{\pi_i, \pi_j} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ij)} \right) \left( \mathbf{T}^{(j)} - \tilde{\mathbb{E}} \mathbf{T}^{(j)} \right) \right\rangle \right|.
 \end{aligned}$$

Next, we use Lemma I.1 and the fact that  $\mathbf{S}_l^{(ij)} \lesssim \mathbf{I}$  to bound the matrix inner products:

$$\begin{aligned}
 I_{11} \leq & \frac{\eta^2 \sigma^2 N D^2 \text{tr}(\mathbf{H})}{2(N!)^{K-2}} \sum_{\substack{i \neq j \\ i, j=1}}^K \sum_{\substack{\pi_1 \dots \pi_K \\ \text{except } \pi_i, \pi_j}} \left( \mathbb{E} \left\| \mathbf{T}^{(i)} - \tilde{\mathbb{E}} \mathbf{T}^{(i)} \right\| + \mathbb{E} \left\| \mathbf{T}^{(j)} - \tilde{\mathbb{E}} \mathbf{T}^{(j)} \right\| \right. \\
 & \left. + \mathbb{E} \left\| \mathbf{T}^{(i)} - \tilde{\mathbb{E}} \mathbf{T}^{(i)} \right\| \left\| \mathbf{T}^{(j)} - \tilde{\mathbb{E}} \mathbf{T}^{(j)} \right\| \right).
 \end{aligned}$$

Notice that Lemma I.2 and Lemma I.5 implies that

$$\begin{aligned}
 \mathbb{E} \left\| \mathbf{T}^{(i)} - \tilde{\mathbb{E}} \mathbf{T}^{(i)} \right\| & \leq (\sqrt{\delta_A \eta^2 N K} + \|\mathbb{E} \mathbf{A}\|)^K - \|\mathbb{E} \mathbf{A}\|^K \\
 & \leq (\sqrt{\delta_A \eta^2 N K} + 1)^K - 1 \\
 & \leq 2K \sqrt{\delta_A \eta^2 N K} \quad \text{when } K \leq \frac{\log 2}{\eta \sqrt{\delta_A T}},
 \end{aligned}$$

where  $\delta_A = \tilde{C} 8e D^4 \log d$  is the constant appeared in Lemma I.4, and  $\tilde{C}$  is some absolute constant. The second inequality uses the fact that  $(\sqrt{\delta_A \eta^2 N K} + \|\mathbb{E} \mathbf{A}\|)^K - \|\mathbb{E} \mathbf{A}\|^K$  monotonously increases with  $\|\mathbb{E} \mathbf{A}\|$ . A similar approach combining Lemma I.2 and Lemma I.6 derives another concentration inequality for  $\mathbf{T}^{(i)}$ :

$$\mathbb{E} \left\| \mathbf{T}^{(i)} - \tilde{\mathbb{E}} \mathbf{T}^{(i)} \right\|^2 \leq \left( 2K \sqrt{\delta_A \eta^2 N K} \right)^2 \quad \text{when } K \leq \frac{\log 2}{\eta \sqrt{\delta_A T}}.$$

Applying Cauchy-Schwarz's inequality and the concentration inequalities for  $(\mathbf{T}^{(i)})_i$ , we get that

$$\begin{aligned} I_{11} &\leq \frac{\eta^2 \sigma^2 N D^2 \text{tr}(\mathbf{H})}{2(N!)^{K-2}} \sum_{\substack{i \neq j \\ i, j=1}}^K \sum_{\substack{\pi_1 \cdots \pi_K \\ \text{except } \pi_i, \pi_j}} \left( \mathbb{E} \left\| \mathbf{T}^{(i)} - \tilde{\mathbb{E}} \mathbf{T}^{(i)} \right\| + \mathbb{E} \left\| \mathbf{T}^{(j)} - \tilde{\mathbb{E}} \mathbf{T}^{(j)} \right\| \right. \\ &\quad \left. + \left( \mathbb{E} \left\| \mathbf{T}^{(i)} - \tilde{\mathbb{E}} \mathbf{T}^{(i)} \right\|^2 \right)^{\frac{1}{2}} \left( \mathbb{E} \left\| \mathbf{T}^{(j)} - \tilde{\mathbb{E}} \mathbf{T}^{(j)} \right\|^2 \right)^{\frac{1}{2}} \right) \\ &\leq \frac{\eta^2 \sigma^2 N D^2 \text{tr}(\mathbf{H})}{2} \sum_{\substack{i \neq j \\ i, j=1}}^K \left( 4K \sqrt{\delta_A \eta^2 N K} + \left( 2K \sqrt{2\delta_A \eta^2 N K} \right)^2 \right). \end{aligned}$$

Our next step is to bound  $I_{12}$ . We first make use of the fact that  $\mathbf{I} - \eta \mathbf{H} \lesssim \mathbf{I}$ , and get that

$$I_{12} \leq \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^2} \sum_{\substack{i \neq j \\ i, j=1}}^K \sum_{\pi_i, \pi_j} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ij)} - \tilde{\mathbb{E}} \mathbf{S}_l^{(ij)} \right) \right\rangle \right|.$$

Recall that for a fixed  $i$ , for all  $m \in [0, N-1]$ , there are  $(N-1)!$  permutations  $\pi_i$  that satisfies  $\pi_i(m) = l$ . So

$$\begin{aligned} I_{12} &\leq \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^2} \sum_{\substack{i \neq j \\ i, j=1}}^K \sum_{l=0}^{N-1} \sum_{m=0}^{N-1} \sum_{n=0}^{N-1} ((N-1)!)^2 \left( \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \mathbf{H} \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} \right. \right. \right. \\ &\quad \left. \left. - \mathbb{E} \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \mathbf{H} \mathbb{E} \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} \right) \right\rangle \right|. \end{aligned}$$

Notice that

$$\begin{aligned} &\mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \mathbf{H} \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} - \mathbb{E} \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \mathbf{H} \mathbb{E} \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} \\ &= \left( \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} - \mathbb{E} \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \right) \mathbf{H} \mathbb{E} \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} + \mathbb{E} \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \mathbf{H} \left( \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} - \mathbb{E} \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} \right) \\ &+ \left( \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} - \mathbb{E} \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \right) \mathbf{H} \left( \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} - \mathbb{E} \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} \right). \end{aligned}$$

Applying [Lemma I.1](#) and using the fact that  $\mathbb{E} \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \lesssim \mathbf{I}$ ,

$$\begin{aligned} I_{12} &\leq \frac{\eta^2 \sigma^2 \text{tr}(\mathbf{H}) \|\mathbf{H}\| N}{2N^2} \mathbb{E} \sum_{\substack{i \neq j \\ i, j=1}}^K \left( \sum_{m=0}^{N-2} \left\| \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} - \mathbb{E} \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \right\| \right. \\ &\quad \left. + \sum_{n=0}^{N-2} \left\| \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} - \mathbb{E} \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} \right\| \right) \end{aligned}$$

$$+ \sum_{m=0}^{N-2} \sum_{n=0}^{N-2} \left\| \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} - \mathbb{E} \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \right\| \left\| \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} - \mathbb{E} \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} \right\| \Bigg).$$

Applying Cauchy-Schwarz inequality and [Lemma I.4](#) gives

$$\begin{aligned} I_{12} &\leq \frac{\eta^2 \sigma^2 \text{tr}(\mathbf{H}) \|\mathbf{H}\| N}{2N^2} \sum_{\substack{i \neq j \\ i, j=1}}^K \left( \sum_{m=0}^{N-2} \mathbb{E} \left\| \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} - \mathbb{E} \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \right\| \right. \\ &\quad + \sum_{n=0}^{N-2} \mathbb{E} \left\| \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} - \mathbb{E} \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} \right\| \\ &\quad \left. + \sum_{m=0}^{N-2} \sum_{n=0}^{N-2} \left( \mathbb{E} \left\| \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} - \mathbb{E} \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \right\|^2 \right)^{\frac{1}{2}} \left( \mathbb{E} \left\| \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} - \mathbb{E} \mathbf{Z}_{n+1 \rightarrow N-1}^{(j)} \right\|^2 \right)^{\frac{1}{2}} \right) \\ &\leq \frac{\eta^2 \sigma^2 \text{tr}(\mathbf{H}) \|\mathbf{H}\| N}{2N^2} \sum_{\substack{i \neq j \\ i, j=1}}^K \left( \sum_{m=0}^{N-2} \left( \sqrt{\delta_A \eta^2 (N-1-m)} \right) + \sum_{n=0}^{N-2} \left( \sqrt{\delta_A \eta^2 (N-1-n)} \right) \right. \\ &\quad \left. + \sum_{m=0}^{N-2} \sum_{n=0}^{N-2} \left( \sqrt{2\delta_A \eta^2 (N-1-m)} \right) \left( \sqrt{2\delta_A \eta^2 (N-1-n)} \right) \right) \\ &\lesssim \eta^3 K^2 \sqrt{N \log d} + \eta^4 K^2 N^2 \log d \quad \text{when } \eta = o\left(\frac{1}{\sqrt{T}}\right). \end{aligned}$$

**Upper bound for  $I_2$ .** We bound  $I_2$  using a similar technique as what we did for  $I_1$ . We first plug in a term that takes pseudo expectation over  $(\mathbf{T}^{(i)})_{i=1}^K$  but does not take pseudo expectation over  $\mathbf{S}_l^{(ii)}$  for every  $l$  and  $i$ , and decompose  $I_2$  into two terms:

$$\begin{aligned} I_2 &\leq \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{i=1}^K \sum_{\substack{\pi_1 \dots \pi_K \\ \text{except } \pi_i}} \mathbf{T}^{(i)} \sum_{\pi_i} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ii)} \right) (\mathbf{T}^{(i)})^\top \right\rangle \right. \\ &\quad - \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{N!} \sum_{i=1}^K (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \sum_{\pi_i} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ii)} \right) (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \right\rangle \Bigg| \\ &\quad + \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{N!} \sum_{i=1}^K (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \sum_{\pi_i} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ii)} \right) (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \right\rangle \right. \\ &\quad \left. - \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{N!} \sum_{i=1}^K (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \left( \sum_{l=0}^{N-1} \sum_{\pi_i} \tilde{\mathbb{E}} \mathbf{S}_l^{(ii)} \right) (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \right\rangle \right| \\ &=: I_{21} + I_{22}. \end{aligned}$$

Next we bound the terms  $I_{21}$  and  $I_{22}$  separately. Notice that

$$\begin{aligned} & \sum_{i=1}^K \sum_{\substack{\pi_1 \cdots \pi_K \\ \text{except } \pi_i}} (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \sum_{\pi_i} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ii)} \right) (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \\ &= (N!)^{K-1} \sum_{i=1}^K (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \sum_{\pi_i} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ii)} \right) (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \end{aligned}$$

because the summands do not depend on the permutations except  $\pi_i$ , we have

$$\begin{aligned} I_{21} &= \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{i=1}^K \sum_{\substack{\pi_1 \cdots \pi_K \\ \text{except } \pi_i}} \mathbf{T}^{(i)} \sum_{\pi_i} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ii)} \right) (\mathbf{T}^{(i)})^\top \right\rangle \right. \\ &\quad \left. - \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{i=1}^K \sum_{\substack{\pi_1 \cdots \pi_K \\ \text{except } \pi_i}} (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \sum_{\pi_i} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ii)} \right) (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \right\rangle \right|. \end{aligned}$$

Then we use the fact that  $\tilde{\mathbb{E}} \mathbf{T}^{(i)} = (\mathbf{I} - \eta \mathbf{H})^{N(K-i)}$  to split  $I_{21}$  into three terms:

$$\begin{aligned} I_{21} &\leq \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{i=1}^K \sum_{\substack{\pi_1 \cdots \pi_K \\ \text{except } \pi_i}} (\mathbf{T}^{(i)} - \tilde{\mathbb{E}} \mathbf{T}^{(i)}) \sum_{\pi_i} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ii)} \right) (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \right\rangle \right| \\ &\quad + \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{i=1}^K \sum_{\substack{\pi_1 \cdots \pi_K \\ \text{except } \pi_i}} (\mathbf{I} - \eta \mathbf{H})^{N(K-i)} \sum_{\pi_i} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ii)} \right) (\mathbf{T}^{(i)} - \tilde{\mathbb{E}} \mathbf{T}^{(i)}) \right\rangle \right| \\ &\quad + \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{i=1}^K \sum_{\substack{\pi_1 \cdots \pi_K \\ \text{except } \pi_i}} (\mathbf{T}^{(i)} - \tilde{\mathbb{E}} \mathbf{T}^{(i)}) \sum_{\pi_i} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ii)} \right) (\mathbf{T}^{(i)} - \tilde{\mathbb{E}} \mathbf{T}^{(i)}) \right\rangle \right|. \end{aligned}$$

Next, we use [Lemma I.1](#) and the fact that  $\mathbf{S}_l^{(ij)} \lesssim \mathbf{I}$  to bound the matrix inner products, and apply the concentration inequalities we derived for  $((\mathbf{T}^{(i)})_i)_i$ :

$$\begin{aligned} I_{21} &\leq \frac{\eta^2 \sigma^2 N D^2 \text{tr}(\mathbf{H})}{2(N!)^{K-1}} \sum_{i=1}^K \sum_{\substack{\pi_1 \cdots \pi_K \\ \text{except } \pi_i}} \left( \mathbb{E} \left\| \mathbf{T}^{(i)} - \tilde{\mathbb{E}} \mathbf{T}^{(i)} \right\| + \mathbb{E} \left\| \mathbf{T}^{(i)} - \tilde{\mathbb{E}} \mathbf{T}^{(i)} \right\|^2 \right) \\ &\leq \frac{\eta^2 \sigma^2 N D^2 \text{tr}(\mathbf{H})}{2} \sum_{i=1}^K \left( 4K \sqrt{\delta_A \eta^2 K N} + \left( 2K \sqrt{2\delta_A \eta^2 K N} \right)^2 \right). \end{aligned}$$



Then we bound  $I_{22}$ . Recall that  $\mathbf{I} - \eta \mathbf{H} \lesssim \mathbf{I}$ , we get

$$I_{22} \leq \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{N!} \sum_{i=1}^K \sum_{\pi_i} \left( \sum_{l=0}^{N-1} \mathbf{s}_l^{(ii)} - \tilde{\mathbb{E}} \mathbf{s}_l^{(ii)} \right) \right\rangle \right|.$$

Recall that for a fixed  $i$ , for all  $m \in [0, N-1]$ , there are  $(N-1)!$  permutations  $\pi_i$  that satisfies  $\pi_i(m) = l$ . So

$$\begin{aligned} I_{22} &\leq \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{N!} \sum_{i=1}^K \sum_{l=0}^{N-1} \sum_{m=0}^{N-1} (N-1)! \left( \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \mathbf{H} \mathbf{Z}_{m+1 \rightarrow N-1}^{(i)} \right. \right. \right. \\ &\quad \left. \left. \left. - \mathbb{E} \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \mathbf{H} \mathbb{E} \mathbf{Z}_{m+1 \rightarrow N-1}^{(i)} \right) \right\rangle \right| \\ &= \left| \frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{N!} \sum_{i=1}^K \sum_{l=0}^{N-1} (N-1)! \sum_{m=0}^{N-2} \left( \left( \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} - \mathbb{E} \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \right) \mathbf{H} \right. \right. \right. \\ &\quad \left. \left. \left. \left( \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} - \mathbb{E} \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \right) \right) \right\rangle \right|. \end{aligned}$$

Using [Lemma I.4](#), we have

$$\begin{aligned} I_{22} &\leq \frac{\eta^2 \sigma^2 \text{tr}(\mathbf{H}) \|\mathbf{H}\| N}{2N} \mathbb{E} \sum_{i=1}^K \left( \sum_{m=0}^{N-2} \left\| \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} - \mathbb{E} \mathbf{Z}_{N-1 \rightarrow m+1}^{(i)} \right\|^2 \right) \\ &\leq \frac{\eta^2 \sigma^2 \text{tr}(\mathbf{H}) \|\mathbf{H}\| N}{2N} \sum_{i=1}^K \sum_{m=0}^{N-2} \left( \sqrt{2\delta_A \eta^2 (N-1-m)} \right)^2 \\ &\lesssim \eta^4 N^2 K \log d \quad \text{when } \eta = o\left(\frac{1}{\sqrt{T}}\right). \end{aligned}$$

Combining all the arguments above, we derive that

$$\begin{aligned} &\left| \tilde{\mathcal{R}}^{\text{var}}(K, N; \eta) - \bar{\mathcal{R}}^{\text{var}}(K, N; \eta) \right| \\ &\leq I_{11} + I_{12} + I_{21} + I_{22} \\ &\leq C \frac{\eta^2 \sigma^2 N D^2 \text{tr}(\mathbf{H})}{2} \sum_{i,j=1}^K \left( 4K \sqrt{\delta_A \eta^2 N K} + \left( 2K \sqrt{\delta_A \eta^2 N K} \right)^2 \right) \\ &\quad + O(\eta^3 K^2 \sqrt{N \log d} + \eta^4 K^2 N^2 \log d) + O(\eta^4 N^2 K \log d) \\ &= O(\eta^3 N^{\frac{3}{2}} K^{\frac{7}{2}} \sqrt{\log d}) \quad \text{when } \eta = o\left(\frac{1}{\sqrt{T}}\right). \end{aligned}$$

The above equation completes the proof.

### F.2.2. BIAS TERM ANALYSIS: PROOF OF LEMMA F.3

For simplicity, and as we did in the proof of Lemma F.2, in this section we omit the superscript "bias" for all the training parameters  $\boldsymbol{\theta}^{\text{bias}}$ . Analogous to the proof of Lemma F.2, we can derive the parameter recursion as

$$\begin{aligned}\boldsymbol{\theta}_{kN} &= (\mathbf{I} - \eta \mathbf{x}_{\pi_k(N-1)} \mathbf{x}_{\pi_k(N-1)}^\top) \boldsymbol{\theta}_{(k-1)N} \\ &= \dots \\ &= (\mathbf{I} - \eta \mathbf{x}_{\pi_k(N-1)} \mathbf{x}_{\pi_k(N-1)}^\top) \dots (\mathbf{I} - \eta \mathbf{x}_{\pi_k(0)} \mathbf{x}_{\pi_k(0)}^\top) \boldsymbol{\theta}_{(k-1)N} \\ &= \mathbf{A}^{(k)} \boldsymbol{\theta}_{(k-1)N}.\end{aligned}$$

For the parameter after  $K$ -epochs updates, we have

$$\boldsymbol{\theta}_{KN} = \mathbf{A}^{(K)} \dots \mathbf{A}^{(1)} \boldsymbol{\theta}_0 = \prod_{l=K}^1 \mathbf{A}^{(l)} \boldsymbol{\theta}_0.$$

We also have the approximation for the bias term

$$\begin{aligned}\bar{\mathcal{R}}^{\text{bias}}(K, N; \eta) &= \frac{1}{2} \langle \mathbf{H}, \mathbb{E} \boldsymbol{\theta}_{KN}^2 \rangle \\ &= \mathbb{E} \frac{1}{2} \boldsymbol{\theta}_{KN}^\top \mathbf{H} \boldsymbol{\theta}_{KN} \\ &= \mathbb{E} \frac{1}{2} \boldsymbol{\theta}_0^\top \left( \prod_{l=K}^1 \mathbf{A}^{(l)} \right)^\top \mathbf{H} \left( \prod_{l=K}^1 \mathbf{A}^{(l)} \right) \boldsymbol{\theta}_0 \\ &\approx \frac{1}{2} \boldsymbol{\theta}_0^\top \left( \prod_{l=K}^1 \mathbb{E} \mathbf{A}^{(l)} \right)^\top \mathbf{H} \left( \prod_{l=K}^1 \mathbb{E} \mathbf{A}^{(l)} \right) \boldsymbol{\theta}_0 \\ &= \underbrace{\frac{1}{2} \boldsymbol{\theta}_0^\top ((\mathbf{I} - \eta \mathbf{H})^{KN}) \mathbf{H} ((\mathbf{I} - \eta \mathbf{H})^{KN}) \boldsymbol{\theta}_0}_{=:\bar{\mathcal{R}}^{\text{var}}(K, N; \eta)}.\end{aligned}$$

The estimate error can be given as

$$\begin{aligned}& \left| \tilde{\mathcal{R}}^{\text{bias}}(K, N; \eta) - \bar{\mathcal{R}}^{\text{bias}}(K, N; \eta) \right| \\ &= \left| \mathbb{E} \frac{1}{2} \boldsymbol{\theta}_0^\top \left( \prod_{l=K}^1 \mathbf{A}^{(l)} \right)^\top \mathbf{H} \left( \prod_{l=K}^1 \mathbf{A}^{(l)} \right) \boldsymbol{\theta}_0 - \frac{1}{2} \boldsymbol{\theta}_0^\top \left( \prod_{l=K}^1 \mathbb{E} \mathbf{A}^{(l)} \right)^\top \mathbf{H} \left( \prod_{l=K}^1 \mathbb{E} \mathbf{A}^{(l)} \right) \boldsymbol{\theta}_0 \right| \\ &= \left| \mathbb{E} \frac{1}{2} \boldsymbol{\theta}_0^\top \left( \prod_{l=K}^1 \mathbf{A}^{(l)} - \|\mathbb{E} \mathbf{A}\|^K \right)^\top \mathbf{H} \left( \prod_{l=K}^1 \mathbf{A}^{(l)} - \|\mathbb{E} \mathbf{A}\|^K \right) \boldsymbol{\theta}_0 \right| \\ &+ 2 \left| \mathbb{E} \frac{1}{2} \boldsymbol{\theta}_0^\top \|\mathbb{E} \mathbf{A}\|^K \mathbf{H} \left( \prod_{l=K}^1 \mathbf{A}^{(l)} - \|\mathbb{E} \mathbf{A}\|^K \right) \boldsymbol{\theta}_0 \right| \\ &\leq \mathbb{E} \frac{1}{2} \|\mathbf{H}\| \|\boldsymbol{\theta}_0\|^2 \left( \left\| \mathbf{A}^K - (\mathbb{E} \mathbf{A})^K \right\|^2 + 2 \|\mathbb{E} \mathbf{A}\|^K \left\| \mathbf{A}^K - (\mathbb{E} \mathbf{A})^K \right\| \right). \tag{11}\end{aligned}$$

where the last equation uses the fact that  $\|\mathbb{E}\mathbf{A}\| \leq 1$ . Next, we discuss the approximation error bound for the bias term in Equation (11), with different categorizations based on the range of  $K$ .

1. Under Assumption 3.1 and  $K = o\left(\frac{N^{\frac{1}{5}}}{(\log N)^{\frac{6}{5}}}\right)$ :

(a)  $\eta \leq \frac{2\log T}{3\lambda_d T}$ . We now verify that  $K = o\left(\frac{\|\mathbb{E}\mathbf{A}\|}{\eta\sqrt{T}}\right)$  under given conditions. We have

$$\begin{aligned} \|\mathbb{E}\mathbf{A}\| &= (1 - \eta\lambda_d)^N = (1 - \eta\lambda_d)^{\frac{T}{K}} \geq (1 - \eta\lambda_d)^{\frac{T}{2}} \\ &\geq \left(1 - \frac{2\log T}{3T}\right)^{\frac{T}{2}} = e^{\frac{T}{2} \log(1 - \frac{2\log T}{3T})} \\ &= e^{-\frac{\log T}{3} + O(\frac{2\log^2 T}{9T})} = \Theta\left(\frac{1}{T^{\frac{1}{3}}}\right). \end{aligned}$$

thus

$$\frac{\|\mathbb{E}\mathbf{A}\|}{\eta\sqrt{T}} = \Omega\left(\frac{T^{\frac{1}{6}}}{\log T}\right).$$

Also, given  $K = o\left(\frac{N^{\frac{1}{5}}}{(\log N)^{\frac{6}{5}}}\right)$ , we obtain that

$$K = o\left(\frac{T^{\frac{1}{6}}}{\log N}\right) = o\left(\frac{T^{\frac{1}{6}}}{\log T}\right).$$

The second equality uses  $\log T = \log N + \log K = \Theta(\log N)$ . Now we use the results in Lemma I.5 and Lemma I.6, and then the estimated distance can be given as

$$\begin{aligned} &\left| \tilde{\mathcal{R}}^{\text{bias}}(K, N; \eta) - \bar{\mathcal{R}}^{\text{bias}}(K, N; \eta) \right| \\ &\leq \frac{1}{2} \|\mathbf{H}\| \|\boldsymbol{\theta}_0\|^2 \|\mathbb{E}\mathbf{A}\|^{2K} \left( \left( \left( \frac{\sqrt{2\delta_A \eta^2 N K}}{\|\mathbb{E}\mathbf{A}\|} + 1 \right)^K - 1 \right)^2 + 2 \left( \left( \frac{\sqrt{2\delta_A \eta^2 N K}}{\|\mathbb{E}\mathbf{A}\|} + 1 \right)^K - 1 \right) \right) \\ &\leq \frac{1}{2} \|\mathbf{H}\| \|\boldsymbol{\theta}_0\|^2 \|\mathbb{E}\mathbf{A}\|^{2K} \left( \frac{8K^2 \delta_A \eta^2 N K}{\|\mathbb{E}\mathbf{A}\|^2} + 4K \frac{\sqrt{2\delta_A \eta^2 N K}}{\|\mathbb{E}\mathbf{A}\|} \right) \\ &= O\left(\|\mathbb{E}\mathbf{A}\|^{2K-1} K \sqrt{\eta^2 N K}\right), \end{aligned}$$

where the second inequality is by Lemma I.2.

(b)  $\eta \geq \frac{2 \log T}{3\lambda_d T}$ . We have

$$\begin{aligned}
 & \left| \tilde{\mathcal{R}}^{\text{bias}}(k, N; \eta) - \bar{\mathcal{R}}^{\text{bias}}(k, N; \eta) \right| \leq \tilde{\mathcal{R}}^{\text{bias}}(k, N; \eta) + \bar{\mathcal{R}}^{\text{bias}}(k, N; \eta) \\
 & \leq \left[ \tilde{\mathcal{R}}^{\text{bias}}(k, N; \eta) + \bar{\mathcal{R}}^{\text{bias}}(k, N; \eta) \right] \Big|_{\eta = \frac{2 \log T}{3\lambda_d T}} \\
 & \leq \left[ \left| \tilde{\mathcal{R}}^{\text{bias}}(k, N; \eta) - \bar{\mathcal{R}}^{\text{bias}}(k, N; \eta) \right| + 2\tilde{\mathcal{R}}^{\text{bias}}(k, N; \eta) \right] \Big|_{\eta = \frac{2 \log T}{3\lambda_d T}} \\
 & \leq \left[ O\left(\|\mathbb{E}\mathbf{A}\|^{2K-1} K \sqrt{\eta^2 K N}\right) + 2 \times \frac{1}{2} \|\mathbf{H}\| \|\boldsymbol{\theta}_0\|^2 \|\mathbb{E}\mathbf{A}\|^{2K} \right] \Big|_{\eta = \frac{2 \log T}{3\lambda_d T}} \\
 & = O\left(\|\mathbb{E}\mathbf{A}\|^{2K}\right) \Big|_{\eta = \frac{2 \log T}{3\lambda_d T}} = O\left(\left(1 - \frac{2 \log T}{3T}\right)^{2KN}\right) \\
 & = O\left(\frac{1}{T^{\frac{4}{3}}}\right) \text{ when } K = o\left(\frac{N^{\frac{1}{5}}}{(\log N)^{\frac{6}{5}}}\right),
 \end{aligned}$$

where the first equality uses the fact that  $K = o\left(\frac{\|\mathbb{E}\mathbf{A}\|}{\eta \sqrt{T}}\right)$  when  $\eta = \frac{2 \log T}{3\lambda_d T}$ .

- For the  $K = 1$  case, which is equivalent to one-pass (OP) SGD, we derive a different upper bound for bias term error. In this scenario, we have the update rule as

$$\boldsymbol{\theta}_t = (\mathbf{I} - \eta \mathbf{x}_t \mathbf{x}_t^\top) \boldsymbol{\theta}_{t-1}.$$

We can denote the covariance as  $\mathbf{B}_t$ , which is

$$\begin{aligned}
 \mathbf{B}_t &:= \mathbb{E} \boldsymbol{\theta}_t \boldsymbol{\theta}_t^\top \\
 &= \mathbb{E} (\mathbf{I} - \eta \mathbf{x}_t \mathbf{x}_t^\top) \boldsymbol{\theta}_{t-1} \boldsymbol{\theta}_{t-1}^\top (\mathbf{I} - \eta \mathbf{x}_t \mathbf{x}_t^\top) \\
 &= \mathbf{B}_{t-1} - \eta \mathbf{H} \mathbf{B}_{t-1} - \eta \mathbf{B}_{t-1} \mathbf{H} + \eta^2 \mathbb{E} \mathbf{x}_t \mathbf{x}_t^\top \boldsymbol{\theta}_{t-1} \boldsymbol{\theta}_{t-1}^\top \mathbf{x}_t \mathbf{x}_t^\top \\
 &= (\mathbf{I} - \eta \mathbf{H}) \mathbf{B}_{t-1} (\mathbf{I} - \eta \mathbf{H}) + \eta^2 \mathbb{E} (\mathbf{x}_t \mathbf{x}_t^\top - \mathbf{H}) \boldsymbol{\theta}_{t-1} \boldsymbol{\theta}_{t-1}^\top (\mathbf{x}_t \mathbf{x}_t^\top - \mathbf{H}). \tag{12}
 \end{aligned}$$

Since the bias term in the excess risk can be represented as

$$\bar{\mathcal{R}}^{\text{bias}}(1, T; \eta) = \frac{1}{2} \langle \mathbf{H}, \mathbf{B}_T \rangle.$$

We then get the lower and upper bounds for  $\mathbf{B}_t$ , and derive the corresponding lower and upper bounds for the bias term in the excess risk.

**Lower bound.** By Equation (12), we get a lower bound of  $\mathbf{B}_t$

$$\begin{aligned}
 \mathbf{B}_T &\succeq (\mathbf{I} - \eta \mathbf{H}) \mathbf{B}_{T-1} (\mathbf{I} - \eta \mathbf{H}) \\
 &\succeq \cdots \succeq (\mathbf{I} - \eta \mathbf{H})^T \mathbf{B}_0 (\mathbf{I} - \eta \mathbf{H})^T
 \end{aligned}$$

and

$$\begin{aligned}
 \bar{\mathcal{R}}^{\text{bias}}(1, T; \eta) &= \frac{1}{2} \langle \mathbf{H}, \mathbf{B}_T \rangle \\
 &\geq \frac{1}{2} \langle \mathbf{H}, (\mathbf{I} - \eta \mathbf{H})^T \mathbf{B}_0 (\mathbf{I} - \eta \mathbf{H})^T \rangle \\
 &= \frac{1}{2} \boldsymbol{\theta}_0^\top ((\mathbf{I} - \eta \mathbf{H})^T) \mathbf{H} ((\mathbf{I} - \eta \mathbf{H})^T) \boldsymbol{\theta}_0.
 \end{aligned}$$

**Upper bound.** By the recursion of  $\mathbf{B}_t$ , we have

$$\begin{aligned}
 \mathbf{B}_t &\preceq (\mathbf{I} - \eta \mathbf{H}) \mathbf{B}_{t-1} (\mathbf{I} - \eta \mathbf{H}) + \eta^2 \mathbb{E}_{\mathbf{x}_{T-1}, \dots, \mathbf{x}_0} \mathbb{E}_{\mathbf{x}_T} (\mathbf{x}_t \mathbf{x}_t^\top - \mathbf{H}) \boldsymbol{\theta}_{t-1} \boldsymbol{\theta}_{t-1}^\top (\mathbf{x}_t \mathbf{x}_t^\top - \mathbf{H}) \\
 &= (\mathbf{I} - \eta \mathbf{H}) \mathbf{B}_{t-1} (\mathbf{I} - \eta \mathbf{H}) + \eta^2 \mathbb{E}_{\mathbf{x}_{T-1}, \dots, \mathbf{x}_0} \left[ \mathbb{E}_{\mathbf{x}_T} \left[ \mathbf{x}_T \mathbf{x}_T^\top \boldsymbol{\theta}_{T-1} \boldsymbol{\theta}_{T-1}^\top \mathbf{x}_T \mathbf{x}_T^\top \right] - \mathbf{H} \boldsymbol{\theta}_{T-1} \boldsymbol{\theta}_{T-1}^\top \mathbf{H} \right] \\
 &\preceq (\mathbf{I} - \eta \mathbf{H}) \mathbf{B}_{t-1} (\mathbf{I} - \eta \mathbf{H}) + \eta^2 \mathbb{E}_{\mathbf{x}_{T-1}, \dots, \mathbf{x}_0} \mathbb{E}_{\mathbf{x}_T} \left[ \mathbf{x}_T \mathbf{x}_T^\top \boldsymbol{\theta}_{T-1} \boldsymbol{\theta}_{T-1}^\top \mathbf{x}_T \mathbf{x}_T^\top \right].
 \end{aligned}$$

Then, combining [Assumption 3.1](#) and [Lemma I.9](#) gives

$$\begin{aligned}
 \mathbf{B}_T &\preceq (\mathbf{I} - \eta \mathbf{H}) \mathbf{B}_{T-1} (\mathbf{I} - \eta \mathbf{H}) + \eta^2 \alpha \mathbb{E}_{\mathbf{x}_{T-1}, \dots, \mathbf{x}_0} \text{tr}(\mathbf{H} \boldsymbol{\theta}_{T-1} \boldsymbol{\theta}_{T-1}^\top) \mathbf{H} \\
 &= (\mathbf{I} - \eta \mathbf{H}) \mathbf{B}_{T-1} (\mathbf{I} - \eta \mathbf{H}) + \eta^2 \alpha \langle \mathbf{H}, \mathbf{B}_{T-1} \rangle \mathbf{H} \\
 &\preceq \dots \\
 &\preceq (\mathbf{I} - \eta \mathbf{H})^T \mathbf{B}_0 (\mathbf{I} - \eta \mathbf{H})^T + \eta^2 \alpha \sum_{i=0}^{T-1} \langle \mathbf{B}_i, \mathbf{H} \rangle (\mathbf{I} - \eta \mathbf{H})^{2(T-i-1)} \mathbf{H},
 \end{aligned}$$

and

$$\langle \mathbf{H}, \mathbf{B}_T \rangle \leq \langle \mathbf{H}, (\mathbf{I} - \eta \mathbf{H})^T \mathbf{B}_0 (\mathbf{I} - \eta \mathbf{H})^T \rangle + \eta^2 \alpha \sum_{i=0}^{T-1} \langle \mathbf{H}, \mathbf{B}_i \rangle \langle (\mathbf{I} - \eta \mathbf{H})^{2(T-i-1)} \mathbf{H}, \mathbf{H} \rangle.$$

We also have

$$\begin{aligned}
 \langle \mathbf{H}, \mathbf{B}_i \rangle &\leq \langle \mathbf{H}, (\mathbf{I} - \eta \mathbf{H}) \mathbf{B}_{i-1} (\mathbf{I} - \eta \mathbf{H}) \rangle + \eta^2 \alpha \text{tr}(\mathbf{H}^2) \langle \mathbf{H}, \mathbf{B}_{i-1} \rangle \\
 &\leq (1 - \eta \lambda_d)^2 \langle \mathbf{H}, \mathbf{B}_{i-1} \rangle + \eta^2 \alpha \text{tr}(\mathbf{H}^2) \langle \mathbf{H}, \mathbf{B}_{i-1} \rangle \\
 &\leq \dots \\
 &\leq [(\lambda_d^2 + \alpha \text{tr}(\mathbf{H}^2)) \eta^2 - 2 \lambda_d \eta + 1]^i \langle \mathbf{H}, \mathbf{B}_0 \rangle \\
 &\leq e^{T \log[(\lambda_d^2 + \alpha \text{tr}(\mathbf{H}^2)) \eta^2 - 2 \lambda_d \eta + 1]} \langle \mathbf{H}, \mathbf{B}_0 \rangle \\
 &= e^{-2 \lambda_d \eta i + O(\eta^2 i)} \langle \mathbf{H}, \mathbf{B}_0 \rangle \\
 &\leq C_1 e^{-2 \lambda_d \eta i} \langle \mathbf{H}, \mathbf{B}_0 \rangle
 \end{aligned}$$

and

$$\begin{aligned}
 \langle (\mathbf{I} - \eta \mathbf{H})^{2(T-i-1)} \mathbf{H}, \mathbf{H} \rangle &= \langle (\mathbf{I} - \eta \mathbf{H})^{2(T-i-1)}, \mathbf{H}^2 \rangle \\
 &\leq \text{tr}(\mathbf{H}^2) (1 - \eta \lambda_d)^{2(T-1-i)} \\
 &\leq \text{tr}(\mathbf{H}^2) e^{2(T-1-i) \log(1 - \eta \lambda_d)} \\
 &= \text{tr}(\mathbf{H}^2) e^{-2(T-1-i) \eta \lambda_d + O(\eta^2(T-1-i))} \\
 &\leq C_2 e^{-2(T-1-i) \eta \lambda_d}
 \end{aligned}$$

So

$$\begin{aligned}\langle \mathbf{H}, \mathbf{B}_i \rangle &\leq \langle \mathbf{H}, (\mathbf{I} - \eta \mathbf{H})^T \mathbf{B}_0 (\mathbf{I} - \eta \mathbf{H})^T \rangle + \eta^2 \alpha \sum_{i=0}^{T-1} C_1 e^{-2\lambda_d \eta i} \langle \mathbf{H}, \mathbf{B}_0 \rangle C_2 e^{-2\lambda_d \eta (T-1-i)} \text{tr}(\mathbf{H}^2) \\ &= \langle \mathbf{H}, (\mathbf{I} - \eta \mathbf{H})^T \mathbf{B}_0 (\mathbf{I} - \eta \mathbf{H})^T \rangle + C_3 \eta^2 T e^{-2\lambda_d \eta T}\end{aligned}$$

And finally we get

$$\left| \bar{\mathcal{R}}^{\text{bias}}(1, T; \eta) - \frac{1}{2} \langle \mathbf{H}, (\mathbf{I} - \eta \mathbf{H})^\top \mathbf{B}_0 (\mathbf{I} - \eta \mathbf{H}) \rangle \right| = O(\eta^2 T e^{-2\lambda_d \eta T}).$$

The above equation completes the proof.

### F.3. Step III: Narrowing the Range for Optimal Learning Rate

We recap that our goal to get the scaling law formula for strongly convex linear regression with multi epoch SGD, and the formula of the effective reuse rate. After we get a risk approximation in Step II, next we start to narrow the range of the optimal learning rate based on the risk approximation. We first give a technical lemma below, as a further simplification of the risk formula.

**Lemma F.4** *Given  $\eta \in \left[ \omega\left(\frac{1}{T}\right), o\left(\frac{1}{\sqrt{T}}\right) \right]$ , and define  $n_d$  to be the number of the minimal eigenvalue  $\lambda_d$  in  $\mathbf{H}$ , then it holds that*

$$\begin{aligned}\sum_{i=1}^d (\mathbf{P}\boldsymbol{\theta}_0)_i^2 \lambda_i (1 - \eta \lambda_i)^{2T} &= \tilde{\theta}_d^2 \lambda_d \exp(-2\lambda_d \eta T) (1 + o(1)), \\ \sum_{i=1}^d \lambda_i (1 - \eta \lambda_i)^{2T} &= n_d \lambda_d \exp(-2\lambda_d \eta T) (1 + o(1)).\end{aligned}$$

*Proof of Lemma F.4.* For the first equation, for any  $\lambda_i > \lambda_d$ , we define  $\rho_i = \frac{\lambda_i}{\lambda_d} > 1$ , then we have

$$\begin{aligned}(1 - \eta \lambda_i)^{2T} &= \exp(2T \log(1 - \eta \lambda_i)) = \exp(2T(-\eta \lambda_i + O(\eta^2 \lambda_i^2))) \\ &= \exp(-2\lambda_i \eta T) \exp(O(\eta^2)) = \exp(-2\lambda_d \rho_i \eta T) (1 + o(1)) \\ &= (\exp(-2\lambda_d \eta T))^{\rho_i} (1 + o(1)) = o(\exp(-2\lambda_d \eta T)).\end{aligned}\tag{13}$$

Since  $\lambda_i \leq D^2$ , we have

$$\sum_{i=1}^{d-n_d} (\mathbf{P}\boldsymbol{\theta}_0)_i^2 \lambda_i (1 - \eta \lambda_i)^{2T} = o(\exp(-2\lambda_d \eta T)),$$

From Equation (13), we can also directly get the second equation, which completes the proof of Lemma F.4.  $\square$

F.3.1. A DESCRIPTION OF THE RANGE OF OPTIMAL LEARNING RATE, SMALL- $K$  CASE

**Lemma F.5** *Under the conditions in Lemma E.4, and when  $K = o(\log N)$ , we have  $\eta^* \in [\frac{\log T}{3\lambda_d T}, \frac{\alpha \log T}{T}]$ , where the constant  $\alpha := \frac{D^2 \text{tr}(\mathbf{H})}{\lambda_d \text{tr}(\mathbf{H}^2)}$ .*

*Proof.* We first prove the upper bound. Given a learning rate  $\eta$ , Equation (7) gives

$$\begin{aligned} \bar{\mathcal{R}}(K, N; \eta) &\geq \bar{\mathcal{R}}^{\text{var}}(K, N; \eta) = \\ &\underbrace{\frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{\substack{i \neq j \\ i, j=1}}^K \sum_{\substack{\pi_1 \dots \pi_K \\ \text{except } \pi_i, \pi_j}} \mathbf{T}^{(i)} \sum_{\pi_i, \pi_j} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ij)} \right) (\mathbf{T}^{(j)})^\top \right\rangle}_{=:\psi_1} \\ &+ \underbrace{\frac{\eta^2 \sigma^2}{2} \mathbb{E} \left\langle \mathbf{H}, \frac{1}{(N!)^K} \sum_{i=1}^K \sum_{\substack{\pi_1 \dots \pi_K \\ \text{except } \pi_i}} \mathbf{T}^{(i)} \sum_{\pi_i} \left( \sum_{l=0}^{N-1} \mathbf{S}_l^{(ii)} \right) (\mathbf{T}^{(i)})^\top \right\rangle}_{=:\psi_2}. \end{aligned}$$

For  $\psi_1$ , using the fact that  $(\mathbf{I} - \eta \mathbf{x} \mathbf{x}^\top) \succeq (\mathbf{I} - \eta D^2 \mathbf{I})$ , we replace all the terms  $(\mathbf{I} - \eta \mathbf{x} \mathbf{x}^\top)$  with  $(\mathbf{I} - \eta D^2 \mathbf{I})$  thus we have a lower bound for  $\psi_1$

$$\begin{aligned} \psi_1 &\geq \frac{\eta^2 \sigma^2}{2} \left\langle \mathbf{H}, \frac{N((N-1)!)^2}{(N!)^K} \sum_{\substack{i \neq j \\ i, j=1}}^K \sum_{\substack{\{\pi_1 \dots \pi_K\} \\ \setminus \{\pi_i, \pi_j\}}} (1 - \eta D^2)^{(2K-i-j)N} \left( \sum_{m,n=0}^{N-1} (1 - \eta D^2)^{2N-2-m-n} \mathbb{E}[\mathbf{x} \mathbf{x}^\top] \right) \right\rangle \\ &= \frac{\eta^2 \sigma^2}{2ND^4} \left\langle \mathbf{H}, \sum_{\substack{i \neq j \\ i, j=1}} (1 - \eta D^2)^{(K-i)N} (1 - \eta D^2)^{(K-j)N} (1 - (1 - \eta D^2)^N)^2 \mathbf{H} \right\rangle \\ &= \frac{\sigma^2}{2ND^4} \text{tr} \left( \mathbf{H}^2 (1 - (1 - \eta D^2)^{KN})^2 \right) - \frac{\sigma^2}{2ND^4} \text{tr} \left( \mathbf{H}^2 \frac{1 - (1 - (1 - \eta D^2)^N)^{2KN}}{1 - (1 - \eta D^2)^{2N}} \right) \\ &= \frac{\sigma^2}{ND^4} \text{tr} \left( \mathbf{H}^2 \frac{1 - (1 - \eta D^2)^{KN}}{1 + (1 - \eta D^2)^N} ((1 - \eta D^2)^N - (1 - \eta D^2)^{KN}) \right). \end{aligned}$$

For  $\psi_2$ , we use a similar argument to get its lower bound

$$\begin{aligned} \psi_2 &\geq \frac{\eta^2 \sigma^2}{2} \left\langle \mathbf{H}, \sum_{i=1}^K (1 - \eta D^2)^{2N(K-i)} \frac{1 - (1 - \eta D^2)^{2N}}{1 - (1 - \eta D^2)^2} \mathbf{H} \right\rangle \\ &= \frac{\eta \sigma^2}{2D^2} \left\langle \mathbf{H}, \frac{1 - (1 - \eta D^2)^{2KN}}{1 - (1 - \eta D^2)^{2N}} \frac{1 - (1 - \eta D^2)^{2N}}{1 - (1 - \eta D^2)^2} \mathbf{H} \right\rangle \\ &= \frac{\eta \sigma^2 \text{tr}(\mathbf{H}^2)}{4D^2} (1 + o(1)). \end{aligned}$$



Notice that from the above lower bound, when  $K = o(\log N)$ , we have

$$\begin{aligned}\bar{\mathcal{R}}(K, N; \eta) &\geq \psi_1 + \psi_2 \\ &\geq O\left(\frac{1}{N}\right) + \frac{\eta\sigma^2\text{tr}(\mathbf{H}^2)}{4D^2} (1 + o(1)) \\ &= \frac{\eta\sigma^2\text{tr}(\mathbf{H}^2)}{4D^2} (1 + o(1)).\end{aligned}\tag{14}$$

Taking  $\eta > \frac{\alpha \log T}{T}$ , and  $\alpha = \frac{D^2\text{tr}(\mathbf{H})}{\lambda_d\text{tr}(\mathbf{H}^2)}$  gives

$$\bar{\mathcal{R}}(K, N; \eta) \geq \frac{\sigma^2\text{tr}(\mathbf{H}) \log T}{4\lambda_d T} (1 + o(1)).$$

Now we recall that

$$\begin{aligned}\bar{\mathcal{R}}^*(K, N) &\leq \bar{\mathcal{R}}(K, N; \eta') = M(K, N; \eta') (1 + o(1)) \\ &= \frac{\sigma^2\text{tr}(\mathbf{H}) \log T}{8\lambda_d T} (1 + o(1)) < \frac{\sigma^2\text{tr}(\mathbf{H}) \log T}{4\lambda_d T} (1 + o(1))\end{aligned}$$

Thus we have that  $\eta^* \leq \frac{\alpha \log T}{T}$ . Next, we give the lower bound of  $\eta^*$ .

When  $\eta < \frac{\log T}{3\lambda_d T}$ , we have that

$$\exp(-2\lambda_d T) = \frac{1}{T^{2/3}} = \omega\left(\frac{\log T}{T}\right) = \omega(\bar{\mathcal{R}}(K, N; \eta')) = \omega(\bar{\mathcal{R}}^*(K, N)).$$

The above equation shows  $\eta^* > \frac{\log T}{3\lambda_d T}$ , which completes the proof.  $\square$

### F.3.2. A DESCRIPTION OF THE RANGE OF OPTIMAL LEARNING RATE, LARGE- $K$ CASE

**Lemma F.6** *Under the conditions in Lemma E.4, and when  $K = \omega(\log N)$ , we have  $\eta^* \in [\frac{\log T}{3\lambda_d T}, o(\frac{1}{N})]$ .*

*Proof.* The proof comprises of three parts. First, we prove that  $\eta^* \geq \frac{\log T}{3\lambda_d T}$  when  $T$  is large. Second, we verify that  $\eta^* \leq \frac{c}{N}$  for sufficiently large  $N$ . Finally, we refine the proof in the second step and justify that  $\eta^* = o(\frac{1}{N})$ . All proofs are carried out by contradiction. The method proceeds as follows: we take a specific  $\eta = \eta'$  and compute its loss, then prove that  $\bar{\mathcal{R}}^*(K, N) > \bar{\mathcal{R}}(K, N; \eta')$  when  $N$  is sufficiently large if  $\eta^*$  does not fall into some interval.

First, by Equation (16), we have

$$\bar{\mathcal{R}}(K, N; \eta') = \frac{\sigma^2 d}{2N} (1 + o(1)).$$

Then we begin our main part of the proof.

*Proof Step I:*  $\eta^* \geq \frac{\log T}{3\lambda_d T}$ .

We assume that  $\eta^* < \frac{\log T}{3\lambda_d T}$ . Observe that  $\bar{\mathcal{R}}^{\text{bias}}(K, N; \eta)$  decreases with  $\eta$ . So

$$\begin{aligned}\bar{\mathcal{R}}^*(K, N) &\geq \bar{\mathcal{R}}^{\text{bias}}(K, N; \eta^*) \geq \bar{\mathcal{R}}^{\text{bias}}(K, N; \eta = \frac{\log T}{3\lambda_d T}) \\ &= \frac{1}{2}(\mathbf{w}_0 - \mathbf{w}^*)^\top (\mathbf{I} - \eta \mathbf{H})^{2T} \mathbf{H}(\mathbf{w}_0 - \mathbf{w}^*) (1 + o(1)) \Big|_{\eta = \frac{\log T}{3\lambda_d T}} \\ &= \left( \frac{1}{2} \tilde{\theta}_d^2 \lambda_d \exp(-2\lambda_d \eta T) \right) (1 + o(1)) \Big|_{\eta = \frac{\log T}{3\lambda_d T}} \\ &= \Theta\left(\frac{1}{T^{\frac{2}{3}}}\right) = \omega\left(\frac{1}{N}\right),\end{aligned}$$

where the first equality is due to [Lemma F.3](#), the second equality is due to [Lemma F.4](#), and the last equality is due to [Assumption 3.1](#).

*Proof Step II:*  $\eta^* \leq \frac{4D^2 d}{\sigma^2 \text{tr}(\mathbf{H}^2) N}$ . We assume that  $\eta^* > \frac{4D^2 d}{\sigma^2 \text{tr}(\mathbf{H}^2) N}$ . By [Equation \(14\)](#), we have

$$\hat{\mathcal{R}}(K, N; \eta) \geq \frac{\eta \sigma^2 \text{tr}(\mathbf{H}^2)}{4D^2} (1 + o(1)) > \frac{\sigma^2 d}{N} (1 + o(1)) > \frac{\sigma^2 d}{2N} (1 + o(1)),$$

which is a contradiction.

A direct corollary is that

$$\bar{\mathcal{R}}^*(K, N) = \hat{\mathcal{R}}(K, N; \eta^*) (1 + o(1)),$$

where

$$\begin{aligned}\hat{\mathcal{R}}(K, N; \eta^*) &= \frac{1}{2}(\mathbf{w}_0 - \mathbf{w}^*)^\top (\mathbf{I} - \eta^* \mathbf{H})^{2T} \mathbf{H}(\mathbf{w}_0 - \mathbf{w}^*) \\ &\quad + \frac{\sigma^2}{N} \text{tr} \left( \frac{(\mathbf{I} - (\mathbf{I} - \eta^* \mathbf{H})^{KN}) ((\mathbf{I} - \eta^* \mathbf{H})^N - (\mathbf{I} - \eta^* \mathbf{H})^{KN})}{\mathbf{I} + (\mathbf{I} - \eta^* \mathbf{H})^N} \right) \\ &\quad + \frac{\eta^* \sigma^2}{2} \langle \mathbf{H}, (\mathbf{I} - (\mathbf{I} - \eta^* \mathbf{H})^{2T}) (2\mathbf{I} - \eta^* \mathbf{H})^{-1} \rangle \\ &= \frac{1}{2} \sum_{i=1}^d (\mathbf{P} \boldsymbol{\theta}_0)_i^2 \lambda_i (1 - \eta^* \lambda_i)^{2T} + \sum_{i=1}^d \frac{\sigma^2}{N} \frac{(1 - \eta^* \lambda_i)^N}{1 + (1 - \eta^* \lambda_i)^N} \\ &\quad + \frac{\eta^* \sigma^2}{4} \text{tr}(\mathbf{H}) - \frac{\eta^* \sigma^2}{4} \sum_{i=1}^d \lambda_i (1 - \eta^* \lambda_i)^{2T} + O((\eta^*)^2) \\ &= \left( \frac{1}{2} \tilde{\theta}_d^2 \lambda_d \exp(-2\lambda_d \eta^* T) + \sum_{i=1}^d \frac{\sigma^2}{N} \frac{e^{-N\eta^* \lambda_i}}{1 + e^{-N\eta^* \lambda_i}} + \frac{\eta^* \sigma^2}{4} \text{tr}(\mathbf{H}) \right) (1 + o(1)).\end{aligned}$$

*Proof Step III:*  $\eta^* = o\left(\frac{1}{N}\right)$ .

We assume that there exists a constant  $\epsilon > 0$  and a sequence  $(N_i)_{i=1}^\infty$  that satisfies  $N_i \rightarrow \infty$  when  $i \rightarrow \infty$  and  $\eta^*(N_i) \geq \frac{\epsilon}{N_i}$  for all  $i$ . As we only conduct our analysis on the sequence  $(N_i)_{i=1}^\infty$ , without loss of generality, we take  $(N_i)_{i=1}^\infty = \mathbb{N}$ .

We define  $f(\delta) = \sum_{i=1}^d \sigma^2 \frac{e^{-\delta\lambda_i}}{1+e^{-\delta\lambda_i}} + \frac{\delta\sigma^2}{4} \text{tr}(\mathbf{H})$ . Then we have

$$f'(\delta) = \frac{\sigma^2}{4} \sum_{i=1}^d \lambda_i - \sum_{i=1}^d \sigma^2 \frac{\lambda_i e^{-\delta\lambda_i}}{(1+e^{-\delta\lambda_i})^2} = \frac{\sigma^2}{4} \sum_{i=1}^d \lambda_i \frac{(1-e^{-\delta\lambda_i})^2}{(1+e^{-\delta\lambda_i})^2} > 0 \text{ when } \delta > 0.$$

So

$$f(\epsilon) > f(0) = \frac{\sigma^2 d}{2N},$$

and

$$\bar{\mathcal{R}}^*(K, N) \geq \frac{1}{N} f(\eta^* N) (1 + o(1)) \geq \frac{1}{N} f(\epsilon) (1 + o(1)) > \frac{\sigma^2 d}{2N} (1 + o(1)) = \bar{\mathcal{R}}(K, N; \eta'),$$

which is a contradiction.  $\square$

### F.3.3. AN APPROXIMATION OF THE EXCESS RISK, SMALL- $K$ CASE

**Lemma F.7** *Let  $\tilde{\theta}_d^2 = \sum_{l=d-n_d+1}^d (\mathbf{P}\boldsymbol{\theta}_0)_l^2$ ,  $\mathbf{H} = \mathbf{P}\mathbf{D}\mathbf{P}^\top$  to be the canonical form under similarity of  $\mathbf{H}$ . Under the conditions in Lemma E.4, for learning rate  $\eta \in \left[\frac{\log KN}{3\lambda_d KN}, \frac{\alpha \log KN}{KN}\right]$  for constant  $\alpha = \frac{D^2 \text{tr}(\mathbf{H})}{\lambda_d \text{tr}(\mathbf{H}^2)}$  and  $K = o(\log N)$ , then we have the approximation of  $\bar{\mathcal{R}}(K, N; \eta)$  as*

$$\begin{aligned} \bar{\mathcal{R}}(K, N; \eta) &= M(K, N; \eta) (1 + o(1)), \\ M(K, N; \eta) &:= \frac{1}{2} \tilde{\theta}_d^2 \lambda_d \exp(-2\lambda_d \eta T) + \frac{\eta \text{tr}(\mathbf{H}) \sigma^2}{4}, \end{aligned}$$

where steps  $T = KN$ .

*Proof.*

From Lemma F.1, we have that  $\bar{\mathcal{R}}(K, N; \eta) = \hat{\mathcal{R}}(K, N; \eta) (1 + o(1))$ , where  $\hat{\mathcal{R}}(K, N; \eta)$  can be written as

$$\begin{aligned} \hat{\mathcal{R}}(K, N; \eta) &= \frac{1}{2} (\mathbf{w}_0 - \mathbf{w}^*)^\top (\mathbf{I} - \eta \mathbf{H})^{2T} \mathbf{H} (\mathbf{w}_0 - \mathbf{w}^*) \\ &\quad + \frac{\sigma^2}{N} \text{tr} \left( \frac{(\mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^{KN}) ((\mathbf{I} - \eta \mathbf{H})^N - (\mathbf{I} - \eta \mathbf{H})^{KN})}{\mathbf{I} + (\mathbf{I} - \eta \mathbf{H})^N} \right) \\ &\quad + \frac{\eta \sigma^2}{2} \langle \mathbf{H}, (\mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^{2T}) (2\mathbf{I} - \eta \mathbf{H})^{-1} \rangle \\ &= \frac{1}{2} \sum_{i=1}^d (\mathbf{P}\boldsymbol{\theta}_0)_i^2 \lambda_i (1 - \eta \lambda_i)^{2T} + \sum_{i=1}^d \frac{\sigma^2}{N} \frac{(1 - \eta \lambda_i)^N}{1 + (1 - \eta \lambda_i)^N} \\ &\quad + \frac{\eta \sigma^2}{4} \text{tr}(\mathbf{H}) - \frac{\eta \sigma^2}{4} \sum_{i=1}^d \lambda_i (1 - \eta \lambda_i)^{2T} + O(\eta^2) \\ &= \underbrace{\left( \frac{1}{2} \tilde{\theta}_d^2 \lambda_d \exp(-2\lambda_d \eta T) + \frac{\eta \sigma^2}{4} \text{tr}(\mathbf{H}) \right)}_{M(K, N; \eta)} (1 + o(1)) + O\left(\frac{1}{N}\right) \end{aligned}$$

$$= \underbrace{\left( \frac{1}{2} \tilde{\theta}_d^2 \lambda_d \exp(-2\lambda_d \eta T) + \frac{\eta \sigma^2}{4} \text{tr}(\mathbf{H}) \right)}_{M(K, N; \eta)} (1 + o(1)), \quad (15)$$

where the second to last equation uses [Lemma F.4](#) and the fact that  $\eta(1 - \eta\lambda_d)^{2T} = o(M(K, N; \eta))$  for  $\eta \in [\frac{\log T}{3\lambda_d T}, \frac{\alpha \log T}{T}]$ , and the last equation uses the fact that when  $K = o(\log N)$ ,  $O(\frac{1}{N}) = o(\frac{\log(N)}{K, N}) = o(M(T; \eta))$ .  $\square$

#### F.3.4. AN APPROXIMATION OF THE EXCESS RISK, LARGE- $K$ CASE

**Lemma F.8** *Under the conditions in [Lemma E.4](#), for  $\eta \in [\frac{\log T}{3\lambda_d T}, o(\frac{1}{N})]$ , and  $K = \omega(\log N)$ , we have*

$$\begin{aligned} \mathbb{E}[\bar{\mathcal{R}}(K, N; \eta)] &= M(K, N; \eta)(1 + o(1)), \\ M(K, N; \eta) &= \frac{1}{2} \tilde{\theta}_d^2 \lambda_d \exp(-2\lambda_d \eta T) + \frac{\eta \text{tr}(\mathbf{H}) \sigma^2}{4} + \frac{\sigma^2 d}{2N}, \end{aligned}$$

where  $\tilde{\theta}_d^2 := \sum_{l=d-n_d+1}^d (\mathbf{P}\boldsymbol{\theta}_0)_l^2$ , and  $\mathbf{P}\mathbf{D}\mathbf{P}^\top$  is the canonical form under similarity of  $\mathbf{H}$ .

*Proof.* Given  $K = O(N^{0.1})$ , one can verify that

$$\lim_{N \rightarrow \infty} K \eta T^{\frac{3}{4}} = \lim_{N \rightarrow \infty} \frac{K^{\frac{7}{4}} N^{\frac{3}{4}}}{N} \eta N = 0.$$

So condition  $K = o(\eta^{-1} T^{-\frac{3}{4}})$  is satisfied, thus by invoking [Lemma F.1](#), we have  $\bar{\mathcal{R}}(K, N; \eta) = \hat{\mathcal{R}}(K, N; \eta)(1 + o(1))$ .

Note that when  $\eta = o(\frac{1}{N})$ , for any  $i \in [1, d]$ , we have

$$(1 - \lambda_i \eta)^N = e^{-\lambda_i \eta N + O(\eta^2 N)} = 1 + o(1).$$

Combining this with [Lemma F.4](#), we have

$$\begin{aligned} \hat{\mathcal{R}}(K, N; \eta) &= \frac{1}{2} (\mathbf{w}_0 - \mathbf{w}^*)^\top (\mathbf{I} - \eta \mathbf{H})^{2T} \mathbf{H} (\mathbf{w}_0 - \mathbf{w}^*) \\ &\quad + \frac{\sigma^2}{N} \text{tr} \left( \frac{(\mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^{KN}) ((\mathbf{I} - \eta \mathbf{H})^N - (\mathbf{I} - \eta \mathbf{H})^{KN})}{\mathbf{I} + (\mathbf{I} - \eta \mathbf{H})^N} \right) \\ &\quad + \frac{\eta \sigma^2}{2} \langle \mathbf{H}, (\mathbf{I} - (\mathbf{I} - \eta \mathbf{H})^{2T}) (2\mathbf{I} - \eta \mathbf{H})^{-1} \rangle \\ &= \frac{1}{2} \sum_{i=1}^d (\mathbf{P}\boldsymbol{\theta}_0)_i^2 \lambda_i (1 - \eta \lambda_i)^{2T} + \sum_{i=1}^d \frac{\sigma^2}{N} \frac{(1 - \eta \lambda_i)^N}{1 + (1 - \eta \lambda_i)^N} \\ &\quad + \frac{\eta \sigma^2}{4} \text{tr}(\mathbf{H}) - \frac{\eta \sigma^2}{4} \sum_{i=1}^d \lambda_i (1 - \eta \lambda_i)^{2T} + O(\eta^2) \\ &= \underbrace{\left( \frac{1}{2} \tilde{\theta}_d^2 \lambda_d \exp(-2\lambda_d \eta T) + \frac{\eta \sigma^2}{4} \text{tr}(\mathbf{H}) \right)}_{M(K, N; \eta)} + \frac{\sigma^2 d}{2N} (1 + o(1)), \quad (16) \end{aligned}$$

which concludes the proof.  $\square$

#### F.4. Step IV: Deriving the Approximately Optimal Learning Rate, Proof of Lemma E.4

After the above preparation, next we present the derivation of approximately optimal learning rate, which is the proof of Lemma E.4. The proof of Lemma E.4 for the small- $K$  case and large- $K$  case follows a similar pattern. First, we minimize the approximate excess risk obtained in Section F.3.3 and Section F.3.4. Then we conduct an error bound analysis and complete the proof.

##### F.4.1. PROOF OF LEMMA E.4, SMALL $K$

##### Part I: Minimizing the Approximation of the Excess Risk

**Lemma F.9** *Under Assumption 3.1 and 3.3, we consider  $K$ -epoch SGD with  $N$  fresh data and learning rate  $\eta$  satisfying  $\eta \in [\frac{\log T}{3\lambda_d T}, \frac{\alpha \log T}{T}]$ , where steps  $T := KN$  and  $\alpha$  is some constant independent of  $T$ , but can depend on  $D$  and  $\lambda_1, \lambda_2, \dots, \lambda_d$ . Then when  $K = o(\log N)$ , the chosen learning rate  $\eta' = \frac{\log \rho T}{2\lambda_d T} = \arg \min_{\eta \in [\frac{\log T}{3\lambda_d T}, \frac{\alpha \log T}{T}]} M(K, N; \eta)$ .*

*Proof.* Given Lemma F.7, we take the derivative of  $M(K, N; \eta)$  with respect to  $\eta$

$$\frac{\partial M}{\partial \eta} = -\tilde{\theta}_d^2 \lambda_d^2 T \exp(-2\lambda_d \eta T) + \frac{\text{tr}(\mathbf{H})\sigma^2}{4}.$$

Define  $\rho := \frac{4\tilde{\theta}_d^2 \lambda_d}{\text{tr}(\mathbf{H})\sigma^2}$ , and we let  $\frac{\partial M}{\partial \eta} = 0$ , then we get

$$\begin{aligned} 0 &= -\rho T \exp(-2\lambda_d \eta T) + 1 \\ \rho T &= \exp(2\lambda_d \eta T) \\ \eta &= \frac{\log \rho T}{2\lambda_d T}. \end{aligned}$$

The above equation completes the proof.  $\square$

##### Part II: Error Bound Analysis

**Lemma F.10** *Consider  $K$ -epoch SGD with  $N$  fresh data and learning rate  $\eta$ . Given a set of learning rate values  $\Gamma$ , and an excess risk estimate that satisfies  $\bar{\mathcal{R}}(K, N; \eta) = M(K, N; \eta)(1 + o(1))$  when  $\eta \in \Gamma$ . Assume that  $\eta' = \arg \min_{\Gamma} M(K, N; \eta)$  and  $\eta^* \in \Gamma$ . Then we have  $\bar{\mathcal{R}}(K, N; \eta'(K, N)) = \bar{\mathcal{R}}^*(K, N)(1 + o(1))$ .*

*Proof.* According to the optimality of  $\eta^*$ , it holds that

$$\bar{\mathcal{R}}^*(K, N) \leq \bar{\mathcal{R}}(K, N; \eta') = M(K, N; \eta')(1 + o(1)).$$

Also, according to the optimality of  $\eta'$ , it holds that

$$M(K, N; \eta')(1 + o(1)) \leq M(K, N; \eta^*)(1 + o(1)) = \bar{\mathcal{R}}^*(K, N)$$

Combining the above two equations gives

$$\bar{\mathcal{R}}(K, N; \eta') = \bar{\mathcal{R}}^*(K, N)(1 + o(1)).$$

$\square$

#### F.4.2. PROOF OF LEMMA E.4, LARGE $K$

##### Part I: Minimizing the Approximation of the Excess Risk

**Lemma F.11** Under Assumption 3.1 and 3.3, we consider  $K$ -epoch SGD with  $N$  fresh data and learning rate  $\eta$  satisfying  $\eta \in [\frac{\log T}{3\lambda_d T}, o(\frac{1}{N})]$ . Then when  $K = \omega(\log N)$ , the chosen learning rate  $\eta' = \frac{\log \rho T}{2\lambda_d T} = \arg \min_{[\frac{\log T}{3\lambda_d T}, o(\frac{1}{N})]} M(K, N; \eta)$ .

*Proof.* Given Lemma F.8, we compute the global minima of  $M(K, N; \eta)$ , we have  $\eta' = \frac{\log T}{2\lambda_d T} + O(\frac{1}{T}) = \arg \min_{\eta \in \mathbb{R}} M(K, N; \eta)$ , which lies in the regime  $[\frac{\log T}{3\lambda_d T}, o(\frac{1}{N})]$  when  $N$  is sufficiently large.  $\square$

**Part II: Error Bound Analysis** The proof of Lemma E.4 concludes directly by applying Lemmas F.6, F.8, F.10 and F.11.

Combine the above two parts and we finish the whole proof.

#### F.5. Proof of Theorem 3.1

Finally, combining the results in the above three steps, we are able to prove our main theorem of the multi-epoch scaling law, as shown in the following.

*Proof.* Notice from Lemma F.1 and Lemma F.4, we have that

$$\begin{aligned} \bar{\mathcal{R}}(K, N; \eta) &= \underbrace{\frac{1}{2} \tilde{\theta}_d^2 \lambda_d (1 - \eta \lambda_d)^{2KN} (1 + o(1))}_{\hat{\mathcal{R}}_1(K, N, \eta)} \\ &\quad + \underbrace{\sum_{i=1}^d \frac{\sigma^2}{N} \frac{(1 - \eta \lambda_i)^N}{1 + (1 - \eta \lambda_i)^N}}_{\hat{\mathcal{R}}_2(K, N, \eta)} \\ &\quad + \underbrace{\frac{\eta \sigma^2}{4} \text{tr}(\mathbf{H}) - \frac{n_d \eta \sigma^2}{4} \lambda_d (1 - \eta \lambda_d)^{2KN} (1 + o(1))}_{\hat{\mathcal{R}}_3(K, N, \eta)} \text{ when } \eta \in \left[ \omega\left(\frac{1}{T}\right), o\left(\frac{1}{T^{\frac{3}{4}}}\right) \right]. \end{aligned}$$

Next, we carefully analyze the magnitude of  $\hat{\mathcal{R}}_1(K, N, \eta)$ ,  $\hat{\mathcal{R}}_2(K, N, \eta)$ , and  $\hat{\mathcal{R}}_3(K, N, \eta)$ , and using these results, we can simplify the excess risk expression.

Now, We take  $\eta = \frac{\log \rho T}{2\lambda_d T} = \frac{\log KN}{2\lambda_d KN} + O(\frac{1}{T})$  in Lemma E.4, then

$$\begin{aligned} (1 - \lambda_d \eta)^{2KN} &= \exp \left( 2KN \log \left( 1 - \frac{\log KN}{2KN} - O\left(\frac{1}{T}\right) \right) \right) \\ &= \exp(-\log KN + O(1)) \\ &= O\left(\frac{1}{T}\right). \end{aligned}$$

Thus

$$\begin{aligned}\widehat{\mathcal{R}}_1(K, N, \eta) &= \frac{1}{2} \tilde{\theta}_d^2 \lambda_d (1 - \lambda_d \eta)^{2KN} \\ &= O\left(\frac{1}{T}\right),\end{aligned}$$

and

$$\begin{aligned}\widehat{\mathcal{R}}_3(K, N, \eta) &= \frac{\sigma^2 \text{tr}(\mathbf{H}) \log T}{8\lambda_d T} - \frac{n_d \sigma^2 \log T}{8\lambda_d T} \lambda_d (1 - \lambda_d \eta)^{2KN} (1 + o(1)) \\ &= \frac{\sigma^2 \text{tr}(\mathbf{H}) \log T}{8\lambda_d T} \left(1 + O\left(\frac{1}{T}\right)\right) \\ &= \frac{\sigma^2 \text{tr}(\mathbf{H}) \log T}{8\lambda_d T} (1 + o(1)) \\ &= \omega(\widehat{\mathcal{R}}_1(K, N, \eta)).\end{aligned}$$

Next, we discuss two scenarios where  $K$  is relatively small and  $K$  is relatively large, to be specific,  $K = o(\log N)$  and  $K = \omega(\log N)$ .

**When  $K = o(\log N)$ ,** We have

$$\begin{aligned}(1 - \lambda_i \eta)^N &= \left(1 - \frac{\log KN}{2KN} \rho_i + O\left(\frac{1}{KN}\right)\right)^N \\ &= \exp\left(N \log\left(1 - \frac{\log KN}{2KN} \rho_i + O\left(\frac{1}{KN}\right)\right)\right) \\ &= \exp\left(-\frac{\log KN}{2K} \rho_i (1 + o(1))\right) \\ &= o(1).\end{aligned}$$

As a consequence,

$$\begin{aligned}\widehat{\mathcal{R}}_2(K, N, \eta) &= \sum_{i=1}^d \frac{\sigma^2}{N} \frac{o(1)}{1 + o(1)} \\ &= o\left(\frac{1}{N}\right).\end{aligned}$$

Meanwhile,

$$\widehat{\mathcal{R}}_3(K, N, \eta) = O\left(\frac{\log KN}{KN}\right) = O\left(\frac{1}{N}\right) = \omega\left(\widehat{\mathcal{R}}_2(K, N, \eta)\right).$$

So

$$\bar{\mathcal{R}}^*(K, N) = \widehat{\mathcal{R}}(K, N; \eta)(1 + o(1)) = \frac{\sigma^2 \text{tr}(\mathbf{H}) \log T}{8\lambda_d T} (1 + o(1)).$$

When  $K = \omega(\log N)$ , we have

$$\begin{aligned}
 (1 - \lambda_i \eta)^N &= \left(1 - \frac{\log KN}{2KN} \rho_i + O\left(\frac{1}{KN}\right)\right)^N \\
 &= \exp\left(N \log\left(1 - \frac{\log KN}{2KN} \rho_i + O\left(\frac{1}{KN}\right)\right)\right) \\
 &= \exp\left(-\frac{\log KN}{2K} \rho_i + O\left(\frac{1}{K}\right)\right) = \exp(o(1)) \\
 &= 1 + o(1).
 \end{aligned}$$

So

$$\begin{aligned}
 \widehat{\mathcal{R}}_2(K, N, \eta) &= \sum_{i=1}^d \frac{\sigma^2}{N} \frac{1 + o(1)}{2 + o(1)} = \frac{\sigma^2 d}{2N} (1 + o(1)) \\
 &= O\left(\frac{1}{N}\right). \\
 \widehat{\mathcal{R}}_3(K, N, \eta) &= O\left(\frac{\log KN}{KN}\right) = o\left(\frac{1}{N}\right) = o\left(\widehat{\mathcal{R}}_2(K, N, \eta)\right).
 \end{aligned}$$

As a result, we have

$$\bar{\mathcal{R}}^*(K, N) = \widehat{\mathcal{R}}(K, N; \eta)(1 + o(1)) = \frac{\sigma^2 d}{2N} (1 + o_N(1)).$$

### F.6. Proof of Theorem 3.2

Now we establish the formulation of  $E(K, N)$  by solving the equation  $\bar{\mathcal{R}}^*(1, T') = \bar{\mathcal{R}}^*(K, N)$ .

When  $K = o(\log N)$ , solving  $\bar{\mathcal{R}}^*(1, T') = \bar{\mathcal{R}}^*(K, N)$ , we get

$$\begin{aligned}
 \frac{\sigma^2 \text{tr}(\mathbf{H}) \log T'}{8\lambda_d T'} (1 + o_{T'}(1)) &= \frac{\sigma^2 \text{tr}(\mathbf{H}) \log T}{8\lambda_d T} (1 + o_T(1)) \\
 \frac{\log T'}{T'} (1 + o_{T'}(1)) &= \frac{\log T}{T} (1 + o_T(1)).
 \end{aligned} \tag{17}$$

According to the definition of the small  $o$  notation, there exists a constant  $\tilde{T}_0$  such that when  $T > \tilde{T}_0$ , the right hand side is smaller than  $\max_{T' \in 1, 2, 3} \frac{\log T'}{T'} (1 + o_{T'}(1))$ . So W.L.O.G, we could assume that  $T' \geq 3$  in the following and use the fact that the function  $\frac{\log x}{x}$  is monotonously decreasing when  $x > 3$ .

**Lemma F.12** *Given  $T'$  and  $N$  satisfying Equation (17), it holds that  $T' \approx T$  when  $T > T_0$  for some constant  $T_0$ .*

*Proof.* Notice that there exists  $T_1$  such that  $|o_T(1)| < \frac{1}{2}$  when  $T > T_1$ , and  $o_{T'}(1)$  is bounded. Furthermore,  $o_{T'}(1) > -1$ , because the left hand side is strictly greater than zero due to the fact



that  $\eta < \frac{1}{D^2}$ . So when  $T > T_1$ , we have

$$c_4 \frac{\log T'}{T'} \leq \frac{3}{2} \frac{\log T}{T} \quad (18)$$

$$c_5 \frac{\log T'}{T'} \geq \frac{1}{2} \frac{\log T}{T} \quad (19)$$

for two constants  $c_4 \leq 1 \leq c_5$ . We claim that  $T' \geq \frac{c_4}{3}T =: \alpha T$  when  $T \geq \frac{1}{\alpha^2}$ ; otherwise,

$$\begin{aligned} c_4 \frac{\log T'}{T'} &\geq c_4 \frac{\log \alpha T}{\alpha T} \\ &= \frac{3 \log \alpha T}{T} \\ &\geq \frac{3 \log T}{2T} \quad \text{when } T \geq \frac{1}{\alpha^2}, \end{aligned}$$

which contradicts Equation (18). We also have  $T' \leq 3c_5 T =: \beta T$  when  $T \geq \beta^2$  by a similar argument; otherwise,

$$\begin{aligned} c_5 \frac{\log T'}{T'} &\leq c_5 \frac{\log \beta T}{\beta T} \\ &= \frac{\log \beta T}{3T} \\ &\leq \frac{\log T}{2T} \quad \text{when } T \geq \beta^2, \end{aligned}$$

which contradicts Equation (19). So  $T' \approx T$  when  $T \geq \min(T_1, \frac{1}{\alpha^2}, \beta^2, \tilde{T}_0) = T_0$ .

Next, we prove the first part in Theorem 3.2, which is  $\mathbb{E}(K, N) = K(1 + o(1))$  when  $K = o(\log N)$ . We define  $F(T) = \frac{\log T}{T}$ ,  $\delta(T) = |o_T(1)|$ , and  $\epsilon(T') = |o_{T'}(1)|$ , so

$$\begin{aligned} F(T')(1 - \epsilon(T')) &\leq F(T)(1 + \delta(T)) \\ F(T')(1 + \epsilon(T')) &\geq F(T)(1 - \delta(T)) \end{aligned}$$

Consequently, we have

$$-F(T)\delta(T) - F(T')\epsilon(T') \leq F(T') - F(T) \leq F(T)\delta(T) + F(T')\epsilon(T'). \quad (20)$$

So due to the convexity of  $F(T)$ ,

$$-\frac{\log T - 1}{T^2}(T' - T) \leq F'(T)(T' - T) \leq F(T') - F(T) \leq F(T)\delta(T) + F(T')\epsilon(T') = \frac{\log T}{T}|o(1)|.$$

Thus we have

$$T' \geq T(1 - o(1)).$$

The above equation completes the proof.  $\square$

Combining Equation (17) and Lemma F.12 gets

$$-\frac{\log T - 1}{T^2}(T - T') \approx -\frac{\log T' - 1}{T'^2}(T - T'). \quad (21)$$

Further using Equation (20),

$$F'(T')(T - T') \leq F(T) - F(T') \leq F(T)\delta(T) + F(T')\epsilon(T') \quad (22)$$

Combining Equation (21) and Equation (22) gives

$$T' \leq T(1 + o(1)).$$

Substituting the definition of  $E(K, N)$  and we get the first part in Theorem 3.2.

**When**  $K = \omega(\log N)$ , solving  $\bar{\mathcal{R}}^*(1, T') = \bar{\mathcal{R}}^*(K, N)$ , we get

$$\frac{\sigma^2 \text{tr}(\mathbf{H}) \log T'}{8\lambda_d T'}(1 + o_{T'}(1)) = \frac{\sigma^2 d}{2N}(1 + o_N(1)). \quad (23)$$

There exists a constant  $\tilde{N}_0$  such that when  $N > \tilde{N}_0$ , the right hand side is smaller than the minimal value of  $\bar{\mathcal{R}}^*(1, T')$  when  $T'$  is finite, that is,  $\min_{T' \in 1, 2, 3} \frac{\sigma^2 \text{tr}(\mathbf{H}) \log T'}{8\lambda_d T'}(1 + o_{T'}(1))$ . So W.L.O.G, we could assume that  $T' \geq 3$  in the following and use the fact that the function  $\frac{\log x}{x}$  is monotonously decreasing when  $x > 3$ .

Now we provide a lemma to give a loose bound of  $T'$  first, and then we apply the lemma to get the formula of  $E(K, N)$ .

**Lemma F.13** *Given  $T'$  and  $N$  satisfying Equation (23). It holds that  $N \leq T' \leq N^{\frac{3}{2}}$  when  $N \geq N_0$  for some constant  $N_0$ .*

*Proof.* Notice that there exists  $N_1$  such that  $|o_N(1)| < \frac{1}{2}$  when  $N > N_1$ , and  $o_{T'}(1)$  is bounded. Furthermore,  $o_{T'}(1) > -1$ , because the left hand side is strictly greater than zero due to the fact that  $\eta < \frac{1}{D^2}$ . So when  $N > N_1$ , for the left side in Equation (23), we have

$$c_6 \frac{\log T'}{T'} \leq \frac{\sigma^2 \text{tr}(\mathbf{H}) \log T'}{8\lambda_d T'}(1 + o_{T'}(1)) \leq c_7 \frac{\log T'}{T'},$$

where  $c_6 < c_7$  are two positive constants. And for the right side,

$$\frac{c_8}{N} \leq \frac{\sigma^2 d}{2N}(1 + o_N(1)) \leq \frac{c_9}{N},$$

where  $c_8 < c_9$  are two positive constants. Then we prove that  $T' \geq N$  when  $N \geq \max\left(e^{\frac{c_9}{c_6}}, 3\right)$ . Otherwise, we have

$$\frac{\sigma^2 \text{tr}(\mathbf{H}) \log T'}{8\lambda_d T'}(1 + o_{T'}(1)) \geq c_6 \frac{\log T'}{T'} \geq c_6 \frac{\log N}{N} \geq \frac{c_9}{N} \geq \frac{\sigma^2 d}{2N}(1 + o_N(1)),$$

which is a contradiction. Then we prove that  $T' \leq N^{\frac{3}{2}}$  when  $N \geq \left(\frac{c_{10}}{c_8}\right)^4$  for some constant  $c_{10}$ . Otherwise, we have

$$\frac{\sigma^2 \text{tr}(\mathbf{H}) \log T'}{8\lambda_d T'} (1 + o_{T'}(1)) \leq c_7 \frac{\log T'}{T'} \leq c_7 \frac{\log N^{\frac{3}{2}}}{N^{\frac{3}{2}}} = \frac{3c_7 \log N}{2 N^{\frac{3}{2}}} \leq \frac{c_{10}}{N^{\frac{5}{4}}} \leq \frac{c_8}{N} \leq \frac{\sigma^2 d}{2N} (1 + o_N(1)),$$

which is another contradiction. The third inequality uses the fact that  $\frac{\log N}{N^{\frac{1}{4}}}$  is bounded. We take

$$N_0 = \max \left( N_1, e^{\frac{c_9}{c_6}}, \left( \frac{c_{10}}{c_8} \right)^4, \tilde{N}_0 \right) \text{ and we prove the claim.} \quad \square$$

Combining [Equation \(23\)](#) and [Lemma F.13](#) gives

$$T' = \Theta(N \log T') = \Theta(N \log N). \quad (24)$$

Again, combining [Equation \(24\)](#) and [Equation \(23\)](#), and we get

$$T' = \frac{\text{tr}(\mathbf{H}) N \log T'}{4\lambda_d d} (1 + o_N(1)) = \frac{\text{tr}(\mathbf{H}) N \log N}{4\lambda_d d} (1 + o_N(1)),$$

and as a direct corollary,

$$E(K, N) = \frac{T'}{N} = \frac{\text{tr}(\mathbf{H}) \log N}{4\lambda_d d} (1 + o_N(1)).$$

The above equation immediately finishes the proof.  $\square$

## Appendix G. Proof Outline for the Solvable Case with Zipf-distributed Data

In this section, we give the proof sketch of [Theorems 4.2](#) and [D.2](#), which characterizes the behavior of  $E(K, N)$  respectively under power spectrum and logarithm power spectrum assumptions. Here we give the proof outline of [Theorem 4.2](#); the proof sketch of [Theorem D.2](#) is similar to that of [Theorem 4.2](#).

Recall the definition of the effective reuse rate

$$E(K, N) := \frac{1}{N} \min\{N' \geq 0 : \bar{\mathcal{R}}^*(1, N') \leq \bar{\mathcal{R}}^*(K, N)\}.$$

So to obtain  $E(K, N)$ , we need to calculate the optimal expected excess risk  $\bar{\mathcal{R}}^*(K, N)$ , which requires us to compute  $\bar{\mathcal{R}}(K, N; \eta)$ . Fortunately, due to the commutativity of the sequential updates and the closed form of all finite-order moments of data, we can obtain a closed-form expression of  $\bar{\mathcal{R}}(K, N; \eta)$  through a direct calculation. The formal statement is given as follows:

**Lemma G.1** *Under [Assumption 4.1](#), the excess risk for  $K$ -epoch training over  $N$  fresh data points, with learning rate  $\eta \geq 0$  can be given by*

$$\bar{\mathcal{R}}(K, N; \eta) = \frac{1}{2} \left\langle \mathbf{P}\Lambda, \left( \mathbf{I} - \mathbf{P} + \mathbf{P}(\mathbf{I} - \eta\Lambda)^{2K} \right)^N \right\rangle,$$

where the expectation is over the randomness of  $\mathbf{w}^*$  and training datasets  $\{\mathbf{x}_i, y_i\}_{i=0}^{N-1}$ .

In the settings of [Theorem 4.2](#),  $\Lambda$  satisfies a power-law spectrum, and we have

$$\bar{\mathcal{R}}(K, N; \eta) = \frac{1}{2} \sum_{i=1}^d \frac{c}{i^a} \left( 1 - \frac{c}{i^{a-b}} \left( 1 - \left( 1 - \frac{\eta}{i^b} \right)^{2K} \right) \right)^N.$$

Recall that our goal is to give a formula of  $\bar{\mathcal{R}}^*(K, N)$ , which can be use to derive the formula for  $E(K, N)$ . To this end, we first simplify the above expression of  $\bar{\mathcal{R}}(K, N; \eta)$ . Specifically, we choose a parameter  $d_1 = \Theta(1)$  such that  $1 - \frac{\eta}{i^b} > 0$  for all  $i > d_1$ . Then  $\bar{\mathcal{R}}(K, N; \eta)$  can be divided into two parts:

$$\begin{aligned} \bar{\mathcal{R}}(K, N; \eta) &= \underbrace{\frac{1}{2} \sum_{i=1}^{d_1} \frac{c}{i^a} \left( 1 - \frac{c}{i^{a-b}} \left( 1 - \left( 1 - \frac{\eta}{i^b} \right)^{2K} \right) \right)^N}_{S_1(K, N; \eta)} \\ &\quad + \underbrace{\frac{1}{2} \sum_{i=d_1+1}^d \frac{c}{i^a} \left( 1 - \frac{c}{i^{a-b}} \left( 1 - \left( 1 - \frac{\eta}{i^b} \right)^{2K} \right) \right)^N}_{S_2(K, N; \eta)}. \end{aligned}$$

First, we derive an approximation for  $S_2(K, N; \eta)$  for the small- $K$  case and large- $K$  case separately. For sufficiently large  $i$ , when  $K$  is small, we can apply a Taylor expansion technique:

$$\begin{aligned} \frac{c}{i^a} \left( 1 - \frac{c}{i^{a-b}} \left( 1 - \left( 1 - \frac{\eta}{i^b} \right)^{2K} \right) \right)^N &\approx \frac{c}{i^a} \left( 1 - \frac{c}{i^{a-b}} \left( 1 - \left( 1 - \frac{2K\eta}{i^b} \right) \right) \right)^N \\ &= \frac{c}{i^a} \left( 1 - \frac{2Kc\eta}{i^a} \right)^N \approx \frac{c}{i^a} e^{-\frac{2KNc\eta}{i^a}}, \end{aligned}$$

where the second approximation uses  $1 - x \approx e^{-x}$  when  $x$  is small. And when  $K$  is large, we have another form of approximation:

$$\frac{c}{i^a} \left( 1 - \frac{c}{i^{a-b}} \left( 1 - \left( 1 - \frac{\eta}{i^b} \right)^{2K} \right) \right)^N \approx \frac{c}{i^a} \left( 1 - \frac{c}{i^{a-b}} \right)^N.$$

These simplified expressions motivate us to give an estimation of  $S_2(K, N, \eta)$  for a large range of learning rate. The formal statements are given in [Lemma H.1](#) and [Lemma H.4](#) respectively.

Unfortunately, we can not use the same way to bound  $S_1(K, N; \eta)$  because Taylor expansion cannot be easily applied for  $i \leq d_1$ . Nevertheless, we can upper bound  $S_1(K, N; \eta)$  by  $d_1 \cdot \max_{i \in [1, d_1]} \frac{c}{i^a} \left( 1 - \frac{c}{i^{a-b}} \left( 1 - \left( 1 - \frac{\eta}{i^b} \right)^{2K} \right) \right)^N$ , and prove that this is an infinitesimal of  $S_2(K, N; \eta)$  for a specific reference learning rate value.

Next we proceed like Step III in [Section E.1](#): we first take a specific reference learning rate  $\eta'$  to narrow down the range of optimal learning rate. This allows us to characterize  $S_2(K, N; \eta^*)$ . Then, notice that

$$\bar{\mathcal{R}}^*(K, N) = S_1(K, N; \eta^*) + S_2(K, N; \eta^*) \leq \bar{\mathcal{R}}(K, N; \eta') = S_1(K, N; \eta') + S_2(K, N; \eta'),$$

and we already have characterizations of  $S_1(K, N; \eta')$ ,  $S_2(K, N; \eta')$  and  $S_2(K, N; \eta^*)$ . Combining these together allows us to give a simple characterization of  $S_1(K, N; \eta^*)$ . Finally, with the characterizations of  $S_1(K, N; \eta^*)$  and  $S_2(K, N; \eta^*)$ , we are able to derive  $E(K, N)$ .

## Appendix H. Proof of Main Results for the Solvable Case with Zipf-distributed Data

### H.1. A Closed Formula for the Excess Risk: Proof of [Lemma G.1](#)

We first write out the update of parameter after  $K$  epochs

$$\begin{aligned} \theta_{KN} &= A^{(K)} \cdots A^{(1)} \theta_0 = \prod_{l=K}^1 A^{(l)} \theta_0 \\ &= \left( \mathbf{I} - \eta \mathbf{x}_{N-1} \mathbf{x}_{N-1}^\top \right)^K \cdots \left( \mathbf{I} - \eta \mathbf{x}_0 \mathbf{x}_0^\top \right)^K \theta_0. \end{aligned}$$

Then we get the excess risk expression as

$$\begin{aligned} \bar{\mathcal{R}}(K, N; \eta) &= \mathbb{E} \frac{1}{2} \theta_{K,N}^\top \mathbf{H} \theta_{K,N} \\ &= \mathbb{E} \frac{1}{2} \theta_0^\top \mathbf{P} \Lambda \left( \mathbf{I} - \eta \mathbf{x}_{N-1} \mathbf{x}_{N-1}^\top \right)^{2K} \cdots \left( \mathbf{I} - \eta \mathbf{x}_0 \mathbf{x}_0^\top \right)^{2K} \theta_0. \end{aligned}$$

[Assumption 4.1](#) gives

$$\begin{aligned} \bar{\mathcal{R}}(K, N; \eta) &= \mathbb{E} \frac{1}{2} \left\langle \theta_0 \theta_0^\top, \mathbf{P} \Lambda \left( \mathbf{I} - \eta \mathbf{x}_{N-1} \mathbf{x}_{N-1}^\top \right)^{2K} \cdots \left( \mathbf{I} - \eta \mathbf{x}_0 \mathbf{x}_0^\top \right)^{2K} \right\rangle \\ &= \frac{1}{2} \left\langle \mathbf{I}, \mathbf{P} \Lambda \left( \mathbb{E} \left( \mathbf{I} - \eta \mathbf{x} \mathbf{x}^\top \right)^{2K} \right)^N \right\rangle. \end{aligned}$$

Direct calculation gives

$$\mathbb{E} \left( \mathbf{x} \mathbf{x}^\top \right)^j = \sum_{i=1}^d \mu_i^{2j-2} p_i \mu_i^2 \mathbf{e}_i \mathbf{e}_i^\top = \mathbf{P} \Lambda^j,$$

and

$$\mathbb{E} \left[ \left( \mathbf{I} - \eta \mathbf{x} \mathbf{x}^\top \right)^{2K} \right] = \mathbf{I} + \sum_{j=1}^{2K} \binom{2K}{j} (-1)^j \eta^j \mathbf{P} \Lambda^j = \mathbf{I} - \mathbf{P} + \mathbf{P} (\mathbf{I} - \eta \Lambda)^{2K}.$$

Then we can write out the excess risk as

$$\bar{\mathcal{R}}(K, N; \eta) = \frac{1}{2} \left\langle \mathbf{P} \Lambda, (\mathbf{I} - \mathbf{P} + \mathbf{P} (\mathbf{I} - \eta \Lambda)^{2K})^N \right\rangle.$$

The above equation completes the proof.

## H.2. Scaling Laws for Power-Law Spectrum: Proof of Theorem 4.1

### H.2.1. PROOF OF THEOREM 4.1: SMALL- $K$ CASE

**An Approximation of  $S_2(K, N; \eta)$ .**

**Lemma H.1** Suppose the assumptions in Theorem 4.2 hold. When  $K = o(N^{\frac{b}{a-b}})$  and  $\eta = \Theta(1)$ , we define the estimator of  $S_2(K, N; \eta)$  as

$$\tilde{S}_2(K, N; \eta) := \frac{1}{2} \sum_{i=d_1+1}^d \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a}}.$$

Then we have  $S_2(K, N; \eta) = \tilde{S}_2(K, N; \eta)(1 + o(1))$ , and  $\tilde{S}_2(K, N; \eta) \asymp \frac{1}{(KN)^{\frac{a-1}{a}}}$ .

*Proof.* By the fact that  $K = o(N^{\frac{b}{a-b}})$ , there exists a constant  $N_2$  such that when  $N \geq N_2$ ,  $K \leq N^{\frac{b}{a-b}}$ . And we define  $F(x) := \frac{c}{x^a} \left( 1 - \frac{c}{x^{a-b}} \left( 1 - \left( 1 - \frac{\eta}{x^b} \right)^{2K} \right) \right)^N$ . Direct observation gives us that under Assumption 4.2,  $\bar{\mathcal{R}}(K, N; \eta) \propto \sum_{i=1}^d F(i)$ . Next we take the derivative of  $F$  and analyze its maximizer.

$$\begin{aligned} F'(x) &= -\frac{ac}{x^{a+1}} \left( 1 - \frac{c}{x^{a-b}} + \frac{c}{x^{a-b}} \left( 1 - \frac{\eta}{x^b} \right)^{2K} \right)^N \\ &\quad + \frac{cN}{x^a} \left( 1 - \frac{c}{x^{a-b}} + \frac{c}{x^{a-b}} \left( 1 - \frac{\eta}{x^b} \right)^{2K} \right)^{N-1} \cdot \Phi(x) \\ &= \frac{c}{x^a} \left( 1 - \frac{c}{x^{a-b}} + \frac{c}{x^{a-b}} \left( 1 - \frac{\eta}{x^b} \right)^{2K} \right)^{N-1} \\ &\quad \left( -\frac{a}{x} \left( 1 - \frac{c}{x^{a-b}} + \frac{c}{x^{a-b}} \left( 1 - \frac{\eta}{x^b} \right)^{2K} \right) + N\Phi(x) \right) \\ &= \frac{c}{x^{2a-b+1}} \left( 1 - \frac{c}{x^{a-b}} + \frac{c}{x^{a-b}} \left( 1 - \frac{\eta}{x^b} \right)^{2K} \right)^{N-1} \cdot G(x). \end{aligned}$$

where we define

$$G(x) := -a \left( x^{a-b} - c + c \left( 1 - \frac{\eta}{x^b} \right)^{2K} \right) + N \left( (a-b)c - (a-b)c \left( 1 - \frac{\eta}{x^b} \right)^{2K} + \frac{2cKb\eta}{x^b} \left( 1 - \frac{\eta}{x^b} \right)^{2K-1} \right),$$

and

$$\Phi(x) := \left( \frac{(a-b)c}{x^{a-b+1}} - \frac{(a-b)c}{x^{a-b+1}} \left( 1 - \frac{\eta}{x^b} \right)^{2K} + \frac{2cKb\eta}{x^{a+1}} \left( 1 - \frac{\eta}{x^b} \right)^{2K-1} \right).$$

We denote the maximizer of  $F(x)$  by  $x_0$ , so  $G(x_0) = 0$ . We claim that:

$$\text{when } N \geq N_2, x_0 \geq \min \left( \left( \frac{KN(a-b)c\eta}{2a} \right)^{\frac{1}{a}}, 6^{\frac{1}{b}} (KN)^{\frac{1}{a}} \right) =: x_1.$$

*Proof of the claim..* Notice that when  $N \geq N_2$ ,

$$\frac{\eta}{x^b} \leq \frac{1}{6(KN)^{\frac{b}{a}}} \leq \frac{1}{6K}.$$

We assume that the claim is wrong, then

$$\begin{aligned} G(x_0) &\geq N \left( (a-b)c - (a-b)c \left( 1 - \frac{\eta}{x^b} \right)^{2K} \right) - ax^{a-b} \\ &\geq \frac{KN(a-b)c\eta - ax^a}{x^b} \\ &\geq \frac{KN(a-b)c\eta}{2x_1^b} \\ &> 0, \end{aligned}$$

which is a contradiction. The third inequality comes from [Lemma I.3](#). □

So  $x_0 = \Omega \left( (KN)^{\frac{1}{a}} \right)$ . Further plugging this into  $G(x_0) = 0$  that

$$G(x_0) = -ax_0^{a-b}(1 + o(1)) + N \left( \frac{2K(a-b)c\eta}{x_0^b}(1 + o(1)) + \frac{2K(a-b)c\eta}{x_0^b}(1 + o(1)) \right) = 0.$$

gives

$$x_0 = \Theta \left( (KN)^{\frac{1}{a}} \right), \quad F(x_0) = \Theta \left( \frac{1}{KN} \right).$$

Then we have

$$\begin{aligned} S_2(K, N; \eta) &= \frac{1}{2} \sum_{i=d_1+1}^{K^{\frac{0.5}{a+0.5b}} (KN)^{\frac{1}{a+0.5b}}} \frac{c}{i^a} \left( 1 - \frac{c}{i^{a-b}} + \frac{c}{i^{a-b}} \left( 1 - \frac{\eta}{i^b} \right)^{2K} \right)^N \\ &\quad + \frac{1}{2} \sum_{K^{\frac{0.5}{a+0.5b}} (KN)^{\frac{1}{a+0.5b}} + 1}^d \frac{c}{i^a} \left( 1 - \frac{c}{i^{a-b}} + \frac{c}{i^{a-b}} \left( 1 - \frac{\eta}{i^b} \right)^{2K} \right)^N \\ &:= J_1 + J_2. \end{aligned}$$

Furthermore, we have

$$J_1 \lesssim K^{\frac{0.5}{a+0.5b}} (KN)^{\frac{1}{a+0.5b}} F(x_0) \lesssim \frac{K^{\frac{0.5}{a+0.5b}} (KN)^{\frac{1}{a+0.5b}}}{KN},$$

and

$$\begin{aligned} J_2 &= \frac{1}{2} \sum_{i=K^{\frac{0.5}{a+0.5b}} (KN)^{\frac{1}{a+0.5b}+1}}^d \frac{c}{i^a} \left( 1 - \frac{c}{i^{a-b}} + \frac{c}{i^{a-b}} \left( 1 - \frac{\eta}{i^b} \right)^{2K} \right)^N \\ &= \frac{1}{2} \sum_{i=K^{\frac{0.5}{a+0.5b}} (KN)^{\frac{1}{a+0.5b}+1}}^d \frac{c}{i^a} \left( 1 - \frac{2Kc\eta}{i^a} + O\left( \frac{K^2}{i^{a+b}} \right) \right)^N \\ &= \frac{1}{2} \sum_{i=K^{\frac{0.5}{a+0.5b}} (KN)^{\frac{1}{a+0.5b}+1}}^d \frac{c}{i^a} e^{N \log \left( 1 - \frac{2Kc\eta}{i^a} + O\left( \frac{K^2}{i^{a+b}} \right) \right)} \\ &= \frac{1}{2} \sum_{i=K^{\frac{0.5}{a+0.5b}} (KN)^{\frac{1}{a+0.5b}+1}}^d \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a} + O\left( \frac{K^2N}{i^{a+b}} \right)} \\ &= \frac{1}{2} \sum_{i=K^{\frac{0.5}{a+0.5b}} (KN)^{\frac{1}{a+0.5b}+1}}^d \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a}} (1 + o(1)). \end{aligned}$$

We define  $K_1(x) = \frac{c}{x^a} e^{\frac{-2KNc\eta}{x^a}}$ . We can derive that  $\arg \max K_1(x) = \Theta \left( (KN)^{\frac{1}{a}} \right)$ , and  $\max K_1(x) = \Theta \left( \frac{1}{KN} \right)$ . So when  $d \geq 3(KN)^{\frac{1}{a}}$ , we have

$$\begin{aligned} J_2 &\geq \frac{1}{2} \sum_{i=(KN)^{\frac{1}{a}}}^{3(KN)^{\frac{1}{a}}} \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a}} (1 + o(1)) \\ &\gtrsim (KN)^{\frac{1}{a}} \times \frac{ce^{-2c\eta}}{KN} \gtrsim \frac{(KN)^{\frac{1}{a}}}{KN}. \end{aligned}$$

We can verify that  $J_1 = o(J_2)$  as a direct consequence. We define

$$\begin{aligned} \tilde{S}_2(K, N; \eta) &= \frac{1}{2} \sum_{i=d_1+1}^d \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a}} \\ &= \frac{1}{2} \sum_{i=d_1+1}^{K^{\frac{0.5}{a+0.5b}} (KN)^{\frac{1}{a+0.5b}}} \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a}} + \frac{1}{2} \sum_{i=K^{\frac{0.5}{a+0.5b}} (KN)^{\frac{1}{a+0.5b}+1}}^d \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a}} \\ &:= \tilde{J}_1 + \tilde{J}_2. \end{aligned}$$



We have  $J_2 = \tilde{J}_2(1 + o(1))$ , and

$$\tilde{J}_1 \leq K^{\frac{0.5}{a+0.5b}} (KN)^{\frac{1}{a+0.5b}} \times \max K_1(x) \lesssim \frac{K^{\frac{0.5}{a+0.5b}} (KN)^{\frac{1}{a+0.5b}}}{KN} = o(\tilde{J}_2).$$

So  $S_2(K, N; \eta) = \tilde{S}_2(K, N; \eta)(1 + o(1))$ .

The matching upper and lower bounds for  $\tilde{S}_2(K, N; \eta)$  comes directly from [Lemma I.7](#).  $\square$

**Additional Technical Lemmas.** By combining the expression of  $\tilde{S}_2(K, N; \eta)$  with [Lemma I.7](#), we get another lemma:

**Lemma H.2** *Suppose the assumptions in [Theorem 4.2](#) hold, and the expression of  $\tilde{S}_2(K, N; \eta)$  is given in [Lemma H.1](#). Then we have  $\frac{\partial}{\partial \eta} \tilde{S}_2(K, N; \eta) \approx -\frac{1}{(KN)^{\frac{a-1}{a}}}$ .*

*Proof.*

$$\begin{aligned} \frac{\partial}{\partial \eta} \tilde{S}_2(K, N; \eta) &= -KN \sum_{i=d_1+1}^d \frac{c}{i^{2a}} e^{\frac{-2KNc\eta}{i^a}} \\ &\approx -\frac{1}{(KN)^{\frac{a-1}{a}}}, \end{aligned}$$

where the second line comes from [Lemma I.7](#).  $\square$

**Lemma H.3** *Suppose the assumptions in [Theorem 4.2](#) hold, and the expression of  $\tilde{S}_2(K, N; \eta)$  is given in [Lemma H.1](#). Consider two learning rate options  $\eta, \eta' = \Theta(1)$  that satisfy  $\eta - \eta' = o(1)$ . Then we have  $\tilde{S}_2(K, N; \eta) = \tilde{S}_2(K, N; \eta')(1 + o(1))$ .*

*Proof.*

$$\begin{aligned} \left| \tilde{S}_2(K, N; \eta) - \tilde{S}_2(K, N; \eta') \right| &= \left| \frac{\partial}{\partial \eta} \tilde{S}_2(K, N; \tilde{\eta}) \right| |(\eta - \eta')| \\ &\approx \frac{1}{(KN)^{\frac{a-1}{a}}} |(\eta - \eta')| \\ &= \tilde{S}_2(K, N; \eta') o(1), \end{aligned}$$

where  $\tilde{\eta} \in [\min(\eta, \eta'), \max(\eta, \eta')] = \Theta(1)$ , and the first line comes from Lagrange's Mean Value Theorem. The second line comes from [Lemma H.2](#), and the last line comes from [Lemma H.1](#).  $\square$

**An Approximation of  $S_1(K, N; \eta)$  under a Reference Learning Rate.** We take  $\eta' = 2 - \frac{(a-1)d_1^{a-b} \log KN}{ac}$ , and we have

$$S_1(K, N; \eta') \leq \frac{d_1 c}{2} \left( 1 - \frac{c}{d_1^{a-b}} + \frac{c}{d_1^{a-b}} \left( 1 - \frac{(a-1)d_1^{a-b} \log KN}{ac} \frac{1}{KN} \right)^{2K} \right)^N.$$

By a Taylor expansion argument, we have

$$\begin{aligned}
 S_1(K, N; \eta') &= \frac{d_1 c}{2} \left( 1 - \frac{2Kc}{d_1^{a-b}} \times \frac{(a-1)d_1^{a-b} \log KN}{ac} (1 + o(1)) \right)^N \\
 &= \frac{d_1 c}{2} \left( 1 - \frac{2(a-1) \log KN}{a} (1 + o(1)) \right)^N \\
 &= \frac{d_1 c}{2} e^{N \log \left( 1 - \frac{2(a-1) \log KN}{a} (1 + o(1)) \right)} \\
 &\approx \frac{1}{(KN)^{\frac{2(a-1)}{a}}} = o(S_2(K, N; \eta')),
 \end{aligned}$$

where the last inequality comes from [Lemma H.1](#).

**The Calculation of  $\bar{\mathcal{R}}^*(K, N)$ .** Under the reference learning rate, we have

$$\begin{aligned}
 \bar{\mathcal{R}}(K, N; \eta') &= S_1(K, N; \eta') + S_2(K, N; \eta') \\
 &= \tilde{S}_2(K, N; \eta') (1 + o(1)) \\
 &= \tilde{S}_2(K, N; 2) (1 + o(1)) \\
 &= \left( \frac{1}{2} \sum_{i=d_1+1}^d \frac{c}{i^a} e^{-\frac{4KNc}{i^a}} \right) (1 + o(1)).
 \end{aligned}$$

Then we prove that  $\eta^* \in [2 - o(1), 2]$ . We prove by contradiction, and assume that there exist a constant  $\epsilon > 0$  and a sequence  $(N_i)_{i=1}^\infty \rightarrow \infty$  such that  $\eta^*(N_i) \leq 2 - \epsilon$  for all  $i \geq 1$ . As we only analyze with respect to the sequence  $(N_i)_{i=1}^\infty$ , without loss of generality, we take  $(N_i)_{i=1}^\infty = \mathbb{N}$ . By [Lemma H.1](#), we have

$$\begin{aligned}
 \bar{\mathcal{R}}^*(K, N) &\geq S_2(K, N; \eta^*) = \tilde{S}_2(K, N; \eta^*) (1 + o(1)) \\
 &\geq \left[ \tilde{S}_2(K, N; 2) + \epsilon \frac{\partial}{\partial \eta} \tilde{S}_2(K, N; 2) \right] (1 + o(1)) > \bar{\mathcal{R}}(K, N; \eta')
 \end{aligned}$$

when  $N$  is sufficiently large, which is a contradiction. So

$$\begin{aligned}
 \bar{\mathcal{R}}^*(K, N) &= S_1(K, N; \eta^*) + S_2(K, N; \eta^*) \\
 &= S_1(K, N; \eta^*) + \tilde{S}_2(K, N; \eta^*) (1 + o(1)) \\
 &= S_1(K, N; \eta^*) + \tilde{S}_2(K, N; 2) (1 + o(1)) \leq \bar{\mathcal{R}}(K, N; \eta').
 \end{aligned}$$

Thus,  $S_1(K, N; \eta^*) = o(\tilde{S}_2(K, N; 2))$ , and  $\bar{\mathcal{R}}^*(K, N) = \tilde{S}_2(K, N; 2) (1 + o(1))$ .

By [Lemma H.1](#) and [Lemma I.7](#), there exist two constants  $C_1$  and  $C_2$  such that  $\bar{\mathcal{R}}^*(K, N) \leq \frac{C_1}{(KN)^{\frac{a-1}{a}}}$  and  $\bar{\mathcal{R}}^*(K, N) \geq \frac{C_2}{(KN)^{\frac{a-1}{a}}}$  when the condition  $d = \Omega(T^{\frac{1}{a}})$  holds. For one-pass case, by [Lemma H.1](#) and [Lemma I.7](#), we have

$$\begin{aligned}
 \bar{\mathcal{R}}^*(1, T') &= \bar{\mathcal{R}}(1, T'; \eta^*(1, T'))|_{d=d} \\
 &\leq \bar{\mathcal{R}}(1, T'; \eta^*(1, T'))|_{d=\infty} \\
 &= \bar{\mathcal{R}}^*(1, T')|_{d=\infty} = \frac{1}{2} \sum_{i=d_1+1}^{\infty} \frac{c}{i^a} e^{-\frac{4KNc}{i^a}} (1 + o(1)) \leq \frac{C_3}{T'^{\frac{a-1}{a}}} \tag{25}
 \end{aligned}$$

and

$$\bar{\mathcal{R}}^*(1, T') = \frac{1}{2} \sum_{i=d_1+1}^d \frac{c}{i^a} e^{-\frac{4KNc}{i^a}} (1 + o(1)) \geq \frac{C_4}{T'^{\frac{a-1}{a}}} \text{ when } d = \Omega\left(T'^{\frac{1}{a}}\right). \quad (26)$$

### H.2.2. PROOF OF [THEOREM 4.1](#): LARGE- $K$ CASE

**An Approximation of  $S_2(K, N; \eta)$ .**

**Lemma H.4** *Suppose the assumptions in [Theorem 4.2](#) hold. When  $K = \omega(N^{\frac{b}{a-b}})$  and  $\eta = \Theta(1)$ , we have  $S_2(K, N; \eta) \approx \frac{1}{N^{\frac{a-1}{a-b}}}$ .*

*Proof.*

There exists  $N_3$  such that when  $N \geq N_3$ , we have  $K \geq N^{\frac{b}{a-b}}$ . Then when  $d \geq 3(KN)^{\frac{1}{a}} \geq 3N^{\frac{1}{a-b}}$ , we give the lower bound of the loss:

$$\begin{aligned} S_2(K, N; \eta) &\geq \frac{1}{2} \sum_{i=N^{\frac{1}{a-b}}}^{3N^{\frac{1}{a-b}}} \frac{c}{i^a} \left(1 - \frac{c}{i^{a-b}}\right)^N \\ &\geq \frac{1}{2} \frac{2N^{\frac{1}{a-b}}}{(3N^{\frac{1}{a-b}})^a} \left(1 - \frac{c}{N}\right)^N \\ &\gtrsim \frac{1}{N^{\frac{a-1}{a-b}}}. \end{aligned}$$

Then we derive the upper bound of the loss:

$$\begin{aligned} S_2(K, N; \eta) &\leq \frac{1}{2} \sum_{i=1}^{\infty} \frac{c}{i^a} \left(1 - \frac{c}{i^{a-b}} + \frac{c}{i^{a-b}} \left(1 - \frac{\eta}{i^b}\right)^{2K}\right)^N \\ &\leq \frac{1}{2} \sum_{i=1}^{N^{\frac{1}{a-b}}} \frac{c}{i^a} \left(1 - \frac{c}{i^{a-b}} + \frac{c}{i^{a-b}} \left(1 - \frac{\eta}{i^b}\right)^{2K}\right)^N + \frac{1}{2} \sum_{i=N^{\frac{1}{a-b}+1}}^{\infty} \frac{c}{i^a}. \end{aligned}$$

When  $K = \omega(N^{\frac{b}{a-b}})$  and  $i \leq N^{\frac{1}{a-b}}$ ,

$$\begin{aligned} \left(1 - \frac{\eta}{i^b}\right)^{2K} &\leq \left(1 - \frac{\eta}{N^{\frac{b}{a-b}}}\right)^{2K} = e^{2K \log\left(1 - \frac{\eta}{N^{\frac{b}{a-b}}}\right)} \\ &\leq e^{-2K \frac{\eta}{N^{\frac{b}{a-b}}}} = o(1). \end{aligned}$$

Then there exists  $N_4$  such that  $(1 - \frac{\eta}{i^b})^{2K} \leq \frac{1}{2}$  when  $N \geq N_4$ . So when  $N \geq \max(N_3, N_4)$ , we have

$$S_2(K, N; \eta) \leq \frac{1}{2} \sum_{i=1}^{N^{\frac{1}{a-b}}} \frac{c}{i^a} \left(1 - \frac{c}{2i^{a-b}}\right)^N + \frac{1}{2} \sum_{i=N^{\frac{1}{a-b}+1}}^{\infty} \frac{c}{i^a}.$$

One can derive that  $\max \frac{c}{i^a} \left(1 - \frac{c}{2i^{a-b}}\right)^N = \Theta\left(\frac{1}{N^{\frac{a-1}{a-b}}}\right)$ . So

$$\begin{aligned} \bar{\mathcal{R}}^*(K, N) &\lesssim \frac{1}{N^{\frac{a-1}{a-b}}} + \frac{1}{N^{\frac{a-1}{a-b}}} \\ &\lesssim \frac{1}{N^{\frac{a-1}{a-b}}}. \end{aligned}$$

And we complete the proof.  $\square$

**An Approximation of  $S_1(K, N; \eta)$  under a Reference Learning Rate.** We take the reference learning rate  $\eta' = 1.5$ , and we have

$$\begin{aligned} S_1(K, N; \eta') &\leq \frac{d_1 c}{2} \left(1 - \frac{c}{d_1^{a-b}} + \frac{c}{d_1^{a-b}} \left(\max\left(0.5, 1 - \frac{1.5}{d_1^b}\right)\right)^{2K}\right)^N \\ &= \frac{d_1 c}{2} (1 - \Theta(1))^N \\ &= o(S_2(K, N; \eta')), \end{aligned}$$

where the last inequality comes from [Lemma H.4](#).

**The Calculation of  $\bar{\mathcal{R}}^*(K, N)$ .** Under the reference learning rate, we have

$$\begin{aligned} \bar{\mathcal{R}}(K, N; \eta') &= S_1(K, N; \eta') + S_2(K, N; \eta') \\ &= S_2(K, N; \eta')(1 + o(1)) \end{aligned}$$

It is obvious that  $\eta^* \in [1, 2]$ . We know that

$$\bar{\mathcal{R}}^*(K, N) = S_1(K, N; \eta^*) + S_2(K, N; \eta^*) \leq \bar{\mathcal{R}}(K, N; \eta') = S_2(K, N; \eta')(1 + o(1)).$$

By [Lemma H.4](#), we have

$$S_2(K, N; \eta^*) = \Theta\left(N^{-\frac{a-1}{a-b}}\right) \quad \text{and} \quad S_2(K, N; \eta') = \Theta\left(N^{-\frac{a-1}{a-b}}\right),$$

which directly implies that

$$S_1(K, N; \eta^*) = O\left(N^{-\frac{a-1}{a-b}}\right), \quad \bar{\mathcal{R}}^*(K, N) = \Theta\left(N^{-\frac{a-1}{a-b}}\right).$$

### H.3. $E(K, N)$ for Power-Law Spectrum: Proof of Theorem 4.2

#### H.3.1. PROOF OF THEOREM 4.2, SMALL- $K$ CASE

Let  $T'$  be defined implicitly by equating the averaged risks at their optimal step sizes:

$$\bar{\mathcal{R}}^*(1, T') = \bar{\mathcal{R}}^*(K, N). \quad (27)$$

We claim that

$$\left(\frac{C_4}{C_1}\right)^{\frac{a}{a-1}} T \leq T' \leq \left(\frac{C_3}{C_2}\right)^{\frac{a}{a-1}} T. \quad (28)$$

*Proof.* We argue by contradiction, considering two exclusive violations of Equation (28).

1. **Case 1:**  $T' > \left(\frac{C_3}{C_2}\right)^{\frac{a}{a-1}} T$ . By the risk bounds encoded by  $(C_2, C_3)$  for one-pass training with  $T'$  fresh data and by  $(C_1, C_4)$  for  $K$ -epoch training with  $N$  fresh data, this inequality forces

$$\bar{\mathcal{R}}^*(1, T') < \bar{\mathcal{R}}^*(K, N),$$

which contradicts the defining equality Equation (27).

2. **Case 2:**  $T' < \left(\frac{C_4}{C_1}\right)^{\frac{a}{a-1}} T$ . given  $d = \Omega(T^{\frac{1}{a}})$  we still have  $d = \Omega((T')^{1/a})$ . The same risk comparisons then yield

$$\bar{\mathcal{R}}^*(1, T') > \bar{\mathcal{R}}^*(K, N),$$

again contradicting Equation (27).

Both contradictions rule out violations; hence Equation (28) holds.  $\square$

Therefore, the desired characterization of  $E(K, N)$  follows directly from Lemma I.8.

#### H.3.2. PROOF OF THEOREM 4.2, LARGE- $K$ CASE

By Theorem 4.1, there exist constants  $C_5, C_6 > 0$  such that, given  $d = \Omega(T^{\frac{1}{a}})$ ,

$$\frac{C_6}{N^{\frac{a-1}{a-b}}} \leq \bar{\mathcal{R}}^*(K, N) \leq \frac{C_5}{N^{\frac{a-1}{a-b}}}. \quad (29)$$

Let  $T'$  be defined by equating the averaged risks at their optimal step sizes:

$$\bar{\mathcal{R}}^*(K, N) = \bar{\mathcal{R}}^*(1, T'). \quad (30)$$

Combining Equation (29), Equation (30) with Equation (25), Equation (26), we claim that

$$\left(\frac{C_4}{C_5}\right)^{\frac{a}{a-1}} N^{\frac{a}{a-b}} \leq T' \leq \left(\frac{C_3}{C_6}\right)^{\frac{a}{a-1}} N^{\frac{a}{a-b}}. \quad (31)$$

*Proof of the claim.* We argue by contradiction.

1. **Upper violation.** If  $T' > \left(\frac{C_3}{C_6}\right)^{\frac{a}{a-1}} N^{\frac{a}{a-b}}$ , then by Equation (25) and Equation (29) (lower bound),

$$\bar{\mathcal{R}}^*(1, T') \leq \frac{C_3}{(T')^{\frac{a-1}{a}}} < \frac{C_6}{N^{\frac{a-1}{a-b}}} \leq \bar{\mathcal{R}}^*(K, N),$$

which contradicts Equation (30).

2. **Lower violation.** If  $T' < (\frac{C_4}{C_5})^{\frac{a}{a-1}} N^{\frac{a}{a-b}}$ , then the condition  $d = \Omega(T'^{\frac{1}{a}})$  gives

$$d = \Omega\left(N^{\frac{1}{a-b}}\right) = \Omega\left((T')^{\frac{1}{a}}\right).$$

Using Equation (26) and Equation (29) (upper bound),

$$\bar{\mathcal{R}}^*(1, T') \geq \frac{C_4}{(T')^{\frac{a-1}{a}}} > \frac{C_5}{N^{\frac{a-1}{a-b}}} \geq \bar{\mathcal{R}}^*(K, N),$$

again contradicting Equation (30).

Both contradictions are impossible; hence Equation (31) holds.  $\square$

The characterization of  $E(K, N)$  follows directly by the claim.

#### H.4. Scaling Laws for Logarithmic Power-Law Spectrum: Proof of Theorem D.1

Similar to the proof of Theorem 4.1, the proof of Theorem D.2 consists of two parts: First part is the case when  $K = o(\log^b N)$ , and the second part is the case when  $K = \omega(\log^b N)$ .

Before we begin our main part of the proof, note that for all  $\eta = \Theta(1)$  and  $\eta \leq 2$ , there exists  $d_2 = \Theta(1) > 0$  such that  $1 - \frac{\eta}{\log^b(i+1)} > 0$  when  $i > d_2$ . Then we divide the loss into two parts:

$$\begin{aligned} \bar{\mathcal{R}}(K, N; \eta) &= \frac{1}{2} \sum_{i=1}^d \frac{c}{i^a} \left( 1 - \frac{c \log^b(i+1)}{i^a} + \frac{c \log^b(i+1)}{i^a} \left( 1 - \left( 1 - \frac{\eta}{\log^b(i+1)} \right)^{2K} \right) \right)^N \\ &= \underbrace{\frac{1}{2} \sum_{i=1}^{d_2} \frac{c}{i^a} \left( 1 - \frac{c \log^b(i+1)}{i^a} + \frac{c \log^b(i+1)}{i^a} \left( 1 - \left( 1 - \frac{\eta}{\log^b(i+1)} \right)^{2K} \right) \right)^N}_{V_1(K, N; \eta)} \\ &\quad + \underbrace{\sum_{d_2+1}^d \frac{c}{i^a} \left( 1 - \frac{c \log^b(i+1)}{i^a} + \frac{c \log^b(i+1)}{i^a} \left( 1 - \left( 1 - \frac{\eta}{\log^b(i+1)} \right)^{2K} \right) \right)^N}_{V_2(K, N; \eta)}. \end{aligned}$$

##### H.4.1. PROOF OF THEOREM D.1: SMALL- $K$ CASE

**An Approximation of  $V_2(K, N; \eta)$ .**

**Lemma H.5** Suppose the assumptions in Theorem D.2 hold. When  $K = o(\log^b N)$ , we define the estimate of  $V(K, N; \eta)$  as

$$\tilde{V}_2(K, N; \eta) := \frac{1}{2} \sum_{i=1}^d \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a}}.$$

Then we have  $V_2(K, N; \eta) = \tilde{V}(K, N; \eta)(1 + o(1))$ , and  $\tilde{V}_2(K, N; \eta) \approx \frac{1}{(KN)^{\frac{a-1}{a}}}$ .

*Proof of Lemma H.5.* We first define a function

$$W(x) := \frac{c}{x^a} \left( 1 - \frac{c \log^b(x+1)}{x^a} \left( 1 - \left( 1 - \frac{\eta}{\log^b(x+1)} \right)^{2K} \right) \right)^N.$$

Direct observation gives us that under [Assumption D.1](#),  $\bar{\mathcal{R}}(K, N; \eta) \propto \sum_{i=1}^d W(i)$ . Similiarly we take the derivative of  $W$ .

$$\begin{aligned} W'(x) &= -\frac{ac}{x^{a+1}} \left( 1 - \frac{c \log^b(x+1)}{x^a} + \frac{c \log^b(x+1)}{x^a} \left( 1 - \frac{\eta}{\log^b(x+1)} \right)^{2K} \right)^N \\ &\quad + \frac{cN}{x^a} \left( 1 - \frac{c \log^b(x+1)}{x^a} + \frac{c \log^b(x+1)}{x^a} \left( 1 - \frac{\eta}{\log^b(x+1)} \right)^{2K} \right)^{N-1} \\ &\quad \left( \left( \frac{ac \log^b(x+1)}{x^{a+1}} - \frac{bc \log^{b-1}(x+1)}{x^a(x+1)} \right) \left( 1 - \left( 1 - \frac{\eta}{\log^b(x+1)} \right)^{2K} \right) \right. \\ &\quad \left. + \frac{2cK \log^b(x+1)}{x^a} \left( 1 - \frac{\eta}{\log^b(x+1)} \right)^{2K-1} \frac{b\eta}{(x+1) \log^{b+1}(x+1)} \right) \\ &= \frac{c}{x^{2a+1}} \left( 1 - \frac{c \log^b(x+1)}{x^a} + \frac{c \log^b(x+1)}{x^a} \left( 1 - \frac{\eta}{\log^b(x+1)} \right)^{2K} \right)^{N-1} \\ &\quad \left( -a \left( x^a - c \log^b(x+1) + c \log^b(x+1) \left( 1 - \frac{\eta}{\log^b(x+1)} \right)^{2K} \right) \right. \\ &\quad \left. + N \left( \left( ac \log^b(x+1) - bc \log^{b-1}(x+1) \frac{x}{x+1} \right) \left( 1 - \left( 1 - \frac{\eta}{\log^b(x+1)} \right)^{2K} \right) \right. \right. \\ &\quad \left. \left. + \frac{2cKb\eta}{\log(x+1)} \left( 1 - \frac{\eta}{\log^b(x+1)} \right)^{2K-1} \frac{x}{x+1} \right) \right). \end{aligned}$$

We define

$$\begin{aligned} G(x) &= -a \left( x^a - c \log^b(x+1) + c \log^b(x+1) \left( 1 - \frac{\eta}{\log^b(x+1)} \right)^{2K} \right) \\ &\quad + N \left( \left( ac \log^b(x+1) - bc \log^{b-1}(x+1) \frac{x}{x+1} \right) \left( 1 - \left( 1 - \frac{\eta}{\log^b(x+1)} \right)^{2K} \right) \right. \\ &\quad \left. + \frac{2cKb\eta}{\log(x+1)} \left( 1 - \frac{\eta}{\log^b(x+1)} \right)^{2K-1} \frac{x}{x+1} \right), \end{aligned}$$

and  $x_0$  is defined to be the maximum of  $W(x)$ , so  $G(x_0) = 0$ .

$$G(x) \geq N \log^b(x+1) \left( ac - \frac{bc}{\log(x+1)} \frac{x}{x+1} \right) \left( 1 - \left( 1 - \frac{\eta}{\log^b x} \right)^{2K} \right) - ax^a$$

$$\begin{aligned}
 &\geq N(a-b)c \log^b(x+1) \left( 1 - \left( 1 - \frac{\eta}{\log^b(x+1)} \right)^{2K} \right) - ax^a \\
 &= N(a-b)c \log^b(x+1) \times \frac{\eta}{\log^b(x+1)} \left( \sum_{i=0}^{2K-1} \left( 1 - \frac{\eta}{\log^b(x+1)} \right)^i \right) - ax^a \\
 &\geq N(a-b)c\eta - ax^a.
 \end{aligned}$$

So  $x_0 = \Omega\left(N^{\frac{1}{a}}\right)$  is an direct conclusion by  $G(x_0) = 0$ . Also , by solving  $G(x_0) = 0$ , we can get the approximation of  $x_0$  as

$$\begin{aligned}
 G(x_0) &= -ax_0^a(1+o(1)) \\
 &\quad + N \left( ac \log^b(x_0+1)(1+o(1)) \times \frac{2K\eta}{\log^b(x_0+1)}(1+o(1)) + O\left(\frac{K}{\log N}\right) \right) \\
 &= -ax_0^a(1+o(1)) + 2KNac\eta(1+o(1)) = 0,
 \end{aligned}$$

thus we have

$$x_0 = \Theta\left((KN)^{\frac{1}{a}}\right), \quad W(x_0) = \Theta\left(\frac{1}{KN}\right).$$

There exists a constant  $N_5$  such that  $K \leq \log^b N$  when  $N \geq N_5$ . So when  $N \geq N_5$  and  $d \geq 3(KN)^{\frac{1}{a}} \geq 3(KN)^{\frac{1}{a}} \left(\frac{K}{\log^b N}\right)^{\frac{1}{2a}}$ , we have

$$\begin{aligned}
 V_2(K, N; \eta) &= \frac{1}{2} \sum_{i=d_2+1}^{(KN)^{\frac{1}{a}} \left(\frac{K}{\log^b N}\right)^{\frac{1}{2a}}} \frac{c}{i^a} \left( 1 - \frac{c \log^b(i+1)}{i^a} \left( 1 - \left( 1 - \frac{\eta}{\log^b(i+1)} \right)^{2K} \right) \right)^N \\
 &\quad + \frac{1}{2} \sum_{(KN)^{\frac{1}{a}} \left(\frac{K}{\log^b N}\right)^{\frac{1}{2a}}}^d \frac{c}{i^a} \left( 1 - \frac{c \log^b(i+1)}{i^a} \left( 1 - \left( 1 - \frac{\eta}{\log^b(i+1)} \right)^{2K} \right) \right)^N \\
 &:= \psi_1 + \psi_2.
 \end{aligned}$$

Furthermore,

$$\psi_1 \lesssim (KN)^{\frac{1}{a}} \left(\frac{K}{\log^b N}\right)^{\frac{1}{2a}} \times W(x_0) \lesssim \frac{(KN)^{\frac{1}{a}} \left(\frac{K}{\log^b N}\right)^{\frac{1}{2a}}}{KN},$$



and

$$\begin{aligned}
 \psi_2 &= \frac{1}{2} \sum_{i=(KN)^{\frac{1}{a}} \left(\frac{K}{\log^b N}\right)^{\frac{1}{2a}}}^d \frac{c}{i^a} \left(1 - \frac{c \log^b(i+1)}{i^a} + \frac{c \log^b(i+1)}{i^a} \left(1 - \frac{\eta}{\log^b(i+1)}\right)^{2K}\right)^N \\
 &= \frac{1}{2} \sum_{i=(KN)^{\frac{1}{a}} \left(\frac{K}{\log^b N}\right)^{\frac{1}{2a}}}^d \frac{c}{i^a} \left(1 - \frac{2Kc\eta}{i^a} + O\left(\frac{K^2}{i^a \log^b(i+1)}\right)\right)^N \\
 &= \frac{1}{2} \sum_{i=(KN)^{\frac{1}{a}} \left(\frac{K}{\log^b N}\right)^{\frac{1}{2a}}}^d \frac{c}{i^a} e^{N \log\left(1 - \frac{2Kc\eta}{i^a} + O\left(\frac{K^2}{i^a \log^b(i+1)}\right)\right)} \\
 &= \frac{1}{2} \sum_{i=(KN)^{\frac{1}{a}} \left(\frac{K}{\log^b N}\right)^{\frac{1}{2a}}}^d \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a} + O\left(\frac{K^2 N}{i^{2a}}\right) + O\left(\frac{K^2 N}{i^a \log^b(i+1)}\right)} \\
 &= \frac{1}{2} \sum_{i=(KN)^{\frac{1}{a}} \left(\frac{K}{\log^b N}\right)^{\frac{1}{2a}}}^d \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a}} (1 + o(1)).
 \end{aligned}$$

We recall  $K_1(x) = \frac{c}{x^a} e^{\frac{-2KNc\eta}{x^a}}$ . We can verify that  $\arg \max K_1(x) = \Theta\left((KN)^{\frac{1}{a}}\right)$  and  $\max K_1(x) = \Theta\left(\frac{1}{KN}\right)$  through a direct calculation. So for  $\psi_2$  we have

$$\begin{aligned}
 \psi_2 &\geq \frac{1}{2} \sum_{i=(KN)^{\frac{1}{a}}}^{3(KN)^{\frac{1}{a}}} \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a}} (1 + o(1)) \\
 &\gtrsim \frac{(KN)^{\frac{1}{a}}}{KN}.
 \end{aligned}$$

We can verify that  $\psi_1 = o(\psi_2)$  as a direct consequence. We define

$$\begin{aligned}
 \tilde{V}_2(K, N; \eta) &= \frac{1}{2} \sum_{i=d_2+1}^d \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a}} \\
 &= \frac{1}{2} \sum_{i=d_2+1}^{(KN)^{\frac{1}{a}} \left(\frac{K}{\log^b N}\right)^{\frac{1}{2a}}} \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a}} + \frac{1}{2} \sum_{i=(KN)^{\frac{1}{a}} \left(\frac{K}{\log^b N}\right)^{\frac{1}{2a}}}^d \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a}} \\
 &:= \tilde{\psi}_1 + \tilde{\psi}_2.
 \end{aligned}$$

We have  $\psi_2 = \tilde{\psi}_2(1 + o(1))$ , and

$$\tilde{\psi}_1 \lesssim \frac{(KN)^{\frac{1}{a}} \left( \frac{K}{\log^b N} \right)^{\frac{1}{2a}}}{KN} = o(\tilde{\psi}_2).$$

So  $V_2(K, N; \eta) = \tilde{V}_2(K, N; \eta)(1 + o(1))$ .

Finally, we derive a matching upper and lower bound for  $\tilde{V}_2(K, N; \eta)$  and conclude the proof:

$$\begin{aligned} \tilde{V}_2(K, N; \eta) &\geq \tilde{J}_2 \gtrsim J_2 \gtrsim \frac{1}{(KN)^{\frac{a-1}{a}}}. \\ \tilde{V}_2(K, N; \eta) &= \frac{1}{2} \sum_{i=d_2+1}^{(KN)^{\frac{1}{a}}} \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a}} + \frac{1}{2} \sum_{i=(KN)^{\frac{1}{a}}+1}^d \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a}} \\ &\leq \frac{1}{2} \sum_{i=1}^{(KN)^{\frac{1}{a}}} \frac{c}{i^a} e^{\frac{-2KNc\eta}{i^a}} + \frac{1}{2} \sum_{i=(KN)^{\frac{1}{a}}+1}^d \frac{c}{i^a} \\ &\lesssim \frac{(KN)^{\frac{1}{a}}}{KN} + \frac{1}{(KN)^{\frac{a-1}{a}}} \lesssim \frac{1}{(KN)^{\frac{a-1}{a}}}. \end{aligned}$$

Then we complete the proof.  $\square$

Notice that  $\tilde{V}_2(K, N; \eta)$  and  $\tilde{S}_2(K, N; \eta)$  are identical to each other, so we can directly apply [Lemma H.2](#) and [Lemma H.3](#) in the remaining proof of [Theorem D.2](#).

**An Approximation of  $V_1(K, N; \eta)$  under a Reference Learning Rate.** We take  $\eta' = 2 \log^b(2) - \epsilon$ , where  $\epsilon := \frac{(a-1)d_2^a \log KN}{ac}$ , and we have

$$\begin{aligned} V_1(K, N; \eta') &\leq \frac{d_2 c}{2} \left( 1 - \frac{c \log^b(2)}{d_2^a} + \frac{c \log^b(2)}{d_2^a} \left( 1 - \frac{\epsilon}{\log^b(2)} \right)^{2K} \right)^N \\ &= \frac{d_2 c}{2} \left( 1 - \frac{2Kc \log^b(2)}{d_2^a} \times \frac{\epsilon}{\log^b(2)} (1 + o(1)) \right)^N \\ &= \frac{d_2 c}{2} \left( 1 - \frac{2(a-1) \log KN}{a} \frac{1}{N} (1 + o(1)) \right)^N \\ &= \frac{d_2 c}{2} e^{N \log \left( 1 - \frac{2(a-1) \log KN}{a} \frac{1}{N} (1 + o(1)) \right)} \\ &\approx \frac{1}{(KN)^{\frac{2(a-1)}{a}}} = o(V_2(K, N; \eta')), \end{aligned}$$

where the last inequality comes from [Lemma H.5](#).

**The Calculation of  $\bar{\mathcal{R}}^*(K, N)$ .** Under the reference learning rate, we have

$$\begin{aligned}\bar{\mathcal{R}}(K, N; \eta') &= V_1(K, N; \eta') + V_2(K, N; \eta') \\ &= \tilde{V}_2(K, N; \eta')(1 + o(1)) \\ &= \tilde{V}_2(K, N; 2)(1 + o(1)) \\ &= \left( \frac{1}{2} \sum_{i=d_1+1}^d \frac{c}{i^a} e^{\frac{-4 \log^b(2) K N c}{i^a}} \right) (1 + o(1)).\end{aligned}$$

Then we prove that  $\eta^* \in [2 \log^b(2) - o(1), 2 \log^b(2)]$ . We prove by contradiction, and assume that there exist a constant  $\epsilon > 0$  and a sequence  $(N_i)_{i=1}^\infty \rightarrow \infty$  such that  $\eta^*(N_i) \leq 2 \log^b(2) - \epsilon$  for all  $i \geq 1$ . As we only analyze with respect to the sequence  $(N_i)_{i=1}^\infty$ , without loss of generality, we take  $(N_i)_{i=1}^\infty = \mathbb{N}$ . By [Lemma H.1](#), we have

$$\begin{aligned}\bar{\mathcal{R}}^*(K, N) &\geq V_2(K, N; \eta^*) = \tilde{V}_2(K, N; \eta^*)(1 + o(1)) \\ &\geq \left[ \tilde{V}_2(K, N; 2) + \epsilon \frac{\partial}{\partial \eta} \tilde{V}_2(K, N; 2 \log^b(2)) \right] (1 + o(1)) > \bar{\mathcal{R}}(K, N; \eta')\end{aligned}$$

when  $N$  is sufficiently large, which is a contradiction. So

$$\begin{aligned}\bar{\mathcal{R}}^*(K, N) &= V_1(K, N; \eta^*) + V_2(K, N; \eta^*) \\ &= V_1(K, N; \eta^*) + \tilde{V}_2(K, N; \eta^*)(1 + o(1)) \\ &= V_1(K, N; \eta^*) + \tilde{V}_2(K, N; 2 \log^b(2))(1 + o(1)) \leq \bar{\mathcal{R}}(K, N; \eta').\end{aligned}$$

So  $V_1(K, N; \eta^*) = o(\tilde{V}_2(K, N; 2 \log^b(2)))$ , and  $\bar{\mathcal{R}}^*(K, N) = \tilde{V}_2(K, N; 2 \log^b(2))(1 + o(1)) \approx \frac{1}{(KN)^{\frac{a-1}{a}}}$ .

#### H.4.2. PROOF OF [THEOREM D.1](#), LARGE- $K$ CASE

**An Approximation of  $V_2(K, N; \eta)$ .**

**Lemma H.6** Suppose the assumptions [Theorem D.2](#) hold. When  $K = \omega(\log^b N)$ , we have  $V_2(K, N; \eta) \approx \frac{1}{(N \log^b N)^{\frac{a-1}{a}}}$ .

*Proof of [Lemma H.6](#).* By  $K = \omega(\log^b N)$ , there exists a constant  $N_6 > 0$  such that  $K > \log^b N$  when  $N \geq N_6$ . We notice that when  $i = \Theta\left((N \log^b N)^{\frac{1}{a}}\right)$ ,  $\log(i+1) = \Theta(\log N)$ . Then, when

$N \geq N_6$  and  $d \geq 3(KN)^{\frac{1}{a}} \geq 3(N \log^b N)^{\frac{1}{a}}$ , we have

$$\begin{aligned} V_2(K, N; \eta) &\geq \frac{1}{2} \sum_{i=(N \log^b N)^{\frac{1}{a}}}^{3(N \log^b N)^{\frac{1}{a}}} \frac{c}{i^a} \left(1 - \frac{c \log^b(i+1)}{i^a}\right)^N \\ &\geq \frac{1}{2} \frac{2(N \log^b N)^{\frac{1}{a}}}{3^a N \log^b N} (1 - \frac{c_{11}}{N})^N \\ &\gtrsim \frac{1}{(N \log^b N)^{\frac{a-1}{a}}}. \end{aligned}$$

For the upper bound, we have

$$\begin{aligned} \bar{\mathcal{R}}(K, N; \eta) &\leq \frac{1}{2} \sum_{i=1}^{\infty} \frac{c}{i^a} \left(1 - \frac{c \log^b(i+1)}{i^a} + \frac{c \log^b(i+1)}{i^a} \left(1 - \frac{\eta}{\log^b(i+1)}\right)^{2K}\right)^N \\ &\leq \frac{1}{2} \sum_{i=1}^{(N \log^b N)^{\frac{1}{a}}} \frac{c}{i^a} \left(1 - \frac{c \log^b(i+1)}{i^a} + \frac{c \log^b(i+1)}{i^a} \left(1 - \frac{\eta}{\log^b(i+1)}\right)^{2K}\right)^N \\ &\quad + \frac{1}{2} \sum_{i=(N \log^b N)^{\frac{1}{a}}+1}^{\infty} \frac{c}{i^a}. \end{aligned}$$

When  $K = \omega(\log^b N)$  and  $i \leq (N \log^b N)^{\frac{1}{a}}$ ,

$$\begin{aligned} \left(1 - \frac{\eta}{\log^b(i+1)}\right)^K &\leq \left(1 - \frac{c_{12}}{\log^b N}\right)^K = e^{K \log\left(1 - \frac{c_{12}}{\log^b N}\right)} \\ &\leq e^{-K \frac{c_{12}}{\log^b N}} = o(1). \end{aligned}$$

So there exists  $N_7$  such that when  $N \geq N_7$ ,  $\left(1 - \frac{\eta}{\log^b(i+1)}\right)^K \leq \frac{1}{2}$ , and when  $N \geq \max(N_6, N_7)$ ,

$$\bar{\mathcal{R}}(K, N; \eta) \leq \frac{1}{2} \sum_{i=1}^{(N \log^b N)^{\frac{1}{a}}} \frac{c}{i^a} \left(1 - \frac{c \log^b(i+1)}{2i^a}\right)^N + \frac{1}{2} \sum_{i=(N \log^b N)^{\frac{1}{a}}+1}^{\infty} \frac{c}{i^a}.$$

One can derive that  $\max_x \frac{c}{x^a} \left(1 - \frac{c \log^b(x+1)}{2x^a}\right)^N = \Theta\left(\frac{1}{N \log^b N}\right)$ .

So finally, we have

$$\begin{aligned} V_2(K, N; \eta) \leq \bar{\mathcal{R}}(K, N; \eta) &\lesssim \frac{1}{(N \log^b N)^{\frac{a-1}{a}}} + \frac{1}{(N \log^b N)^{\frac{a-1}{a}}} \\ &\lesssim \frac{1}{(N \log^b N)^{\frac{a-1}{a}}}, \end{aligned}$$

and we get the result.  $\square$

**An Approximation of  $V_1(K, N; \eta)$  under a Reference Learning Rate.** We take  $\eta' = 1.5 \log^b(2)$ , and we have

$$\begin{aligned} V_1(K, N; \eta') &\leq \frac{d_2 c}{2} \left( 1 - \frac{c \log^b(2)}{d_2^a} + \frac{c \log^b(2)}{d_2^a} \max \left( 0.5, 1 - \frac{1.5 \log^b(2)}{\log^b(d_2 + 1)} \right)^{2K} \right)^N \\ &= \frac{d_1 c}{2} (1 - \Theta(1))^N \\ &= o(V_2(K, N; \eta')), \end{aligned}$$

where the last inequality comes from [Lemma H.4](#).

**The Calculation of  $\bar{\mathcal{R}}^*(K, N)$ .** Under the reference learning rate, we have

$$\begin{aligned} \bar{\mathcal{R}}(K, N; \eta') &= V_1(K, N; \eta') + V_2(K, N; \eta') \\ &= \tilde{V}_2(K, N; \eta')(1 + o(1)) \end{aligned}$$

It is obvious that  $\eta^* \in [\log^b(2), 2 \log^b(2)]$ . We know that

$$\begin{aligned} \bar{\mathcal{R}}^*(K, N) &= V_1(K, N; \eta^*) + V_2(K, N; \eta^*) \leq \bar{\mathcal{R}}(K, N; \eta') = V_2(K, N; \eta')(1 + o(1)) \\ &\lesssim \frac{1}{(N \log^b N)^{\frac{a-1}{a}}}. \end{aligned}$$

## H.5. $E(K, N)$ for Logarithmic Power-Law Spectrum: Proof of [Theorem D.2](#)

### H.5.1. PROOF OF [THEOREM D.2](#), SMALL- $K$ CASE

The proof here is almost a reproduction of the proof in [Appendix H.2.1](#).

### H.5.2. PROOF OF [THEOREM D.2](#), LARGE- $K$ CASE

Consider the multi-epoch training setting with  $d = \Omega\left((KN)^{\frac{1}{a+b}}\right)$ . By [Lemmas H.6](#) and [I.7](#), there exist constants  $C_7, C_8 > 0$  such that

$$\frac{C_8}{N(\log N)^b} \leq \bar{\mathcal{R}}^*(K, N) \leq \frac{C_7}{N(\log N)^b}. \quad (32)$$

Let  $T'$  be defined by matching the expected risks:

$$\bar{\mathcal{R}}^*(K, N) = \bar{\mathcal{R}}^*(1, T'). \quad (33)$$

In the one-pass case, we use the constants  $C_3, C_4 > 0$  (as defined in the proof of [Theorem 4.2](#)) to control  $\bar{\mathcal{R}}^*(1, T')$ .

We claim that

$$\left( \frac{C_4}{C_7} \right)^{\frac{a}{a-1}} N(\log N)^b \leq T' \leq \left( \frac{C_3}{C_8} \right)^{\frac{a}{a-1}} N(\log N)^b. \quad (34)$$

*Proof of the claim.* We argue by contradiction.

1. **Upper bound violation.** If  $T' > \left(\frac{C_3}{C_8}\right)^{\frac{a}{a-1}} N(\log N)^b$ , then the one-pass upper bound together with Equation (32) (multi-epoch lower bound) imply

$$\bar{\mathcal{R}}^*(K, N) < \bar{\mathcal{R}}^*(1, T'),$$

which contradicts the defining equality Equation (33).

2. **Lower bound violation.** If  $T' < \left(\frac{C_4}{C_7}\right)^{\frac{a}{a-1}} N(\log N)^b$ , then  $d = \Omega\left((KN)^{\frac{1}{a+b}}\right)$  yields

$$d = \Omega\left((N(\log N)^b)^{1/a}\right) = \Omega\left((T')^{1/a}\right),$$

so the one-pass lower bound together with Equation (32) (multi-epoch upper bound) give

$$\bar{\mathcal{R}}^*(K, N) > \bar{\mathcal{R}}^*(1, T'),$$

again contradicting Equation (33).

Both violations are impossible; hence Equation (34) holds.  $\square$

Thus, in the large- $K$  multi-epoch regime, the matched one-epoch training time satisfies  $T' = \Theta(N(\log N)^b)$  up to fixed constants. Therefore, the desired characterization of  $E(K, N)$  follows directly.

## Appendix I. Additional Technical lemmas

**Lemma I.1** *For any PSD matrix  $\mathbf{A}$ , it holds that*

$$\langle \mathbf{H}, \mathbf{A} \rangle \leq \text{tr}(\mathbf{H}) \|\mathbf{A}\|.$$

*Proof.* We denote the PSD decomposition of  $\mathbf{H}$  by

$$\mathbf{H} = \sum_{i=1}^d \lambda_i q_i q_i^\top$$

where  $\lambda_i$  and  $q_i$  are the eigenvalues and corresponding eigenvectors of  $\mathbf{H}$ . So we get

$$\begin{aligned} \langle \mathbf{H}, \mathbf{A} \rangle &= \left\langle \sum_{i=1}^d \lambda_i q_i q_i^\top, \mathbf{A} \right\rangle \\ &= \sum_{i=1}^d \lambda_i q_i^\top \mathbf{A} q_i \\ &\leq \sum_{i=1}^d \lambda_i \|\mathbf{A}\| \\ &= \text{tr}(\mathbf{H}) \|\mathbf{A}\|, \end{aligned}$$

which completes the proof.  $\square$

**Lemma I.2** *When  $l \geq 1$ , we have*

$$(1+x)^l \leq 1 + 2lx, \quad x \in [0, \frac{\log 2}{l}]$$

*Proof.* We define  $f(x) := (1+x)^l - (1+2lx)$ . Calculating the derivative and notice the fact that  $2^x - 1 \geq (\log 2)x$ , we obtain

$$\begin{aligned} f'(x) &= l(1+x)^{l-1} - 2l \\ &\leq l(1+2^{\frac{1}{l}} - 1)^{l-1} - 2l \\ &\leq l \times 2^{\frac{l-1}{l}} - 2l \leq 0. \end{aligned}$$

The above equation completes the proof.  $\square$

**Lemma I.3** *When  $l \geq 1$ , we have*

$$(1-x)^{2l} \leq 1 - lx, \quad x \in [0, \frac{1}{6l}]$$

.

*Proof.* We define  $g(x) := (1-x)^{2l} - (1-lx)$ . Calculating the derivative, we obtain

$$g'(x) = -2l(1-x)^{2l-1} + l \leq 0 \quad \text{when } x \in [0, 1 - 2^{-\frac{1}{2l-1}}].$$

Notice that  $h(x) = 2^x$  is convex, so for  $x \in [0, 1]$ , we have

$$h(-x + 0 \times (1-x)) \leq xh(-1) + (1-x)h(0),$$

that is

$$2^{-x} \leq 1 - \frac{x}{2} \quad \text{when } x \in [0, 1].$$

So

$$\begin{aligned} 1 - 2^{-\frac{1}{2l-1}} &\geq 1 - \left(1 - \frac{1}{2(2l-1)}\right) \\ &= \frac{1}{2(2l-1)} \geq \frac{1}{6l} \quad \text{when } l \geq 1, \end{aligned}$$

which concludes the proof.  $\square$

**Lemma I.4** *Given  $N$  data points such that  $\mathbf{x}_0, \dots, \mathbf{x}_{n-1} \stackrel{i.i.d}{\sim} \mathcal{N}(0, \mathbf{H})$ , and define  $\mathbf{A} = (\mathbf{I} - \eta \mathbf{x}_{N-1} \mathbf{x}_{N-1}^\top) \cdots (\mathbf{I} - \eta \mathbf{x}_0 \mathbf{x}_0^\top)$ . Then we have*

$$\mathbb{E} \|\mathbf{A} - \mathbb{E} \mathbf{A}\|^l \leq \left( \sqrt{\delta_A \eta^2 N l} \right)^l,$$

where  $\delta_A := \tilde{C} 8eD^4 \log d$  for some absolute constant  $\tilde{C} > 0$ .

*Proof.* We define  $\mathbf{Q} := \mathbf{A} - \mathbb{E} \mathbf{A}$  for convenience. We can obtain a concentration inequality for  $\|\mathbf{Q}\|$  due to the boundedness of  $\mathbf{x}$  according to Theorem 7.1 in Huang et al. [19].

We define

$$\mathbf{Y}_i := \mathbf{I} - \eta \mathbf{x}_i \mathbf{x}_i^\top$$

For any  $1 \leq i \leq N$ , we can choose  $m_i = 1$ , and we have

$$\|\mathbf{Y}_i - \mathbb{E} \mathbf{Y}_i\| = \|\eta(\mathbf{H} - \mathbf{x}_i \mathbf{x}_i^\top)\| \leq 2D^2\eta := \sigma_i$$

So we know that  $M_A = 1, v_A = 4D^4\eta^2N$ , and

$$\mathbb{P}\{\|\mathbf{Q}\| \geq t\} \leq de^{-\frac{t^2}{2ev_A}} = de^{-\frac{t^2}{8eD^4\eta^2N}} \quad \text{when } t^2 \geq 8eD^4\eta^2N.$$

Furthermore, we have

$$\mathbb{P}\{\|\mathbf{Q}\| \geq t\} \leq e^{-\frac{t^2}{16eD^4\eta^2N}} \quad \text{when } t^2 \geq 16eD^4\eta^2N \log d.$$

So there exists a non-negative sub-Gaussian random variable  $Z$ , s.t

$$\mathbb{P}\{\|\mathbf{Q}\| \geq t\} \leq \mathbb{P}\{Z \geq t\} \leq e^{-\frac{t^2}{16eD^4\eta^2N}} \quad \text{when } t^2 \geq 16eD^4\eta^2N \log d.$$



Then for all  $l \geq 1$ , we can get

$$\begin{aligned}
 \mathbb{E}\|\mathbf{Q}\|^l &= \mathbb{E}\|\mathbf{Q}\|^l (\mathbb{1}_{\{\|\mathbf{Q}\| \leq \sqrt{16eD^4\eta^2 N \log d}\}} + \mathbb{1}_{\{\|\mathbf{Q}\| > \sqrt{16eD^4\eta^2 N \log d}\}}) \\
 &\leq \left(\sqrt{16eD^4\eta^2 N \log d}\right)^l + \mathbb{E}\|\mathbf{Q}\|^l \mathbb{1}_{\{\|\mathbf{Q}\| > \sqrt{16eD^4\eta^2 N \log d}\}} \\
 &\leq \left(\sqrt{16eD^4\eta^2 N \log d}\right)^l + \int_{\sqrt{16eD^4\eta^2 N \log d}}^{+\infty} \mathbb{P}\{\|\mathbf{Q}\| \geq t\} l t^{l-1} dt \\
 &\leq \left(\sqrt{16eD^4\eta^2 N \log d}\right)^l + \int_0^{+\infty} \mathbb{P}\{Z \geq t\} l t^{l-1} dt \\
 &\leq \left(\sqrt{16eD^4\eta^2 N \log d}\right)^l + \mathbb{E}Z^l \\
 &\leq \left(\sqrt{16eD^4\eta^2 N \log d}\right)^l + (\sqrt{C16eD^4\eta^2 N l \log d})^l \\
 &\leq \left(\sqrt{\tilde{C}8eD^4\eta^2 N l \log d}\right)^l.
 \end{aligned}$$

where  $C$  and  $\tilde{C}$  are absolute constants, the fifth inequality is due to Proposition 2.5.2 in [44].  $\square$

**Lemma I.5** For any  $l \leq K$ , we have

$$\mathbb{E} \left\| \prod_{k=1}^l \mathbf{A}^{(k)} - (\mathbb{E}\mathbf{A})^l \right\| \leq \left( \sqrt{\delta_A \eta^2 N l} + \|\mathbb{E}\mathbf{A}\| \right)^l - \|\mathbb{E}\mathbf{A}\|^l,$$

where  $\delta_A$  is the same positive constant appearing in Lemma I.4.

*Proof.* Let  $a = \|\mathbb{E}\mathbf{A}\|$  and  $c_l = \sqrt{\tilde{C}8eD^4\eta^2 N l \log d}$ . Define the perturbation  $\mathbf{Q}^{(k)} = \mathbf{A}^{(k)} - \mathbb{E}\mathbf{A}$ . Expanding the product as

$$\prod_{k=1}^l \mathbf{A}^{(k)} = \prod_{k=1}^l (\mathbf{Q}^{(k)} + \mathbb{E}\mathbf{A}) = \sum_{m=0}^l \sum_{\mathcal{S} \in \binom{[l]}{m}} P_{\mathcal{S}},$$

where  $P_{\mathcal{S}}$  is the matrix product with  $\mathbf{Q}^{(k)}$  at positions  $k \in \mathcal{S}$  and  $\mathbb{E}\mathbf{A}$  elsewhere, preserving order. The difference is

$$\prod_{k=1}^l \mathbf{A}^{(k)} - (\mathbb{E}\mathbf{A})^l = \sum_{m=1}^l \sum_{\mathcal{S} \in \binom{[l]}{m}} P_{\mathcal{S}}.$$

By the triangle inequality and linearity of expectation:

$$\mathbb{E} \left\| \prod_{k=1}^l \mathbf{A}^{(k)} - (\mathbb{E}\mathbf{A})^l \right\| \leq \sum_{m=1}^l \sum_{\mathcal{S} \in \binom{[l]}{m}} \mathbb{E}\|P_{\mathcal{S}}\|.$$

For each  $\mathcal{S}$ , decompose into  $t$  maximal consecutive blocks  $\mathcal{B}_1, \dots, \mathcal{B}_t$  with sizes  $s_1, \dots, s_t$  ( $\sum s_i = m$ ). By Folland's Hölder inequality and Lemma I.4:

$$\mathbb{E}\|P_{\mathcal{S}}\| \leq a^{l-m} \mathbb{E} \prod_{i=1}^t \prod_{j \in \mathcal{B}_i} \|\mathbf{Q}^{(j)}\| \leq a^{l-m} \prod_{i=1}^t \prod_{j \in \mathcal{B}_i} \left( \mathbb{E} \|\mathbf{Q}^{(j)}\|^{s_i} \right)^{\frac{1}{s_i}} \leq a^{l-m} \prod_{i=1}^t c_{s_i}^{s_i}.$$

Since  $c_s^s = \left( \sqrt{\tilde{C}8eD^4\eta^2Ns\log d} \right)^s$  is increasing in  $s$  and  $s_i \leq l$ :

$$c_{s_i}^{s_i} \leq c_l^{s_i} \Rightarrow \mathbb{E}\|\mathbf{P}_S\| \leq a^{l-m}c_l^m.$$

Summing over all  $S$  with  $|S| = m$ :

$$\sum_{S \in \binom{[l]}{m}} \mathbb{E}\|\mathbf{P}_S\| \leq \binom{l}{m} a^{l-m}c_l^m.$$

Thus the total bound is:

$$\sum_{m=1}^l \binom{l}{m} a^{l-m}c_l^m = (a + c_l)^l - a^l,$$

completing the proof.  $\square$

**Lemma I.6** For any  $l \leq K$ , it holds that

$$\mathbb{E} \left\| \prod_{k=1}^l \mathbf{A}^{(k)} - (\mathbb{E}\mathbf{A})^l \right\|^2 \leq \left[ \left( \sqrt{2\delta_A\eta^2Nl} + \|\mathbb{E}\mathbf{A}\| \right)^l - \|\mathbb{E}\mathbf{A}\|^l \right]^2,$$

where  $\delta_A$  is the same positive constant appearing in [Lemma I.4](#).

*Proof.* Set  $a = \|\mathbb{E}\mathbf{A}\|_2$  and  $c_l = \sqrt{\tilde{C}16eD^4\eta^2Nl\log d}$ . Define the perturbation  $\mathbf{Q}^{(k)} = \mathbf{A}^{(k)} - \mathbb{E}\mathbf{A}$ . Expand the matrix product as:

$$\prod_{k=1}^l \mathbf{A}^{(k)} = \prod_{k=1}^l (\mathbf{Q}^{(k)} + \mathbb{E}\mathbf{A}) = \sum_{m=0}^l \sum_{S \in \binom{[l]}{m}} \mathbf{P}_S,$$

where  $\mathbf{P}_S$  denotes the ordered matrix product with  $\mathbf{Q}^{(k)}$  at positions  $k \in S$  and  $\mathbb{E}\mathbf{A}$  elsewhere. The target difference is:

$$\prod_{k=1}^l \mathbf{A}^{(k)} - (\mathbb{E}\mathbf{A})^l = \sum_{m=1}^l \sum_{S \in \binom{[l]}{m}} \mathbf{P}_S.$$

For the squared spectral norm, we have:

$$\begin{aligned} \mathbb{E} \left\| \sum_{m=1}^l \sum_S \mathbf{P}_S \right\|^2 &\leq \mathbb{E} \left( \sum_{m=1}^l \sum_S \|\mathbf{P}_S\| \right)^2 \\ &= \sum_{m=1}^l \sum_{n=1}^l \sum_{S_m} \sum_{S_n} \mathbb{E} [\|\mathbf{P}_{S_m}\| \|\mathbf{P}_{S_n}\|], \end{aligned}$$

where  $S_m$  and  $S_n$  range over all subsets of  $[l]$  with sizes  $m$  and  $n$ , respectively. For each pair  $(S_m, S_n)$ , decompose the union  $\mathcal{U} = S_m \cup S_n$  into  $t$  maximal consecutive blocks  $\mathcal{B}_1, \dots, \mathcal{B}_t$  with sizes  $s_i = |\mathcal{B}_i|$  ( $\sum_{i=1}^t s_i = |\mathcal{U}| = m + n$ ). By [Folland's](#) Hölder inequality and [Lemma I.4](#):

$$\begin{aligned}
 \mathbb{E} [\|P_{S_m}\| \|P_{S_n}\|] &\leq a^{2l-m-n} \mathbb{E} \prod_{i=1}^t \prod_{j \in \mathcal{B}_i} \|Q_j\| \\
 &\leq a^{2l-m-n} \prod_{i=1}^t \prod_{j \in \mathcal{B}_i} \mathbb{E} (\|Q_j\|^{m+n})^{\frac{1}{m+n}} \\
 &\leq a^{2l-m-n} \left( \sqrt{\tilde{C} 8e D^4 \eta^2 N(m+n) \log d} \right)^{m+n} \\
 &\leq a^{2l-m-n} c_l^{m+n}.
 \end{aligned}$$

The combinatorial count satisfies:

$$\sum_{S_m} \sum_{S_n} 1 = \binom{l}{m} \binom{l}{n}.$$

Combining all terms:

$$\mathbb{E} \left\| \prod_{k=1}^l A^{(k)} - a^l \right\|^2 \leq \sum_{m=1}^l \sum_{n=1}^l \binom{l}{m} \binom{l}{n} a^{2l-m-n} c_l^{m+n} = \left[ (a + c_l)^l - a^l \right]^2,$$

where the last equality follows from the binomial theorem applied to  $(a + c_l)^{2l}$ .  $\square$

**Lemma I.7** Consider a function of training time  $T$  given by

$$\mathcal{L}(T) = \frac{1}{2} \sum_{i=d_1+1}^d \frac{c}{i^l} e^{-\frac{2Tc\eta}{i^a}},$$

where  $c, l$  are some absolute constants,  $d_1 = \Theta(1)$ , and  $l > 1$ . Then we have:

1.  $\mathcal{L}(T) \lesssim \frac{1}{T^{\frac{1}{l-1}}}$ ;
2. Given  $d = \Theta\left((KN)^{\frac{1}{a}}\right)$ ,  $\mathcal{L}(T) \gtrsim \frac{1}{T^{\frac{1}{l-1}}}$ .

*Proof.* Computing the derivative of  $f(x) = \frac{c}{x^l} e^{-\frac{2Tc\eta}{x^a}}$ , we have

$$\begin{aligned}
 \arg \max_x f(x) &= \Theta\left((KN)^{\frac{1}{a}}\right), \\
 \max_x f(x) &= \Theta\left(\frac{1}{(KN)^{\frac{l}{a}}}\right).
 \end{aligned}$$

Then

1. For the upper bound, we have

$$\begin{aligned}\mathcal{L}(T) &\leq \frac{1}{2} \sum_{i=d_1+1}^{\infty} \frac{c}{i^l} e^{-\frac{2Tc\eta}{i^a}} \leq \frac{1}{2} \sum_{i=d_1+1}^{(KN)^{\frac{1}{a}}} \frac{c}{i^l} e^{-\frac{2Tc\eta}{i^a}} + \frac{1}{2} \sum_{i=(KN)^{\frac{1}{a}}+1}^{\infty} \frac{c}{i^l} \\ &\lesssim (KN)^{\frac{1}{a}} \times \frac{1}{(KN)^{\frac{l}{a}}} + \frac{1}{(KN)^{\frac{l-1}{a}}} \lesssim \frac{1}{(KN)^{\frac{l-1}{a}}}.\end{aligned}$$

2. For the lower bound, when  $d \geq 3T^{\frac{1}{a}}$ , we have

$$\mathcal{L}(T) \geq \frac{1}{2} \sum_{i=(KN)^{\frac{1}{a}}}^{3(KN)^{\frac{1}{a}}} \frac{c}{i^l} e^{-\frac{2Tc\eta}{i^a}} \geq \frac{1}{2} \frac{c}{3^l (KN)^{\frac{l}{a}}} e^{-2c\eta} \times 2(KN)^{\frac{1}{a}} \gtrsim \frac{1}{(KN)^{\frac{l-1}{a}}}.$$

The above equation completes the proof.  $\square$

**Lemma I.8** *Given an estimator of the excess risk for ME and OP cases*

$$\tilde{S}_2(K, N; \eta) = \frac{1}{2} \sum_{i=d_1+1}^d \frac{c}{i^a} e^{-\frac{2KNc\eta}{i^a}},$$

and

$$\tilde{S}_2(1, T'; \eta) = \frac{1}{2} \sum_{i=d_1+1}^d \frac{c}{i^a} e^{-\frac{2T'c\eta}{i^a}}$$

for some  $d_1 = \Theta(1)$ . If the ME excess risk and OP excess risk satisfy that

$$\begin{aligned}\bar{\mathcal{R}}(K, N; \eta) &= \tilde{S}_2(K, N; \eta)(1 + o(1)) \\ \bar{\mathcal{R}}(1, T'; \eta) &= \tilde{S}_2(1, T'; \eta)(1 + o(1)),\end{aligned}$$

then give  $d = \Omega(T^{\frac{1}{a}})$  and when  $T' \asymp T$ , it holds that

$$E(K, N) \in [K(1 - o(1)), K(1 + o(1))].$$

*Proof.* We define  $H(T) = \tilde{S}_2(K, N; \eta)$  and  $\alpha = \frac{T'}{T}$ . By definition of  $E(K, N)$ , we have  $T' = E(K, N)N$ . Our goal is to prove that  $\alpha \in [1 - o(1), 1 + o(1)]$ .

Solving  $\bar{\mathcal{R}}(K, N; \eta) = \bar{\mathcal{R}}(1, T'; \eta)$ , we can get  $H(T)(1 + o_N(1)) = H(T')(1 + o_{T'}(1))$ . We define  $\delta(K, N) = \frac{\bar{\mathcal{R}}(K, N; \eta) - \tilde{S}_2(K, N; \eta)}{\tilde{S}_2(K, N; \eta)} = o(1)$ , and  $\delta(1, T') = \frac{\bar{\mathcal{R}}(1, T'; \eta) - \tilde{S}_2(1, T'; \eta)}{\tilde{S}_2(1, T'; \eta)} = o(1)$ . Then we can derive that

$$\begin{aligned}H(T')(1 - \delta(1, T')) &\leq H(T)(1 + \delta(K, N)) \\ H(T')(1 + \delta(1, T')) &\geq H(T)(1 - \delta(K, N))\end{aligned}$$

which indicates that

$$-\delta(1, T')H(T') - \delta(K, N)H(T) \leq H(T') - H(T) \leq \delta(1, T')H(T') + \delta(K, N)H(T).$$

Notice that  $H(T)$  is strongly convex, and we have  $H(T) \approx \frac{1}{(KN)^{\frac{a-1}{a}}}$  and  $H'(T) = \frac{1}{2} \sum_{i=1}^d \frac{c}{i^{2a}} e^{\frac{-2KNc\eta}{i^a}} \approx \frac{1}{(KN)^{\frac{2a-1}{a}}}$  by [Lemma I.7](#). We are now ready to prove that  $\alpha \in [1 - o(1), 1 + o(1)]$ .

$$\begin{aligned} -\frac{1}{T^{(2-\frac{1}{a})}}(T' - T) &\lesssim H'(T)(T' - T) \leq H(T') - H(T) \leq H'(T')(T' - T) \lesssim -\frac{1}{T'^{(2-\frac{1}{a})}}(T' - T) \\ \delta(1, T')H(T') + \delta(K, N)H(T) &\lesssim \frac{\delta(1, T')}{T'^{(1-\frac{1}{a})}} + \frac{\delta(K, N)}{T^{(1-\frac{1}{a})}} \lesssim \frac{o(1)}{T^{(1-\frac{1}{a})}}. \end{aligned}$$

So

$$\begin{aligned} \frac{T - T'}{T^{(1-\frac{1}{a})}} &\lesssim \frac{o(1)}{T^{(1-\frac{1}{a})}} \\ -\frac{o(1)}{T^{(1-\frac{1}{a})}} &\lesssim -\frac{1}{T'^{(1-\frac{1}{a})}}(T' - T). \end{aligned}$$

Direct calculation yields the result.  $\square$

**Lemma I.9 (Hyper-Contractivity)** *Given  $d$ -dimension random vector  $\mathbf{x} \sim \mathcal{D}$  satisfying that  $\|\mathbf{x}\| \leq D$  for some constant  $D$ , and the covariance matrix  $\mathbf{H} := \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\mathbf{x}\mathbf{x}^\top] = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$ , where  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq c$  for some constant  $c > 0$ , then the following holds:*

$$\mathbb{E} [\mathbf{x}\mathbf{x}^\top \mathbf{P}\mathbf{x}\mathbf{x}^\top] \leq \alpha \text{tr}(\mathbf{H}\mathbf{P})\mathbf{H}$$

for some constant  $\alpha > 0$  independent of  $\mathbf{P}$ .

*Proof.* By Dieuleveut et al. [10], the above lemma holds for data distributions with a bounded kurtosis along every direction, i.e., there exists a constant  $\kappa > 0$  such that

$$\text{for every } \mathbf{v} \in \mathbb{R}^d, \quad \mathbb{E} [\langle \mathbf{v}, \mathbf{x} \rangle^4] \leq \kappa \langle \mathbf{v}, \mathbf{H}\mathbf{v} \rangle^2.$$

So that it suffices to verify the above inequality. Since  $\lambda_d \geq c$ , we have

$$\langle \mathbf{v}, \mathbf{H}\mathbf{v} \rangle^2 \geq c^2 \|\mathbf{v}\|^4.$$

For the left side, by the triangle inequality and that  $\|\mathbf{x}\|$  is bounded

$$\langle \mathbf{v}, \mathbf{x} \rangle^4 \leq \|\mathbf{v}\|^4 \|\mathbf{x}\|^4 \leq D^4 \|\mathbf{v}\|^4.$$

Combining the above two inequalities gives

$$\mathbb{E} [\langle \mathbf{v}, \mathbf{x} \rangle^4] \leq \frac{D^4}{c^2} \langle \mathbf{v}, \mathbf{H}\mathbf{v} \rangle^2.$$

Now setting  $\kappa = \frac{D^4}{c^2}$  completes the proof.  $\square$