# Learning under Quantization for High-Dimensional Linear Regression

Dechen Zhang[*]    Junwei Su[†]    Difan Zou[‡]

October 22, 2025

**Abstract**

The use of low-bit quantization has emerged as an indispensable technique for enabling the efficient training of large-scale models. Despite its widespread empirical success, a rigorous theoretical understanding of its impact on learning performance remains notably absent, even in the simplest linear regression setting. We present the first systematic theoretical study of this fundamental question, analyzing finite-step stochastic gradient descent (SGD) for high-dimensional linear regression under a comprehensive range of quantization targets: data, labels, parameters, activations, and gradients. Our novel analytical framework establishes precise algorithm-dependent and data-dependent excess risk bounds that characterize how different quantization affects learning: parameter, activation, and gradient quantization amplify noise during training; data quantization distorts the data spectrum; and data and label quantization introduce additional approximation and quantized error. Crucially, we prove that for multiplicative quantization (with input-dependent quantization step), this spectral distortion can be eliminated, and for additive quantization (with constant quantization step), a beneficial scaling effect with batch size emerges. Furthermore, for common polynomial-decay data spectra, we quantitatively compare the risks of multiplicative and additive quantization, drawing a parallel to the comparison between FP and integer quantization methods. Our theory provides a powerful lens to characterize how quantization shapes the learning dynamics of optimization algorithms, paving the way to further explore learning theory under practical hardware constraints.

## 1 Introduction

Quantization has garnered widespread attention as an essential technique for deploying large-scale deep learning models, particularly large language models (LLMs) (Lang et al., 2024; Shen et al., 2024). In line with this low-precision paradigm, a new frontier of research has emerged: quantization scaling laws, which seek to formalize the trade-offs between model size, dataset size, and computational bit-width. Seminal work by Kumar et al. (2024) treated bit-width as a discrete measure of precision. This was extended by Sun et al. (2025), who established a more comprehensive scaling law for floating-point (FP) quantization (Kuzmin et al., 2022) by separately accounting for the distinct roles of exponent and mantissa bits. Going further, Chen et al. (2025) proposed a unified scaling law that models quantized error as a function of model size, training data volume, and quantization group size. Collectively, these studies identify a crucial insight: for large-scale model training, strategic low-bit quantization can drastically reduce memory, computation, and communication overhead. This efficiency enables the training of significantly larger models on more extensive datasets under a fixed memory budget, all without sacrificing final model performance.

[*]Institute of Data Science, The University of Hong Kong. Email: `dechenzhang`@connect.hku.hk

[†]School of Computing & Data Science, The University of Hong Kong. Email: `junweisu`@connect.hku.hk

[‡]School of Computing & Data Science and Institute of Data Science, The University of Hong Kong. Email: `dzou`@hku.hk

The practical deployment of low-precision training has advanced rapidly, yet a significant theory-practice gap persists. Theoretical research remains predominantly restricted to analyzing *convergence guarantees on the training loss* for quantized optimizers (Nadiradze et al., 2021; Liu et al., 2023; Xin et al., 2025; Markov et al., 2023). For example, Markov et al. (2023) proves convergence guarantee for the communication-efficient variant of Fully-Shared Data-Parallel distributed training under parameter and gradient quantization. While these studies offer crucial insights into optimization, they overlook a more fundamental question: *how does quantization affect the model's learning performance?* Specifically, a rigorous characterization of the interplay between quantization, model dimension, dataset size, and their joint effect on the *population risk* remains largely unexplored. A notable step in this direction is Zhang et al. (2022), which analyzes the generalization of quantized two-layer networks through the lens of neural tangent kernel (NTK). However, their work is limited in three key aspects: it only considers parameter quantization; its analysis is confined to the lazy-training regime; and it fails to provide explicit generalization bounds in terms of core parameters like sample size, dimension, and quantization error. These limitations restrict its applicability to modern low-precision training practices.

Motivated by recent theoretical advances in scaling laws (Lin et al., 2024, 2025), we analyze the learning performance of quantized training using a high-dimensional linear model. This model serves as a powerful and well-established testbed for isolating phenomena like learning rate and batch size effects (Kunstner and Bach, 2025; Luo et al., 2025; Zhang et al., 2024b; Xiao, 2024; Ren et al., 2025; Bordelon et al., 2025). Its simplicity provides the analytical flexibility necessary to derive precise relationships between generalization error and critical parameters such as dimension, sample size, and quantization error (or bit-width).

**Our Setting.** In this paper, we consider SGD for linear regression under quantization. We first iterate the standard linear regression problem as follows:

$$\min_{\mathbf{w}} L(\mathbf{w}), \text{ where } L(\mathbf{w}) = \frac{1}{2}\mathbb{E}_{\mathbf{x},y}\left[(y - \langle \mathbf{w}, \mathbf{x}\rangle)^2\right],$$

where $\mathbf{x} \in \mathcal{H}$, is the feature vector, where, $\mathcal{H}$ is some (finite $d$-dimensional or countably infinite dimensional) Hilbert space; $y \in \mathbb{R}$ is the response; $\mathcal{D}$ is an unknown distribution over $\mathbf{x}$ and $y$; and $\mathbf{w} \in \mathcal{H}$ is the weight vector to be optimized. We consider the constant stepsize SGD under quantization as follows: at each iteration $t$, an i.i.d. batch (with batchsize $B$) of examples $(\mathbf{X}_t, \mathbf{y}_t) \in \mathbb{R}^{B \times d} \times \mathbb{R}^B$ is observed, and the weight $\mathbf{w}_t \in \mathbb{R}^d$ is updated according to SGD as follows.

$$\mathbf{w}_t = \mathbf{w}_{t-1} + \gamma \frac{1}{B}\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_o\Big(\mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_a\big(\mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1})\big)\Big), \quad t = 1, ..., N, \quad \text{(Quantized SGD)}$$

where $\gamma > 0$ is a constant stepsize, $N$ is the number of sample batches observed, the master weights be initialized at $\mathbf{w}_0$, and $\mathcal{Q}_d, \mathcal{Q}_l, \mathcal{Q}_p, \mathcal{Q}_a, \mathcal{Q}_o$ are independent general quantization operations for data features, labels, model parameters, activations and output gradients respectively. Notably, for theoretical simplicity, we assume all matrix operations (e.g., addition and multiplication) are computed in full precision, with quantization applied subsequently to obtain low-precision values. Then, we consider the iterate average as the algorithm output, i.e., $\overline{\mathbf{w}}_N := \frac{1}{N}\sum_{t=0}^{N-1}\mathbf{w}_t$.

The goal of this work is to characterize the learning performance of the quantized SGD via evaluating the population risk $L(\overline{\mathbf{w}}_N)$, and more importantly, its relationship with the quantization error. Let $\mathbf{w}^* = \arg\min L(\mathbf{w})$, we define the following excess risk as a surrogate of the population risk:

$$\mathcal{E}(\overline{\mathbf{w}}_N) = L(\overline{\mathbf{w}}_N) - L(\mathbf{w}^*). \tag{Excess Risk}$$

**Our Contributions.** We develop a novel theoretical study on the learnability of the quantized SGD algorithm for high-dimensional linear regression problems. Our contributions are as follows:

- By systematically analyzing a class of quantization techniques, we establish a theoretical bound for the excess risk in quantized SGD. This bound is explicitly formulated as a function of the full eigenspectrum of the quantized data feature covariance matrix, sample size, and quantization errors (see Theorem 4.1 for details). Our results precisely reveal how quantization applied to different model components impacts learning performance: data quantization distorts the data spectrum; the quantization of both data and labels introduces additional approximation and quantized error; while the quantization of parameters, activations, and output gradients amplifies noise throughout the training process on the quantized domain.

- We analyze two standard quantization error models: additive and multiplicative, which conceptually relate to the integer and FP quantization techniques. Our theoretical result shows that multiplicative quantization eliminates spectrum distortion and subsumes the additional quantized error into dominated terms (see Theorem 4.2 for details). For additive quantization, our theoretical bound suggests that the impacts of activation and gradient quantization diminish as batch size scales up (see Corollary 4.1 for details).

- We further derive the conditions on the quantization errors such that the learning performance of the full-precision SGD can be maintained (in orders). Our results indicate that compared with multiplicative quantization, additive quantization imposes more strict spectrum-related requirements on data quantization but weaker batchsize-related requirements on activation and parameter quantization (see Corollary 4.2 for details). By extending the comparison to data spectra with polynomial decay, we show that in high-dimensional settings, multiplicative quantization is applicable while additive quantization is not (see Corollary 4.3 for details). These simplified theoretical results also draw implications for comparing integer and FP quantization, allowing us to identify the conditions under which each type is likely to yield superior performance.

## 2 Related Works

**High-dimensional Linear Regression via SGD.** Theoretical guarantees for the generalization property have garnered significant attention in machine learning and deep learning. Seminal work by Bartlett et al. (2020); Tsigler and Bartlett (2023) derived nearly tight upper and lower excess risk bounds in linear (ridge) regression for general regularization schemes. With regards to the classical underparameterized regime, a large number of works studied the learnability of iterate averaged SGD in linear regression (Polyak and Juditsky, 1992; Défossez and Bach, 2015; Bach and Moulines, 2013; Dieuleveut et al., 2017; Jain et al., 2018, 2017; Zou et al., 2021). With regards to modern overparameterized setting, one-pass SGD in linear regression has also been extensively studied (Dieuleveut and Bach, 2015; Berthier et al., 2020; Varre et al., 2021; Zou et al., 2023; Wu et al., 2022a,b; Zhang et al., 2024a), providing a framework to characterize how the optimization algorithm affects the generalization performance for various data distributions. Another line of work analyzed the behavior of multi-pass SGD on a high-dimensional $\ell^2$-regularized least-squares problem, characterizing excess risk bounds (Lei et al., 2021; Zou et al., 2022) and the exact dynamics of excess risk (Paquette et al., 2024a). From a technical perspective, our work builds on the sharp finite-sample and dimension-free analysis of SGD developed by Zou et al. (2023). However, these works did not concern the practical quantization operations. It remains unclear how quantization error affects the learning behavior of SGD for linear regression.

**Theoretical Analysis for Quantization.** As a powerful technique for deploying large-scale deep learning models, quantization has attracted significant attention. From the theoretical perspective, a line of works focus on the convergence guarantee in both quantized training (SGD) algorithms (De Sa et al., 2015; Alistarh et al., 2017; Faghri et al., 2020; Gorbunov et al., 2020;

Gandikota et al., 2021; Markov et al., 2023; Xin et al., 2025) and post-training quantization methods (Lybrand and Saab, 2021; Zhang and Saab, 2023; Zhang et al., 2023, 2025). For low-precision SGD, De Sa et al. (2015) was the first to consider the convergence guarantees. Assuming unbiased stochastic quantization, convexity, and gradient sparsity, they gave upper bounds on the error probability of SGD. Alistarh et al. (2017) refined these results by focusing on the trade-off between communication and convergence and proposed Quantized SGD (QSGD). Faghri et al. (2020) extended the fixed quantization scheme (Alistarh et al., 2017) to two adaptive quantization schemes, providing a more general convergence guarantee for quantized training. For post-training quantization, Lybrand and Saab (2021) derived an error bound for ternary weight quantization under independent Gaussian data distribution. Zhang et al. (2023) extended this results to more general quantization grids and a wider range of data distributions using a different proof technique. More recently, Zhang et al. (2025) presented the first quantitative error bounds for OPTQ post-training algorithm framework. However, no prior work provides explicit generalization bounds.

**Linear Models for Theory of Scaling Law.** Several recent studies have sought to formalize and explain the empirical scaling laws using conceptually simplified linear models (Bahri et al., 2024; Atanasov et al., 2024; Paquette et al., 2024b; Bordelon et al., 2024; Lin et al., 2024, 2025). Among them, Bahri et al. (2024) considered a linear teacher-student model with power-law spectrum and showed that the test loss of the ordinary least square estimator decreases following a power law in sample size $N$ (or model size $M$) when the other parameter goes to infinity. Bordelon et al. (2024) analyzed the test error of the solution found by gradient flow in a linear random feature model and established power-law scaling in one of $N$, $M$ and training time $T$ while the other two parameters go to infinity. Building on the technique in Zou et al. (2023), Lin et al. (2024) analyzed the test error of the last iterate of one-pass SGD in a sketched linear model. They presented the first systematic study to establish a finite-sample joint scaling law (in $M$ and $N$) for linear models that aligns with empirical observations (Kaplan et al., 2020). More recently, Lin et al. (2025) extended the scaling law analysis to the setting with data reuse (i.e., multi-pass SGD) in data-constrained regimes.

## 3 Preliminary

### 3.1 Quantization operations

For all quantization operations in (Quantized SGD), we employ the stochastic quantization method (Markov et al., 2023), which unbiasedly rounds values using randomly adjusted probabilities. We summarize this in the following assumption.

**Assumption 3.1.** *Let $\mathcal{Q}_i, i \in \{d, l, p, a, o\}$ be the coordinate-wise quantization operation for data feature, label, model parameters, activations, and output gradients, respectively. Then for any $\mathbf{u}$, the quantization operation is unbiased:*
$$\mathbb{E}\left[\mathcal{Q}_i(\mathbf{u})|\mathbf{u}\right] = \mathbf{u}.$$

Furthermore, to better uncover the effect of quantization, we consider the following two types of quantization error: multiplicative quantization and additive quantization, which are motivated by abstracting the behavior of prevalent numerical formats used in practice.

**Definition 3.1.** *Let $\mathcal{Q}$ be an unbiased quantization operation. We categorize it based on the structure of its error variance:*
- ***Multiplicative quantization.*** *We call the quantization is $\epsilon$-multiplicative if the conditional second moment of quantization error is proportional to the outer product of raw data itself, i.e.,*

$$\mathbb{E}\left[\left(\mathcal{Q}(\mathbf{x}) - \mathbf{x}\right)\left(\mathcal{Q}(\mathbf{x}) - \mathbf{x}\right)^\top \Big| \mathbf{x}\right] = \epsilon \mathbf{x}\mathbf{x}^\top.$$

4

- ***Additive quantization.*** *We call the quantization is $\epsilon$-additive if the conditional second moment of quantization error is proportional to identity, i.e.,*

$$\mathbb{E}\left[\left(\mathcal{Q}(\mathbf{x}) - \mathbf{x}\right)\left(\mathcal{Q}(\mathbf{x}) - \mathbf{x}\right)^{\top} \middle| \mathbf{x}\right] = \epsilon\mathbf{I}.$$

This theoretical distinction is grounded in practical quantization schemes. For instance, integer quantization (e.g., INT8, INT16) uses a fixed bin length, resulting in an error that is largely independent of the value's magnitude. This characteristic aligns with our definition of additive quantization, where the error variance is uniform across coordinates. Conversely, floating-point quantization (e.g., FP8, FP32) employs a value-aware bin length via its exponent and mantissa bits (e.g., the E4M3 format in FP8). This structure causes the quantization error to scale with the magnitude of the value itself, corresponding to the model of multiplicative quantization.

To precisely capture the quantization error, we further introduce some relevant notations on quantization errors during the training. Denote the activation and output gradient at time $t$ as

$$\mathbf{a}_t = \mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1}), \quad \mathbf{o}_t = \mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_a\left(\mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1})\right).$$

Then we are ready to define quantization errors.

**Definition 3.2.** *The quantization error on data $\boldsymbol{\epsilon}^{(d)}$, on label $\boldsymbol{\epsilon}^{(l)}$, on parameter $\boldsymbol{\epsilon}_t^{(p)}$ at time $t$, on activation $\boldsymbol{\epsilon}_t^{(a)}$ at time $t$ and on output gradient $\boldsymbol{\epsilon}_t^{(o)}$ at time $t$ are defined as follows.*

$$\boldsymbol{\epsilon}^{(d)} := \mathcal{Q}_d(\mathbf{x}) - \mathbf{x}, \; \boldsymbol{\epsilon}^{(l)} := \mathcal{Q}_l(y) - y, \; \boldsymbol{\epsilon}_t^{(p)} := \mathcal{Q}_p(\mathbf{w}_t) - \mathbf{w}_t, \; \boldsymbol{\epsilon}_t^{(a)} := \mathcal{Q}_a(\mathbf{a}_t) - \mathbf{a}_t, \; \boldsymbol{\epsilon}_t^{(o)} := \mathcal{Q}_o(\mathbf{o}_t) - \mathbf{o}_t.$$

### 3.2 Data model

We then state the regularity assumptions on the data distribution, which align with those common in prior works Zou et al. (2023); Lin et al. (2024). A key distinction in our setting is that all training is performed on quantized data, $\mathcal{Q}_d(\mathbf{x})$ and $\mathcal{Q}_l(y)$. Consequently, we formulate these assumptions directly on the quantized data rather than the full-precision versions.

**Assumption 3.2** (Data covariance). *Let $\mathbf{H} = \mathbb{E}[\mathbf{x}\mathbf{x}^{\top}]$ be the data covariance matrix and*

$$\mathbf{H}^{(q)} := \mathbb{E}[\mathcal{Q}_d(\mathbf{x})\mathcal{Q}_d(\mathbf{x})^{\top}], \quad \mathbf{D} := \mathbb{E}[(\mathcal{Q}_d(\mathbf{x}) - \mathbf{x})(\mathcal{Q}_d(\mathbf{x}) - \mathbf{x})^{\top}],$$

*be the covariance matrices of the quantized data feature and quantization error, respectively. Then we assume that $\mathrm{tr}(\mathbf{H})$ and $\mathrm{tr}(\mathbf{H}^{(q)})$ are finite.*

Further let $\mathbf{H} = \sum_i \lambda_i \mathbf{v}_i \mathbf{v}_i^{\top}$ be the eigen-decomposition of $\mathbf{H}$, where $\{\lambda_i\}_{i=1}^{\infty}$ are the eigenvalues of $\mathbf{H}$ sorted in non-increasing order and $\mathbf{v}_i$ are the corresponding eigenvectors. As in Zou et al. (2023), we denote

$$\mathbf{H}_{0:k} := \sum_{i=1}^{k} \lambda_i \mathbf{v}_i \mathbf{v}_i^{\top}, \quad \mathbf{H}_{k:\infty} := \sum_{i>k} \lambda_i \mathbf{v}_i \mathbf{v}_i^{\top}, \quad \mathbf{I}_{0:k} := \sum_{i=1}^{k} \mathbf{v}_i \mathbf{v}_i^{\top}, \quad \mathbf{I}_{k:\infty} := \sum_{i>k} \mathbf{v}_i \mathbf{v}_i^{\top}.$$

Similarly, we denote the eigendecomposition of $\mathbf{H}^{(q)}$ as $\mathbf{H}^{(q)} = \sum_i \lambda_i^{(q)} \mathbf{v}_i^{(q)} \mathbf{v}_i^{(q)^{\top}}$ and correspondingly obtain $\mathbf{H}_{0:k}^{(q)}, \mathbf{H}_{k:\infty}^{(q)}, \mathbf{I}_{0:k}^{(q)}, \mathbf{I}_{k:\infty}^{(q)}$.

**Assumption 3.3** (Fourth-order moment). *Let $\mathbf{x}^{(q)} = \mathcal{Q}_d(\mathbf{x})$. Then for any PSD matrix $\mathbf{A}$, there exists a constant $\alpha_B > 0$ such that*

$$\mathbb{E}\left[\mathbf{x}^{(q)}\mathbf{x}^{(q)^{\top}}\mathbf{A}\mathbf{x}^{(q)}\mathbf{x}^{(q)^{\top}}\right] \preceq \alpha_B \, \mathrm{tr}(\mathbf{H}^{(q)}\mathbf{A})\mathbf{H}^{(q)}.$$

We note that the above assumptions are adopted primarily to simplify the exposition of our final theoretical results. In practice, for specific quantization mechanisms (e.g., the multiplicative or additive schemes in Definition 3.1), Assumption 3.2 is naturally satisfied by combining standard regularity conditions on the full-precision data (Assumptions 2.1 and 2.2 in Zou et al. (2023)) with the specific properties of the quantization error[1].

Furthermore, to extend the model noise assumption to the quantization setting, we define the optimal model weights regarding the quantized data features and labels:

$$\mathbf{w}^{(q)^*} = \operatorname{argmin}_{\mathbf{w}} \mathbb{E}_{\mathbf{x},y \sim \mathcal{D}} \left[ (\mathcal{Q}_l(y) - \langle \mathbf{w}, \mathcal{Q}_d(\mathbf{x}) \rangle)^2 \right].$$

Accordingly, we make the following assumption on the model noise $\xi := \mathcal{Q}_l(y) - \langle \mathbf{w}^{(q)^*}, \mathcal{Q}_d(\mathbf{x}) \rangle$ based on the optimun under quantization.

**Assumption 3.4.** *Denote* $\xi := \mathcal{Q}_l(y) - \langle \mathbf{w}^{(q)^*}, \mathcal{Q}_d(\mathbf{x}) \rangle$. *Assume there exists a positive constant* $\sigma > 0$ *such that*
$$\mathbb{E}\left[ \xi^2 \mathcal{Q}_d(\mathbf{x}) \mathcal{Q}_d(\mathbf{x})^\top \right] \preceq \sigma^2 \mathbf{H}^{(q)}.$$

In fact, Assumptions 3.3 and 3.4 can be directly inferred from the standard assumptions on the full-precision data under specific quantization regimes. We defer the discussion to Section G.

# 4 Main Theoretical Results

We first derive excess risk upper bounds for Quantized SGD in Section 4.1, then compare these rates with the full-precision SGD (in orders) in Section 4.2 and perform specific case study in Section 4.3.

## 4.1 Excess Risk Bounds

We now provide excess risk bounds under general quantization, multiplicative quantization and additive quantization. Denote the effective dimension for $\mathbf{H}^q$: $k^* = \max \left\{ k : \lambda_k^{(q)} \geq \frac{1}{N\gamma} \right\}$.

**Theorem 4.1.** *Consider the general data quantization, let* $\mathbf{D}_1^H = \mathbf{H}(\mathbf{H} + \mathbf{D})^{-1} \mathbf{D}(\mathbf{H} + \mathbf{D})^{-1} \mathbf{H}$ *and* $\mathbf{D}_2^H = \mathbf{D}(\mathbf{H} + \mathbf{D})^{-1}\mathbf{H}(\mathbf{H} + \mathbf{D})^{-1}\mathbf{D}$, *and consider the zero initialization* $\mathbf{w}_0 = \mathbf{0}$. *Under Assumption 3.1, 3.2, 3.3, and 3.4, if the stepsize* $\gamma < \frac{1}{\gamma \alpha_B \operatorname{tr}(\mathbf{H}+\mathbf{D})}$, *then it holds that,*

$$\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] \leq \text{VarErr} + \text{BiasErr} + \text{ApproxErr} + \text{QuantizedErr},$$

*where*

$$\text{VarErr} = \frac{\sigma_G^{(q)2} + \frac{2\alpha_B}{N\gamma}\left( \|\mathbf{w}^{(q)^*}\|^2_{\mathbf{I}_{0:k^*}^{(q)}} + N\gamma\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}_{k^*:\infty}^{(q)}} \right)}{1 - \gamma\alpha_B \operatorname{tr}(\mathbf{H}+\mathbf{D})} \left( \frac{k^*}{N} + N\gamma^2 \sum_{i>k^*} \lambda_i^{(q)2} \right),$$

$$\text{BiasErr} = \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^{(q)^*}\|^2_{(\mathbf{H}_{0:k^*}^{(q)})^{-1}} + \|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}_{k^*:\infty}^{(q)}},$$

$$\text{ApproxErr} = \mathbb{E}\left[ (\epsilon^{(l)})^2 \right] + \frac{3}{2}\|\mathbf{w}^*\|^2_{\mathbf{D}_1^H} + \frac{1}{2}\|\mathbf{w}^*\|^2_{\mathbf{D}_2^H},$$

---

[1]We remark that the multiplicative quantization regime will require no dimensional constraints, making the results applicable even to infinite-dimensional settings. In contrast, additive quantization necessitates a finite dimension to prevent the variance of the quantization error, $\operatorname{tr}(\epsilon\mathbf{I})$, from becoming infinite.

$$\text{QuantizedErr} = 2\|\mathbf{D}\| \left[ \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^{(q)^*}\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-2}} + \|\mathbf{w}^{(q)^*}\|^2_{\mathbf{I}^{(q)}_{k^*:\infty}} \right]$$

$$+ 2\|\mathbf{D}\| \frac{\sigma_G^{(q)2} + \frac{2\alpha_B}{N\gamma} \left( \|\mathbf{w}^{(q)^*}\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} \right)}{1 - \gamma\alpha_B \text{tr}\,(\mathbf{H} + \mathbf{D})} \left( \sum_{i \le k^*} \frac{1}{N\lambda_i^{(q)}} + N\gamma^2 \sum_{i > k^*} \lambda_i^{(q)} \right),$$

with [2] $\sigma_G^{(q)2} = \frac{\sigma^2 + \sup_t \left\| \mathbb{E}\left[ \boldsymbol{\epsilon}_t^{(o)} \boldsymbol{\epsilon}_t^{(o)\top} | \mathbf{o}_t \right] + \mathbb{E}\left[ \boldsymbol{\epsilon}_t^{(a)} \boldsymbol{\epsilon}_t^{(a)\top} | \mathbf{a}_t \right] \right\|}{B} + \alpha_B \mathbb{E}\left[ \text{tr}\left( \mathbf{H}^{(q)} \boldsymbol{\epsilon}_{t-1}^{(p)} \boldsymbol{\epsilon}_{t-1}^{(p)\top} \right) \right].$

Theorem 4.1 establishes an excess risk bound for quantized SGD under a general quantization paradigm, which is decomposed into four terms: variance error, bias error, approximation error, and quantized error. In particular, the variance and bias errors resemble those for the full-precision SGD (Zou et al., 2023) and can be equivalent by setting the quantization error to be zero. The key role that quantization plays relies is two-fold: data quantization significantly influences the effective data Hessian $\mathbf{H}^{(q)}$, while activation, output gradient and parameter quantization affect the effective noise variance $\sigma_G^{(q)}$ (which will be further characterized in the subsequent theorems when given specific quantization type). Specifically, the quantized Hessian arises from performing SGD in quantized data space and the quantized noise variance corresponds to additional quantization error introduced in the parameter update rule.

The additional two error terms, i.e., approximation error and quantized error, can be interpreted as follows. The approximation error, resulting from quantization of both data and labels, corresponds precisely to the discrepancy between the optimal solution in non-quantized data space and quantized data space, i.e., $\frac{1}{2}\mathbb{E}\left[ (\mathcal{Q}_l(y) - \langle \mathbf{w}^{(q)^*}, \mathcal{Q}_d(\mathbf{x})\rangle)^2 \right] - \frac{1}{2}\mathbb{E}\left[ (y - \langle \mathbf{w}^*, \mathbf{x}\rangle)^2 \right]$. The quantized error originates from the risk associated with applying the quantized averaged SGD iteration $\overline{\mathbf{w}}_N$ to the discrepancy between the quantized data $\mathcal{Q}_d(\mathbf{x})$ and the raw data $\mathbf{x}$, which takes the form $\langle \mathbf{D}, \mathbb{E}\left[ \overline{\mathbf{w}}_N \otimes \overline{\mathbf{w}}_N \right] \rangle$. This expression resembles the quantized bias and quantized variance, but includes an additional factor accounting for data quantization error.

Moreover, in the absence of quantization, our bound exactly reduces to the standard results presented in Zou et al. (2023). It is also worth noting that under the unbiased quantization assumption, parameter, gradient and activation quantization will not affect the BiasErr term [3] To better uncover the effects of quantization, we consider two specific quantization regimes: multiplicative quantization and additive quantization.

**Theorem 4.2 (Multiplicative quantization).** *Under Assumption 3.1, 3.2, 3.3, and 3.4, if there exist $\epsilon_d, \epsilon_l, \epsilon_p, \epsilon_a$ and $\epsilon_o$ such that for any $i \in \{d, l, p, a, o\}$, quantization $\mathcal{Q}_i$ is $\epsilon_i$-multiplicative with $\epsilon_i \le O(1)$ and the stepsize satisfies $\gamma < \frac{1}{\alpha_B C_\epsilon \text{tr}(\mathbf{H})}$ [4], then the excess risk can be upper bounded as follows.*

$$\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] \lesssim \text{ApproxErr} + \text{VarErr} + \text{BiasErr},$$

*where*

$$\text{ApproxErr} \lesssim \|\mathbf{w}^*\|^2_{\mathbf{H}}\, \epsilon_d + \epsilon_l, \quad \text{BiasErr} \lesssim \frac{1}{\gamma^2 N^2} \|\mathbf{w}^*\|^2_{\mathbf{H}^{-1}_{0:k^*}} + \|\mathbf{w}^*\|^2_{\mathbf{H}_{k^*:\infty}},$$

---

[2] In Theorem 4.1, $\|\cdot\|$ denotes the spectral norm.

[3] In absence of the unbiased assumption, the conditional expectation for $\boldsymbol{\eta}_t := \mathbf{w}_t - \mathbf{w}^*$ (Eq. (D.2)) involves additional terms related to quantization expectations, thereby introducing extra terms (related to parameter, output gradient and activation quantization) into bias. Our framework can easily extend to this case. The unbiased assumption is applied for theoretical simplicity.

[4] $C_\epsilon = (1 + \epsilon_d)(1 + 2\epsilon_p + 4\epsilon_o(1 + \epsilon_a)(1 + \epsilon_p) + 2\epsilon_a(1 + \epsilon_p)) \le O(1)$.

$$\text{VarErr} \lesssim \left( \frac{k^*}{N} + N\gamma^2 \sum_{i>k^*} \lambda_i^2 \right) \left( \frac{\frac{\sigma^2}{B} + \alpha_B \left( (\epsilon_o + \epsilon_a + \epsilon_p) \|\mathbf{w}^*\|_{\mathbf{H}}^2 + \frac{\|\mathbf{w}^*\|_{\mathbf{I}_{0:k^*}}^2}{N\gamma} + \|\mathbf{w}^*\|_{\mathbf{H}_{k^*:\infty}}^2 \right)}{1 - \gamma\alpha_B C_\epsilon \mathrm{tr}\,(\mathbf{H})} \right).$$

We would like to remark that, compared with Theorem 4.1, the discrepancy between $\mathbf{H}^{(q)}$ and $\mathbf{H}$ can be incorporated into $\mathbf{H}$ under multiplicative quantization, eliminating additional error due to the spectral gap. Though the multiplicative nature of the quantization introduces additional complexity into the iteration update rule, this additional error is merged when $\epsilon_a, \epsilon_o, \epsilon_p$ are at most constants (see Lemma D.3 for details). Regarding additive quantization, the excess risk bound can be directly adapted from Theorem 4.1, which is summarized in the following corollary.

**Corollary 4.1 (Additive quantization).** *Denote* $\sigma_A^{(q)2} = \frac{\epsilon_o + \epsilon_a}{B} + \alpha_B \epsilon_p \mathrm{tr}\,(\mathbf{H} + \epsilon_d \mathbf{I}) + \frac{\sigma^2}{B}$. *Under Assumption 3.1, 3.2, 3.3, and 3.4, if there exist $\epsilon_d, \epsilon_l, \epsilon_p, \epsilon_a$ and $\epsilon_o$ such that for any $i \in \{d, l, p, a, o\}$, quantization $\mathcal{Q}_i$ is $\epsilon_i$-additive, and the stepsize satisfies $\gamma < \frac{1}{\gamma \alpha_B \mathrm{tr}(\mathbf{H}+\epsilon_d \mathbf{I})}$, then*

$$\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] \lesssim \text{ApproxErr} + \text{VarErr} + \text{BiasErr},$$

*where*

$$\text{ApproxErr} \lesssim \epsilon_l + \epsilon_d \|\mathbf{w}^*\|^2, \quad \text{BiasErr} \lesssim \frac{1}{\gamma^2 N^2} \|\mathbf{w}^*\|_{\mathbf{H}_{0:k^*}^{(q)}}^2 {}^{-1} + \|\mathbf{w}^*\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2,$$

$$\text{VarErr} \lesssim \frac{\sigma_A^{(q)2} + \frac{2\alpha_B}{N\gamma} \left( \|\mathbf{w}^*\|_{\mathbf{I}_{0:k^*}^{(q)}}^2 + N\gamma \|\mathbf{w}^*\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2 \right)}{1 - \gamma\alpha_B \mathrm{tr}\,(\mathbf{H} + \epsilon_d \mathbf{I})} \left( \frac{k^*}{N} + N\gamma^2 \sum_{i>k^*} (\lambda_i + \epsilon_d)^2 \right).$$

A key point to emphasize under additive quantization is that, compared with multiplicative quantization, the contribution of activation and output gradient quantization error to the effective noise variance $\sigma_A^{(q)2}$ is scaled by a factor of $\frac{1}{B}$. The interpretation is that additive quantization provides a constant-level conditional second moment for the quantization error, which diminishes the underlying effect of the (quantized) data. Specifically, this reduction manifests as a change from a dependence on the fourth moment of the data (i.e., $\frac{1}{B^2} \mathbb{E}[\mathbf{X}^{q\top} \mathbf{X}^q \mathbf{X}^{q\top} \mathbf{X}^q]$) to a dependence on the second moment (i.e., $\frac{1}{B^2} \mathbb{E}[\mathbf{X}^{q\top} \mathbf{X}^q]$), consequently introducing an extra factor of $1/B$ in the output gradient and activation quantization error (see Lemma D.2 and Lemma D.3 for details). This leads to a distinction in how quantization affects different components: the influence of $\epsilon_a$ and $\epsilon_o$ diminishes as batchsize $B$ increases, while the effect of $\epsilon_p$ remains independent of batchsize $B$.

## 4.2 Comparisons with Standard Excess Risk Bound

In this part, we will provide a detailed comparison with standard excess risk bounds and identify the conditions on the quantization error such that the excess risk bound will not be largely affected. First, let $k_0^* = \max\{k : \lambda_k \geq \frac{1}{N\gamma}\}$, we recall the standard excess risk bound (Zou et al., 2023):

$$R_0 = \left( \frac{k_0^*}{N} + N\gamma^2 \cdot \sum_{i>k_0^*} \lambda_i^2 \right) \frac{4\alpha_B \left( \|\mathbf{w}^*\|_{\mathbf{I}_{0:k_0^*}}^2 + N\gamma \|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}}^2 \right)}{N\gamma \left[ 1 - \gamma\alpha_B \mathrm{tr}\,(\mathbf{H}) \right]}$$

$$+ \frac{1}{B} \left( \frac{k_0^*}{N} + N\gamma^2 \cdot \sum_{i>k_0^*} \lambda_i^2 \right) \frac{\sigma^2}{1 - \gamma\alpha_B \mathrm{tr}\,(\mathbf{H})} + \frac{2}{\gamma^2 N^2} \cdot \|\mathbf{w}^*\|_{(\mathbf{H}_{0:k_0^*})^{-1}}^2 + 2\|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}}^2.$$

The following corollary derives the conditions on the quantization errors such that the learning performance of the full-precision SGD can be maintained (in orders).

8

**Corollary 4.2.** *To ensure that $\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] \lesssim R_0$, conditions on the quantization error are as follows:*
- *For Multiplicative quantization, under the Assumptions in Theorem 4.2, we require*

$$\epsilon_l \lesssim R_0, \quad \epsilon_p, \epsilon_a, \epsilon_o \lesssim \frac{\sigma^2}{B\|\mathbf{w}^*\|_{\mathbf{H}}^2} \wedge 1, \quad \epsilon_d \lesssim \frac{R_0}{\|\mathbf{w}^*\|_{\mathbf{H}}^2} \wedge 1.$$

- *For Additive quantization, under the Assumptions in Corollary 4.1, we require*

$$\epsilon_l \lesssim R_0, \quad \epsilon_o + \epsilon_a \lesssim \sigma^2, \quad \epsilon_p \lesssim \frac{\sigma^2}{B\mathrm{tr}(\mathbf{H} + \epsilon_d \mathbf{I})},$$

$$\epsilon_d \lesssim \frac{R_0}{\|\mathbf{w}^*\|^2} \wedge \sqrt{\frac{\sum_{i>k_0^*} \lambda_i^2}{d - k_0^*}} \wedge \frac{\|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}}^2 + \frac{1}{N^2\gamma^2}\|\mathbf{w}^*\|_{\mathbf{H}_{0:k_0^*}^{-1}}^2}{\|\mathbf{w}^*\|_{\mathbf{I}_{k_0^*:\infty}}^2}.$$

Corollary 4.2 establishes explicit conditions under which the quantized population risk matches its non-quantized counterpart $R_0$ up to a constant factor. Our theoretical results indicate that compared with multiplicative quantization, additive quantization imposes more strict requirements on data quantization $\epsilon_d$, with an extra term related to data spectrum, but weaker requirements on activation quantization $\epsilon_a$ and output gradient quantization $\epsilon_o$ (weakened by a factor of $\frac{1}{B}$). These theoretical results align well with our insights: (1) multiplicative data quantization diminishes the spectral gap between $\mathbf{H}$ and $\mathbf{H}^{(q)}$; (2) additive activation and gradient quantization leads to a beneficial scaling with the batch size $B$.

## 4.3 Case Study on Data Distribution with Polynomially-decay Spectrum

Following Lin et al. (2024, 2025), we study the excess risk bound assuming constant level optimal parameter (i.e., $\|\mathbf{w}_i^*\|^2 = \Theta(1)$, $\forall i > 0$) and the power-law spectrum. In particular,

**Assumption 4.1.** *There exists $a > 1$ such that the eigenvalues of $\mathbf{H}$ satisfy $\lambda_i \sim i^{-a}$, $i > 0$.*

**Corollary 4.3.** *Under Assumption 4.1, we have:*
- *For multiplicative quantization, under the Assumptions of Theorem 4.2, if we further assume $\|\mathbf{w}^*\|_{\mathbf{H}}^2 \lesssim \sigma^2$, let $d_{eff}^{(M)} = [N\gamma(1 + \epsilon_d)]^{\frac{1}{a}}$, then*
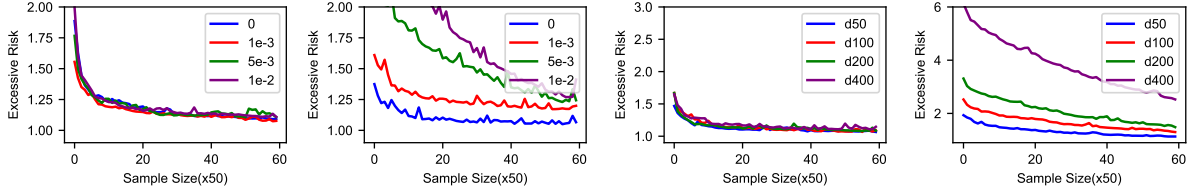
$$\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] \lesssim \epsilon_d + \epsilon_l + \frac{d_{eff}^{(M)} \wedge d}{N\gamma} + \frac{d_{eff}^{(M)} \wedge d}{N}\left(\frac{\sigma^2}{B} + \epsilon_p + \epsilon_o + \epsilon_a + \frac{d_{eff}^{(M)} \wedge d}{N\gamma}\right).$$

- *For additive quantization, under the Assumptions of Corollary 4.1, it holds*

$$\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] \lesssim \epsilon_d d + \epsilon_l + \frac{d_{eff}^{(A)} \wedge d}{N\gamma} + \frac{d_{eff}^{(A)} \wedge d}{N}\left(\frac{\sigma^2}{B} + (1 + d\epsilon_d)\epsilon_p + \frac{\epsilon_o + \epsilon_a}{B} + \frac{d_{eff}^{(A)} \wedge d}{N\gamma}\right),$$

*where $d_{eff}^{(A)} = \left(d - \max\left\{d^{-a}, \frac{1}{N\gamma} - \epsilon_d\right\}^{-\frac{1}{a}}\right)\epsilon_d N\gamma + \max\left\{d^{-a}, \frac{1}{N\gamma} - \epsilon_d\right\}^{-\frac{1}{a}}.$*

Our findings in polynomial-decay data spectrum scenarios reveal distinct scaling behaviors under multiplicative and additive quantization. Specifically, under additive quantization, the impacts of both data quantization and parameter quantization scale with the data dimension $d$. In contrast, under multiplicative quantization, data quantization remains independent of $d$. This difference arises because multiplicative data quantization is applied coordinate-wise, causing its behavior to

(a) **Multiplicative** (FP-like) (b) **Additive** (INT-like) (c) **Multiplicative** (FP-like) (d) **Additive** (INT-like)

Figure 1: **Generalization under quantization.** Test risk for SGD with iterate averaging under multiplicative (FP-like) vs. additive (INT-like) quantization. (a) and (b): vary the quantization level at fixed dimension. (c) and (d): vary dimension at fixed quantization level.

depend solely on the original data spectrum—making it applicable even in infinite-dimensional settings. Additive data quantization, however, employs uniform quantization strength across all dimensions, leading to a strong dependence on the data dimension $d$. We also note that under additive quantization, the influence of the spectral gap is captured by the effective dimension term $d_{eff}^{(A)}$, which includes an additional term $\left[d - \max\{d^{-a}, \frac{1}{N\gamma} - \epsilon_d\}^{-\frac{1}{a}}\right]\epsilon_d N\gamma$ that depends significantly on $\epsilon_d$.

**Implications to integer and FP quantization.** In practical integer quantization with bit-width $b$ and FP quantization with mantissa bit-width $m$, the quantization stepsize for a value $x$ are approximately $\delta(x) = 2^{-b}$ and $\delta(x) = 2^{\lfloor \log_2 x \rfloor - m}$ [5], respectively. Consequently, the conditional second moment of quantization error $\mathbb{E}[(\mathcal{Q}(x) - x)^2|x]$ is roughly proportional to the square of the quantization stepsize $\delta(x)^2$. This implies a fundamental correspondence: multiplicative quantization exhibits the characteristic of FP quantization (input-dependent quantization stepsize), whereas additive quantization characterizes integer quantization (constant quantization stepsize).

A practical takeaway is that, given specific FP and integer quantization with bit-width $b$ and mantissa bit-width $m$, practitioners can directly apply Corollary 4.3 to determine which quantization scheme is more suitable under specific scenarios. A notable observation is the distinct role of dimension $d$: for data quantization, FP quantization becomes preferable when $m_d \geq b_d - \frac{1}{2}\log_2 d$ whereas integer quantization is favored when $b_d \geq m_d + \frac{1}{2}\log_2 d$ [6]. This means FP quantization can outperform integer quantization even when its mantissa bit-width is smaller than the integer bit-width by $\frac{1}{2}\log_2 d$, highlighting the advantage of FP quantization in high-dimensional settings.

**Numerical experiments.** We evaluate constant–stepsize SGD with iterate averaging on a Gaussian least–squares model. The feature distribution has covariance matrix with eigenvalues $\lambda_i = i^{-2}$. The ground–truth parameter is $\mathbf{w}^*$ with entries $\mathbf{w}^*[i] = 1$, and the observation noise variance is $\sigma^2 = 1$. This study answers two questions: **Q1**: How do *additive* vs. *multiplicative* quantization errors affect learning? **Q2**: How does *dimension* $d$ interact with these two quantization types?

**Q1 (quantization level).** We fix $d = 200$ and $B = 1$, and vary the quantization error level $\varepsilon \in \{0.001, 0.005, 0.01\}$ for each scheme. Results are shown in Fig. 1(a,b). This empirically validates our theory: additive errors distort the data Hessian spectrum, increasing risk, whereas multiplicative errors diminish the spectral gap, maintaining risk constant despite higher error levels.

**Q2 (dimension).** We fix the quantization level at $\varepsilon = 0.01$ and $B = 1$, and vary $d \in \{50, 100, 200, 400\}$. Results are shown in Fig. 1(c,d). These empirical results align with our theoretical finding: additive quantization leads to a dramatic increase in excess risk with larger $d$,

---

[5]We assume the exponent bits in FP quantization can cover the scaling of $x$.

[6]$b_d$ and $m_d$ are the bit-width for integer data quantization and the mantissa bit-width for FP data quantization respectively.

while multiplicative quantization maintains stable performance even with high-dimensional data.

# 5    Conclusion and Limitations

In this work, we analyze the excess risk of quantized SGD for high-dimensional linear regression. Our novel theoretical framework characterizes the distinct impacts of various quantization types on learnability: data quantization distorts the data spectrum (eliminated by multiplicative quantization); parameter, activation and gradient quantization amplify noise (mitigated by additive quantization); data and label quantization introduce additional error (scale with dimension in additive quantization yet are dimension-independent in multiplicative quantization). Our theory establishes the conditions on quantization errors required to maintain full-precision SGD performance, and it identifies the scenarios under which FP and integer quantization are each likely to yield superior performance.

Our limitations are twofold: (i) we only establish excess risk upper bounds without a corresponding lower-bound analysis, and (ii) our analysis is confined to one-pass SGD, leaving multi-pass SGD and algorithms with momentum as open problems.

# References

ALISTARH, D., GRUBIC, D., LI, J., TOMIOKA, R. and VOJNOVIC, M. (2017). Qsgd: Communication-efficient sgd via gradient quantization and encoding. *Advances in neural information processing systems* **30**. 3, 4

ATANASOV, A., ZAVATONE-VETH, J. A. and PEHLEVAN, C. (2024). Scaling and renormalization in high-dimensional regression. *arXiv preprint arXiv:2405.00592* . 4

BACH, F. and MOULINES, E. (2013). Non-strongly-convex smooth stochastic approximation with convergence rate o $(1/n)$. *Advances in neural information processing systems* **26**. 3

BAHRI, Y., DYER, E., KAPLAN, J., LEE, J. and SHARMA, U. (2024). Explaining neural scaling laws. *Proceedings of the National Academy of Sciences* **121** e2311878121. 4

BARTLETT, P. L., LONG, P. M., LUGOSI, G. and TSIGLER, A. (2020). Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences* **117** 30063–30070. 3

BERTHIER, R., BACH, F. and GAILLARD, P. (2020). Tight nonparametric convergence rates for stochastic gradient descent under the noiseless linear model. *Advances in Neural Information Processing Systems* **33** 2576–2586. 3

BORDELON, B., ATANASOV, A. and PEHLEVAN, C. (2024). A dynamical model of neural scaling laws. *arXiv preprint arXiv:2402.01092* . 4

BORDELON, B., ATANASOV, A. and PEHLEVAN, C. (2025). How feature learning can improve neural scaling laws. *Journal of Statistical Mechanics: Theory and Experiment* **2025** 084002. 2

CHEN, M., ZHANG, C., LIU, J., ZENG, Y., XUE, Z., LIU, Z., LI, Y., MA, J., HUANG, J., ZHOU, X. ET AL. (2025). Scaling law for quantization-aware training. *arXiv preprint arXiv:2505.14302* . 1

DE SA, C. M., ZHANG, C., OLUKOTUN, K. and RÉ, C. (2015). Taming the wild: A unified analysis of hogwild-style algorithms. *Advances in neural information processing systems* **28**. 3, 4

DÉFOSSEZ, A. and BACH, F. (2015). Averaged least-mean-squares: Bias-variance trade-offs and optimal sampling distributions. In *Artificial Intelligence and Statistics*. PMLR. 3

DIEULEVEUT, A. and BACH, F. (2015). Non-parametric stochastic approximation with large step sizes. *Annals of Statistics* **44**. 3

DIEULEVEUT, A., FLAMMARION, N. and BACH, F. (2017). Harder, better, faster, stronger convergence rates for least-squares regression. *Journal of Machine Learning Research* **18** 1–51. 3

FAGHRI, F., TABRIZIAN, I., MARKOV, I., ALISTARH, D., ROY, D. M. and RAMEZANI-KEBRYA, A. (2020). Adaptive gradient quantization for data-parallel sgd. *Advances in neural information processing systems* **33** 3174–3185. 3, 4

GANDIKOTA, V., KANE, D., MAITY, R. K. and MAZUMDAR, A. (2021). vqsgd: Vector quantized stochastic gradient descent. In *International Conference on Artificial Intelligence and Statistics*. PMLR. 4

GORBUNOV, E., HANZELY, F. and RICHTÁRIK, P. (2020). A unified theory of sgd: Variance reduction, sampling, quantization and coordinate descent. In *International Conference on Artificial Intelligence and Statistics*. PMLR. 3

JAIN, P., KAKADE, S. M., KIDAMBI, R., NETRAPALLI, P., PILLUTLA, V. K. and SIDFORD, A. (2017). A markov chain theory approach to characterizing the minimax optimality of stochastic gradient descent (for least squares). *arXiv preprint arXiv:1710.09430* . 3

JAIN, P., KAKADE, S. M., KIDAMBI, R., NETRAPALLI, P. and SIDFORD, A. (2018). Parallelizing stochastic gradient descent for least squares regression: mini-batching, averaging, and model misspecification. *Journal of machine learning research* **18** 1–42. 3

KAPLAN, J., MCCANDLISH, S., HENIGHAN, T., BROWN, T. B., CHESS, B., CHILD, R., GRAY, S., RADFORD, A., WU, J. and AMODEI, D. (2020). Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361* . 4

KUMAR, T., ANKNER, Z., SPECTOR, B. F., BORDELON, B., MUENNIGHOFF, N., PAUL, M., PEHLEVAN, C., RÉ, C. and RAGHUNATHAN, A. (2024). Scaling laws for precision. *arXiv preprint arXiv:2411.04330* . 1

KUNSTNER, F. and BACH, F. (2025). Scaling laws for gradient descent and sign descent for linear bigram models under zipf's law. *arXiv preprint arXiv:2505.19227* . 2

KUZMIN, A., VAN BAALEN, M., REN, Y., NAGEL, M., PETERS, J. and BLANKEVOORT, T. (2022). Fp8 quantization: The power of the exponent. *Advances in Neural Information Processing Systems* **35** 14651–14662. 1

LANG, J., GUO, Z. and HUANG, S. (2024). A comprehensive study on quantization techniques for large language models. In *2024 4th International Conference on Artificial Intelligence, Robotics, and Communication (ICAIRC)*. IEEE. 1

LEI, Y., HU, T. and TANG, K. (2021). Generalization performance of multi-pass stochastic gradient descent with convex loss functions. *Journal of Machine Learning Research* **22** 1–41. 3

LIN, L., WU, J. and BARTLETT, P. L. (2025). Improved scaling laws in linear regression via data reuse. *arXiv preprint arXiv:2506.08415* . 2, 4, 9

LIN, L., WU, J., KAKADE, S. M., BARTLETT, P. L. and LEE, J. D. (2024). Scaling laws in linear regression: Compute, parameters, and data. *Advances in Neural Information Processing Systems* **37** 60556–60606. 2, 4, 5, 9

LIU, L., ZHANG, J., SONG, S. and LETAIEF, K. B. (2023). Hierarchical federated learning with quantization: Convergence analysis and system design. *IEEE Transactions on Wireless Communications* **22** 2–18. 2

LUO, K., WEN, H., HU, S., SUN, Z., LIU, Z., SUN, M., LYU, K. and CHEN, W. (2025). A multi-power law for loss curve prediction across learning rate schedules. In *The Thirteenth International Conference on Learning Representations.* 2

LYBRAND, E. and SAAB, R. (2021). A greedy algorithm for quantizing neural networks. *Journal of Machine Learning Research* **22** 1–38. 4

MARKOV, I., VLADU, A., GUO, Q. and ALISTARH, D. (2023). Quantized distributed training of large models with convergence guarantees. In *International Conference on Machine Learning.* PMLR. 2, 4

NADIRADZE, G., SABOUR, A., DAVIES, P., LI, S. and ALISTARH, D. (2021). Asynchronous decentralized sgd with quantized and local updates. *Advances in Neural Information Processing Systems* **34** 6829–6842. 2

PAQUETTE, C., PAQUETTE, E., ADLAM, B. and PENNINGTON, J. (2024a). Homogenization of sgd in high-dimensions: Exact dynamics and generalization properties. *Mathematical Programming* 1–90. 3

PAQUETTE, E., PAQUETTE, C., XIAO, L. and PENNINGTON, J. (2024b). 4+ 3 phases of compute-optimal neural scaling laws. *Advances in Neural Information Processing Systems* **37** 16459–16537. 4

POLYAK, B. T. and JUDITSKY, A. B. (1992). Acceleration of stochastic approximation by averaging. *SIAM journal on control and optimization* **30** 838–855. 3

REN, Y., NICHANI, E., WU, D. and LEE, J. D. (2025). Emergence and scaling laws in sgd learning of shallow neural networks. *arXiv preprint arXiv:2504.19983* . 2

SHEN, A., LAI, Z. and LI, D. (2024). Exploring quantization techniques for large-scale language models: Methods, challenges and future directions. In *Proceedings of the 2024 9th International Conference on Cyber Security and Information Engineering.* 1

SUN, X., LI, S., XIE, R., HAN, W., WU, K., YANG, Z., LI, Y., WANG, A., LI, S., XUE, J. ET AL. (2025). Scaling laws for floating point quantization training. *arXiv preprint arXiv:2501.02423* . 1

TSIGLER, A. and BARTLETT, P. L. (2023). Benign overfitting in ridge regression. *Journal of Machine Learning Research* **24** 1–76. 3

VARRE, A. V., PILLAUD-VIVIEN, L. and FLAMMARION, N. (2021). Last iterate convergence of sgd for least-squares in the interpolation regime. *Advances in Neural Information Processing Systems* **34** 21581–21591. 3

WU, J., ZOU, D., BRAVERMAN, V., GU, Q. and KAKADE, S. (2022a). Last iterate risk bounds of sgd with decaying stepsize for overparameterized linear regression. In *International conference on machine learning.* PMLR. 3

WU, J., ZOU, D., BRAVERMAN, V., GU, Q. and KAKADE, S. (2022b). The power and limitation of pretraining-finetuning for linear regression under covariate shift. *Advances in Neural Information Processing Systems* **35** 33041–33053. 3

XIAO, L. (2024). Rethinking conventional wisdom in machine learning: From generalization to scaling. *arXiv preprint arXiv:2409.15156* . 2

XIN, J., CANINI, M., RICHTÁRIK, P. and HORVÁTH, S. (2025). Global-qsgd: Allreduce-compatible quantization for distributed learning with theoretical guarantees. In *Proceedings of the 5th Workshop on Machine Learning and Systems.* 2, 4

ZHANG, H., LIU, Y., CHEN, Q. and FANG, C. (2024a). The optimality of (accelerated) sgd for high-dimensional quadratic optimization. *arXiv preprint arXiv:2409.09745* . 3

ZHANG, H., MORWANI, D., VYAS, N., WU, J., ZOU, D., GHAI, U., FOSTER, D. and KAKADE, S. (2024b). How does critical batch size scale in pre-training? *arXiv preprint arXiv:2410.21676* . 2

ZHANG, H., ZHANG, S., COLBERT, I. and SAAB, R. (2025). Provable post-training quantization: Theoretical analysis of optq and qronos. *arXiv preprint arXiv:2508.04853* . 4

ZHANG, J. and SAAB, R. (2023). Spfq: A stochastic algorithm and its error analysis for neural network quantization. *arXiv preprint arXiv:2309.10975* . 4

ZHANG, J., ZHOU, Y. and SAAB, R. (2023). Post-training quantization for neural networks with provable guarantees. *SIAM journal on mathematics of data science* **5** 373–399. 4

ZHANG, K., YIN, M. and WANG, Y.-X. (2022). Why quantization improves generalization: Ntk of binary weight neural networks. *arXiv preprint arXiv:2206.05916* . 2

ZOU, D., WU, J., BRAVERMAN, V., GU, Q., FOSTER, D. P. and KAKADE, S. (2021). The benefits of implicit regularization from sgd in least squares problems. *Advances in neural information processing systems* **34** 5456–5468. 3

ZOU, D., WU, J., BRAVERMAN, V., GU, Q. and KAKADE, S. (2022). Risk bounds of multi-pass sgd for least squares in the interpolation regime. *Advances in Neural Information Processing Systems* **35** 12909–12920. 3

ZOU, D., WU, J., BRAVERMAN, V., GU, Q. and KAKADE, S. M. (2023). Benign overfitting of constant-stepsize sgd for linear regression. *Journal of Machine Learning Research* **24** 1–58. 3, 4, 5, 6, 7, 8, 15, 17, 21, 27, 29, 30, 37, 48, 54, 56
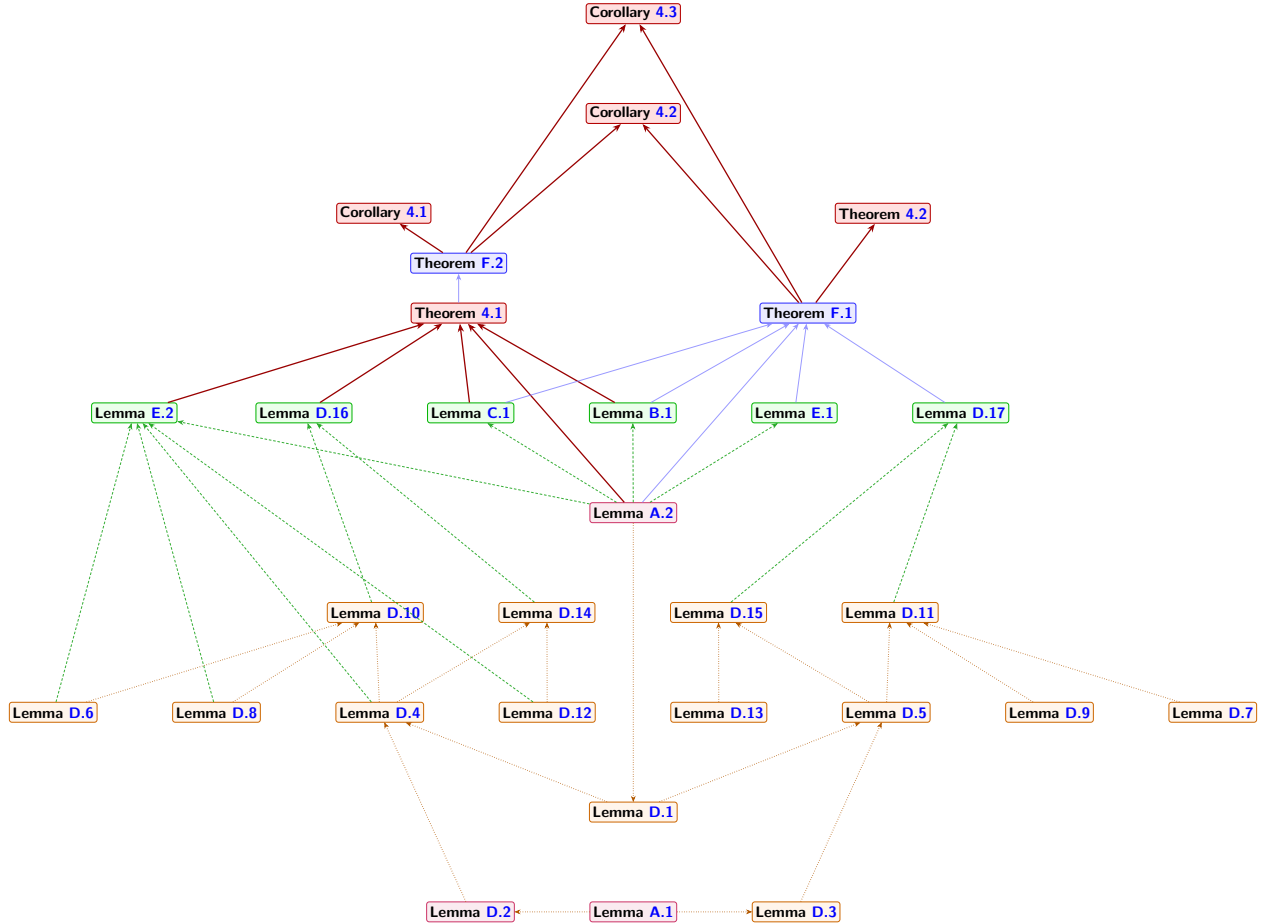
# Appendix

The appendix is organized as follows. In Section A, we begin the analysis of excess risk bounds for the iteratively averaged quantized SGD by first deriving the update rule for the parameter deviation $\mathbf{w}_t - \mathbf{w}^{(q)*}$ (detailed in Section A.1) and second performing an excess risk decomposition (detailed in Section A.2):

$$\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] = R_1 + R_2 + R_3 + R_4.$$

We then conduct a refined analysis of $R_4$ and $R_3$ in Sections B and C, respectively. For $R_2$, we extend techniques from Zou et al. (2023) in Section D. In particular, we first introduce useful notations in Section D.1 and then present a comprehensive analysis of the update rule for $\mathbb{E}[\boldsymbol{\eta}_t \boldsymbol{\eta}_t^\top]$ in Section D.2. This analysis is crucial for adapting previous proof techniques to the quantized SGD setting. Based on these results, we perform a bias–variance decomposition in Section D.3, and analyze the bias and variance errors separately in Section D.4 and D.5. Bounds for $R_1$ are derived in Section E.

Finally, we provide detailed proofs of Theorem 4.1 in Section F.1, Theorem 4.2 in Section F.2, Corollary 4.1 in Section F.3, Corollary 4.2 in Sections F.4 and F.5, and Corollary 4.3 in Sections F.6 and F.7. In addition, we discuss Assumptions 3.3 and 3.4 in Section G.

The following proof dependency graph visually encapsulates the logical structure and organizational architecture of the theoretical results in our paper. In particular, the arrow from element $X$ to element $Y$ means the proof of $Y$ relies on $X$.

# Appendix Contents

# A Initial Study

For simplicity, we denote $y^{(q)} = \mathcal{Q}_l(y)$, $\mathbf{w}_t^{(q)} = \mathcal{Q}_p(\mathbf{w}_t)$, $\mathbf{x}^{(q)} = \mathcal{Q}_d(\mathbf{x})$. For convenience, we assume that $\mathbf{H}$ is strictly positive definite and that $L(\mathbf{w})$ admits a unique global optimum as Zou et al. (2023). We first recall the definition of the global minima $\mathbf{w}^*$ and $\mathbf{w}^{(q)^*}$:

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \mathbb{E}_{\mathbf{x},y}\left[(y - \langle \mathbf{w}, \mathbf{x}\rangle)^2\right], \quad \mathbf{w}^{(q)^*} = \operatorname{argmin}_{\mathbf{w}} \mathbb{E}_{\mathbf{x},y}\left[(\mathcal{Q}_l(y) - \langle \mathbf{w}, \mathcal{Q}_d(\mathbf{x})\rangle)^2\right].$$

The first order optimality shows that

$$\mathbb{E}_{(\mathbf{x},y)\sim\mathcal{D}}[(y - \langle \mathbf{w}^*, \mathbf{x}\rangle)\mathbf{x}] = \mathbf{0}, \quad \mathbb{E}_{(\mathbf{x},y)\sim\mathcal{D}}[(\mathcal{Q}_l(y) - \langle \mathbf{w}^{(q)^*}, \mathcal{Q}_d(\mathbf{x})\rangle)\mathcal{Q}_d(\mathbf{x})] = \mathbf{0}, \tag{A.1}$$

which implies

$$\mathbf{w}^* = \mathbf{H}^{-1}\mathbb{E}_{(\mathbf{x},y)\sim\mathcal{D}}[y\mathbf{x}], \quad \mathbf{w}^{(q)^*} = (\mathbf{H}^{(q)})^{-1}\mathbb{E}\left[\mathcal{Q}_l(y)\mathcal{Q}_d(\mathbf{x})\right] = (\mathbf{H}^{(q)})^{-1}\mathbb{E}_{(\mathbf{x},y)\sim\mathcal{D}}[y\mathbf{x}].$$

Hence, by $\mathbf{H}^{(q)} = \mathbf{H} + \mathbf{D}$, we can characterize the difference between $\mathbf{w}^{(q)^*}$ and $\mathbf{w}^*$ as:

$$\begin{aligned}
\mathbf{w}^{(q)^*} - \mathbf{w}^* &= \left[(\mathbf{H}^{(q)})^{-1} - \mathbf{H}^{-1}\right]\mathbb{E}_{(\mathbf{x},y)\sim\mathcal{D}}[y\mathbf{x}] \\
&= (\mathbf{H}^{(q)})^{-1}\left(\mathbf{H} - \mathbf{H}^{(q)}\right)\mathbf{H}^{-1}\mathbb{E}_{(\mathbf{x},y)\sim\mathcal{D}}[y\mathbf{x}] \\
&= (\mathbf{H}^{(q)})^{-1}\left(\mathbf{H} - \mathbf{H}^{(q)}\right)\mathbf{w}^* \\
&= -(\mathbf{H}^{(q)})^{-1}\mathbf{D}\mathbf{w}^* \\
&= -(\mathbf{H} + \mathbf{D})^{-1}\mathbf{D}\mathbf{w}^*.
\end{aligned} \tag{A.2}$$

## A.1 Deviation of the Update Rule

In this section, we derive the evolution of parameter deviation $\boldsymbol{\eta}_t := \mathbf{w}_t - \mathbf{w}^{(q)^*}$.

**Lemma A.1.** (Error Propagation)

$$\boldsymbol{\eta}_t = \left(\mathbf{I} - \frac{1}{B}\gamma\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right)\boldsymbol{\eta}_{t-1} + \gamma\frac{1}{B}\mathcal{Q}_d(\mathbf{X}_t)^\top\left[\boldsymbol{\xi}_t + \boldsymbol{\epsilon}_t^{(o)} - \boldsymbol{\epsilon}_t^{(a)} - \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\epsilon}_{t-1}^{(p)}\right],$$

*where the quantization errors are*

$$\begin{aligned}
\boldsymbol{\epsilon}_t^{(o)} &:= \mathcal{Q}_o\left(\mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_a\left(\mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1})\right)\right) - \left[\mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_a\left(\mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1})\right)\right], \\
\boldsymbol{\epsilon}_t^{(a)} &:= \mathcal{Q}_a\left(\mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1})\right) - \mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1}), \\
\boldsymbol{\epsilon}_{t-1}^{(p)} &:= \mathcal{Q}_p(\mathbf{w}_{t-1}) - \mathbf{w}_{t-1}, \\
\boldsymbol{\xi}_t &:= \mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_d(\mathbf{X}_t)\mathbf{w}^{(q)^*}.
\end{aligned}$$

*Proof.* The lemma can be proved directly by the parameter update rule. By definition and the update rule of $\mathbf{w}_t$ (Quantized SGD),

$$\begin{aligned}
\boldsymbol{\eta}_t &= \mathbf{w}_t - \mathbf{w}^{(q)^*} \\
&= \mathbf{w}_{t-1} - \mathbf{w}^{(q)^*} + \gamma\frac{1}{B}\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_o\left(\mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_a\left(\mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1})\right)\right) \\
&= \boldsymbol{\eta}_{t-1} + \gamma\frac{1}{B}\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_o\left(\mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_a\left(\mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1})\right)\right).
\end{aligned}$$

We then introduce quantization errors to better characterize each quantization operation $\mathcal{Q}(\cdot)$. In particular, define quantization erros:

$$\boldsymbol{\epsilon}_t^{(o)} := \mathcal{Q}_o\left(\mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_a\left(\mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1})\right)\right) - \left[\mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_a\left(\mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1})\right)\right],$$

$$\boldsymbol{\epsilon}_t^{(a)} := \mathcal{Q}_a\left(\mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1})\right) - \mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1}),$$

$$\boldsymbol{\epsilon}_{t-1}^{(p)} := \mathcal{Q}_p(\mathbf{w}_{t-1}) - \mathbf{w}_{t-1},$$

$$\boldsymbol{\xi}_t := \mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_d(\mathbf{X}_t)\mathbf{w}^{(q)^*}.$$

Then the update rule for the parameter deviation can be expressed as:

$$\begin{aligned}
\boldsymbol{\eta}_t =& \boldsymbol{\eta}_{t-1} + \gamma\frac{1}{B}\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_o\left(\mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_a\left(\mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1})\right)\right) \\
=& \boldsymbol{\eta}_{t-1} + \gamma\mathcal{Q}_d(\mathbf{X}_t)^\top\frac{1}{B}\left[\mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_a\left(\mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1})\right)\right] + \gamma\frac{1}{B}\mathcal{Q}_d(\mathbf{X}_t)^\top\boldsymbol{\epsilon}_t^{(o)} \\
=& \boldsymbol{\eta}_{t-1} + \gamma\mathcal{Q}_d(\mathbf{X}_t)^\top\frac{1}{B}\left[\mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1})\right] + \gamma\frac{1}{B}\mathcal{Q}_d(\mathbf{X}_t)^\top(\boldsymbol{\epsilon}_t^{(o)} - \boldsymbol{\epsilon}_t^{(a)}) \\
=& \boldsymbol{\eta}_{t-1} + \gamma\frac{1}{B}\mathcal{Q}_d(\mathbf{X}_t)^\top(\boldsymbol{\epsilon}_t^{(o)} - \boldsymbol{\epsilon}_t^{(a)}) + \gamma\mathcal{Q}_d(\mathbf{X}_t)^\top\frac{1}{B} \\
& \left[\mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_d(\mathbf{X}_t)\mathbf{w}^{(q)^*} - \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\eta}_{t-1} - \mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1}) + \mathcal{Q}_d(\mathbf{X}_t)\mathbf{w}_{t-1}\right] \\
=& \boldsymbol{\eta}_{t-1} + \gamma\mathcal{Q}_d(\mathbf{X}_t)^\top(\boldsymbol{\epsilon}_t^{(o)} - \boldsymbol{\epsilon}_t^{(a)}) + \gamma\mathcal{Q}_d(\mathbf{X}_t)^\top\frac{1}{B} \\
& \left[\mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_d(\mathbf{X}_t)\mathbf{w}^{(q)^*} - \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\eta}_{t-1} - \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\epsilon}_{t-1}^{(p)}\right] \\
=& \boldsymbol{\eta}_{t-1} + \gamma\frac{1}{B}\mathcal{Q}_d(\mathbf{X}_t)^\top(\boldsymbol{\epsilon}_t^{(o)} - \boldsymbol{\epsilon}_t^{(a)} + \boldsymbol{\xi}_t) - \gamma\mathcal{Q}_d(\mathbf{X}_t)^\top\frac{1}{B}\left[\mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\eta}_{t-1} + \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\epsilon}_{t-1}^{(p)}\right] \\
=& \left(\mathbf{I} - \frac{1}{B}\gamma\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right)\boldsymbol{\eta}_{t-1} + \gamma\frac{1}{B}\mathcal{Q}_d(\mathbf{X}_t)^\top\left[\boldsymbol{\xi}_t + \boldsymbol{\epsilon}_t^{(o)} - \boldsymbol{\epsilon}_t^{(a)} - \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\epsilon}_{t-1}^{(p)}\right].
\end{aligned}$$

$\square$

## A.2 Decomposition of the Excess Risk

In this section, we take the initial step to analyze the excess risk of averaged SGD iterate $\overline{\mathbf{w}}_N$. In particular, we define the deviation of the averaged SGD iterate as $\overline{\boldsymbol{\eta}}_N := \frac{1}{N}\sum_{t=0}^{N-1}\boldsymbol{\eta}_t$. We decompose the excess risk as follows.

**Lemma A.2.** (Excess Risk Decomposition) *Under Assumption 3.1 and Assumption 3.2,*

$$\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] = R_1 + R_2 + R_3 + R_4,$$

*where*

$$\begin{aligned}
R_1 =& \frac{1}{2}\mathbb{E}\left[(y - \mathcal{Q}_l(y))^2\right] + \frac{1}{2}\mathbb{E}\left[\langle\overline{\mathbf{w}}_N, \mathcal{Q}_d(\mathbf{x}) - \mathbf{x}\rangle^2\right], \\
R_2 =& \frac{1}{2}\langle\mathbf{H}^{(q)}, \mathbb{E}[\overline{\boldsymbol{\eta}}_N \otimes \overline{\boldsymbol{\eta}}_N]\rangle, \\
R_3 =& \frac{1}{2}\mathbb{E}\left[(\mathcal{Q}_l(y) - y)^2\right] + \frac{1}{2}\mathbb{E}\left[\langle\mathbf{w}^{(q)^*}, \mathcal{Q}_d(\mathbf{x}) - \mathbf{x}\rangle^2\right], \\
R_4 =& \frac{1}{2}\langle\mathbf{H}, \mathbb{E}[(\mathbf{w}^* - \mathbf{w}^{(q)^*}) \otimes (\mathbf{w}^* - \mathbf{w}^{(q)^*})]\rangle.
\end{aligned}$$

*Proof.* By the definition of the excess risk ([Excess Risk](#)),

$$
\begin{aligned}
\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] =& \frac{1}{2}\mathbb{E}\left[\mathbb{E}_{\mathbf{x},y}\left[(y-\langle\overline{\mathbf{w}}_N,\mathbf{x}\rangle)^2\right]\right] - \frac{1}{2}\mathbb{E}_{\mathbf{x},y}\left[(y-\langle\mathbf{w}^*,\mathbf{x}\rangle)^2\right] \\
=& \underbrace{\frac{1}{2}\mathbb{E}\left[\mathbb{E}_{\mathbf{x},y}\left[(y-\langle\overline{\mathbf{w}}_N,\mathbf{x}\rangle)^2\right]\right] - \frac{1}{2}\mathbb{E}\left[(\mathcal{Q}_l(y)-\langle\overline{\mathbf{w}}_N,\mathcal{Q}_d(\mathbf{x})\rangle)^2\right]}_{R_1} \\
& + \underbrace{\frac{1}{2}\mathbb{E}\left[(\mathcal{Q}_l(y)-\langle\overline{\mathbf{w}}_N,\mathcal{Q}_d(\mathbf{x})\rangle)^2\right] - \frac{1}{2}\mathbb{E}\left[(\mathcal{Q}_l(y)-\langle\mathbf{w}^{(q)^*},\mathcal{Q}_d(\mathbf{x})\rangle)^2\right]}_{R_2} \\
& + \underbrace{\frac{1}{2}\mathbb{E}\left[(\mathcal{Q}_l(y)-\langle\mathbf{w}^{(q)^*},\mathcal{Q}_d(\mathbf{x})\rangle)^2\right] - \frac{1}{2}\mathbb{E}\left[\left(y-\langle\mathbf{w}^{(q)^*},\mathbf{x}\rangle\right)^2\right]}_{R_3} \\
& + \underbrace{\frac{1}{2}\mathbb{E}\left[\left(y-\langle\mathbf{w}^{(q)^*},\mathbf{x}\rangle\right)^2\right] - \frac{1}{2}\mathbb{E}_{\mathbf{x},y}\left[(y-\langle\mathbf{w}^*,\mathbf{x}\rangle)^2\right]}_{R_4},
\end{aligned}
$$

where $R_1$ captures the gap of the averaged SGD iterate between the full-precision and quantized domains, $R_2$ characterizes the distance from the averaged SGD iterate to the quantized optimal solution within the quantized domain, $R_3$ represents the mismatch of the quantized optimal solution in full-precision data space and quantized data space and $R_4$ defines the discrepancy between the averaged SGD iterate and the quantized optimal solution in the full-precision domain. Next, we compute $R_1, R_2, R_3$ and $R_4$ respectively. These computations are mainly based on the first order optimality condition (A.1) and the unbiased quantization Assumption 3.1. For $R_4$,

$$
\begin{aligned}
R_4 =& \frac{1}{2}\mathbb{E}\left[\mathbb{E}_{\mathbf{x},y}\left[\left(y-\langle\mathbf{w}^{(q)^*},\mathbf{x}\rangle\right)^2\right]\right] - \frac{1}{2}\mathbb{E}_{\mathbf{x},y}\left[(y-\langle\mathbf{w}^*,\mathbf{x}\rangle)^2\right] \\
=& \frac{1}{2}\mathbb{E}\left[\mathbb{E}_{\mathbf{x},y}\left[\langle\mathbf{w}^*-\mathbf{w}^{(q)^*},\mathbf{x}\rangle\cdot\left(2y-\langle\mathbf{w}^*+\mathbf{w}^{(q)^*},\mathbf{x}\rangle\right)\right]\right] \\
=& \frac{1}{2}\mathbb{E}\left[\mathbb{E}_{\mathbf{x}}\left[\langle\mathbf{w}^*-\mathbf{w}^{(q)^*},\mathbf{x}\rangle^2\right]\right] \\
=& \frac{1}{2}\langle\mathbf{H},\mathbb{E}[(\mathbf{w}^*-\mathbf{w}^{(q)^*})\otimes(\mathbf{w}^*-\mathbf{w}^{(q)^*})]\rangle,
\end{aligned}
$$

where the third equality uses the first order optimality condition that $\mathbb{E}_{(\mathbf{x},y)\sim\mathcal{D}}[(y-\langle\mathbf{w}^*,\mathbf{x}\rangle)\mathbf{x}]=\mathbf{0}$.

For $R_2$, similarly by the first order optimality condition (A.1) with respect to $\mathbf{w}^{(q)^*}$, it holds

$$
\begin{aligned}
R_2 =& \frac{1}{2}\mathbb{E}\left[(\mathcal{Q}_l(y)-\langle\overline{\mathbf{w}}_N,\mathcal{Q}_d(\mathbf{x})\rangle)^2\right] - \frac{1}{2}\mathbb{E}\left[(\mathcal{Q}_l(y)-\langle\mathbf{w}^{(q)^*},\mathcal{Q}_d(\mathbf{x})\rangle)^2\right] \\
=& \frac{1}{2}\mathbb{E}\left[\langle\mathbf{w}^{(q)^*}-\overline{\mathbf{w}}_N,\mathcal{Q}_d(\mathbf{x})\rangle\cdot\left(2\mathcal{Q}_l(y)-\langle\mathbf{w}^{(q)^*}+\overline{\mathbf{w}}_N,\mathcal{Q}_d(\mathbf{x})\rangle\right)\right] \\
=& \frac{1}{2}\mathbb{E}\left[\mathbb{E}_{\mathbf{x}}\left[\langle\mathbf{w}^{(q)^*}-\overline{\mathbf{w}}_N,\mathcal{Q}_d(\mathbf{x})\rangle^2\right]\right] \\
=& \frac{1}{2}\langle\mathbf{H}^{(q)},\mathbb{E}[\bar{\boldsymbol{\eta}}_N\otimes\bar{\boldsymbol{\eta}}_N]\rangle.
\end{aligned}
$$

For $R_3$,

$$
\begin{aligned}
R_3 =& \frac{1}{2}\mathbb{E}\left[(\mathcal{Q}_l(y) - \langle \mathbf{w}^{(q)*}, \mathcal{Q}_d(\mathbf{x})\rangle)^2\right] - \frac{1}{2}\mathbb{E}\left[\left(y - \langle \mathbf{w}^{(q)*}, \mathbf{x}\rangle\right)^2\right] \\
=& \frac{1}{2}\mathbb{E}\left[\left(\mathcal{Q}_l(y) - y - \langle \mathbf{w}^{(q)*}, \mathcal{Q}_d(\mathbf{x}) - \mathbf{x}\rangle\right) \cdot \left(\mathcal{Q}_l(y) + y - \langle \mathbf{w}^{(q)*}, \mathcal{Q}_d(\mathbf{x}) + \mathbf{x}\rangle\right)\right] \\
=& \frac{1}{2}\mathbb{E}\left[\mathcal{Q}_l(y)^2 - y^2 + \langle \mathbf{w}^{(q)*}, \mathcal{Q}_d(\mathbf{x}) - \mathbf{x}\rangle\langle \mathbf{w}^{(q)*}, \mathcal{Q}_d(\mathbf{x}) + \mathbf{x}\rangle\right] \\
=& \frac{1}{2}\mathbb{E}\left[(\mathcal{Q}_l(y) - y)^2\right] + \frac{1}{2}\mathbb{E}\left[\langle \mathbf{w}^{(q)*}, \mathcal{Q}_d(\mathbf{x}) - \mathbf{x}\rangle^2\right],
\end{aligned}
$$

where the third and the last equality utilize the unbiased quantization Assumption 3.1.

For $R_1$, similarly by the unbiased quantization Assumption 3.1, it holds

$$
\begin{aligned}
R_1 =& \frac{1}{2}\mathbb{E}\left[\mathbb{E}_{\mathbf{x},y}\left[(y - \langle \overline{\mathbf{w}}_N, \mathbf{x}\rangle)^2\right]\right] - \frac{1}{2}\mathbb{E}\left[(\mathcal{Q}_l(y) - \langle \overline{\mathbf{w}}_N, \mathcal{Q}_d(\mathbf{x})\rangle)^2\right] \\
=& \frac{1}{2}\mathbb{E}\left[(y - \mathcal{Q}_l(y) - \langle \overline{\mathbf{w}}_N, \mathbf{x} - \mathcal{Q}_d(\mathbf{x})\rangle) \cdot (y + \mathcal{Q}_l(y) - \langle \overline{\mathbf{w}}_N, \mathbf{x} + \mathcal{Q}_d(\mathbf{x})\rangle)\right] \\
=& \frac{1}{2}\mathbb{E}\left[(y - \mathcal{Q}_l(y))^2\right] + \frac{1}{2}\mathbb{E}\left[\langle \overline{\mathbf{w}}_N, \mathcal{Q}_d(\mathbf{x}) - \mathbf{x}\rangle^2\right].
\end{aligned}
$$

$\square$

In the following Section B-Section E, we provide analysis for $R_1, R_2, R_3$ and $R_4$ respectively. Building on this bounds, we can derive the excess risk bounds and prove our main theorems.

# B  Analysis of $R_4$

**Lemma B.1.** *Under assumptions and notations in Lemma A.2,*

$$
R_4 = \frac{1}{2}\|\mathbf{w}^*\|_{\mathbf{D}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{H}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{D}}.
$$

*Proof.* Recall that by Lemma A.2,

$$
R_4 = \frac{1}{2}\langle \mathbf{H}, \mathbb{E}[(\mathbf{w}^* - \mathbf{w}^{(q)*}) \otimes (\mathbf{w}^* - \mathbf{w}^{(q)*})]\rangle,
$$

and further note that by Equation (A.2),

$$
\mathbb{E}[(\mathbf{w}^* - \mathbf{w}^{(q)*}) \otimes (\mathbf{w}^* - \mathbf{w}^{(q)*})] = \mathbb{E}\left[(\mathbf{H}+\mathbf{D})^{-1}\mathbf{D}\mathbf{w}^*\mathbf{w}^{*\top}\mathbf{D}(\mathbf{H}+\mathbf{D})^{-1}\right],
$$

we have

$$
\begin{aligned}
R_4 =& \frac{1}{2}\mathbb{E}\left[\mathrm{tr}\left(\mathbf{H}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{D}\mathbf{w}^*\mathbf{w}^{*\top}\mathbf{D}(\mathbf{H}+\mathbf{D})^{-1}\right)\right] \\
=& \frac{1}{2}\mathbb{E}\left[\mathrm{tr}\left(\mathbf{w}^{*\top}\mathbf{D}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{H}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{D}\mathbf{w}^*\right)\right] \\
=& \frac{1}{2}\|\mathbf{w}^*\|_{\mathbf{D}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{H}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{D}}.
\end{aligned}
$$

This immediately completes the proof. $\square$

# C   Analysis of $R_3$

**Lemma C.1.** *Under assumptions and notations in Lemma A.2,*

$$R_3 = \frac{\mathbb{E}\left[(\epsilon^{(l)})^2\right]}{2} + \frac{1}{2}\|\mathbf{w}^*\|^2_{\mathbf{H}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{D}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{H}}.$$

*Proof.* By the definition of $R_3$ in Lemma A.2,

$$
\begin{aligned}
R_3 &= \frac{1}{2}\mathbb{E}\left[(\mathcal{Q}_l(y) - y)^2\right] + \frac{1}{2}\mathbb{E}\left[\langle \mathbf{w}^{(q)*}, \mathcal{Q}_d(\mathbf{x}) - \mathbf{x}\rangle^2\right] \\
&= \frac{1}{2}\mathbb{E}\left[(\mathcal{Q}_l(y) - y)^2\right] + \frac{1}{2}\mathbf{w}^{(q)*\top}\mathbf{D}\mathbf{w}^{(q)*} \\
&= \frac{1}{2}\mathbb{E}\left[(\mathcal{Q}_l(y) - y)^2\right] + \frac{1}{2}\mathbf{w}^{*\top}\mathbf{H}\left(\mathbf{H}+\mathbf{D}\right)^{-1}\mathbf{D}\left(\mathbf{H}+\mathbf{D}\right)^{-1}\mathbf{H}\mathbf{w}^* \\
&= \frac{\mathbb{E}\left[(\epsilon^{(l)})^2\right]}{2} + \frac{1}{2}\|\mathbf{w}^*\|^2_{\mathbf{H}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{D}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{H}},
\end{aligned}
$$

where the second equality uses the definition of $\mathbf{D} = \mathbb{E}[(\mathcal{Q}_d(\mathbf{x}) - \mathbf{x})(\mathcal{Q}_d(\mathbf{x}) - \mathbf{x})^\top]$, the third equality uses the expression of $\mathbf{w}^{(q)*}$ (A.1), and the last equality utilizes the definition of $\epsilon^{(l)} = \mathcal{Q}_l(y) - y$. $\quad\square$

# D   Analysis of $R_2$

## D.1   Preliminary

We first define the following linear operators as in Zou et al. (2023):

$$\mathcal{I} = \mathbf{I} \otimes \mathbf{I}, \quad \mathcal{M}^{(q)} = \mathbb{E}[\mathbf{x}^{(q)} \otimes \mathbf{x}^{(q)} \otimes \mathbf{x}^{(q)} \otimes \mathbf{x}^{(q)}], \quad \widetilde{\mathcal{M}}^{(q)} = \mathbf{H}^{(q)} \otimes \mathbf{H}^{(q)},$$

$$\mathcal{T}^{(q)} = \mathbf{H}^{(q)} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^{(q)} - \gamma\mathcal{M}^{(q)}, \quad \widetilde{\mathcal{T}}^{(q)} = \mathbf{H}^{(q)} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}^{(q)} - \gamma\mathbf{H}^{(q)} \otimes \mathbf{H}^{(q)}.$$

For a symmetric matrix $\mathbf{A}$, the above definitions result in:

$$\mathcal{I} \circ \mathbf{A} = \mathbf{A}, \quad \mathcal{M}^{(q)} \circ \mathbf{A} = \mathbb{E}[(\mathbf{x}^{(q)\top}\mathbf{A}\mathbf{x}^{(q)})\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}], \quad \widetilde{\mathcal{M}}^{(q)} \circ \mathbf{A} = \mathbf{H}^{(q)}\mathbf{A}\mathbf{H}^{(q)},$$

$$(\mathcal{I} - \gamma\mathcal{T}^{(q)}) \circ \mathbf{A} = \mathbb{E}[(\mathbf{I} - \gamma\mathbf{x}^{(q)}\mathbf{x}^{(q)\top})\mathbf{A}(\mathbf{I} - \gamma\mathbf{x}^{(q)}\mathbf{x}^{(q)\top})], \quad (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)}) \circ \mathbf{A} = (\mathbf{I} - \gamma\mathbf{H}^{(q)})\mathbf{A}(\mathbf{I} - \gamma\mathbf{H}^{(q)}).$$

Further, we generalize the linear operators from Zou et al. (2023) to account for batch size effects. For a symmetric matrix $\mathbf{A}$, we define

$$
\begin{aligned}
\mathcal{M}_B^{(q)} \circ \mathbf{A} &= \mathbb{E}\left[\frac{1}{B^2}\mathbf{X}^{(q)\top}\mathbf{X}^{(q)}\mathbf{A}\mathbf{X}^{(q)\top}\mathbf{X}^{(q)}\right], \\
(\mathcal{I} - \gamma\mathcal{T}_B^{(q)}) \circ \mathbf{A} &= \mathbb{E}\left[\left(\mathbf{I} - \gamma\frac{1}{B}\mathbf{X}^{(q)\top}\mathbf{X}^{(q)}\right)\mathbf{A}\left(\mathbf{I} - \gamma\frac{1}{B}\mathbf{X}^{(q)\top}\mathbf{X}^{(q)}\right)\right].
\end{aligned}
$$

## D.2   Initial Study of $R_2$

To analyze $R_2$, we first substitute $\bar{\boldsymbol{\eta}}_N$ with the summation of $\boldsymbol{\eta}_t$. This step mainly based on the expansion of Lemma A.1,

$$\mathbb{E}\left[\boldsymbol{\eta}_t|\boldsymbol{\eta}_{t-1}\right] = \left(\mathbf{I} - \gamma\mathbf{H}^{(q)}\right)\boldsymbol{\eta}_{t-1},$$

which holds under the unbiased quantization Assumption 3.1 and the first order optimality condition (A.1). We summarize as the following lemma.

**Lemma D.1.** *Under assumptions and notation in Lemma A.2,*

$$R_2 \le \frac{1}{N^2} \cdot \sum_{t=0}^{N-1} \sum_{k=t}^{N-1} \left\langle (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{k-t} \mathbf{H}^{(q)}, \mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t] \right\rangle.$$

*Proof.* By Lemma A.2,

$$R_2 = \frac{1}{2} \langle \mathbf{H}^{(q)}, \mathbb{E}[\bar{\boldsymbol{\eta}}_N \otimes \bar{\boldsymbol{\eta}}_N] \rangle.$$

Then we focus on $\mathbb{E}[\bar{\boldsymbol{\eta}}_N \otimes \bar{\boldsymbol{\eta}}_N]$. By definition $\bar{\boldsymbol{\eta}}_N = \frac{1}{N} \sum_{t=0}^{N-1} \boldsymbol{\eta}_t$,

$$\begin{aligned}
\mathbb{E}[\bar{\boldsymbol{\eta}}_N \otimes \bar{\boldsymbol{\eta}}_N] =& \frac{1}{N^2} \cdot \left( \sum_{0 \le k \le t \le N-1} \mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_k] + \sum_{0 \le t < k \le N-1} \mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_k] \right) \\
\preceq& \frac{1}{N^2} \cdot \left( \sum_{0 \le k \le t \le N-1} \mathbb{E}\left[ \mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_k | \boldsymbol{\eta}_k] \right] + \sum_{0 \le t \le k \le N-1} \mathbb{E}\left[ \mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_k | \boldsymbol{\eta}_t] \right] \right).
\end{aligned} \tag{D.1}$$

Note that by the unbiased Assumption 3.1,

$$\mathbb{E}\left[ \gamma \mathcal{Q}_d(\mathbf{X}_t)^\top \left( \boldsymbol{\epsilon}_t^{(o)} - \boldsymbol{\epsilon}_t^{(a)} - \mathcal{Q}_d(\mathbf{X}_t) \boldsymbol{\epsilon}_{t-1}^{(p)} \right) \Big| \boldsymbol{\eta}_{t-1} \right] = \mathbf{0}.$$

Further, by the optimality (A.1),

$$\mathbb{E}\left[ \gamma \mathcal{Q}_d(\mathbf{X}_t)^\top \boldsymbol{\xi}_t \Big| \boldsymbol{\eta}_{t-1} \right] = \mathbb{E}\left[ \gamma \mathcal{Q}_d(\mathbf{X}_t)^\top \left[ \mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_d(\mathbf{X}_t) \mathbf{w}^{(q)*} \right] \Big| \boldsymbol{\eta}_{t-1} \right] = \mathbf{0}.$$

Hence, by Lemma A.1,

$$\mathbb{E}\left[ \boldsymbol{\eta}_t | \boldsymbol{\eta}_{t-1} \right] = \left( \mathbf{I} - \gamma \mathbf{H}^{(q)} \right) \boldsymbol{\eta}_{t-1}. \tag{D.2}$$

Therefore, by (D.1) and (D.2),

$$\begin{aligned}
& \mathbb{E}[\bar{\boldsymbol{\eta}}_N \otimes \bar{\boldsymbol{\eta}}_N] \\
\preceq& \frac{1}{N^2} \cdot \left( \sum_{0 \le k \le t \le N-1} \mathbb{E}\left[ \mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_k | \boldsymbol{\eta}_k] \right] + \sum_{0 \le t \le k \le N-1} \mathbb{E}\left[ \mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_k | \boldsymbol{\eta}_t] \right] \right) \\
=& \frac{1}{N^2} \cdot \left( \sum_{0 \le k \le t \le N-1} (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{t-k} \mathbb{E}[\boldsymbol{\eta}_k \otimes \boldsymbol{\eta}_k] + \sum_{0 \le t \le k \le N-1} \mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t] (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{k-t} \right) \\
=& \frac{1}{N^2} \cdot \sum_{t=0}^{N-1} \sum_{k=t}^{N-1} \left( (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{k-t} \mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t] + \mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t] (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{k-t} \right).
\end{aligned} \tag{D.3}$$

Applying (D.3) into $R_2$, we have

$$\begin{aligned}
R_2 =& \frac{1}{2} \langle \mathbf{H}^{(q)}, \mathbb{E}[\bar{\boldsymbol{\eta}}_N \otimes \bar{\boldsymbol{\eta}}_N] \rangle \\
\le& \frac{1}{2N^2} \cdot \sum_{t=0}^{N-1} \sum_{k=t}^{N-1} \left\langle \mathbf{H}^{(q)}, (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{k-t} \mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t] + \mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t] (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{k-t} \right\rangle \\
=& \frac{1}{N^2} \cdot \sum_{t=0}^{N-1} \sum_{k=t}^{N-1} \left\langle (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{k-t} \mathbf{H}^{(q)}, \mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t] \right\rangle,
\end{aligned}$$

where the last equality holds since $\mathbf{H}^{(q)}$ and $(\mathbf{I} - \gamma \mathbf{H}^{(q)})^{k-t}$ commute. This completes the proof. $\square$

Lemma D.1 implies that, to bound $R_2$, the main goal is to bound $\mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t]$. Next lemma provides an update rule for iteration $\mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t]$.

**Lemma D.2.** (Update Rule) *Under Assumption 3.1,3.2,3.3,3.4,*

$$
\begin{aligned}
\mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t] =& \mathbb{E}\left[\left(\mathbf{I} - \gamma\frac{1}{B}\mathcal{Q}_d(\mathbf{X})^\top \mathcal{Q}_d(\mathbf{X})\right)\mathbb{E}[\boldsymbol{\eta}_{t-1} \otimes \boldsymbol{\eta}_{t-1}]\left(\mathbf{I} - \gamma\frac{1}{B}\mathcal{Q}_d(\mathbf{X})^\top \mathcal{Q}_d(\mathbf{X})\right)\right] + \gamma^2 \boldsymbol{\Sigma}_t \\
\preceq& \mathbb{E}\left[\left(\mathbf{I} - \gamma\frac{1}{B}\mathcal{Q}_d(\mathbf{X})^\top \mathcal{Q}_d(\mathbf{X})\right)\mathbb{E}[\boldsymbol{\eta}_{t-1} \otimes \boldsymbol{\eta}_{t-1}]\left(\mathbf{I} - \gamma\frac{1}{B}\mathcal{Q}_d(\mathbf{X})^\top \mathcal{Q}_d(\mathbf{X})\right)\right] + \gamma^2 \sigma_G^{(q)2}\mathbf{H}^{(q)},
\end{aligned}
$$

*where*

$$
\begin{aligned}
\boldsymbol{\Sigma}_t :=& \frac{1}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top\left[\boldsymbol{\xi}_t + \boldsymbol{\epsilon}_t^{(o)} - \boldsymbol{\epsilon}_t^{(a)} - \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\epsilon}_{t-1}^{(p)}\right]\left[\boldsymbol{\xi}_t + \boldsymbol{\epsilon}_t^{(o)} - \boldsymbol{\epsilon}_t^{(a)} - \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\epsilon}_{t-1}^{(p)}\right]^\top \mathcal{Q}_d(\mathbf{X}_t)\right] \\
=& \underbrace{\frac{1}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top\boldsymbol{\xi}_t\boldsymbol{\xi}_t^\top \mathcal{Q}_d(\mathbf{X}_t)\right]}_{\boldsymbol{\Sigma}_t^\xi} + \underbrace{\frac{1}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top\boldsymbol{\epsilon}_t^{(o)}\boldsymbol{\epsilon}_t^{(o)\top} \mathcal{Q}_d(\mathbf{X}_t)\right]}_{\boldsymbol{\Sigma}_t^{\epsilon(o)}} \\
& + \underbrace{\frac{1}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top\boldsymbol{\epsilon}_t^{(a)}\boldsymbol{\epsilon}_t^{(a)\top} \mathcal{Q}_d(\mathbf{X}_t)\right]}_{\boldsymbol{\Sigma}_t^{\epsilon(a)}} + \underbrace{\frac{1}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\epsilon}_{t-1}^{(p)}\boldsymbol{\epsilon}_{t-1}^{(p)\top} \mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right]}_{\boldsymbol{\Sigma}_t^{\epsilon(p)}},
\end{aligned}
$$

*and*

$$
\sigma_G^{(q)2} = \frac{\sup_t\left\|\mathbb{E}\left[\boldsymbol{\epsilon}_t^{(o)}\boldsymbol{\epsilon}_t^{(o)\top}|\mathbf{o}_t\right] + \mathbb{E}\left[\boldsymbol{\epsilon}_t^{(a)}\boldsymbol{\epsilon}_t^{(a)\top}|\mathbf{a}_t\right]\right\|}{B} + \alpha_B\mathbb{E}_{\mathbf{w}_{t-1}}\left[\operatorname{tr}\left(\mathbf{H}^{(q)}\mathbb{E}\left[\boldsymbol{\epsilon}_{t-1}^{(p)}\boldsymbol{\epsilon}_{t-1}^{(p)\top}|\mathbf{w}_{t-1}\right]\right)\right] + \frac{\sigma^2}{B},
$$

*with* $\mathbf{a}_t = \mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1})$, $\mathbf{o}_t = \mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_a\left(\mathcal{Q}_d(\mathbf{X}_t)\mathcal{Q}_p(\mathbf{w}_{t-1})\right)$ *and* $\|\cdot\|$ *denoting the spectral norm.*

*Proof.* By Lemma A.1,

$$
\boldsymbol{\eta}_t = \left(\mathbf{I} - \frac{1}{B}\gamma\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\right)\boldsymbol{\eta}_{t-1} + \gamma\frac{1}{B}\mathcal{Q}_d(\mathbf{X}_t)^\top\left[\boldsymbol{\xi}_t + \boldsymbol{\epsilon}_t^{(o)} - \boldsymbol{\epsilon}_t^{(a)} - \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\epsilon}_{t-1}^{(p)}\right].
$$

Hence, by the unbiased quantization Assumption 3.1,

$$
\mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t] = \mathbb{E}\left[\left(\mathbf{I} - \frac{1}{B}\gamma\mathcal{Q}_d(\mathbf{X})^\top \mathcal{Q}_d(\mathbf{X})\right)\mathbb{E}[\boldsymbol{\eta}_{t-1} \otimes \boldsymbol{\eta}_{t-1}]\left(\mathbf{I} - \frac{1}{B}\gamma\mathcal{Q}_d(\mathbf{X})^\top \mathcal{Q}_d(\mathbf{X})\right)\right] + \gamma^2 \boldsymbol{\Sigma}_t,
\tag{D.4}
$$

where we complete the proof for the first statement.

Next, we cope with each term in $\boldsymbol{\Sigma}_t$ to provide an upper bound. For $\boldsymbol{\Sigma}_t^{\epsilon(p)}$,

$$
\begin{aligned}
\boldsymbol{\Sigma}_t^{\epsilon(p)} =& \frac{1}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\epsilon}_{t-1}^{(p)}\boldsymbol{\epsilon}_{t-1}^{(p)\top}\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right] \\
=& \frac{1}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\mathbb{E}\left[\boldsymbol{\epsilon}_{t-1}^{(p)}\boldsymbol{\epsilon}_{t-1}^{(p)\top}|\mathbf{w}_{t-1}\right]\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right] \\
\preceq& \alpha_B\mathbb{E}_{\mathbf{w}_{t-1}}\left[\operatorname{tr}\left(\mathbf{H}^{(q)}\mathbb{E}\left[\boldsymbol{\epsilon}_{t-1}^{(p)}\boldsymbol{\epsilon}_{t-1}^{(p)\top}|\mathbf{w}_{t-1}\right]\right)\right]\mathbf{H}^{(q)},
\end{aligned}
$$

where the inequality holds by Assumption 3.3.

23

For $\mathbf{\Sigma}_t^\xi$,

$$
\begin{aligned}
\mathbf{\Sigma}_t^\xi &= \frac{1}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top\boldsymbol{\xi}_t\boldsymbol{\xi}_t^\top\mathcal{Q}_d(\mathbf{X}_t)\right]\\
&= \frac{1}{B^2}\mathbb{E}\left[\sum_{i=1}^B\sum_{j=1}^B\mathcal{Q}_d(\mathbf{X}_t)^{i\top}\boldsymbol{\xi}_t^i\left(\mathcal{Q}_d(\mathbf{X}_t)^{j\top}\boldsymbol{\xi}_t^j\right)^\top\right]\\
&= \frac{1}{B^2}\sum_{i=1}^B\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^{i\top}\boldsymbol{\xi}_t^i\left(\mathcal{Q}_d(\mathbf{X}_t)^{i\top}\boldsymbol{\xi}_t^i\right)^\top\right]\\
&= \frac{1}{B}\cdot\mathbb{E}\left[\mathcal{Q}_d(\mathbf{x})\xi\left(\mathcal{Q}_d(\mathbf{x})\xi\right)^\top\right]\\
&= \frac{1}{B}\cdot\mathbb{E}\left[\xi^2\mathcal{Q}_d(\mathbf{x})\mathcal{Q}_d(\mathbf{x})^\top\right]\\
&\preceq \frac{\sigma^2}{B}\cdot\mathbf{H}^{(q)},
\end{aligned}
\tag{D.5}
$$

where the second equality holds as samples are independent and data quantization is applied to each sample independently, the inequality holds by Assumption 3.4.

For $\mathbf{\Sigma}_t^{\epsilon^{(o)}}+\mathbf{\Sigma}_t^{\epsilon^{(a)}}$,

$$
\begin{aligned}
\mathbf{\Sigma}_t^{\epsilon^{(o)}}+\mathbf{\Sigma}_t^{\epsilon^{(a)}} &= \frac{1}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top(\boldsymbol{\epsilon}_t^{(o)}\boldsymbol{\epsilon}_t^{(o)\top}+\boldsymbol{\epsilon}_t^{(a)}\boldsymbol{\epsilon}_t^{(a)\top})\mathcal{Q}_d(\mathbf{X}_t)\right]\\
&= \frac{1}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top\left(\mathbb{E}\left[\boldsymbol{\epsilon}_t^{(o)}\boldsymbol{\epsilon}_t^{(o)\top}|\mathbf{o}_t\right]+\mathbb{E}\left[\boldsymbol{\epsilon}_t^{(a)}\boldsymbol{\epsilon}_t^{(a)\top}|\mathbf{a}_t\right]\right)\mathcal{Q}_d(\mathbf{X}_t)\right]\\
&\preceq \frac{1}{B^2}\mathbb{E}\left[\left(\left\|\mathbb{E}\left[\boldsymbol{\epsilon}_t^{(o)}\boldsymbol{\epsilon}_t^{(o)\top}|\mathbf{o}_t\right]+\mathbb{E}\left[\boldsymbol{\epsilon}_t^{(a)}\boldsymbol{\epsilon}_t^{(a)\top}|\mathbf{a}_t\right]\right\|\right)\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right]\\
&\preceq \frac{1}{B^2}\sup_t\left\|\mathbb{E}\left[\boldsymbol{\epsilon}_t^{(o)}\boldsymbol{\epsilon}_t^{(o)\top}|\mathbf{o}_t\right]+\mathbb{E}\left[\boldsymbol{\epsilon}_t^{(a)}\boldsymbol{\epsilon}_t^{(a)\top}|\mathbf{a}_t\right]\right\|\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right]\\
&= \frac{1}{B}\sup_t\left\|\mathbb{E}\left[\boldsymbol{\epsilon}_t^{(o)}\boldsymbol{\epsilon}_t^{(o)\top}|\mathbf{o}_t\right]+\mathbb{E}\left[\boldsymbol{\epsilon}_t^{(a)}\boldsymbol{\epsilon}_t^{(a)\top}|\mathbf{a}_t\right]\right\|\mathbf{H}^{(q)},
\end{aligned}
$$

where $\|\cdot\|$ represents the matrix spectral norm.

Combining the upper bounds for $\mathbf{\Sigma}_t^{\epsilon^{(p)}}$, $\mathbf{\Sigma}_t^\xi$ and $\mathbf{\Sigma}_t^{\epsilon^{(o)}}+\mathbf{\Sigma}_t^{\epsilon^{(a)}}$ immediately completes the proof. $\qquad\square$

For multiplicative quantization, the explicit dependence of the conditional expectations on $\mathbf{w}_t$ renders Lemma D.2 inapplicable to the update rule for $\mathbb{E}[\boldsymbol{\eta}_t\otimes\boldsymbol{\eta}_t]$. We thus propose the following alternative update rule.

**Lemma D.3.** (Update Rule under Multiplicative Quantization) *If there exist $\epsilon_d,\epsilon_l,\epsilon_p,\epsilon_a$ and $\epsilon_o$ such that for any $i\in\{d,l,p,a,o\}$, quantization $\mathcal{Q}_i$ is $\epsilon_i$-multiplicative, then under Assumption 3.1,3.2,3.3,3.4, it holds*

$$
\begin{aligned}
\mathbb{E}[\boldsymbol{\eta}_t\otimes\boldsymbol{\eta}_t] &= \mathbb{E}\left[\left(\mathbf{I}-\gamma\frac{1}{B}\mathcal{Q}_d(\mathbf{X})^\top\mathcal{Q}_d(\mathbf{X})\right)\mathbb{E}[\boldsymbol{\eta}_{t-1}\otimes\boldsymbol{\eta}_{t-1}]\left(\mathbf{I}-\gamma\frac{1}{B}\mathcal{Q}_d(\mathbf{X})^\top\mathcal{Q}_d(\mathbf{X})\right)\right]+\gamma^2\mathbf{\Sigma}_t\\
&\preceq \mathbb{E}\left[\left(\mathbf{I}-\frac{1}{B}\gamma\mathcal{Q}_d(\mathbf{X})^\top\mathcal{Q}_d(\mathbf{X})\right)\mathbb{E}[\boldsymbol{\eta}_{t-1}\otimes\boldsymbol{\eta}_{t-1}]\left(\mathbf{I}-\frac{1}{B}\gamma\mathcal{Q}_d(\mathbf{X})^\top\mathcal{Q}_d(\mathbf{X})\right)\right]\\
&\quad+\tilde{\epsilon}\mathbb{E}\left[\frac{\gamma}{B}\mathcal{Q}_d(\mathbf{X})^\top\mathcal{Q}_d(\mathbf{X})\mathbb{E}[\boldsymbol{\eta}_{t-1}\otimes\boldsymbol{\eta}_{t-1}]\frac{\gamma}{B}\mathcal{Q}_d(\mathbf{X})^\top\mathcal{Q}_d(\mathbf{X})\right]+\gamma^2\sigma_M^{(q)2}\mathbf{H}^{(q)},
\end{aligned}
$$

*where $\mathbf{\Sigma}_t$ is the same as defined in Lemma D.2 and*

$$
\tilde{\epsilon} = 2\epsilon_p+4\epsilon_o(1+\epsilon_a)(1+\epsilon_p)+2\epsilon_a(1+\epsilon_p),
$$

24

$$\sigma_M^{(q)2} = \frac{(1+4\epsilon_o)\sigma^2}{B} + \frac{\|\mathbf{w}^*\|_{\mathbf{H}}^2}{1+\epsilon_d}\alpha_B\left(4\epsilon_o[(1+\epsilon_a)(1+\epsilon_p)+1]+2\epsilon_a(1+\epsilon_p)+2\epsilon_p\right).$$

*Proof.* The first statement has been proved in Lemma D.2. To complete the proof, we merely need to derive the upper bound for $\mathbf{\Sigma}_t = \mathbf{\Sigma}_t^\xi + \mathbf{\Sigma}_t^{(a)} + \mathbf{\Sigma}_t^{(o)} + \mathbf{\Sigma}_t^{(p)}$, where $\mathbf{\Sigma}_t^\xi$, $\mathbf{\Sigma}_t^{(a)}$, $\mathbf{\Sigma}_t^{(o)}$ and $\mathbf{\Sigma}_t^{(p)}$ are defined in Lemma D.2. Regarding $\mathbf{\Sigma}_t^\xi$, by the computation in the proof of Lemma D.2, i.e., (D.5),

$$\mathbf{\Sigma}_t^\xi \preceq \frac{\sigma^2}{B}\mathbf{H}^{(q)}. \tag{D.6}$$

Regarding $\mathbf{\Sigma}_t^{\epsilon^{(p)}}$,

$$
\begin{aligned}
\mathbf{\Sigma}_t^{\epsilon^{(p)}} =& \frac{1}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\mathbb{E}\left[\boldsymbol{\epsilon}_{t-1}^{(p)}\boldsymbol{\epsilon}_{t-1}^{(p)}{}^\top\Big|\mathbf{w}_{t-1}\right]\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\right]\\
=& \frac{\epsilon_p}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\mathbf{w}_{t-1}\mathbf{w}_{t-1}^\top \mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\right]\\
\preceq& \frac{2\epsilon_p}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\eta}_{t-1}\boldsymbol{\eta}_{t-1}{}^\top \mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\right]\\
+& \frac{2\epsilon_p}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\mathbf{w}^{(q)*}\mathbf{w}^{(q)*}{}^\top \mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\right],
\end{aligned}
\tag{D.7}
$$

where the second equality utilizes the definition of multiplicative quantization and the inequality uses the definition of $\boldsymbol{\eta}_t$. Building on the upper bound for $\mathbf{\Sigma}_t^{(p)}$, we can derive upper bounds for $\mathbf{\Sigma}_t^{\epsilon^{(a)}}$ and $\mathbf{\Sigma}_t^{\epsilon^{(o)}}$.

Regarding $\mathbf{\Sigma}_t^{\epsilon^{(a)}}$,

$$
\begin{aligned}
\mathbf{\Sigma}_t^{\epsilon^{(a)}} =& \frac{1}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \boldsymbol{\epsilon}_t^{(a)}\boldsymbol{\epsilon}_t^{(a)}{}^\top \mathcal{Q}_d(\mathbf{X}_t)\right]\\
=& \frac{\epsilon_a}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\mathbf{w}_{t-1}^{(q)}\mathbf{w}_{t-1}^{(q)}{}^\top \mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\right]\\
=& \frac{\epsilon_a}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\mathbf{w}_{t-1}\mathbf{w}_{t-1}{}^\top \mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\right]\\
+& \frac{\epsilon_a}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\epsilon}_{t-1}^{(p)}\boldsymbol{\epsilon}_{t-1}^{(p)}{}^\top \mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\right]\\
=& \frac{(1+\epsilon_p)\epsilon_a}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\mathbf{w}_{t-1}\mathbf{w}_{t-1}{}^\top \mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\right]\\
\preceq& \frac{2(1+\epsilon_p)\epsilon_a}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\eta}_{t-1}\boldsymbol{\eta}_{t-1}{}^\top \mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\right]\\
+& \frac{2(1+\epsilon_p)\epsilon_a}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\mathbf{w}^{(q)*}\mathbf{w}^{(q)*}{}^\top \mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\right],
\end{aligned}
\tag{D.8}
$$

where the second and the fourth equality holds by the definition of multiplicative quantization, the third equality holds by the definition of the quantization error $\boldsymbol{\epsilon}_t^{(p)}$ and the inequality uses the bound for $\mathbf{\Sigma}_t^{(p)}$ (D.7).

Regarding $\boldsymbol{\Sigma}_t^{\epsilon^{(o)}}$, similar to $\boldsymbol{\Sigma}_t^{\epsilon^{(a)}}$, it holds

$$
\begin{aligned}
\boldsymbol{\Sigma}_t^{\epsilon^{(o)}} =& \frac{1}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \boldsymbol{\epsilon}_t^{(o)}\boldsymbol{\epsilon}_t^{(o)^\top}\mathcal{Q}_d(\mathbf{X}_t)\right]\\
=& \frac{\epsilon_o}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathbf{o}_t\mathbf{o}_t^\top\mathcal{Q}_d(\mathbf{X}_t)\right]\\
\preceq& \frac{2\epsilon_o}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_l(\mathbf{y}_t)\mathcal{Q}_l(\mathbf{y}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right] + \frac{2\epsilon_o}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_a(\mathbf{a}_t)\mathcal{Q}_a(\mathbf{a}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right]\\
\preceq& \frac{2\epsilon_o}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_l(\mathbf{y}_t)\mathcal{Q}_l(\mathbf{y}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right] + \frac{2(1+\epsilon_a)\epsilon_o}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathbf{a}_t\mathbf{a}_t^\top\mathcal{Q}_d(\mathbf{X}_t)\right]\\
\preceq& \frac{2\epsilon_o}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_l(\mathbf{y}_t)\mathcal{Q}_l(\mathbf{y}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right]\\
&+ \frac{4(1+\epsilon_p)(1+\epsilon_a)\epsilon_o}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\eta}_{t-1}\boldsymbol{\eta}_{t-1}^\top\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right]\\
&+ \frac{4(1+\epsilon_p)(1+\epsilon_a)\epsilon_o}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\mathbf{w}^{(q)*}\mathbf{w}^{(q)*^\top}\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right]\\
\preceq& \frac{4\epsilon_o}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \boldsymbol{\xi}_t\boldsymbol{\xi}_t^\top\mathcal{Q}_d(\mathbf{X}_t)\right]\\
&+ \frac{4(1+\epsilon_p)(1+\epsilon_a)\epsilon_o}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\eta}_{t-1}\boldsymbol{\eta}_{t-1}^\top\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right]\\
&+ \frac{4[(1+\epsilon_p)(1+\epsilon_a)\epsilon_o+1]}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\mathbf{w}^{(q)*}\mathbf{w}^{(q)*^\top}\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right],
\end{aligned}
$$

where the last inequality holds by the fact that $\boldsymbol{\xi}_t = \mathcal{Q}_l(\mathbf{y}_t) - \mathcal{Q}_d(\mathbf{X}_t)\mathbf{w}^{(q)*}$. Further, by the bound for $\boldsymbol{\Sigma}_t^\xi$ (D.5),

$$
\frac{1}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \boldsymbol{\xi}_t\boldsymbol{\xi}_t^\top\mathcal{Q}_d(\mathbf{X}_t)\right] \preceq \frac{\sigma^2}{B}\mathbf{H}^{(q)},
$$

we have

$$
\begin{aligned}
\boldsymbol{\Sigma}_t^{\epsilon^{(o)}} \preceq& \frac{4\epsilon_o\sigma^2}{B}\mathbf{H}^{(q)} + \frac{4(1+\epsilon_p)(1+\epsilon_a)\epsilon_o}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\eta}_{t-1}\boldsymbol{\eta}_{t-1}^\top\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right]\\
&+ \frac{4[(1+\epsilon_p)(1+\epsilon_a)\epsilon_o+1]}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\mathbf{w}^{(q)*}\mathbf{w}^{(q)*^\top}\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right].
\end{aligned} \tag{D.9}
$$

Further, by Assumption 3.3,

$$
\frac{1}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\mathbf{w}^{(q)*}\mathbf{w}^{(q)*^\top}\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right] \preceq \alpha_B\mathrm{tr}\left(\mathbf{H}^{(q)}\mathbf{w}^{(q)*}\mathbf{w}^{(q)*^\top}\right)\mathbf{H}^{(q)},
$$

then together with (D.6), (D.7), (D.8) and (D.9) it holds

$$
\begin{aligned}
\boldsymbol{\Sigma}_t \preceq& \frac{(1+4\epsilon_o)\sigma^2}{B}\mathbf{H}^{(q)} + \alpha_B\left(4\epsilon_o[(1+\epsilon_a)(1+\epsilon_p)+1] + 2\epsilon_a(1+\epsilon_p)+2\epsilon_p\right)\mathrm{tr}\left(\mathbf{H}^{(q)}\mathbf{w}^{(q)*}\mathbf{w}^{(q)*^\top}\right)\mathbf{H}^{(q)}\\
&+ \frac{2\epsilon_p + 4\epsilon_o(1+\epsilon_a)(1+\epsilon_p)+2\epsilon_a(1+\epsilon_p)}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\eta}_{t-1}\boldsymbol{\eta}_{t-1}^\top\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right].
\end{aligned}
$$

Note that by the definition of multiplicative quantization,

$$
\mathrm{tr}\left(\mathbf{H}^{(q)}\mathbf{w}^{(q)*}\mathbf{w}^{(q)*^\top}\right) = \frac{\|\mathbf{w}^*\|_{\mathbf{H}}^2}{1+\epsilon_d},
$$

then

$$
\begin{aligned}
\boldsymbol{\Sigma}_t \preceq& \left[\frac{(1+4\epsilon_o)\sigma^2}{B} + \frac{\|\mathbf{w}^*\|_{\mathbf{H}}^2}{1+\epsilon_d}\alpha_B\left(4\epsilon_o[(1+\epsilon_a)(1+\epsilon_p)+1]+2\epsilon_a(1+\epsilon_p)+2\epsilon_p\right)\right]\mathbf{H}^{(q)}\\
&+ \frac{2\epsilon_p + 4\epsilon_o(1+\epsilon_a)(1+\epsilon_p)+2\epsilon_a(1+\epsilon_p)}{B^2}\mathbb{E}\left[\mathcal{Q}_d(\mathbf{X}_t)^\top \mathcal{Q}_d(\mathbf{X}_t)\boldsymbol{\eta}_{t-1}\boldsymbol{\eta}_{t-1}^\top\mathcal{Q}_d(\mathbf{X}_t)^\top\mathcal{Q}_d(\mathbf{X}_t)\right].
\end{aligned} \tag{D.10}
$$

Hence, by (D.10) and (D.4)

$$\mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t] \leq \mathbb{E}\left[\left(\mathbf{I} - \frac{1}{B}\gamma \mathcal{Q}_d(\mathbf{X})^\top \mathcal{Q}_d(\mathbf{X})\right)\mathbb{E}[\boldsymbol{\eta}_{t-1} \otimes \boldsymbol{\eta}_{t-1}]\left(\mathbf{I} - \frac{1}{B}\gamma \mathcal{Q}_d(\mathbf{X})^\top \mathcal{Q}_d(\mathbf{X})\right)\right]$$

$$+ [2\epsilon_p + 4\epsilon_o(1+\epsilon_a)(1+\epsilon_p) + 2\epsilon_a(1+\epsilon_p)]\,\mathbb{E}\left[\frac{\gamma}{B}\mathcal{Q}_d(\mathbf{X})^\top \mathcal{Q}_d(\mathbf{X})\mathbb{E}[\boldsymbol{\eta}_{t-1} \otimes \boldsymbol{\eta}_{t-1}]\frac{\gamma}{B}\mathcal{Q}_d(\mathbf{X})^\top \mathcal{Q}_d(\mathbf{X})\right]$$

$$+\gamma^2\left[\frac{(1+4\epsilon_o)\sigma^2}{B} + \frac{\|\mathbf{w}^*\|_{\mathbf{H}}^2}{1+\epsilon_d}\alpha_B\left(4\epsilon_o[(1+\epsilon_a)(1+\epsilon_p)+1] + 2\epsilon_a(1+\epsilon_p) + 2\epsilon_p\right)\right]\mathbf{H}^{(q)}.$$

$\square$

Equipped with Lemma D.1, Lemma D.2 and Lemma D.3, we are ready to derive bounds for $R_2$. As shown in Zou et al. (2023), we first perform bias-variance decomposition.

## D.3 Bias-Variance Decomposition

For simplicity, denote $\mathbf{A}_t := \mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t]$. Then under general quantization, Lemma D.2 shows

$$\mathbf{A}_t = (\mathcal{I} - \gamma \mathcal{T}_B^{(q)}) \circ \mathbf{A}_{t-1} + \gamma^2 \boldsymbol{\Sigma}_t \preceq (\mathcal{I} - \gamma \mathcal{T}_B^{(q)}) \circ \mathbf{A}_{t-1} + \gamma^2 \sigma_G^{(q)^2}\mathbf{H}^{(q)}. \tag{D.11}$$

Under multiplicative quantization, Lemma D.3 shows

$$\mathbf{A}_t = (\mathcal{I} - \gamma \mathcal{T}_B^{(q)}) \circ \mathbf{A}_{t-1} + \gamma^2 \boldsymbol{\Sigma}_t \preceq (\mathcal{I} - \gamma \mathcal{T}_B^{(q)} + \tilde{\epsilon}\gamma^2 \mathcal{M}_B^{(q)}) \circ \mathbf{A}_{t-1} + \gamma^2 \sigma_M^{(q)^2}\mathbf{H}^{(q)}. \tag{D.12}$$

As in Zou et al. (2023), we perform bias-variance for excess risk, which is summarized as the following lemma.

**Lemma D.4.** (Bias-Variance Decomposition) *Under Assumption 3.1,3.2,3.3,3.4,*

$$R_2 \leq \underbrace{\frac{1}{N^2} \cdot \sum_{t=0}^{N-1}\sum_{k=t}^{N-1}\left\langle(\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k-t}\mathbf{H}^{(q)}, \mathbf{B}_t\right\rangle}_{\text{bias}} + \underbrace{\frac{1}{N^2} \cdot \sum_{t=0}^{N-1}\sum_{k=t}^{N-1}\left\langle(\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k-t}\mathbf{H}^{(q)}, \mathbf{C}_t\right\rangle}_{\text{variance}},$$

*where*

$$\mathbf{B}_t := (\mathcal{I} - \gamma\mathcal{T}_B^{(q)})^t \circ \mathbf{B}_0, \quad \mathbf{B}_0 = \mathbb{E}\left[\boldsymbol{\eta}_0 \otimes \boldsymbol{\eta}_0\right].$$

$$\mathbf{C}_t := (\mathcal{I} - \gamma\mathcal{T}_B^{(q)})\mathbf{C}_{t-1} + \gamma^2\sigma_G^{(q)^2}\mathbf{H}^{(q)}, \quad \mathbf{C}_0 = \mathbf{0}.$$

*Proof.* By Lemma D.1,

$$R_2 \leq \frac{1}{N^2} \cdot \sum_{t=0}^{N-1}\sum_{k=t}^{N-1}\left\langle(\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k-t}\mathbf{H}^{(q)}, \mathbf{A}_t\right\rangle.$$

The proof is immediately completed owing to

$$\mathbf{A}_t = \mathbf{B}_t + \mathbf{C}_t,$$

where $\mathbf{A}_t$ is defined in (D.11). $\square$

For multiplicative quantization, we can directly deduce Lemma D.4 by the update rule under multiplicative quantization (D.12).

**Lemma D.5.** (Bias-Variance Decomposition under Multiplicative Quantization) *Under Assumption 3.1,3.2,3.3,3.4, if there exist $\epsilon_d, \epsilon_l, \epsilon_p, \epsilon_a$ and $\epsilon_o$ such that for any $i \in \{d, l, p, a, o\}$, quantization $\mathcal{Q}_i$ is $\epsilon_i$-multiplicative, then*

$$R_2 \leq \underbrace{\frac{1}{N^2} \cdot \sum_{t=0}^{N-1} \sum_{k=t}^{N-1} \left\langle (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{k-t} \mathbf{H}^{(q)}, \mathbf{B}_t^{(M)} \right\rangle}_{\text{bias}} + \underbrace{\frac{1}{N^2} \cdot \sum_{t=0}^{N-1} \sum_{k=t}^{N-1} \left\langle (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{k-t} \mathbf{H}^{(q)}, \mathbf{C}_t^{(M)} \right\rangle}_{\text{variance}},$$

*where*

$$\mathbf{B}_t^{(M)} := (\mathcal{I} - \gamma \mathcal{T}_B^{(q)} + \tilde{\epsilon} \gamma^2 \mathcal{M}_B^{(q)})^t \circ \mathbf{B}_0^{(M)}, \quad \mathbf{B}_0^{(M)} = \mathbb{E}\left[\boldsymbol{\eta}_0 \otimes \boldsymbol{\eta}_0\right].$$

$$\mathbf{C}_t^{(M)} := (\mathcal{I} - \gamma \mathcal{T}_B^{(q)} + \tilde{\epsilon} \gamma^2 \mathcal{M}_B^{(q)}) \mathbf{C}_{t-1}^{(M)} + \gamma^2 {\sigma_M^{(q)}}^2 \mathbf{H}^{(q)}, \quad \mathbf{C}_0^{(M)} = 0.$$

## D.4   Bounding the Bias Error

By Lemma D.4,

$$
\begin{aligned}
\text{bias} &= \frac{1}{N^2} \sum_{t=0}^{N-1} \sum_{k=t}^{N-1} \left\langle (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{k-t} \mathbf{H}^{(q)}, \mathbf{B}_t \right\rangle \\
&= \frac{1}{\gamma N^2} \sum_{t=0}^{N-1} \left\langle \mathbf{I} - (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{N-t}, \mathbf{B}_t \right\rangle \\
&\leq \frac{1}{\gamma N^2} \langle \mathbf{I} - (\mathbf{I} - \gamma \mathbf{H}^{(q)})^N, \sum_{t=0}^{N-1} \mathbf{B}_t \rangle.
\end{aligned}
\tag{D.13}
$$

For $1 \leq n \leq N$, let $\mathbf{S}_n = \sum_{t=0}^{n-1} \mathbf{B}_t$, $\mathbf{S}_n^{(M)} = \sum_{t=0}^{n-1} \mathbf{B}_t^{(M)}$, then we only need to bound $\mathbf{S}_N$ and $\mathbf{S}_N^{(M)}$ to bound bias term under general quantization and multiplicative quantization, respectively. We first derive the update rule for $\mathbf{S}_t$ and $\mathbf{S}_t^{(M)}$.

**Lemma D.6.** (Initial Study of $\mathbf{S}_t$) *For $1 \leq t \leq N$,*

$$\mathbf{S}_t \preceq (\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)}) \circ \mathbf{S}_{t-1} + \gamma^2 \mathcal{M}_B^{(q)} \circ \mathbf{S}_N + \mathbf{B}_0.$$

*Proof.* By definition,

$$
\begin{aligned}
\mathbf{S}_t &= \sum_{k=0}^{t-1} (\mathcal{I} - \gamma \mathcal{T}_B^{(q)})^k \circ \mathbf{B}_0 \\
&= (\mathcal{I} - \gamma \mathcal{T}_B^{(q)}) \circ \left( \sum_{k=1}^{t-1} (\mathcal{I} - \gamma \mathcal{T}_B^{(q)})^{k-1} \circ \mathbf{B}_0 \right) + \mathbf{B}_0 \\
&= (\mathcal{I} - \gamma \mathcal{T}_B^{(q)}) \circ \mathbf{S}_{t-1} + \mathbf{B}_0.
\end{aligned}
\tag{D.14}
$$

Then we convert $\mathcal{T}_B^{(q)}$ to $\mathcal{T}_B^{(q)}$. By (D.14),

$$
\begin{aligned}
\mathbf{S}_t &= (\mathcal{I} - \gamma \mathcal{T}_B^{(q)}) \circ \mathbf{S}_{t-1} + \mathbf{B}_0 \\
&= (\mathcal{I} - \gamma \widetilde{\mathcal{T}}^{(q)}) \circ \mathbf{S}_{t-1} + \gamma (\widetilde{\mathcal{T}}^{(q)} - \mathcal{T}_B^{(q)}) \circ \mathbf{S}_{t-1} + \mathbf{B}_0 \\
&= (\mathcal{I} - \gamma \widetilde{\mathcal{T}}^{(q)}) \circ \mathbf{S}_{t-1} + \gamma^2 (\mathcal{M}_B^{(q)} - \widetilde{\mathcal{M}}^{(q)}) \circ \mathbf{S}_{t-1} + \mathbf{B}_0 \\
&\preceq (\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)}) \circ \mathbf{S}_{t-1} + \gamma^2 \mathcal{M}_B^{(q)} \circ \mathbf{S}_N + \mathbf{B}_0,
\end{aligned}
$$

where the third equality holds by the definition of linear operators. $\qquad \square$

**Lemma D.7.** (Initial Study of $\mathbf{S}_t^{(M)}$) *For* $1 \le t \le N$,

$$\mathbf{S}_t^{(M)} \preceq (\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)}) \circ \mathbf{S}_{t-1}^{(M)} + (1 + \tilde{\epsilon}) \gamma^2 \mathcal{M}_B^{(q)} \circ \mathbf{S}_N^{(M)} + \mathbf{B}_0.$$

*Proof.* The proof is similar to the proof for Lemma D.6.

$$\begin{aligned}
\mathbf{S}_t^{(M)} &= (\mathcal{I} - \gamma \mathcal{T}_B^{(q)} + \tilde{\epsilon} \gamma^2 \mathcal{M}_B^{(q)}) \circ \mathbf{S}_{t-1}^{(M)} + \mathbf{B}_0 \\
&= (\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)}) \circ \mathbf{S}_{t-1} + \gamma (\tilde{\mathcal{T}}^{(q)} - \mathcal{T}_B^{(q)}) \circ \mathbf{S}_{t-1}^{(M)} + \tilde{\epsilon} \gamma^2 \mathcal{M}_B^{(q)} \circ \mathbf{S}_{t-1}^{(M)} + \mathbf{B}_0 \\
&= (\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)}) \circ \mathbf{S}_{t-1}^{(M)} + \gamma^2 ((1 + \tilde{\epsilon}) \mathcal{M}_B^{(q)} - \widetilde{\mathcal{M}}^{(q)}) \circ \mathbf{S}_{t-1}^{(M)} + \mathbf{B}_0 \\
&\preceq (\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)}) \circ \mathbf{S}_{t-1}^{(M)} + (1 + \tilde{\epsilon}) \gamma^2 \mathcal{M}_B^{(q)} \circ \mathbf{S}_N^{(M)} + \mathbf{B}_0.
\end{aligned}$$

$\square$

**Lemma D.8.** (A Bound for $\mathcal{M}_B^{(q)} \circ \mathbf{S}_t$) *For* $1 \le t \le N$, *under Assumption 3.1,3.2,3.3,3.4, if* $\gamma < \frac{1}{\alpha_B \mathrm{tr}(\mathbf{H}^{(q)})}$, *then*

$$\mathcal{M}_B^{(q)} \circ \mathbf{S}_t \preceq \frac{\alpha_B \cdot \mathrm{tr}\left(\left[\mathcal{I} - (\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)})^t\right] \circ \mathbf{B}_0\right)}{\gamma(1 - \gamma \alpha_B \mathrm{tr}(\mathbf{H}^{(q)}))} \cdot \mathbf{H}^{(q)}.$$

*Proof.* The first step is to derive a crude bound for $\mathbf{S}_t$. Take summation via the update rule, we have

$$\mathbf{S}_t = \sum_{k=0}^{t-1} (\mathcal{I} - \gamma \mathcal{T}_B^{(q)})^k \circ \mathbf{B}_0 = \gamma^{-1} {\mathcal{T}_B^{(q)}}^{-1} \circ \left[\mathcal{I} - (\mathcal{I} - \gamma \mathcal{T}_B^{(q)})^t\right] \circ \mathbf{B}_0.$$

Note that

$$\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)} \preceq \mathcal{I} - \gamma \mathcal{T}_B^{(q)}, \quad (\mathcal{I} - (\mathcal{I} - \gamma \mathcal{T}_B^{(q)})^t) \preceq (\mathcal{I} - (\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)})^t),$$

and further note that ${\mathcal{T}_B^{(q)}}^{-1}$ is a PSD mapping [7], and $[\mathcal{I} - (\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)})^t] \circ \mathbf{B}_0$ is a PSD matrix, we obtain

$$\mathbf{S}_t \preceq \gamma^{-1} {\mathcal{T}_B^{(q)}}^{-1} \circ (\mathcal{I} - (\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)})^t) \circ \mathbf{B}_0.$$

For simplicity, we denote $\mathbf{A} := (\mathcal{I} - (\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)})^t) \circ \mathbf{B}_0$. We then tackle ${\mathcal{T}_B^{(q)}}^{-1} \circ \mathbf{A}$. To be specific, we apply $\tilde{\mathcal{T}}^{(q)}$.

$$\begin{aligned}
\tilde{\mathcal{T}}^{(q)} \circ {\mathcal{T}_B^{(q)}}^{-1} \circ \mathbf{A} &= \gamma \mathcal{M}_B^{(q)} \circ {\mathcal{T}_B^{(q)}}^{-1} \circ \mathbf{A} + \mathbf{A} - \gamma \mathbf{H}^{(q)} ({\mathcal{T}_B^{(q)}}^{-1} \circ \mathbf{A}) \mathbf{H}^{(q)} \\
&\preceq \gamma \mathcal{M}_B^{(q)} \circ {\mathcal{T}_B^{(q)}}^{-1} \circ \mathbf{A} + \mathbf{A}.
\end{aligned}$$

Therefore,

$${\mathcal{T}_B^{(q)}}^{-1} \circ \mathbf{A} \preceq \gamma (\tilde{\mathcal{T}}^{(q)})^{-1} \circ \mathcal{M}_B^{(q)} \circ {\mathcal{T}_B^{(q)}}^{-1} \circ \mathbf{A} + (\tilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{A}.$$

Then we undertake the second step, applying $\mathcal{M}_B^{(q)}$ on both sides.

$$\begin{aligned}
\mathcal{M}_B^{(q)} \circ ({\mathcal{T}_B^{(q)}}^{-1} \circ \mathbf{A}) &\preceq \mathcal{M}_B^{(q)} \circ \gamma (\tilde{\mathcal{T}}^{(q)})^{-1} \circ \mathcal{M}_B^{(q)} \circ {\mathcal{T}_B^{(q)}}^{-1} \circ \mathbf{A} + \mathcal{M}_B^{(q)} \circ (\tilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{A} \\
&\preceq \sum_{t=0}^{\infty} (\gamma \mathcal{M}_B^{(q)} \circ (\tilde{\mathcal{T}}^{(q)})^{-1})^t \circ (\mathcal{M}_B^{(q)} \circ (\tilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{A}) \text{ (By recursion).}
\end{aligned} \tag{D.15}$$

---

[7] ${\mathcal{T}_B^{(q)}}^{-1}$ is a PSD mapping under the condition that $\gamma < \frac{1}{\alpha_B \mathrm{tr}(\mathbf{H}^{(q)})}$, which can be directly deduced by Lemma B.1 in Zou et al. (2023). We omit the proof here for simplicity.

By Assumption 3.3,

$$\mathcal{M}_B^{(q)} \circ (\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{A} \preceq \alpha_B \operatorname{tr}(\mathbf{H}^{(q)}(\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{A})\mathbf{H}^{(q)}$$

$$= \alpha_B \gamma \operatorname{tr}\left(\sum_{t=0}^{\infty} \mathbf{H}^{(q)}(\mathbf{I} - \gamma\mathbf{H}^{(q)})^t \mathbf{A}(\mathbf{I} - \gamma\mathbf{H}^{(q)})^t\right)\mathbf{H}^{(q)}$$

$$= \alpha_B \operatorname{tr}\left(\mathbf{H}^{(q)}(2\mathbf{H}^{(q)} - \gamma(\mathbf{H}^{(q)})^2)^{-1}\mathbf{A}\right)\mathbf{H}^{(q)}$$

$$\preceq \alpha_B \operatorname{tr}(\mathbf{A})\mathbf{H}^{(q)},$$

where the first equality holds by the definition of $\widetilde{\mathcal{T}}^{(q)}$ and the last inequality requires the condition that $\gamma < \frac{1}{\alpha_B \operatorname{tr}(\mathbf{H}^{(q)})}$. Hence, by (D.15), and further by $(\widetilde{\mathcal{T}}^{(q)})^{-1}\mathbf{H}^{(q)} \preceq \mathbf{I}$ and $\mathcal{M}_B^{(q)} \circ \mathbf{I} \preceq \alpha_B \operatorname{tr}(\mathbf{H}^{(q)})\mathbf{H}^{(q)}$, we obtain

$$\mathcal{M}_B^{(q)} \circ (\mathcal{T}_B^{(q)^{-1}} \circ \mathbf{A}) \preceq \sum_{t=0}^{\infty} (\gamma\mathcal{M}_B^{(q)} \circ (\widetilde{\mathcal{T}}^{(q)})^{-1})^t \circ (\mathcal{M}_B^{(q)} \circ (\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{A})$$

$$\preceq \alpha_B \operatorname{tr}(\mathbf{A})\sum_{t=0}^{\infty}(\gamma\alpha_B \operatorname{tr}(\mathbf{H}^{(q)}))^t \mathbf{H}^{(q)}$$

$$\preceq \frac{\alpha_B \operatorname{tr}(\mathbf{A})}{1 - \gamma\alpha_B \operatorname{tr}(\mathbf{H}^{(q)})} \cdot \mathbf{H}^{(q)}.$$

Therefore,

$$\mathcal{M}_B^{(q)} \circ \mathbf{S}_t \preceq \gamma^{-1}\frac{\alpha_B \operatorname{tr}(\mathbf{A})}{1 - \gamma\alpha_B \operatorname{tr}(\mathbf{H}^{(q)})} \cdot \mathbf{H}^{(q)} = \frac{\alpha_B \cdot \operatorname{tr}\left(\left[\mathcal{I} - (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)})^t\right] \circ \mathbf{B}_0\right)}{\gamma(1 - \gamma\alpha_B \operatorname{tr}(\mathbf{H}^{(q)}))} \cdot \mathbf{H}^{(q)}.$$

$\square$

**Lemma D.9.** (A Bound for $\mathcal{M}_B^{(q)} \circ \mathbf{S}_t^{(M)}$) *For $1 \leq t \leq N$, under Assumption 3.1,3.2,3.3,3.4, if $\gamma < \frac{1}{(1+\tilde{\epsilon})\alpha_B \operatorname{tr}(\mathbf{H}^{(q)})}$,*

$$\mathcal{M}_B^{(q)} \circ \mathbf{S}_t^{(M)} \preceq \frac{\alpha_B \cdot \operatorname{tr}\left(\left[\mathcal{I} - (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)})^t\right] \circ \mathbf{B}_0\right)}{\gamma(1 - (1 + \tilde{\epsilon})\gamma\alpha_B \operatorname{tr}(\mathbf{H}^{(q)}))} \cdot \mathbf{H}^{(q)}.$$

*Proof.* The first step is to derive a crude bound for $\mathbf{S}_t^{(M)}$. Take summation via the update rule, we have [8]

$$\mathbf{S}_t^{(M)} = \sum_{k=0}^{t-1}(\mathcal{I} - \gamma\mathcal{T}_B^{(q)} + \tilde{\epsilon}\gamma^2\mathcal{M}_B^{(q)})^k \circ \mathbf{B}_0 = \gamma^{-1}(\mathcal{T}_B^{(q)} - \tilde{\epsilon}\gamma\mathcal{M}_B^{(q)})^{-1} \circ \left[\mathcal{I} - (\mathcal{I} - \gamma\mathcal{T}_B^{(q)} + \tilde{\epsilon}\gamma^2\mathcal{M}_B^{(q)})^t\right] \circ \mathbf{B}_0.$$

Note that

$$\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)} \preceq \mathcal{I} - \gamma\mathcal{T}_B^{(q)}, \quad (\mathcal{I} - (\mathcal{I} - \gamma\mathcal{T}_B^{(q)} + \tilde{\epsilon}\gamma^2\mathcal{M}_B^{(q)})^t) \preceq (\mathcal{I} - (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)} + \tilde{\epsilon}\gamma^2\mathcal{M}_B^{(q)})^t),$$

we obtain

$$\mathbf{S}_t^{(M)} \preceq \gamma^{-1}(\mathcal{T}_B^{(q)} - \tilde{\epsilon}\gamma\mathcal{M}_B^{(q)})^{-1} \circ (\mathcal{I} - (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)} + \tilde{\epsilon}\gamma^2\mathcal{M}_B^{(q)})^t) \circ \mathbf{B}_0.$$

---

[8] $(\mathcal{T}_B^{(q)} - \tilde{\epsilon}\gamma\mathcal{M}_B^{(q)})^{-1}$ is a PSD mapping under the condition that $\gamma < \frac{1}{(1+\tilde{\epsilon})\alpha_B \operatorname{tr}(\mathbf{H}^{(q)})}$, which can be directly deduced by Lemma B.1 in Zou et al. (2023). We omit the proof here for simplicity.

Denote $\mathbf{A} := (\mathcal{I} - (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)} + \tilde{\epsilon}\gamma^2\mathcal{M}_B^{(q)})^t) \circ \mathbf{B}_0$, then

$$\widetilde{\mathcal{T}}^{(q)} \circ (\mathcal{T}_B^{(q)} - \tilde{\epsilon}\gamma\mathcal{M}_B^{(q)})^{-1} \circ \mathbf{A} \preceq (1+\tilde{\epsilon})\gamma\mathcal{M}_B^{(q)} \circ (\mathcal{T}_B^{(q)} - \tilde{\epsilon}\gamma\mathcal{M}_B^{(q)})^{-1} \circ \mathbf{A} + \mathbf{A}.$$

Therefore

$$(\mathcal{T}_B^{(q)} - \tilde{\epsilon}\gamma\mathcal{M}_B^{(q)})^{-1} \circ \mathbf{A} \preceq (1+\tilde{\epsilon})\gamma(\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathcal{M}_B^{(q)} \circ (\mathcal{T}_B^{(q)} - \tilde{\epsilon}\gamma\mathcal{M}_B^{(q)})^{-1} \circ \mathbf{A} + (\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{A}.$$

Then we undertake the second step, applying $\mathcal{M}_B^{(q)}$ on both sides.

$$\mathcal{M}_B^{(q)} \circ (\mathcal{T}_B^{(q)} - \tilde{\epsilon}\gamma\mathcal{M}_B^{(q)})^{-1} \circ \mathbf{A} \preceq \sum_{t=0}^{\infty}((1+\tilde{\epsilon})\gamma\mathcal{M}_B^{(q)} \circ (\widetilde{\mathcal{T}}^{(q)})^{-1})^t \circ (\mathcal{M}_B^{(q)} \circ (\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{A}). \quad \text{(D.16)}$$

By Assumption 3.3,

$$\begin{aligned}
\mathcal{M}_B^{(q)} \circ (\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{A} &\preceq \alpha_B \operatorname{tr}(\mathbf{H}^{(q)}(\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{A})\mathbf{H}^{(q)} \\
&= \alpha_B\gamma \operatorname{tr}\left(\sum_{t=0}^{\infty}\mathbf{H}^{(q)}(\mathbf{I}-\gamma\mathbf{H}^{(q)})^t\mathbf{A}(\mathbf{I}-\gamma\mathbf{H}^{(q)})^t\right)\mathbf{H}^{(q)} \\
&= \alpha_B\operatorname{tr}\left(\mathbf{H}^{(q)}(2\mathbf{H}^{(q)} - \gamma(\mathbf{H}^{(q)})^2)^{-1}\mathbf{A}\right)\mathbf{H}^{(q)} \\
&\preceq \alpha_B\operatorname{tr}(\mathbf{A})\mathbf{H}^{(q)},
\end{aligned} \quad \text{(D.17)}$$

where the last inequality requires the condition that $\gamma < \frac{1}{\alpha_B\operatorname{tr}(\mathbf{H}^{(q)})}$. Hence, by (D.16), (D.17), and further by $(\widetilde{\mathcal{T}}^{(q)})^{-1}\mathbf{H}^{(q)} \preceq \mathbf{I}$ and $\mathcal{M}_B^{(q)} \circ \mathbf{I} \preceq \alpha_B\operatorname{tr}(\mathbf{H}^{(q)})\mathbf{H}^{(q)}$, we obtain

$$\begin{aligned}
\mathcal{M}_B^{(q)} \circ ((\mathcal{T}_B^{(q)} - \tilde{\epsilon}\gamma\mathcal{M}_B^{(q)})^{-1} \circ \mathbf{A}) &\preceq \sum_{t=0}^{\infty}((1+\tilde{\epsilon})\gamma\mathcal{M}_B^{(q)} \circ (\widetilde{\mathcal{T}}^{(q)})^{-1})^t \circ (\mathcal{M}_B^{(q)} \circ (\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{A}) \\
&\preceq \alpha_B\operatorname{tr}(\mathbf{A})\sum_{t=0}^{\infty}((1+\tilde{\epsilon})\gamma\alpha_B\operatorname{tr}(\mathbf{H}^{(q)}))^t\mathbf{H}^{(q)} \\
&\preceq \frac{\alpha_B\operatorname{tr}(\mathbf{A})}{1-(1+\tilde{\epsilon})\gamma\alpha_B\operatorname{tr}(\mathbf{H}^{(q)})} \cdot \mathbf{H}^{(q)}.
\end{aligned}$$

Therefore,

$$\mathcal{M}_B^{(q)} \circ \mathbf{S}_t^{(M)} \preceq \gamma^{-1}\frac{\alpha_B\operatorname{tr}(\mathbf{A})}{1-(1+\tilde{\epsilon})\gamma\alpha_B\operatorname{tr}(\mathbf{H}^{(q)})} \cdot \mathbf{H}^{(q)} \preceq \frac{\alpha_B \cdot \operatorname{tr}\left(\left[\mathcal{I} - (\mathcal{I}-\gamma\widetilde{\mathcal{T}}^{(q)})^t\right] \circ \mathbf{B}_0\right)}{\gamma(1-(1+\tilde{\epsilon})\gamma\alpha_B\operatorname{tr}(\mathbf{H}^{(q)}))} \cdot \mathbf{H}^{(q)}.$$

$\square$

By Lemma D.6, Lemma D.7, Lemma D.8 and Lemma D.9, we can provide a refined bound for $\mathbf{S}_t$ and $\mathbf{S}_t^{(M)}$. Then we are ready to bound the bias error.

**Lemma D.10.** (A Bound for bias) *Under Assumption 3.1,3.2,3.3,3.4, if the stepsize satisfies* $\gamma < \frac{1}{\alpha_B\operatorname{tr}(\mathbf{H}^{(q)})}$, *then*

$$\begin{aligned}
\text{bias} \leq &\frac{2\alpha_B\left(\|\mathbf{w}_0 - \mathbf{w}^{(q)^*}\|_{\mathbf{I}_{0:k^*}^{(q)}}^2 + N\gamma\|\mathbf{w}_0 - \mathbf{w}^{(q)^*}\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2\right)}{N\gamma(1-\gamma\alpha_B\operatorname{tr}(\mathbf{H}^{(q)}))} \cdot \left(\frac{k^*}{N} + N\gamma^2\sum_{i>k^*}(\lambda_i^{(q)})^2\right) \\
&+ \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}_0 - \mathbf{w}^{(q)^*}\|_{(\mathbf{H}_{0:k^*}^{(q)})^{-1}}^2 + \|\mathbf{w}_0 - \mathbf{w}^{(q)^*}\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2.
\end{aligned}$$

*Proof.* Recalling Lemma D.6, we can derive a refined upper bound for $\mathbf{S}_t$ by Lemma D.8:

$$
\begin{aligned}
\mathbf{S}_t \preceq & (\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)}) \circ \mathbf{S}_{t-1} + \gamma^2 \mathcal{M}_B^{(q)} \circ \mathbf{S}_N + \mathbf{B}_0 \\
\preceq & (\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)}) \circ \mathbf{S}_{t-1} + \frac{\gamma \alpha_B \cdot \operatorname{tr}\left(\left[\mathcal{I} - (\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)})^N\right] \circ \mathbf{B}_0\right)}{(1 - \gamma \alpha_B \operatorname{tr}(\mathbf{H}^{(q)}))} \cdot \mathbf{H}^{(q)} + \mathbf{B}_0 \\
= & \sum_{k=0}^{t-1} (\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)})^k \left(\frac{\gamma \alpha_B \cdot \operatorname{tr}\left(\left[\mathcal{I} - (\mathcal{I} - \gamma \tilde{\mathcal{T}}^{(q)})^N\right] \circ \mathbf{B}_0\right)}{(1 - \gamma \alpha_B \operatorname{tr}(\mathbf{H}^{(q)}))} \cdot \mathbf{H}^{(q)} + \mathbf{B}_0\right) \\
= & \sum_{k=0}^{t-1} (\mathbf{I} - \gamma \mathbf{H}^{(q)})^k \left(\frac{\gamma \alpha_B \cdot \operatorname{tr}\left(\mathbf{B}_0 - (\mathbf{I} - \gamma \mathbf{H}^{(q)})^N \mathbf{B}_0 (\mathbf{I} - \gamma \mathbf{H}^{(q)})^N\right)}{(1 - \gamma \alpha_B \operatorname{tr}(\mathbf{H}^{(q)}))} \cdot \mathbf{H}^{(q)} + \mathbf{B}_0\right) (\mathbf{I} - \gamma \mathbf{H}^{(q)})^k.
\end{aligned}
$$

$$(\text{D.18})$$

Before providing our upper bound for the bias error, we denote

$$
\mathbf{B}_{a,b} := \mathbf{B}_a - (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{b-a} \mathbf{B}_a (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{b-a}.
$$

Then by (D.13) and (D.18),

$$
\begin{aligned}
\text{bias} \leq & \frac{1}{\gamma N^2} \langle \mathbf{I} - (\mathbf{I} - \gamma \mathbf{H}^{(q)})^N, \sum_{t=0}^{N-1} \mathbf{B}_t \rangle \\
\leq & \frac{1}{\gamma N^2} \sum_{k=0}^{N-1} \left\langle \mathbf{I} - (\mathbf{I} - \gamma \mathbf{H}^{(q)})^N, (\mathbf{I} - \gamma \mathbf{H}^{(q)})^k \left(\frac{\gamma \alpha_B \cdot \operatorname{tr}(\mathbf{B}_{0,N})}{1 - \gamma \alpha_B \operatorname{tr}(\mathbf{H}^{(q)})} \cdot \mathbf{H}^{(q)} + \mathbf{B}_0\right) (\mathbf{I} - \gamma \mathbf{H}^{(q)})^k \right\rangle \\
= & \frac{1}{\gamma N^2} \sum_{k=0}^{N-1} \left\langle (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{2k} - (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{N+2k}, \left(\frac{\gamma \alpha_B \cdot \operatorname{tr}(\mathbf{B}_{0,N})}{1 - \gamma \alpha_B \operatorname{tr}(\mathbf{H}^{(q)})} \cdot \mathbf{H}^{(q)} + \mathbf{B}_0\right) \right\rangle.
\end{aligned}
$$

Note that

$$
\begin{aligned}
(\mathbf{I} - \gamma \mathbf{H}^{(q)})^{2k} - (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{N+2k} &= \left(\mathbf{I} - \gamma \mathbf{H}^{(q)}\right)^k \left(\left(\mathbf{I} - \gamma \mathbf{H}^{(q)}\right)^k - \left(\mathbf{I} - \gamma \mathbf{H}^{(q)}\right)^{N+k}\right) \\
&\preceq (\mathbf{I} - \gamma \mathbf{H}^{(q)})^k - (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{N+k},
\end{aligned}
$$

we obtain

$$
\text{bias} \leq \frac{1}{\gamma N^2} \sum_{k=0}^{N-1} \left\langle (\mathbf{I} - \gamma \mathbf{H}^{(q)})^k - (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{N+k}, \frac{\gamma \alpha_B \cdot \operatorname{tr}(\mathbf{B}_{0,N})}{1 - \gamma \alpha_B \operatorname{tr}(\mathbf{H}^{(q)})} \cdot \mathbf{H}^{(q)} + \mathbf{B}_0 \right\rangle.
$$

Therefore, it suffices to upper bound the following two terms

$$
\begin{aligned}
I_1 &= \frac{\alpha_B \operatorname{tr}(\mathbf{B}_{0,N})}{N^2 (1 - \gamma \alpha \operatorname{tr}(\mathbf{H}^{(q)}))} \sum_{k=0}^{N-1} \left\langle (\mathbf{I} - \gamma \mathbf{H}^{(q)})^k - (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{N+k}, \mathbf{H}^{(q)} \right\rangle \\
I_2 &= \frac{1}{\gamma N^2} \sum_{k=0}^{N-1} \left\langle (\mathbf{I} - \gamma \mathbf{H}^{(q)})^k - (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{N+k}, \mathbf{B}_0 \right\rangle.
\end{aligned}
$$

Regarding $I_1$, since $\mathbf{H}^{(q)}$ and $\mathbf{I} - \gamma\mathbf{H}^{(q)}$ can be diagonalized simultaneously,

$$
\begin{aligned}
I_1 &= \frac{\alpha_B \operatorname{tr}(\mathbf{B}_{0,N})}{N^2(1 - \gamma\alpha_B \operatorname{tr}(\mathbf{H}^{(q)}))} \sum_{k=0}^{N-1} \sum_i \left[ (1 - \gamma\lambda_i^{(q)})^k - (1 - \gamma\lambda_i^{(q)})^{N+k} \right] \lambda_i^{(q)} \\
&= \frac{\alpha_B \operatorname{tr}(\mathbf{B}_{0,N})}{\gamma N^2(1 - \gamma\alpha_B \operatorname{tr}(\mathbf{H}^{(q)}))} \sum_i \left[ 1 - (1 - \gamma\lambda_i^{(q)})^N \right]^2 \\
&\leq \frac{\alpha_B \operatorname{tr}(\mathbf{B}_{0,N})}{\gamma N^2(1 - \gamma\alpha_B \operatorname{tr}(\mathbf{H}^{(q)}))} \sum_i \min\left\{ 1, \gamma^2 N^2 (\lambda_i^{(q)})^2 \right\} \\
&\leq \frac{\alpha_B \operatorname{tr}(\mathbf{B}_{0,N})}{\gamma(1 - \gamma\alpha_B \operatorname{tr}(\mathbf{H}^{(q)}))} \cdot \left( \frac{k^*}{N^2} + \gamma^2 \sum_{i>k^*} (\lambda_i^{(q)})^2 \right),
\end{aligned}
$$

where $k^* = \max\{k : \lambda_k^{(q)} \geq \frac{1}{N\gamma}\}$. Then we tackle $\operatorname{tr}(\mathbf{B}_{0,N})$.

$$
\begin{aligned}
\operatorname{tr}(\mathbf{B}_{0,N}) &= \operatorname{tr}\left( \mathbf{B}_0 - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^N \mathbf{B}_0 (\mathbf{I} - \gamma\mathbf{H}^{(q)})^N \right) \\
&= \sum_i \left( 1 - (1 - \gamma\lambda_i^{(q)})^{2N} \right) \cdot \left( \langle \mathbf{w}_0 - \mathbf{w}^{(q)*}, \mathbf{v}_i^{(q)} \rangle \right)^2 \\
&\leq 2 \sum_i \min\{1, N\gamma\lambda_i^{(q)}\} \left( \langle \mathbf{w}_0 - \mathbf{w}^{(q)*}, \mathbf{v}_i^{(q)} \rangle \right)^2 \\
&\leq 2 \left( \| \mathbf{w}_0 - \mathbf{w}^{(q)*} \|_{\mathbf{I}_{0:k^*}}^2 + N\gamma \| \mathbf{w}_0 - \mathbf{w}^{(q)*} \|_{\mathbf{H}_{k^*:\infty}}^2 \right).
\end{aligned}
\tag{D.19}
$$

Hence,

$$
I_1 \leq \frac{2\alpha_B \left( \| \mathbf{w}_0 - \mathbf{w}^{(q)*} \|_{\mathbf{I}_{0:k^*}^{(q)}}^2 + N\gamma \| \mathbf{w}_0 - \mathbf{w}^{(q)*} \|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2 \right)}{N\gamma(1 - \gamma\alpha_B \operatorname{tr}(\mathbf{H}^{(q)}))} \cdot \left( \frac{k^*}{N} + N\gamma^2 \sum_{i>k^*} (\lambda_i^{(q)})^2 \right).
$$

Regarding $I_2$, decompose $\mathbf{H}^{(q)} = \mathbf{V}^{(q)} \mathbf{\Lambda}^{(q)} \mathbf{V}^{(q)\top}$, then

$$
I_2 = \frac{1}{\gamma N^2} \sum_{k=0}^{N-1} \langle (\mathbf{I} - \gamma\mathbf{\Lambda}^{(q)})^k - (\mathbf{I} - \gamma\mathbf{\Lambda}^{(q)})^{N+k}, \mathbf{V}^{(q)\top} \mathbf{B}_0 \mathbf{V}^{(q)} \rangle.
$$

Note that $\mathbf{B}_0 = \boldsymbol{\eta}_0 \boldsymbol{\eta}_0^\top$, it can be shown that the diagonal entries of $\mathbf{V}^{(q)\top} \mathbf{B}_0 \mathbf{V}^{(q)}$ are $\omega_1^2, \ldots,$ where $\omega_i = \mathbf{v}_i^{(q)\top} \boldsymbol{\eta}_0 = \mathbf{v}_i^{(q)\top} (\mathbf{w}_0 - \mathbf{w}^{(q)*})$. Hence,

$$
\begin{aligned}
I_2 &= \frac{1}{\gamma N^2} \sum_{k=0}^{N-1} \sum_i \left[ (1 - \gamma\lambda_i^{(q)})^k - (1 - \gamma\lambda_i^{(q)})^{N+k} \right] \omega_i^2 \\
&= \frac{1}{\gamma^2 N^2} \sum_i \frac{\omega_i^2}{\lambda_i^{(q)}} \left[ 1 - (1 - \gamma\lambda_i^{(q)})^N \right]^2 \\
&\leq \frac{1}{\gamma^2 N^2} \sum_i \frac{\omega_i^2}{\lambda_i^{(q)}} \min\left\{ 1, \gamma^2 N^2 (\lambda_i^{(q)})^2 \right\} \\
&\leq \frac{1}{\gamma^2 N^2} \cdot \sum_{i \leq k^*} \frac{\omega_i^2}{\lambda_i^{(q)}} + \sum_{i>k^*} \lambda_i^{(q)} \omega_i^2 \\
&= \frac{1}{\gamma^2 N^2} \cdot \| \mathbf{w}_0 - \mathbf{w}^{(q)*} \|_{(\mathbf{H}_{0:k^*}^{(q)})^{-1}}^2 + \| \mathbf{w}_0 - \mathbf{w}^{(q)*} \|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2.
\end{aligned}
$$

In conclusion, if the stepsize satisfies $\gamma < \frac{1}{\alpha_B \text{tr}(\mathbf{H}^{(q)})}$,

$$\text{bias} \leq I_1 + I_2$$

$$\leq \frac{2\alpha_B \left( \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} \right)}{N\gamma(1 - \gamma\alpha_B \text{tr}(\mathbf{H}^{(q)}))} \cdot \left( \frac{k^*}{N} + N\gamma^2 \sum_{i>k^*} (\lambda_i^{(q)})^2 \right)$$

$$+ \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-1}} + \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}}.$$

$\square$

**Lemma D.11.** (A Bound for bias under Multiplicative Quantization) *Under Assumption 3.1,3.2,3.3,3.4, if the stepsize satisfies* $\gamma < \frac{1}{(1+\tilde{\epsilon})\alpha_B \text{tr}(\mathbf{H}^{(q)})}$, *if there exist* $\epsilon_d, \epsilon_l, \epsilon_p, \epsilon_a$ *and* $\epsilon_o$ *such that for any* $i \in \{d, l, p, a, o\}$, *quantization* $\mathcal{Q}_i$ *is* $\epsilon_i$-*multiplicative, then*

$$\text{bias} \leq \frac{2(1+\tilde{\epsilon})\alpha_B \left( \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} \right)}{N\gamma(1 - (1+\tilde{\epsilon})\gamma\alpha_B \text{tr}(\mathbf{H}^{(q)}))} \cdot \left( \frac{k^*}{N} + N\gamma^2 \sum_{i>k^*} (\lambda_i^{(q)})^2 \right)$$

$$+ \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-1}} + \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}}.$$

*Proof.* Recalling Lemma D.7, we can derive an upper bound for $\mathbf{S}_t$ by Lemma D.9:

$$\mathbf{S}_t \preceq (\mathcal{I} - \gamma\tilde{\mathcal{T}}^{(q)}) \circ \mathbf{S}_{t-1} + (1+\tilde{\epsilon})\gamma^2 \mathcal{M}_B^{(q)} \circ \mathbf{S}_N + \mathbf{B}_0$$

$$\preceq (\mathcal{I} - \gamma\tilde{\mathcal{T}}^{(q)}) \circ \mathbf{S}_{t-1} + \frac{(1+\tilde{\epsilon})\gamma\alpha_B \cdot \text{tr}\left( \left[ \mathcal{I} - (\mathcal{I} - \gamma\tilde{\mathcal{T}}^{(q)})^N \right] \circ \mathbf{B}_0 \right)}{(1 - (1+\tilde{\epsilon})\gamma\alpha_B \text{tr}(\mathbf{H}^{(q)}))} \cdot \mathbf{H}^{(q)} + \mathbf{B}_0$$

$$= \sum_{k=0}^{t-1} (\mathcal{I} - \gamma\tilde{\mathcal{T}}^{(q)})^k \left( \frac{(1+\tilde{\epsilon})\gamma\alpha_B \cdot \text{tr}\left( \left[ \mathcal{I} - (\mathcal{I} - \gamma\tilde{\mathcal{T}}^{(q)})^N \right] \circ \mathbf{B}_0 \right)}{(1 - (1+\tilde{\epsilon})\gamma\alpha_B \text{tr}(\mathbf{H}^{(q)}))} \cdot \mathbf{H}^{(q)} + \mathbf{B}_0 \right)$$

$$\preceq \sum_{k=0}^{t-1} (\mathcal{I} - \gamma\tilde{\mathcal{T}}^{(q)})^k \left( \frac{(1+\tilde{\epsilon})\gamma\alpha_B \cdot \text{tr}\left( \left[ \mathcal{I} - (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)})^N \right] \circ \mathbf{B}_0 \right)}{(1 - (1+\tilde{\epsilon})\gamma\alpha_B \text{tr}(\mathbf{H}^{(q)}))} \cdot \mathbf{H}^{(q)} + \mathbf{B}_0 \right)$$

$$= \sum_{k=0}^{t-1} (\mathbf{I} - \gamma\mathbf{H}^{(q)})^k \left( \frac{(1+\tilde{\epsilon})\gamma\alpha_B \cdot \text{tr}\left( \mathbf{B}_0 - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^N \mathbf{B}_0 (\mathbf{I} - \gamma\mathbf{H}^{(q)})^N \right)}{(1 - (1+\tilde{\epsilon})\gamma\alpha_B \text{tr}(\mathbf{H}^{(q)}))} \cdot \mathbf{H}^{(q)} + \mathbf{B}_0 \right) (\mathbf{I} - \gamma\mathbf{H}^{(q)})^k.$$

Repeat the same computation in the proof of Lemma D.10, we obtain

$$\text{bias} \leq \frac{1}{\gamma N^2} \sum_{k=0}^{N-1} \left\langle (\mathbf{I} - \gamma\mathbf{H}^{(q)})^k - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{N+k}, \frac{(1+\tilde{\epsilon})\gamma\alpha_B \cdot \text{tr}(\mathbf{B}_{0,N})}{1 - (1+\tilde{\epsilon})\gamma\alpha_B \text{tr}(\mathbf{H}^{(q)})} \cdot \mathbf{H}^{(q)} + \mathbf{B}_0 \right\rangle.$$

Therefore, it suffices to upper bound the following two terms

$$I_1 = \frac{(1+\tilde{\epsilon})\alpha_B \text{tr}(\mathbf{B}_{0,N})}{N^2(1 - (1+\tilde{\epsilon})\gamma\alpha \text{tr}(\mathbf{H}^{(q)}))} \sum_{k=0}^{N-1} \left\langle (\mathbf{I} - \gamma\mathbf{H}^{(q)})^k - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{N+k}, \mathbf{H}^{(q)} \right\rangle$$

$$I_2 = \frac{1}{\gamma N^2} \sum_{k=0}^{N-1} \left\langle (\mathbf{I} - \gamma\mathbf{H}^{(q)})^k - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{N+k}, \mathbf{B}_0 \right\rangle.$$

34

Repeating the computation in the proof of Lemma D.10,

$$I_1 \leq \frac{2(1+\tilde{\epsilon})\alpha_B \left( \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} \right)}{N\gamma(1 - (1+\tilde{\epsilon})\gamma\alpha_B \operatorname{tr}(\mathbf{H}^{(q)}))} \cdot \left( \frac{k^*}{N} + N\gamma^2 \sum_{i>k^*}(\lambda_i^{(q)})^2 \right).$$

$$I_2 \leq \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-1}} + \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}}.$$

In conclusion, if the stepsize satisfies $\gamma < \frac{1}{(1+\tilde{\epsilon})\alpha_B \operatorname{tr}(\mathbf{H}^{(q)})}$,

$$\text{bias} \leq \frac{2(1+\tilde{\epsilon})\alpha_B \left( \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} \right)}{N\gamma(1 - (1+\tilde{\epsilon})\gamma\alpha_B \operatorname{tr}(\mathbf{H}^{(q)}))} \cdot \left( \frac{k^*}{N} + N\gamma^2 \sum_{i>k^*}(\lambda_i^{(q)})^2 \right)$$

$$+ \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-1}} + \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}}.$$

$\square$

## D.5  Bounding the Variance Error

Recalling Lemma D.4 and Lemma D.5, the key part of bounding the variance error is to derive an upper bound for $\mathbf{C}_t$ and $\mathbf{C}_t^{(M)}$, where

$$\mathbf{C}_t := (\mathcal{I} - \gamma\mathcal{T}_B^{(q)})\mathbf{C}_{t-1} + \gamma^2\sigma_G^{(q)^2}\mathbf{H}^{(q)}, \quad \mathbf{C}_0 = \mathbf{0}.$$

$$\mathbf{C}_t^{(M)} := (\mathcal{I} - \gamma\mathcal{T}_B^{(q)} + \tilde{\epsilon}\gamma^2\mathcal{M}_B^{(q)})\mathbf{C}_{t-1}^{(M)} + \gamma^2\sigma_M^{(q)^2}\mathbf{H}^{(q)}, \quad \mathbf{C}_0^{(M)} = 0.$$

We first estimate $\mathbf{C}_t$ by converting $\mathcal{T}_B^{(q)}$ to $\widetilde{\mathcal{T}}^{(q)}$.

$$\begin{aligned}
\mathbf{C}_t &= (\mathcal{I} - \gamma\mathcal{T}_B^{(q)}) \circ \mathbf{C}_{t-1} + \gamma^2\sigma_G^{(q)^2}\mathbf{H}^{(q)} \\
&= (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)}) \circ \mathbf{C}_{t-1} + \gamma(\widetilde{\mathcal{T}}^{(q)} - \mathcal{T}_B^{(q)}) \circ \mathbf{C}_{t-1} + \gamma^2\sigma_G^{(q)^2}\mathbf{H}^{(q)} \\
&= (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)}) \circ \mathbf{C}_{t-1} + \gamma^2(\mathcal{M}_B^{(q)} - \widetilde{\mathcal{M}}^{(q)}) \circ \mathbf{C}_{t-1} + \gamma^2\sigma_G^{(q)^2}\mathbf{H}^{(q)} \\
&\preceq (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)}) \circ \mathbf{C}_{t-1} + \gamma^2\mathcal{M}_B^{(q)} \circ \mathbf{C}_{t-1} + \gamma^2\sigma_G^{(q)^2}\mathbf{H}^{(q)}.
\end{aligned} \tag{D.20}$$

Similarly,

$$\begin{aligned}
\mathbf{C}_t^{(M)} &= (\mathcal{I} - \gamma\mathcal{T}_B^{(q)} + \tilde{\epsilon}\gamma^2\mathcal{M}_B^{(q)})\mathbf{C}_{t-1}^{(M)} + \gamma^2\sigma_M^{(q)^2}\mathbf{H}^{(q)} \\
&= (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)}) \circ \mathbf{C}_{t-1} + \gamma(\widetilde{\mathcal{T}}^{(q)} - \mathcal{T}_B^{(q)} + \tilde{\epsilon}\gamma\mathcal{M}_B^{(q)}) \circ \mathbf{C}_{t-1}^{(M)} + \gamma^2\sigma_M^{(q)^2}\mathbf{H}^{(q)} \\
&= (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)}) \circ \mathbf{C}_{t-1}^{(M)} + \gamma^2(\mathcal{M}_B^{(q)} - \widetilde{\mathcal{M}}^{(q)} + \tilde{\epsilon}\gamma\mathcal{M}_B^{(q)}) \circ \mathbf{C}_{t-1}^{(M)} + \gamma^2\sigma_M^{(q)^2}\mathbf{H}^{(q)} \\
&\preceq (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)}) \circ \mathbf{C}_{t-1}^{(M)} + \gamma^2(1+\tilde{\epsilon})\mathcal{M}_B^{(q)} \circ \mathbf{C}_{t-1}^{(M)} + \gamma^2\sigma_M^{(q)^2}\mathbf{H}^{(q)}.
\end{aligned} \tag{D.21}$$

The following two lemmas provide upper bounds for $\mathcal{M}_B^{(q)} \circ \mathbf{C}_t$ and $\mathcal{M}_B^{(q)} \circ \mathbf{C}_t^{(M)}$.

**Lemma D.12.** (A Bound for $\mathcal{M}_B^{(q)} \circ \mathbf{C}_t$) *For $t \geq 1$, under Assumption 3.1,3.2,3.3,3.4, if the stepsize $\gamma \leq \frac{1}{\alpha_B \operatorname{tr}(\mathbf{H}^{(q)})}$, then*

$$\mathcal{M}_B^{(q)} \circ \mathbf{C}_t \preceq \frac{\alpha_B \operatorname{tr}(\mathbf{H}^{(q)})\gamma\sigma_G^{(q)^2}}{1 - \gamma\alpha_B \operatorname{tr}(\mathbf{H}^{(q)})}\mathbf{H}^{(q)}.$$

*Proof.* The main goal is to derive a crude upper bound for $\mathbf{C}_t$. Denote $\mathbf{\Sigma} = \sigma_G^{(q)^2} \mathbf{H}^{(q)}$.

**Step 1: $\mathbf{C}_t$ is increasing.** By definition,

$$
\begin{aligned}
\mathbf{C}_t &= (\mathcal{I} - \gamma\mathcal{T}_B^{(q)}) \circ \mathbf{C}_{t-1} + \gamma^2\mathbf{\Sigma} \\
&= \gamma^2 \sum_{k=0}^{t-1} (\mathcal{I} - \gamma\mathcal{T}_B^{(q)})^k \circ \mathbf{\Sigma} \quad \text{(solving the recursion)} \\
&= \mathbf{C}_{t-1} + \gamma^2 (\mathcal{I} - \gamma\mathcal{T}_B^{(q)})^{t-1} \circ \mathbf{\Sigma} \\
&\succeq \mathbf{C}_{t-1}. \quad \text{(since } \mathcal{I} - \gamma\mathcal{T}_B^{(q)} \text{ is a PSD mapping)}.
\end{aligned}
$$

**Step 2: $\mathbf{C}_\infty$ exists.** It suffices to show that $\operatorname{tr}(\mathbf{C}_t)$ is uniformly upper bounded. To be specific, for any $t \geq 1$,

$$
\mathbf{C}_t = \gamma^2 \sum_{k=0}^{t-1} (\mathcal{I} - \gamma\mathcal{T}_B^{(q)})^k \circ \mathbf{\Sigma} \preceq \gamma^2 \sum_{t=0}^{\infty} (\mathcal{I} - \gamma\mathcal{T}_B^{(q)})^t \circ \mathbf{\Sigma}.
$$

Then

$$
\operatorname{tr}(\mathbf{C}_t) \leq \gamma^2 \sum_{t=0}^{\infty} \operatorname{tr}\left((\mathcal{I} - \gamma\mathcal{T}_B^{(q)})^t \circ \mathbf{\Sigma}\right) := \gamma^2 \sum_{t=0}^{\infty} \operatorname{tr}(\mathbf{E}_t) \leq \frac{\gamma \operatorname{tr}(\mathbf{\Sigma})}{\lambda_d^{(q)}} < \infty,
$$

where the second inequality holds by the iteration:

$$
\begin{aligned}
\operatorname{tr}(\mathbf{E}_t) &= \operatorname{tr}(\mathbf{E}_{t-1}) - 2\gamma\operatorname{tr}(\mathbf{H}^{(q)}\mathbf{E}_{t-1}) + \gamma^2\operatorname{tr}\left(\mathbf{E}_{t-1}\mathbb{E}\left[\frac{1}{B^2}\mathbf{X}^{(q)\top}\mathbf{X}^{(q)}\mathbf{X}^{(q)\top}\mathbf{X}^{(q)}\right]\right) \\
&\leq \operatorname{tr}(\mathbf{E}_{t-1}) - (2\gamma - \gamma^2\alpha_B\operatorname{tr}(\mathbf{H}^{(q)}))\operatorname{tr}(\mathbf{H}^{(q)}\mathbf{E}_{t-1}) \\
&\leq \operatorname{tr}\left((\mathbf{I} - \gamma\mathbf{H}^{(q)})\mathbf{E}_{t-1}\right) \\
&\leq (1 - \gamma\lambda_d^{(q)})\operatorname{tr}(\mathbf{E}_{t-1}),
\end{aligned}
$$

where the first inequality holds by Assumption 3.3 and the second inequality holds if $\gamma \leq \frac{1}{\alpha_B\operatorname{tr}(\mathbf{H}^{(q)})}$.

**Step 3: upper bound $\mathbf{C}_\infty$.** By the update rule for $\mathbf{C}_t$,

$$
\mathbf{C}_\infty = (\mathcal{I} - \gamma\mathcal{T}_B^{(q)}) \circ \mathbf{C}_\infty + \gamma^2\mathbf{\Sigma},
$$

which immediately implies

$$
\mathbf{C}_\infty = \gamma\mathcal{T}_B^{(q)^{-1}} \circ \mathbf{\Sigma}. \tag{D.22}
$$

We provide the upper bound by applying $\widetilde{\mathcal{T}}^{(q)}$.

$$
\begin{aligned}
\widetilde{\mathcal{T}}^{(q)} \circ \mathbf{C}_\infty &= \mathcal{T}_B^{(q)} \circ \mathbf{C}_\infty + \gamma\mathcal{M}_B^{(q)} \circ \mathbf{C}_\infty - \gamma\widetilde{\mathcal{M}}^{(q)} \circ \mathbf{C}_\infty \\
&= \gamma\mathbf{\Sigma} + \gamma\mathcal{M}_B^{(q)} \circ \mathbf{C}_\infty - \gamma\widetilde{\mathcal{M}}^{(q)} \circ \mathbf{C}_\infty \\
&\preceq \gamma\mathbf{\Sigma} + \gamma\mathcal{M}_B^{(q)} \circ \mathbf{C}_\infty,
\end{aligned}
$$

where the first equality holds by the definition of $\mathcal{T}_B^{(q)}$ and $\widetilde{\mathcal{T}}^{(q)}$ and the second equality holds by (D.22). Hence,

$$
\widetilde{\mathcal{T}}^{(q)} \circ \mathbf{C}_\infty \preceq \gamma\sigma_G^{(q)^2}\mathbf{H}^{(q)} + \gamma\mathcal{M}_B^{(q)} \circ \mathbf{C}_\infty.
$$

Therefore, by applying $(\widetilde{\mathcal{T}}^{(q)})^{-1}$ we have

$$
\begin{aligned}
\mathbf{C}_\infty &\preceq \gamma\sigma_G^{(q)^2} \cdot (\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{H}^{(q)} + \gamma(\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathcal{M}_B^{(q)} \circ \mathbf{C}_\infty \\
&\preceq \gamma\sigma_G^{(q)^2} \cdot \sum_{t=0}^{\infty} \left(\gamma(\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathcal{M}_B^{(q)}\right)^t \circ (\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{H}^{(q)}. \quad \text{(solving the recursion)}
\end{aligned} \tag{D.23}
$$

We first deal with $(\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{H}^{(q)}$.

$$
\begin{aligned}
(\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{H}^{(q)} &= \gamma \sum_{t=0}^{\infty} (\mathcal{I} - \gamma \widetilde{\mathcal{T}}^{(q)})^t \circ \mathbf{H}^{(q)} \\
&= \gamma \sum_{t=0}^{\infty} (\mathbf{I} - \gamma \mathbf{H}^{(q)})^t \mathbf{H}^{(q)} (\mathbf{I} - \gamma \mathbf{H}^{(q)})^t \\
&\preceq \gamma \sum_{t=0}^{\infty} (\mathbf{I} - \gamma \mathbf{H}^{(q)})^t \mathbf{H}^{(q)} \\
&= \mathbf{I},
\end{aligned}
\tag{D.24}
$$

where the second equality uses the definition of $\widetilde{\mathcal{T}}^{(q)}$. Hence, by (D.23) and (D.24),

$$
\begin{aligned}
\mathbf{C}_{\infty} &\preceq \gamma \sigma_G^{(q)\,2} \cdot \sum_{t=0}^{\infty} (\gamma (\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathcal{M}_B^{(q)})^t \circ \mathbf{I} \\
&= \gamma \sigma_G^{(q)\,2} \cdot \sum_{t=0}^{\infty} (\gamma (\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathcal{M}_B^{(q)})^{t-1} \gamma (\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathcal{M}_B^{(q)} \circ \mathbf{I} \\
&\preceq \gamma \sigma_G^{(q)\,2} \cdot \sum_{t=0}^{\infty} (\gamma (\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathcal{M}_B^{(q)})^{t-1} \circ \gamma \alpha_B \mathrm{tr}(\mathbf{H}^{(q)}) \mathbf{I} \\
&\preceq \gamma \sigma_G^{(q)\,2} \cdot \sum_{t=0}^{\infty} \left( \gamma \alpha_B \mathrm{tr}(\mathbf{H}^{(q)}) \right)^t \mathbf{I} \\
&= \frac{\gamma \sigma_G^{(q)\,2}}{1 - \gamma \alpha_B \mathrm{tr}(\mathbf{H}^{(q)})} \mathbf{I},
\end{aligned}
$$

where the second inequality holds by the fact that $\mathcal{M}_B^{(q)} \circ \mathbf{I} \preceq \alpha_B \mathrm{tr}(\mathbf{H}^{(q)}) \mathbf{H}^{(q)}$.

Here we complete deriving a crude upper bound for $\mathbf{C}_t$:

$$
\mathbf{C}_t \preceq \mathbf{C}_{\infty} \preceq \frac{\gamma \sigma_G^{(q)\,2}}{1 - \gamma \alpha_B \mathrm{tr}(\mathbf{H}^{(q)})} \mathbf{I}.
$$

Then by $\mathcal{M}_B^{(q)} \circ \mathbf{I} \preceq \alpha_B \mathrm{tr}(\mathbf{H}^{(q)}) \mathbf{H}^{(q)}$ again,

$$
\mathcal{M}_B^{(q)} \circ \mathbf{C}_t \preceq \frac{\alpha_B \mathrm{tr}(\mathbf{H}^{(q)}) \gamma \sigma_G^{(q)\,2}}{1 - \gamma \alpha_B \mathrm{tr}(\mathbf{H}^{(q)})} \mathbf{H}^{(q)}.
$$

$\square$

**Lemma D.13.** (A Bound for $\mathcal{M}_B^{(q)} \circ \mathbf{C}_t^{(M)}$) *For $t \geq 1$, under Assumption 3.1,3.2,3.3,3.4, if the stepsize $\gamma \leq \frac{1}{(1+\tilde{\epsilon})\alpha_B \mathrm{tr}(\mathbf{H}^{(q)})}$, then*

$$
\mathcal{M}_B^{(q)} \circ \mathbf{C}_t^{(M)} \preceq \frac{\alpha_B \mathrm{tr}(\mathbf{H}^{(q)}) \gamma \sigma_M^{(q)\,2}}{1 - (1 + \tilde{\epsilon}) \gamma \alpha_B \mathrm{tr}(\mathbf{H}^{(q)})} \mathbf{H}^{(q)}.
$$

*Proof.* The proof idea is similar to the proof of Lemma D.12 while the main goal is to derive a crude upper bound for $\mathbf{C}_t^{(M)}$. We deduce from the proof of Lemma D.12 that [9]

$$
\mathbf{C}_t^{(M)} \preceq \mathbf{C}_{\infty}^{(M)} = \gamma (\mathcal{T}_B^{(q)} - \tilde{\epsilon} \gamma \mathcal{M}_B^{(q)})^{-1} \circ \sigma_M^{(q)\,2} \mathbf{H}^{(q)}.
\tag{D.25}
$$

---

[9] $(\mathcal{T}_B^{(q)} - \tilde{\epsilon} \gamma \mathcal{M}_B^{(q)})^{-1}$ exists under the condition that $\gamma < \frac{1}{(1+\tilde{\epsilon})\alpha_B \mathrm{tr}(\mathbf{H}^{(q)})}$, which can be directly deduced by Lemma B.1 in Zou et al. (2023). We omit the proof here for simplicity.

We provide the upper bound for $\mathbf{C}_\infty^{(M)}$ by applying $\widetilde{\mathcal{T}}^{(q)}$.

$$\widetilde{\mathcal{T}}^{(q)} \circ \mathbf{C}_\infty^{(M)} = (\mathcal{T}_B^{(q)} - \tilde{\epsilon}\gamma\mathcal{M}_B^{(q)}) \circ \mathbf{C}_\infty^{(M)} + (1+\tilde{\epsilon})\gamma\mathcal{M}_B^{(q)} \circ \mathbf{C}_\infty^{(M)} - \gamma\widetilde{\mathcal{M}}^{(q)} \circ \mathbf{C}_\infty^{(M)}$$

$$= \gamma\sigma_M^{(q)2}\mathbf{H}^{(q)} + (1+\tilde{\epsilon})\gamma\mathcal{M}_B^{(q)} \circ \mathbf{C}_\infty^{(M)} - \gamma\widetilde{\mathcal{M}}^{(q)} \circ \mathbf{C}_\infty^{(M)}$$

$$\preceq \gamma\sigma_M^{(q)2}\mathbf{H}^{(q)} + (1+\tilde{\epsilon})\gamma\mathcal{M}_B^{(q)} \circ \mathbf{C}_\infty^{(M)},$$

where the first equality holds by the definition of $\widetilde{\mathcal{T}}^{(q)}$ and the second equality holds by the definition of $\mathbf{C}_\infty^{(M)}$ (D.25). Therefore, applying $(\widetilde{\mathcal{T}}^{(q)})^{-1}$ we have

$$\mathbf{C}_\infty^{(M)} \preceq \gamma\sigma_M^{(q)2} \cdot (\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{H}^{(q)} + (1+\tilde{\epsilon})\gamma(\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathcal{M}_B^{(q)} \circ \mathbf{C}_\infty^{(M)}$$

$$\preceq \gamma\sigma_M^{(q)2} \cdot \sum_{t=0}^\infty ((1+\tilde{\epsilon})\gamma(\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathcal{M}_B^{(q)})^t \circ (\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{H}^{(q)}. \quad \text{(solving the recursion)} \tag{D.26}$$

By the computation (D.24) in the proof for Lemma D.12,

$$(\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathbf{H}^{(q)} \preceq \mathbf{I}. \tag{D.27}$$

Hence, by (D.26) and (D.27),

$$\mathbf{C}_\infty^{(M)} \preceq \gamma\sigma_M^{(q)2} \cdot \sum_{t=0}^\infty ((1+\tilde{\epsilon})\gamma(\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathcal{M}_B^{(q)})^t \circ \mathbf{I}$$

$$= \gamma\sigma_M^{(q)2} \cdot \sum_{t=0}^\infty ((1+\tilde{\epsilon})\gamma(\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathcal{M}_B^{(q)})^{t-1}(1+\tilde{\epsilon})\gamma(\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathcal{M}_B^{(q)} \circ \mathbf{I}$$

$$\preceq \gamma\sigma_M^{(q)2} \cdot \sum_{t=0}^\infty ((1+\tilde{\epsilon})\gamma(\widetilde{\mathcal{T}}^{(q)})^{-1} \circ \mathcal{M}_B^{(q)})^{t-1} \circ (1+\tilde{\epsilon})\gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})\mathbf{I}$$

$$\preceq \gamma\sigma_M^{(q)2} \cdot \sum_{t=0}^\infty \left((1+\tilde{\epsilon})\gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})\right)^t \mathbf{I}$$

$$= \frac{\gamma\sigma_M^{(q)2}}{1 - (1+\tilde{\epsilon})\gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})}\mathbf{I},$$

where the second inequality holds by the fact that $\mathcal{M}_B^{(q)} \circ \mathbf{I} \preceq \alpha_B\,\mathrm{tr}(\mathbf{H}^{(q)})\mathbf{H}^{(q)}$.

Therefore, we complete the proof by

$$\mathcal{M}_B^{(q)} \circ \mathbf{C}_t^{(M)} \preceq \frac{\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})\gamma\sigma_M^{(q)2}}{1 - (1+\tilde{\epsilon})\gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})}\mathbf{H}^{(q)}.$$

$\square$

By (D.20), (D.21), Lemma D.12 and Lemma D.13, we can provide a refined bound for $\mathbf{C}_t$ and $\mathbf{C}_t^{(M)}$. Then we are ready to bound the variance error.

**Lemma D.14.** (A Bound for variance) *Under Assumption 3.1,3.2,3.3,3.4, if the stepsize satisfies $\gamma < \frac{1}{\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})}$, then*

$$\text{variance} \leq \frac{\sigma_G^{(q)2}}{1 - \gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})}\left(\frac{k^*}{N} + N\gamma^2 \cdot \sum_{i>k^*}(\lambda_i^{(q)})^2\right).$$

38

*Proof.* We first provide a refined upper bound for $\mathbf{C}_t$. By (D.20),

$$\mathbf{C}_t \preceq (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)}) \circ \mathbf{C}_{t-1} + \gamma^2\mathcal{M}_B^{(q)} \circ \mathbf{C}_{t-1} + \gamma^2\sigma_G^{(q)2}\mathbf{H}^{(q)}$$

$$\preceq (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)}) \circ \mathbf{C}_{t-1} + \frac{\gamma^2\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})\gamma\sigma_G^{(q)2}}{1 - \gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})}\mathbf{H}^{(q)} + \gamma^2\sigma_G^{(q)2}\mathbf{H}^{(q)}$$

$$= (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)}) \circ \mathbf{C}_{t-1} + \frac{\gamma^2\sigma_G^{(q)2}}{1 - \gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})}\mathbf{H}^{(q)}$$

$$\preceq \frac{\gamma^2\sigma_G^{(q)2}}{1 - \gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})} \cdot \sum_{k=0}^{t-1}(\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)})^k \circ \mathbf{H}^{(q)} \quad \text{(solving the recursion)} \tag{D.28}$$

$$= \frac{\gamma^2\sigma_G^{(q)2}}{1 - \gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})} \cdot \sum_{k=0}^{t-1}(\mathbf{I} - \gamma\mathbf{H}^{(q)})^k\mathbf{H}^{(q)}(\mathbf{I} - \gamma\mathbf{H}^{(q)})^k$$

$$\preceq \frac{\gamma^2\sigma_G^{(q)2}}{1 - \gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})} \cdot \sum_{k=0}^{t-1}(\mathbf{I} - \gamma\mathbf{H}^{(q)})^k\mathbf{H}^{(q)}$$

$$= \frac{\gamma\sigma_G^{(q)2}}{1 - \gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})} \cdot \left(\mathbf{I} - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^t\right),$$

where the second inequality holds by Lemma D.12 and the second equality holds by the definition of $\widetilde{\mathcal{T}}^{(q)}$.

After providing a refined bound for $\mathbf{C}_t$, we are ready to bound the variance. By Lemma D.4,

$$\mathrm{variance} = \frac{1}{N^2} \cdot \sum_{t=0}^{N-1}\sum_{k=t}^{N-1}\left\langle (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k-t}\mathbf{H}^{(q)}, \mathbf{C}_t \right\rangle$$

$$= \frac{1}{\gamma N^2}\sum_{t=0}^{N-1}\left\langle \mathbf{I} - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{N-t}, \mathbf{C}_t \right\rangle$$

$$\leq \frac{1}{\gamma^2 N^2}\frac{\gamma^2\sigma_G^{(q)2}}{1 - \gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})}\sum_{t=0}^{N-1}\left\langle \mathbf{I} - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{N-t}, \mathbf{I} - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^t \right\rangle$$

$$= \frac{1}{\gamma^2 N^2}\frac{\gamma^2\sigma_G^{(q)2}}{1 - \gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})}\sum_i\sum_{t=0}^{N-1}\left[1 - (1 - \gamma\lambda_i^{(q)})^{N-t}\right]\left[1 - (1 - \gamma\lambda_i^{(q)})^t\right]$$

$$\leq \frac{1}{\gamma^2 N^2}\frac{\gamma^2\sigma_G^{(q)2}}{1 - \gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})}\sum_i\sum_{t=0}^{N-1}\left[1 - (1 - \gamma\lambda_i^{(q)})^N\right]\left[1 - (1 - \gamma\lambda_i^{(q)})^N\right]$$

$$= \frac{1}{\gamma^2 N}\frac{\gamma^2\sigma_G^{(q)2}}{1 - \gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})}\sum_i\left[1 - (1 - \gamma\lambda_i^{(q)})^N\right]^2$$

$$\leq \frac{1}{\gamma^2 N}\frac{\gamma^2\sigma_G^{(q)2}}{1 - \gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})}\sum_i\min\left\{1, \gamma^2 N^2(\lambda_i^{(q)})^2\right\}$$

$$\leq \frac{1}{\gamma^2 N}\frac{\gamma^2\sigma_G^{(q)2}}{1 - \gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})}\left(k^* + N^2\gamma^2 \cdot \sum_{i>k^*}(\lambda_i^{(q)})^2\right)$$

$$= \frac{\sigma_G^{(q)2}}{1 - \gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})}\left(\frac{k^*}{N} + N\gamma^2 \cdot \sum_{i>k^*}(\lambda_i^{(q)})^2\right),$$

39

where the first inequality holds by (D.28) and the last inequality holds by the definition of $k^* = \max\left\{k : \lambda_k^{(q)} \geq \frac{1}{N\gamma}\right\}$. This immediately completes the proof.

$\square$

**Lemma D.15.** (A Bound for variance under Multiplicative Quantization) *Under Assumption 3.1,3.2,3.3,3.4, if the stepsize satisfies $\gamma < \frac{1}{(1+\tilde{\epsilon})\alpha_B \mathrm{tr}(\mathbf{H}^{(q)})}$, if there exist $\epsilon_d, \epsilon_l, \epsilon_p, \epsilon_a$ and $\epsilon_o$ such that for any $i \in \{d, l, p, a, o\}$, quantization $\mathcal{Q}_i$ is $\epsilon_i$-multiplicative, then*

$$\mathrm{variance} \leq \frac{\sigma_M^{(q)2}}{1 - (1+\tilde{\epsilon})\gamma\alpha_B \mathrm{tr}(\mathbf{H}^{(q)})}\left(\frac{k^*}{N} + N\gamma^2 \cdot \sum_{i>k^*}(\lambda_i^{(q)})^2\right).$$

*Proof.* Applying (D.21), and repeating the computation in the proof of Lemma D.14,

$$\mathbf{C}_t^{(M)} \preceq (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)}) \circ \mathbf{C}_{t-1}^{(M)} + \gamma^2(1+\tilde{\epsilon})\mathcal{M}_B^{(q)} \circ \mathbf{C}_{t-1}^{(M)} + \gamma^2\sigma_M^{(q)2}\mathbf{H}^{(q)}$$

$$\preceq (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)}) \circ \mathbf{C}_{t-1}^{(M)} + \gamma^2(1+\tilde{\epsilon})\frac{\alpha_B \mathrm{tr}(\mathbf{H}^{(q)})\gamma\sigma_M^{(q)2}}{1 - (1+\tilde{\epsilon})\gamma\alpha_B \mathrm{tr}(\mathbf{H}^{(q)})}\mathbf{H}^{(q)} + \gamma^2\sigma_M^{(q)2}\mathbf{H}^{(q)}$$

$$= (\mathcal{I} - \gamma\widetilde{\mathcal{T}}^{(q)}) \circ \mathbf{C}_{t-1}^{(M)} + \frac{\gamma^2\sigma_M^{(q)2}\mathbf{H}^{(q)}}{1 - (1+\tilde{\epsilon})\gamma\alpha_B \mathrm{tr}(\mathbf{H}^{(q)})}\mathbf{H}^{(q)}$$

$$\preceq \frac{\gamma\sigma_M^{(q)2}}{1 - (1+\tilde{\epsilon})\gamma\alpha_B \mathrm{tr}(\mathbf{H}^{(q)})} \cdot \left(\mathbf{I} - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^t\right),$$

where the second inequality holds by Lemma D.13 and the last inequality repeats the proof in (D.28).

Therefore, repeating the procedure in the proof for Lemma D.14, we directly deduce that

$$\mathrm{variance} \leq \frac{\sigma_M^{(q)2}}{1 - (1+\tilde{\epsilon})\gamma\alpha_B \mathrm{tr}(\mathbf{H}^{(q)})}\left(\frac{k^*}{N} + N\gamma^2 \cdot \sum_{i>k^*}(\lambda_i^{(q)})^2\right),$$

which immediately completes the proof.

$\square$

**Lemma D.16.** (A Bound for $R_2$) *Under Assumption 3.1,3.2,3.3,3.4, if the stepsize satisfies $\gamma < \frac{1}{\alpha_B \mathrm{tr}(\mathbf{H}^{(q)})}$, then*

$$R_2 \leq \frac{2\alpha_B\left(\|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|_{\mathbf{I}_{0:k^*}^{(q)}}^2 + N\gamma\|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2\right)}{N\gamma(1 - \gamma\alpha_B \mathrm{tr}(\mathbf{H}^{(q)}))} \cdot \left(\frac{k^*}{N} + N\gamma^2\sum_{i>k^*}(\lambda_i^{(q)})^2\right)$$

$$+ \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|_{(\mathbf{H}_{0:k^*}^{(q)})^{-1}}^2 + \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2$$

$$+ \frac{\sigma_G^{(q)2}}{1 - \gamma\alpha_B \mathrm{tr}(\mathbf{H}^{(q)})}\left(\frac{k^*}{N} + N\gamma^2 \cdot \sum_{i>k^*}(\lambda_i^{(q)})^2\right).$$

*Proof.* The proof is immediately completed by Lemma D.10 and Lemma D.14.

$\square$

**Lemma D.17.** (A Bound for $R_2$ under Multiplicative Quantization) *Under Assumption 3.1,3.2,3.3,3.4, if the stepsize satisfies $\gamma < \frac{1}{(1+\tilde{\epsilon})\alpha_B \mathrm{tr}(\mathbf{H}^{(q)})}$, if there exist $\epsilon_d, \epsilon_l, \epsilon_p, \epsilon_a$ and $\epsilon_o$ such that for any*

$i \in \{d, l, p, a, o\}$, *quantization $\mathcal{Q}_i$ is $\epsilon_i$-multiplicative, then*

$$
\begin{aligned}
R_2 \leq & \frac{2(1+\tilde{\epsilon})\alpha_B \left( \|\mathbf{w}_0 - \mathbf{w}^{(q)^*}\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}_0 - \mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} \right)}{N\gamma(1 - (1+\tilde{\epsilon})\gamma\alpha_B\,\mathrm{tr}(\mathbf{H}^{(q)}))} \cdot \left( \frac{k^*}{N} + N\gamma^2 \sum_{i > k^*} (\lambda_i^{(q)})^2 \right) \\
& + \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}_0 - \mathbf{w}^{(q)^*}\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-1}} + \|\mathbf{w}_0 - \mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} \\
& + \frac{\sigma_M^{(q)^2}}{1 - (1+\tilde{\epsilon})\gamma\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})} \left( \frac{k^*}{N} + N\gamma^2 \cdot \sum_{i > k^*} (\lambda_i^{(q)})^2 \right).
\end{aligned}
$$

*Proof.* The proof is immediately completed by Lemma D.11 and Lemma D.15. □

# E   Analysis of $R_1$

Recalling the excess risk decomposition Lemma A.2, we analyze $R_1$ in this section. Utilizing the definition of the covariance matrix of data quantization error $\mathbf{D}$, we have

$$
\begin{aligned}
R_1 = & \frac{1}{2}\mathbb{E}\left[ (y - \mathcal{Q}_l(y))^2 \right] + \frac{1}{2}\mathbb{E}\left[ \langle \overline{\mathbf{w}}_N, \mathcal{Q}_d(\mathbf{x}) - \mathbf{x} \rangle^2 \right] \\
= & \frac{1}{2}\mathbb{E}\left[ (y - \mathcal{Q}_l(y))^2 \right] + \frac{1}{2}\mathbb{E}\left[ \overline{\mathbf{w}}_N^\top \mathbf{D}\overline{\mathbf{w}}_N \right] \\
= & \frac{1}{2}\mathbb{E}\left[ (y - \mathcal{Q}_l(y))^2 \right] + \frac{1}{2}\mathbb{E}\left[ (\overline{\boldsymbol{\eta}}_N + \mathbf{w}^{(q)^*})^\top \mathbf{D}(\overline{\boldsymbol{\eta}}_N + \mathbf{w}^{(q)^*}) \right] \\
\leq & \frac{1}{2}\mathbb{E}\left[ (y - \mathcal{Q}_l(y))^2 \right] + \mathbb{E}\left[ \overline{\boldsymbol{\eta}}_N^\top \mathbf{D}\overline{\boldsymbol{\eta}}_N \right] + \mathbb{E}\left[ \mathbf{w}^{(q)^{*\top}} \mathbf{D}\mathbf{w}^{(q)^*} \right] \\
= & \frac{1}{2}\mathbb{E}\left[ (y - \mathcal{Q}_l(y))^2 \right] + \langle \mathbf{D}, \mathbb{E}\left[ \overline{\boldsymbol{\eta}}_N \otimes \overline{\boldsymbol{\eta}}_N \right] \rangle + \mathbb{E}\left[ \mathbf{w}^{(q)^{*\top}} \mathbf{D}\mathbf{w}^{(q)^*} \right].
\end{aligned}
$$

The main goal in the following part is to derive bounds for $\langle \mathbf{D}, \mathbb{E}\left[ \overline{\boldsymbol{\eta}}_N \otimes \overline{\boldsymbol{\eta}}_N \right] \rangle$.

## E.1   Multiplicative Data Quantization

Under multiplicative quantization, $\mathbf{D}$ is proportional to $\mathbf{H}$. Hence, the term $\langle \mathbf{D}, \mathbb{E}\left[ \overline{\boldsymbol{\eta}}_N \otimes \overline{\boldsymbol{\eta}}_N \right] \rangle$ can be directly bounded by our analysis for $R_2$, i.e., Lemma D.17.

**Lemma E.1.** (A Bound for $R_1$ under Multiplicative Quantization) *Under Assumption 3.1,3.2,3.3,3.4, if the stepsize satisfies $\gamma < \frac{1}{(1+\tilde{\epsilon})\alpha_B\mathrm{tr}(\mathbf{H}^{(q)})}$, if there exist $\epsilon_d, \epsilon_l, \epsilon_p, \epsilon_a$ and $\epsilon_o$ such that for any $i \in \{d, l, p, a, o\}$, quantization $\mathcal{Q}_i$ is $\epsilon_i$-multiplicative, then*

$$
R_1 \leq \frac{\mathbb{E}\left[ (\epsilon^{(l)})^2 \right]}{2} + \frac{2\epsilon_d}{1 + \epsilon_d} R_2 + \frac{\epsilon_d}{(1 + \epsilon_d)^2} \|\mathbf{w}^*\|_{\mathbf{H}}^2,
$$

*where $R_2$ is bounded by Lemma D.17.*

*Proof.* The proof is easily completed by

$$
\mathbf{D} = \frac{\epsilon_d}{1 + \epsilon_d}\mathbf{H}^{(q)}, \quad \text{and} \quad \mathbb{E}\left[ \mathbf{w}^{(q)^{*\top}} \mathbf{D}\mathbf{w}^{(q)^*} \right] = \frac{\epsilon_d}{(1 + \epsilon_d)^2}\|\mathbf{w}^*\|_{\mathbf{H}}^2.
$$

□

## E.2 General Data Quantization

Under general data quantization, as $\mathbf{D}$ is not proportional to $\mathbf{H}$, we need to extract $\mathbf{D}$ from $\langle \mathbf{D}, \mathbb{E}\left[\overline{\boldsymbol{\eta}}_N \otimes \overline{\boldsymbol{\eta}}_N\right]\rangle$, which is summarized as the following lemma.

**Lemma E.2.** (A Bound for $R_1$) *Under Assumption 3.1,3.2,3.3,3.4, if the stepsize satisfies $\gamma < \frac{1}{\alpha_B \mathrm{tr}(\mathbf{H}^{(q)})}$ then*

$$R_1 \leq \frac{\mathbb{E}\left[\left(\epsilon^{(l)}\right)^2\right]}{2} + 2\|\mathbf{D}\|(\mathrm{B} + \mathrm{V}) + \|\mathbf{w}^*\|^2_{\mathbf{H}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{D}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{H}},$$

*where*

$$\mathrm{B} \leq \frac{2\alpha_B \left( \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} \right)}{N\gamma(1 - \gamma\alpha_B \,\mathrm{tr}(\mathbf{H}^{(q)}))} \left( \sum_{i \leq k^*} \frac{1}{N\lambda_i^{(q)}} + N\gamma^2 \sum_{i > k^*} \lambda_i^{(q)} \right)$$

$$+ \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-2}} + \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{I}^{(q)}_{k^*:\infty}},$$

$$\mathrm{V} \leq \frac{\sigma_G^{(q)2}}{1 - \gamma\alpha_B \mathrm{tr}(\mathbf{H}^{(q)})} \left( \sum_{i \leq k^*} \frac{1}{N\lambda_i^{(q)}} + \sum_{i > k^*} \gamma^2 N\lambda_i^{(q)} \right).$$

*Proof.* It is easy to verify that

$$\mathbb{E}\left[\mathbf{w}^{(q)*\top}\mathbf{D}\mathbf{w}^{(q)*}\right] = \|\mathbf{w}^*\|^2_{\mathbf{H}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{D}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{H}}.$$

Hence, it suffices to prove

$$\langle \mathbf{D}, \mathbb{E}\left[\overline{\boldsymbol{\eta}}_N \otimes \overline{\boldsymbol{\eta}}_N\right]\rangle \leq \|\mathbf{D}\|\mathrm{tr}\left(\mathbb{E}\left[\overline{\boldsymbol{\eta}}_N \otimes \overline{\boldsymbol{\eta}}_N\right]\right) \leq 2\|\mathbf{D}\|(\mathrm{B} + \mathrm{V}).$$

In the following part, we analyze $\mathbb{E}\left[\overline{\boldsymbol{\eta}}_N \otimes \overline{\boldsymbol{\eta}}_N\right]$ utilizing the technique in the analysis of $R_2$. From the computation in the analysis of $R_2$ (D.3), we have

$$\mathbb{E}\left[\overline{\boldsymbol{\eta}}_N \otimes \overline{\boldsymbol{\eta}}_N\right] \preceq \frac{1}{N^2} \cdot \sum_{t=0}^{N-1}\sum_{k=t}^{N-1} \left( (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k-t}\mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t] + \mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t](\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k-t} \right).$$

Hence,

$$\frac{1}{2}\mathrm{tr}\left(\mathbb{E}\left[\overline{\boldsymbol{\eta}}_N \otimes \overline{\boldsymbol{\eta}}_N\right]\right) \leq \frac{1}{2N^2} \cdot \sum_{t=0}^{N-1}\sum_{k=t}^{N-1} \mathrm{tr}\left( (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k-t}\mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t] + \mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t](\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k-t} \right)$$

$$= \frac{1}{N^2} \cdot \sum_{t=0}^{N-1}\sum_{k=t}^{N-1} \mathrm{tr}\left( (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k-t}\mathbb{E}[\boldsymbol{\eta}_t \otimes \boldsymbol{\eta}_t] \right)$$

$$\leq \underbrace{\frac{1}{N^2} \cdot \sum_{t=0}^{N-1}\sum_{k=t}^{N-1} \mathrm{tr}\left( (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k-t}\mathbf{B}_t \right)}_{\mathrm{B}} + \underbrace{\frac{1}{N^2} \cdot \sum_{t=0}^{N-1}\sum_{k=t}^{N-1} \mathrm{tr}\left( (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k-t}\mathbf{C}_t \right)}_{\mathrm{V}},$$

where the last inequality holds by the bias-variance decomposition Lemma D.4 with

$$\mathbf{B}_t := (\mathcal{I} - \gamma\mathcal{T}_B^{(q)})^t \circ \mathbf{B}_0, \quad \mathbf{B}_0 = \mathbb{E}\left[\boldsymbol{\eta}_0 \otimes \boldsymbol{\eta}_0\right],$$

and
$$\mathbf{C}_t := (\mathcal{I} - \gamma\mathcal{T}_B^{(q)})\mathbf{C}_{t-1} + \gamma^2\sigma_G^{(q)^2}\mathbf{H}^{(q)}, \quad \mathbf{C}_0 = \mathbf{0}.$$

We then focus on bounding B and V using the same idea in the analysis of $R_2$. Regarding B,

$$\begin{aligned}
\text{B} &= \frac{1}{N^2} \cdot \sum_{t=0}^{N-1}\sum_{k=t}^{N-1} \text{tr}\left((\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k-t}\mathbf{B}_t\right) \\
&= \frac{1}{N^2} \cdot \sum_{t=0}^{N-1}\sum_{k=t}^{N-1} \left\langle(\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k-t}, \mathbf{B}_t\right\rangle \\
&= \frac{1}{\gamma N^2} \cdot \sum_{t=0}^{N-1} \left\langle\left(\mathbf{I} - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{N-t}\right)(\mathbf{H}^{(q)})^{-1}, \mathbf{B}_t\right\rangle \\
&\leq \frac{1}{\gamma N^2} \cdot \left\langle\left(\mathbf{I} - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{N}\right)(\mathbf{H}^{(q)})^{-1}, \sum_{t=0}^{N-1}\mathbf{B}_t\right\rangle \\
&\leq \frac{1}{\gamma N^2} \sum_{k=0}^{N-1} \left\langle\left(\mathbf{I} - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{N}\right)(\mathbf{H}^{(q)})^{-1}, (\mathbf{I} - \gamma\mathbf{H}^{(q)})^k\left(\frac{\gamma\alpha_B \cdot \text{tr}\left(\mathbf{B}_{0,N}\right)}{1 - \gamma\alpha_B\,\text{tr}(\mathbf{H}^{(q)})}\cdot\mathbf{H}^{(q)} + \mathbf{B}_0\right)(\mathbf{I} - \gamma\mathbf{H}^{(q)})^k\right\rangle \\
&= \frac{1}{\gamma N^2} \sum_{k=0}^{N-1} \left\langle\left((\mathbf{I} - \gamma\mathbf{H}^{(q)})^{2k} - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{N+2k}\right)(\mathbf{H}^{(q)})^{-1}, \left(\frac{\gamma\alpha_B \cdot \text{tr}\left(\mathbf{B}_{0,N}\right)}{1 - \gamma\alpha_B\,\text{tr}(\mathbf{H}^{(q)})}\cdot\mathbf{H}^{(q)} + \mathbf{B}_0\right)\right\rangle \\
&\leq \frac{1}{\gamma N^2} \sum_{k=0}^{N-1} \left\langle\left((\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k} - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{N+k}\right)(\mathbf{H}^{(q)})^{-1}, \left(\frac{\gamma\alpha_B \cdot \text{tr}\left(\mathbf{B}_{0,N}\right)}{1 - \gamma\alpha_B\,\text{tr}(\mathbf{H}^{(q)})}\cdot\mathbf{H}^{(q)} + \mathbf{B}_0\right)\right\rangle,
\end{aligned}$$

where the second inequality holds by the bound for $\mathbf{S}_N = \sum_{t=0}^{N-1}\mathbf{B}_t$, i.e., (D.18). Denote

$$I_3 = \frac{1}{\gamma N^2} \sum_{k=0}^{N-1} \left\langle\left((\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k} - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{N+k}\right), \frac{\gamma\alpha_B \cdot \text{tr}\left(\mathbf{B}_{0,N}\right)}{1 - \gamma\alpha_B\,\text{tr}(\mathbf{H}^{(q)})}\cdot\mathbf{I}\right\rangle,$$

$$I_4 = \frac{1}{\gamma N^2} \sum_{k=0}^{N-1} \left\langle\left((\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k} - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{N+k}\right)(\mathbf{H}^{(q)})^{-1}, \mathbf{B}_0\right\rangle,$$

then

$$\text{B} \leq I_3 + I_4.$$

We then respectively analyze $I_3$ and $I_4$.

Regarding $I_3$, repeating the analysis of $I_1$,

$$\begin{aligned}
I_3 &= \frac{\alpha_B \cdot \text{tr}\left(\mathbf{B}_{0,N}\right)}{N^2(1 - \gamma\alpha_B\,\text{tr}(\mathbf{H}^{(q)}))} \sum_{k=0}^{N-1}\sum_{i} \left[(1 - \gamma\lambda_i^{(q)})^k - (1 - \gamma\lambda_i^{(q)})^{N+k}\right] \\
&= \frac{\alpha_B \cdot \text{tr}\left(\mathbf{B}_{0,N}\right)}{N^2(1 - \gamma\alpha_B\,\text{tr}(\mathbf{H}^{(q)}))} \sum_{i} \left[1 - (1 - \gamma\lambda_i^{(q)})^N\right]\sum_{k=0}^{N-1}(1 - \gamma\lambda_i^{(q)})^k \\
&= \frac{\alpha_B \cdot \text{tr}\left(\mathbf{B}_{0,N}\right)}{\gamma N^2(1 - \gamma\alpha_B\,\text{tr}(\mathbf{H}^{(q)}))} \sum_{i} \left[1 - (1 - \gamma\lambda_i^{(q)})^N\right]^2(\lambda_i^{(q)})^{-1} \\
&\leq \frac{\alpha_B \cdot \text{tr}\left(\mathbf{B}_{0,N}\right)}{\gamma N^2(1 - \gamma\alpha_B\,\text{tr}(\mathbf{H}^{(q)}))} \sum_{i} \min\left\{(\lambda_i^{(q)})^{-1}, \gamma^2 N^2\lambda_i^{(q)}\right\} \\
&\leq \frac{\alpha_B\,\text{tr}(\mathbf{B}_{0,N})}{\gamma(1 - \gamma\alpha_B\,\text{tr}(\mathbf{H}^{(q)}))} \cdot \left(\sum_{i\leq k^*}\frac{1}{N^2\lambda_i^{(q)}} + \gamma^2\sum_{i>k^*}\lambda_i^{(q)}\right).
\end{aligned}$$

Note that by (D.19),

$$\operatorname{tr}(\mathbf{B}_{0,N}) \le 2(\|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{I}_{0:k^*}} + N\gamma\|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{H}_{k^*:\infty}}),$$

we then obtain

$$I_3 \le \frac{2\alpha_B \left( \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{I}_{0:k^*}} + N\gamma\|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{H}_{k^*:\infty}} \right)}{N\gamma(1 - \gamma\alpha_B \operatorname{tr}(\mathbf{H}^{(q)}))} \left( \sum_{i \le k^*} \frac{1}{N\lambda_i^{(q)}} + N\gamma^2 \sum_{i > k^*} \lambda_i^{(q)} \right).$$

Regarding $I_4$, decompose $\mathbf{H}^{(q)} = \mathbf{V}^{(q)}\mathbf{\Lambda}^{(q)}\mathbf{V}^{(q)\top}$, note that $\mathbf{B}_0 = \boldsymbol{\eta}_0\boldsymbol{\eta}_0^\top$, it can be shown that the diagonal entries of $\mathbf{V}^{(q)\top}\mathbf{B}_0\mathbf{V}^{(q)}$ are $\omega_1^2, \ldots,$ where $\omega_i = \mathbf{v}_i^{(q)\top}\boldsymbol{\eta}_0 = \mathbf{v}_i^{(q)\top}(\mathbf{w}_0 - \mathbf{w}^{(q)*})$. Hence,

$$
\begin{aligned}
I_4 =& \frac{1}{\gamma N^2} \sum_{k=0}^{N-1} \sum_i \left[ (1 - \gamma\lambda_i^{(q)})^k - (1 - \gamma\lambda_i^{(q)})^{N+k} \right] (\lambda_i^{(q)})^{-1}\omega_i^2 \\
=& \frac{1}{\gamma N^2} \sum_i \left[ 1 - (1 - \gamma\lambda_i^{(q)})^N \right] (\lambda_i^{(q)})^{-1}\omega_i^2 \sum_{k=0}^{N-1}(1 - \gamma\lambda_i^{(q)})^k \\
=& \frac{1}{\gamma^2 N^2} \sum_i \frac{\omega_i^2}{(\lambda_i^{(q)})^2} \left[ 1 - (1 - \gamma\lambda_i^{(q)})^N \right]^2 \\
\le& \frac{1}{\gamma^2 N^2} \cdot \sum_{i \le k^*} \frac{\omega_i^2}{(\lambda_i^{(q)})^2} + \sum_{i > k^*} \omega_i^2 \\
=& \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{(\mathbf{H}_{0:k^*}^{(q)})^{-2}} + \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{I}_{k^*:\infty}^{(q)}}.
\end{aligned}
$$

Therefore,

$$
\begin{aligned}
\mathrm{B} \le& I_3 + I_4 \\
\le& \frac{2\alpha_B \left( \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{I}_{0:k^*}} + N\gamma\|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{H}_{k^*:\infty}} \right)}{N\gamma(1 - \gamma\alpha_B \operatorname{tr}(\mathbf{H}^{(q)}))} \left( \sum_{i \le k^*} \frac{1}{N\lambda_i^{(q)}} + N\gamma^2 \sum_{i > k^*} \lambda_i^{(q)} \right) \\
& + \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{(\mathbf{H}_{0:k^*}^{(q)})^{-2}} + \|\mathbf{w}_0 - \mathbf{w}^{(q)*}\|^2_{\mathbf{I}_{k^*:\infty}^{(q)}}.
\end{aligned}
$$

Regarding V,

$$
\begin{aligned}
\mathrm{V} =& \frac{1}{N^2} \cdot \sum_{t=0}^{N-1} \sum_{k=t}^{N-1} \operatorname{tr}\left( (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k-t}\mathbf{C}_t \right) \\
=& \frac{1}{N^2} \cdot \sum_{t=0}^{N-1} \sum_{k=t}^{N-1} \left\langle (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{k-t}, \mathbf{C}_t \right\rangle \\
=& \frac{1}{\gamma N^2} \cdot \sum_{t=0}^{N-1} \left\langle \left[ \mathbf{I} - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^{N-t} \right] (\mathbf{H}^{(q)})^{-1}, \mathbf{C}_t \right\rangle.
\end{aligned}
$$

By the upper bound for $\mathbf{C}_t$, i.e., (D.28),

$$\mathbf{C}_t \preceq \frac{\gamma\sigma_G^{(q)2}}{1 - \gamma\alpha_B\operatorname{tr}(\mathbf{H}^{(q)})} \cdot \left( \mathbf{I} - (\mathbf{I} - \gamma\mathbf{H}^{(q)})^t \right),$$

44

we have

$$
\begin{aligned}
\mathrm{V} \leq & \frac{1}{\gamma^2 N^2} \frac{\gamma^2 \sigma_G^{(q)^2}}{1 - \gamma \alpha_B \mathrm{tr}(\mathbf{H}^{(q)})} \sum_{t=0}^{N-1} \left\langle \left[ \mathbf{I} - (\mathbf{I} - \gamma \mathbf{H}^{(q)})^{N-t} \right] (\mathbf{H}^{(q)})^{-1}, \mathbf{I} - (\mathbf{I} - \gamma \mathbf{H}^{(q)})^t \right\rangle \\
= & \frac{1}{\gamma^2 N^2} \frac{\gamma^2 \sigma_G^{(q)^2}}{1 - \gamma \alpha_B \mathrm{tr}(\mathbf{H}^{(q)})} \sum_i \sum_{t=0}^{N-1} \left[ 1 - (1 - \gamma \lambda_i^{(q)})^{N-t} \right] (\lambda_i^{(q)})^{-1} \left[ 1 - (1 - \gamma \lambda_i^{(q)})^t \right] \\
\leq & \frac{1}{\gamma^2 N^2} \frac{\gamma^2 \sigma_G^{(q)^2}}{1 - \gamma \alpha_B \mathrm{tr}(\mathbf{H}^{(q)})} \sum_i \sum_{t=0}^{N-1} \left[ 1 - (1 - \gamma \lambda_i^{(q)})^N \right]^2 (\lambda_i^{(q)})^{-1} \\
= & \frac{1}{\gamma^2 N} \frac{\gamma^2 \sigma_G^{(q)^2}}{1 - \gamma \alpha_B \mathrm{tr}(\mathbf{H}^{(q)})} \sum_i \left[ 1 - (1 - \gamma \lambda_i^{(q)})^N \right]^2 (\lambda_i^{(q)})^{-1} \\
= & \frac{1}{\gamma^2 N} \frac{\gamma^2 \sigma_G^{(q)^2}}{1 - \gamma \alpha_B \mathrm{tr}(\mathbf{H}^{(q)})} \sum_i \min \left\{ (\lambda_i^{(q)})^{-1}, \gamma^2 N^2 \lambda_i^{(q)} \right\} \\
= & \frac{1}{\gamma^2} \frac{\gamma^2 \sigma_G^{(q)^2}}{1 - \gamma \alpha_B \mathrm{tr}(\mathbf{H}^{(q)})} \left( \sum_{i \leq k^*} \frac{1}{N \lambda_i^{(q)}} + \sum_{i > k^*} \gamma^2 N \lambda_i^{(q)} \right) \\
\leq & \frac{\sigma_G^{(q)^2}}{1 - \gamma \alpha_B \mathrm{tr}(\mathbf{H}^{(q)})} \left( \sum_{i \leq k^*} \frac{1}{N \lambda_i^{(q)}} + \sum_{i > k^*} \gamma^2 N \lambda_i^{(q)} \right).
\end{aligned}
$$

Combining the bound for B and V, we complete the proof. $\qquad\square$

# F  Deferring Proofs

## F.1  Proof for Theorem 4.1

*Proof.* By the fact that $\mathbf{H}^{(q)} = \mathbf{H} + \mathbf{D}$, Theorem 4.1 can be directly proved by Lemma A.2, Lemma E.2, Lemma D.16, Lemma C.1 and Lemma B.1. $\qquad\square$

## F.2  Proof for Theorem 4.2

We present another Theorem F.1 to provide precise excess risk bound rather than only in order and prove Theorem F.1 in this section. By Theorem F.1, Theorem 4.2 can be immediately proved by the fact that $\mathbf{w}^{(q)^*} = \mathbf{H}^{(q)^{-1}} \mathbf{H} \mathbf{w}^*$ and $\epsilon_i \leq O(1)$.

**Theorem F.1.** *Under Assumption 3.1,3.2,3.3,3.4 and notations in Theorem 4.1, if there exist $\epsilon_d, \epsilon_l, \epsilon_p, \epsilon_a$ and $\epsilon_o$ such that for any $i \in \{d, l, p, a, o\}$, quantization $\mathcal{Q}_i$ is $\epsilon_i$-multiplicative and the stepsize satisfies $\gamma < \frac{1}{\alpha_B (1 + \epsilon_d)(1 + \tilde{\epsilon}) \mathrm{tr}(\mathbf{H})}$, then the excess risk can be upper bounded as follows.*

$$
\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] \leq \mathrm{ApproxErr} + \frac{1 + 3\epsilon_d}{1 + \epsilon_d} \mathrm{VarErr} + \frac{1 + 3\epsilon_d}{1 + \epsilon_d} \mathrm{BiasErr},
$$

*where*

$$
\mathrm{ApproxErr} = \|\mathbf{w}^*\|_{\mathbf{H}}^2 \cdot \frac{\left( \frac{3}{2} + \frac{1}{2}\epsilon_d \right) \epsilon_d}{(1 + \epsilon_d)^2} + \epsilon_l \mathbb{E}[y^2],
$$

$$
\mathrm{BiasErr} = \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^{(q)^*}\|_{(\mathbf{H}_{0:k^*}^{(q)})^{-1}}^2 + \|\mathbf{w}^{(q)^*}\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2,
$$

$$\text{VarErr} = \left( \frac{k^*}{N} + N\gamma^2(1+\epsilon_d)^2 \cdot \sum_{i>k^*} \lambda_i^2 \right) \frac{\sigma_M^{(q)2}}{1 - (1+\tilde{\epsilon})\gamma\alpha_B(1+\epsilon_d)\text{tr}\,(\mathbf{H})}$$

$$+ \left( \frac{k^*}{N} + N\gamma^2(1+\epsilon_d)^2 \cdot \sum_{i>k^*} \lambda_i^2 \right) \frac{2(1+\tilde{\epsilon})\alpha_B \left( \|\mathbf{w}^{(q)*}\|_{\mathbf{I}_{0:k^*}^{(q)}}^2 + N\gamma\|\mathbf{w}^{(q)*}\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2 \right)}{N\gamma\,[1 - (1+\tilde{\epsilon})\gamma\alpha_B(1+\epsilon_d)\text{tr}\,(\mathbf{H})]},$$

*with*

$$\sigma_M^{(q)2} := \frac{(1+4\epsilon_o)\sigma^2}{B} + \frac{\|\mathbf{w}^*\|_{\mathbf{H}}^2}{1+\epsilon_d}[4\epsilon_o[(1+\epsilon_a)(1+\epsilon_p)+1]\alpha_B + 2\epsilon_a(1+\epsilon_p)\alpha_B + 2\epsilon_p\alpha_B],$$

$$\tilde{\epsilon} := 2\epsilon_p + 4\epsilon_o(1+\epsilon_a)(1+\epsilon_p) + 2\epsilon_a(1+\epsilon_p).$$

*Proof.* By Lemma A.2, Lemma E.1, Lemma D.17, Lemma C.1 and Lemma B.1, applying the multiplicative condition, Theorem F.1 is immediately proved. □

We would like to remark that, the multiplicative nature introduces additional complexity into the update rule, resulting in an additional parameter $\tilde{\epsilon}$ in VarErr. It is worth noting that when $\epsilon_p, \epsilon_a, \epsilon_o$ are at most constant level, $\tilde{\epsilon}$ is at most constant level and can therefore be merged.

## F.3 Proof for Corollary 4.1

We present Theorem F.2 to provide precise excess risk bound and prove Theorem F.2 in this section. By Theorem F.2, Corollary 4.1 can be immediately proved by the fact that $\mathbf{w}^{(q)*} = \mathbf{H}^{(q)-1}\mathbf{H}\mathbf{w}^*$.

**Theorem F.2.** *Under Assumption 3.1,3.2,3.3,3.4 and notations in Theorem 4.1, if there exist $\epsilon_d, \epsilon_l, \epsilon_p, \epsilon_a$ and $\epsilon_o$ such that for any $i \in \{d,l,p,a,o\}$, quantization $\mathcal{Q}_i$ is $\epsilon_i$-additive and the stepsize satisfies $\gamma < \frac{1}{\gamma\alpha_B\text{tr}(\mathbf{H}+\epsilon_d\mathbf{I})}$, then*

$$\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] \leq \text{ApproxErr} + 2\text{VarErr} + 2\text{BiasErr},$$

*where*

$$\text{ApproxErr} = \epsilon_l + \frac{3\epsilon_d}{2}\|\mathbf{w}^*\|_{\mathbf{H}(\mathbf{H}+\epsilon_d\mathbf{I})^{-1}(\mathbf{H}+\epsilon_d\mathbf{I})^{-1}\mathbf{H}}^2 + \frac{\epsilon_d^2}{2}\|\mathbf{w}^*\|_{(\mathbf{H}+\epsilon_d\mathbf{I})^{-1}\mathbf{H}(\mathbf{H}+\epsilon_d\mathbf{I})^{-1}}^2,$$

$$\text{BiasErr} = \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^{(q)*}\|_{(\mathbf{H}_{0:k^*}^{(q)})^{-1}}^2 + \|\mathbf{w}^{(q)*}\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2,$$

$$\text{VarErr} = \frac{\sigma_A^{(q)2} + \frac{2\alpha_B}{N\gamma} \left( \|\mathbf{w}^{(q)*}\|_{\mathbf{I}_{0:k^*}^{(q)}}^2 + N\gamma\|\mathbf{w}^{(q)*}\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2 \right)}{1 - \gamma\alpha_B\text{tr}\,(\mathbf{H}+\epsilon_d\mathbf{I})} \left( \frac{k^*}{N} + N\gamma^2 \sum_{i>k^*}(\lambda_i+\epsilon_d)^2 \right),$$

*with $\sigma_A^{(q)2} = \frac{\epsilon_o+\epsilon_a}{B} + \alpha_B\epsilon_p\text{tr}\,(\mathbf{H}+\epsilon_d\mathbf{I}) + \frac{\sigma^2}{B}$.*

*Proof.* By Theorem 4.1 and the additive condition, Theorem F.2 is immediately proved. □

## F.4 Proof for the Multiplicative Statement in Corollary 4.2

*Proof.* We prove by applying Theorem F.1. Denote $k_0^* = \max\left\{ k : \lambda_k \geq \frac{1}{N\gamma} \right\}$, the key of the proof is to convert $k^*$ into $k_0^*$. The key property is that for $k_0^* < i \leq k^*$, $\frac{1}{N\gamma} \leq \lambda_i^{(q)}$. We first handle

$\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}}$ in VarErr.

$$\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}}$$

$$=\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{I}^{(q)}_{0:k_0^*}} + \|\mathbf{w}^{(q)^*}\|^2_{\mathbf{I}^{(q)}_{k_0^*:k^*}} - N\gamma\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k_0^*:k^*}} + N\gamma\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k_0^*:\infty}}$$

$$\leq\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{I}^{(q)}_{0:k_0^*}} + N\gamma\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k_0^*:\infty}} \tag{F.1}$$

$$=\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{I}_{0:k_0^*}} + (1+\epsilon_d)N\gamma\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}_{k_0^*:\infty}}$$

$$=\frac{1}{(1+\epsilon_d)^2}\|\mathbf{w}^*\|^2_{\mathbf{I}_{0:k_0^*}} + \frac{1}{(1+\epsilon_d)}N\gamma\|\mathbf{w}^*\|^2_{\mathbf{H}_{k_0^*:\infty}},$$

where the first inequality holds by $\frac{1}{N\gamma} \leq \lambda_i^{(q)}, k_0^* < i \leq k^*$ and the last inequality holds by $\mathbf{w}^{(q)^*} = \mathbf{H}^{(q)^{-1}}\mathbf{H}\mathbf{w}^* = \frac{1}{1+\epsilon_d}\mathbf{w}^*$.

We then analyze $\frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^{(q)^*}\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-1}} + \|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}}$ in BiasErr.

$$\frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^{(q)^*}\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-1}} + \|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}}$$

$$=\frac{1}{\gamma^2 N^2} \cdot \left(\|\mathbf{w}^{(q)^*}\|^2_{(\mathbf{H}^{(q)}_{0:k_0^*})^{-1}} + \|\mathbf{w}^{(q)^*}\|^2_{(\mathbf{H}^{(q)}_{k_0^*:k^*})^{-1}}\right) - \|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k_0^*:k^*}} + \|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k_0^*:\infty}}$$

$$\leq\frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^{(q)^*}\|^2_{(\mathbf{H}^{(q)}_{0:k_0^*})^{-1}} + \|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k_0^*:\infty}} \tag{F.2}$$

$$=\frac{1}{(1+\epsilon_d)\gamma^2 N^2} \cdot \|\mathbf{w}^{(q)^*}\|^2_{(\mathbf{H}_{0:k_0^*})^{-1}} + (1+\epsilon_d)\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}_{k_0^*:\infty}}$$

$$=\frac{1}{(1+\epsilon_d)^3\gamma^2 N^2} \cdot \|\mathbf{w}^*\|^2_{(\mathbf{H}_{0:k_0^*})^{-1}} + \frac{1}{(1+\epsilon_d)}\|\mathbf{w}^*\|^2_{\mathbf{H}_{k_0^*:\infty}},$$

where the first inequality holds by $\frac{1}{N\gamma} \leq \lambda_i^{(q)}, k_0^* < i \leq k^*$ and the last inequality holds by $\mathbf{w}^{(q)^*} = \mathbf{H}^{(q)^{-1}}\mathbf{H}\mathbf{w}^* = \frac{1}{1+\epsilon_d}\mathbf{w}^*$.

We next handle $\frac{k^*}{N} + N\gamma^2(1+\epsilon_d)^2 \cdot \sum_{i>k^*} \lambda_i^2$. Note that for $k_0^* < i \leq k^*$, $\frac{1}{N\gamma(1+\epsilon_d)} \leq \lambda_i < \frac{1}{N\gamma}$, we have

$$\frac{k^*}{N} + N\gamma^2(1+\epsilon_d)^2 \cdot \sum_{i>k^*} \lambda_i^2$$

$$=\frac{k_0^*}{N} + \frac{k^* - k_0^*}{N} - N\gamma^2(1+\epsilon_d)^2 \cdot \sum_{k_0^*<i\leq k^*} \lambda_i^2 + N\gamma^2(1+\epsilon_d)^2 \cdot \sum_{i>k_0^*} \lambda_i^2$$

$$\leq\frac{k_0^*}{N} + \frac{k^* - k_0^*}{N} - N\gamma^2(1+\epsilon_d)^2(k^* - k_0^*)\frac{1}{N^2\gamma^2(1+\epsilon_d)^2} + N\gamma^2(1+\epsilon_d)^2 \cdot \sum_{i>k_0^*} \lambda_i^2 \tag{F.3}$$

$$=\frac{k_0^*}{N} + N\gamma^2(1+\epsilon_d)^2 \cdot \sum_{i>k_0^*} \lambda_i^2.$$

Therefore, applying (F.1), (F.2) and (F.3) into Theorem F.1, we have

$$
\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)]
$$
$$
\leq \|\mathbf{w}^*\|_{\mathbf{H}}^2 \cdot \frac{\left(\frac{3}{2} + \frac{1}{2}\epsilon_d\right)\epsilon_d}{(1+\epsilon_d)^2} + \epsilon_l \mathbb{E}[y^2]
$$
$$
+ \left(\frac{k_0^*}{N} + N\gamma^2(1+\epsilon_d)^2 \cdot \sum_{i>k_0^*}\lambda_i^2\right) \cdot \frac{1+3\epsilon_d}{1+\epsilon_d} \cdot \frac{\sigma_M^{(q)^2}}{1-(1+\tilde{\epsilon})\gamma\alpha_B(1+\epsilon_d)\mathrm{tr}\,(\mathbf{H})}
$$
$$
+ \left(\frac{k_0^*}{N} + N\gamma^2(1+\epsilon_d)^2 \cdot \sum_{i>k_0^*}\lambda_i^2\right) \cdot \frac{1+3\epsilon_d}{1+\epsilon_d} \cdot \frac{2(1+\tilde{\epsilon})\alpha_B\left(\frac{1}{(1+\epsilon_d)^2}\|\mathbf{w}^*\|_{\mathbf{I}_{0:k_0^*}}^2 + \frac{1}{(1+\epsilon_d)}N\gamma\|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}}^2\right)}{N\gamma\left[1-(1+\tilde{\epsilon})\gamma\alpha_B(1+\epsilon_d)\mathrm{tr}\,(\mathbf{H})\right]}
$$
$$
+ \frac{1}{(1+\epsilon_d)^3\gamma^2N^2} \cdot \|\mathbf{w}^*\|_{(\mathbf{H}_{0:k_0^*})^{-1}}^2 + \frac{1}{(1+\epsilon_d)}\|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}}^2.
$$

Denote the standard excess risk bound in full-precision setting (Zou et al., 2023)

$$
R_0 = \text{EffectiveVarI} + \text{EffectiveVarI} + \text{EffectiveBias},
$$

where

$$
\text{EffectiveVarI} = \left(\frac{k_0^*}{N} + N\gamma^2 \cdot \sum_{i>k_0^*}\lambda_i^2\right)\frac{4\alpha_B\left(\|\mathbf{w}_0 - \mathbf{w}^*\|_{\mathbf{I}_{0:k_0^*}}^2 + N\gamma\|\mathbf{w}_0 - \mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}}^2\right)}{N\gamma\left[1-\gamma\alpha_B\mathrm{tr}\,(\mathbf{H})\right]},
$$
$$
\text{EffectiveVarII} = \left(\frac{k_0^*}{N} + N\gamma^2 \cdot \sum_{i>k_0^*}\lambda_i^2\right)\frac{1}{B}\frac{\sigma^2}{1-\gamma\alpha_B\mathrm{tr}\,(\mathbf{H})},
$$
$$
\text{EffectiveBias} = \frac{2}{\gamma^2N^2} \cdot \|\mathbf{w}_0 - \mathbf{w}^*\|_{(\mathbf{H}_{0:k_0^*})^{-1}}^2 + 2\|\mathbf{w}_0 - \mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}}^2,
$$

with $k_0^* = \max\left\{k : \lambda_k \geq \frac{1}{N\gamma}\right\}$. If $\epsilon_i \leq O(1), \forall i \in \{d, p, a, o\}$, then

$$
\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] \lesssim \|\mathbf{w}^*\|_{\mathbf{H}}^2\,\epsilon_d + \epsilon_l + \text{EffectiveVarI} + \frac{\sigma_M^{(q)^2}}{\sigma^2/B}\text{EffectiveVarII} + \text{EffectiveBias}.
$$

Therefore, if

$$
\epsilon_l \leq O\left(R_0\right),\ \ \epsilon_p, \epsilon_a, \epsilon_o \leq O\left(\frac{\sigma^2}{B\|\mathbf{w}^*\|_{\mathbf{H}}^2} \wedge 1\right),\ \ \epsilon_d \leq O\left(\frac{R_0}{\|\mathbf{w}^*\|_{\mathbf{H}}^2} \wedge 1\right).
$$

then

$$
\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] \leq O(R_0).
$$

$\square$

## F.5  Proof for the Additive Statement in Corollary 4.2

*Proof.* We prove by applying Theorem F.2. The proof idea is similar to the proof in Section F.4, which mainly relies on the fact that $\lambda_i\,(\mathbf{H} + \mathbf{D}) \geq \frac{1}{N\gamma}, k_0^* < i \leq k^*$. We begin by handling

$\frac{k^*}{N} + N\gamma^2 \cdot \sum_{i>k^*} \lambda_i \left(\mathbf{H} + \mathbf{D}\right)^2$ in VarErr.

$$
\begin{aligned}
\frac{k^*}{N} + N\gamma^2 \sum_{i>k^*} \lambda_i^2 \left(\mathbf{H} + \mathbf{D}\right) &= \frac{k_0^*}{N} + \frac{k^* - k_0^*}{N} - N\gamma^2 \sum_{k_0^* < i \leq k^*} \lambda_i^2 \left(\mathbf{H} + \mathbf{D}\right) + N\gamma^2 \sum_{i>k_0^*} \lambda_i^2 \left(\mathbf{H} + \mathbf{D}\right) \\
&\leq \frac{k_0^*}{N} + N\gamma^2 \sum_{i>k_0^*} \lambda_i^2 \left(\mathbf{H} + \mathbf{D}\right).
\end{aligned}
\tag{F.4}
$$

We second handle $\|\mathbf{w}^{(q)*}\|_{\mathbf{I}_{0:k^*}^{(q)}}^2 + N\gamma\|\mathbf{w}^{(q)*}\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2$ in VarErr. A key observation is that for $\mathbf{D} = \epsilon_d \mathbf{I}$, $\mathbf{H}^{(q)} = \mathbf{H} + \mathbf{D}$ has same eigenvectors as $\mathbf{H}$. This immediately implies that

$$
\begin{aligned}
&\|\mathbf{w}^{(q)*}\|_{\mathbf{I}_{0:k^*}^{(q)}}^2 + N\gamma\|\mathbf{w}^{(q)*}\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2 \\
&= \mathbf{w}^{*\top}\mathbf{H}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{I}_{0:k^*}^{(q)}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{H}\mathbf{w}^* + N\gamma\mathbf{w}^{*\top}\mathbf{H}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{H}_{k^*:\infty}^{(q)}(\mathbf{H}+\mathbf{D})^{-1}\mathbf{H}\mathbf{w}^* \\
&= \sum_{i\leq k^*} \frac{\lambda_i^2}{(\lambda_i^{(q)})^2}(\mathbf{w}^{*\top}\mathbf{v}_i)^2 + \sum_{i>k^*} \frac{N\gamma\lambda_i^2}{\lambda_i^{(q)}}(\mathbf{w}^{*\top}\mathbf{v}_i)^2 \\
&\leq \|\mathbf{w}^*\|_{\mathbf{I}_{0:k^*}^{(q)}}^2 + N\gamma\|\mathbf{w}^*\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2 \\
&= \|\mathbf{w}^*\|_{\mathbf{I}_{0:k_0^*}^{(q)}}^2 + \|\mathbf{w}^*\|_{\mathbf{I}_{k_0^*:k^*}^{(q)}}^2 - N\gamma\|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:k^*}^{(q)}}^2 + N\gamma\|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}^{(q)}}^2 \\
&\leq \|\mathbf{w}^*\|_{\mathbf{I}_{0:k_0^*}^{(q)}}^2 + N\gamma\|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}^{(q)}}^2,
\end{aligned}
\tag{F.5}
$$

where the last inequality holds by $\lambda_i \left(\mathbf{H} + \mathbf{D}\right) \geq \frac{1}{N\gamma}, k_0^* < i \leq k^*$. Further note that

$$
\begin{aligned}
N\gamma\|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}^{(q)}}^2 &= N\gamma \sum_{i>k_0^*}(\mathbf{w}^{*\top}\mathbf{v}_i)^2(\lambda_i + \epsilon_d) \\
&= N\gamma \sum_{i>k_0^*}(\mathbf{w}^{*\top}\mathbf{v}_i)^2\lambda_i + N\gamma \sum_{i>k_0^*}(\mathbf{w}^{*\top}\mathbf{v}_i)^2\epsilon_d \\
&= N\gamma\|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}}^2 + \epsilon_d N\gamma\|\mathbf{w}^*\|_{\mathbf{I}_{k_0^*:\infty}}^2,
\end{aligned}
\tag{F.6}
$$

hence by (F.5) and (F.6) we have

$$
\|\mathbf{w}^{(q)*}\|_{\mathbf{I}_{0:k^*}^{(q)}}^2 + N\gamma\|\mathbf{w}^{(q)*}\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2 \leq \|\mathbf{w}^*\|_{\mathbf{I}_{0:k_0^*}}^2 + N\gamma\|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}}^2 + \epsilon_d N\gamma\|\mathbf{w}^*\|_{\mathbf{I}_{k_0^*:\infty}}^2.
\tag{F.7}
$$

We next cope with $\frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^{(q)*}\|_{(\mathbf{H}_{0:k^*}^{(q)})^{-1}}^2 + \|\mathbf{w}^{(q)*}\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2$ in BiasErr.

$$
\begin{aligned}
&\frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^{(q)*}\|_{(\mathbf{H}_{0:k^*}^{(q)})^{-1}}^2 + \|\mathbf{w}^{(q)*}\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2 \\
&= \frac{1}{\gamma^2 N^2} \sum_{i\leq k^*} \frac{\lambda_i^2}{(\lambda_i^{(q)})^2}\frac{1}{\lambda_i^{(q)}}(\mathbf{w}^{*\top}\mathbf{v}_i)^2 + \sum_{i>k^*} \frac{\lambda_i^2}{(\lambda_i^{(q)})^2}\lambda_i^{(q)}(\mathbf{w}^{*\top}\mathbf{v}_i)^2 \\
&\leq \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^*\|_{(\mathbf{H}_{0:k^*}^{(q)})^{-1}}^2 + \|\mathbf{w}^*\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2 \\
&= \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^*\|_{(\mathbf{H}_{0:k_0^*}^{(q)})^{-1}}^2 + \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^*\|_{(\mathbf{H}_{k_0^*:k^*}^{(q)})^{-1}}^2 - \|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:k^*}^{(q)}}^2 + \|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}^{(q)}}^2 \\
&\leq \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^*\|_{(\mathbf{H}_{0:k_0^*}^{(q)})^{-1}}^2 + \|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}^{(q)}}^2,
\end{aligned}
\tag{F.8}
$$

where the first equality holds by $\mathbf{w}^{(q)*} = \mathbf{H}^{(q)-1}\mathbf{H}\mathbf{w}^*$ and the last inequality holds by $\lambda_i(\mathbf{H} + \mathbf{D}) \geq \frac{1}{N\gamma}, k_0^* < i \leq k^*$. Further note that

$$\|\mathbf{w}^*\|_{(\mathbf{H}_{0:k_0^*}^{(q)})^{-1}}^2 = \sum_{i \leq k_0^*} \frac{(\mathbf{w}^{*\top}\mathbf{v}_i)^2}{\lambda_i + \epsilon_d} \leq \sum_{i \leq k_0^*} \frac{(\mathbf{w}^{*\top}\mathbf{v}_i)^2}{\lambda_i} = \|\mathbf{w}^*\|_{(\mathbf{H}_{0:k_0^*})^{-1}}^2, \tag{F.9}$$

hence by (F.8), (F.6) and (F.9), we have

$$\frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^{(q)*}\|_{(\mathbf{H}_{0:k^*}^{(q)})^{-1}}^2 + \|\mathbf{w}^{(q)*}\|_{\mathbf{H}_{k^*:\infty}^{(q)}}^2 \leq \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^*\|_{(\mathbf{H}_{0:k_0^*})^{-1}}^2 + \|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}}^2 + \epsilon_d \|\mathbf{w}^*\|_{\mathbf{I}_{k_0^*:\infty}}^2. \tag{F.10}$$

Finally, applying (F.4), (F.7) and (F.10) into Theorem F.2, we have

$$\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)]$$
$$\leq \epsilon_l + \frac{3\epsilon_d}{2}\|\mathbf{w}^*\|_{\mathbf{H}(\mathbf{H}+\epsilon_d\mathbf{I})^{-1}(\mathbf{H}+\epsilon_d\mathbf{I})^{-1}\mathbf{H}}^2 + \frac{\epsilon_d^2}{2}\|\mathbf{w}^*\|_{(\mathbf{H}+\epsilon_d\mathbf{I})^{-1}\mathbf{H}(\mathbf{H}+\epsilon_d\mathbf{I})^{-1}}^2$$
$$+ \frac{\sigma_A^{(q)2} + \frac{2\alpha_B}{N\gamma}\left(\|\mathbf{w}^*\|_{\mathbf{I}_{0:k_0^*}}^2 + N\gamma\|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}}^2 + \epsilon_d N\gamma\|\mathbf{w}^*\|_{\mathbf{I}_{k_0^*:\infty}}^2\right)}{1 - \gamma\alpha_B\mathrm{tr}(\mathbf{H}+\epsilon_d\mathbf{I})}\left(\frac{k_0^*}{N} + N\gamma^2\sum_{i>k_0^*}(\lambda_i + \epsilon_d)^2\right)$$
$$+ \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^*\|_{(\mathbf{H}_{0:k_0^*})^{-1}}^2 + \|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}}^2 + \epsilon_d\|\mathbf{w}^*\|_{\mathbf{I}_{k_0^*:\infty}}^2.$$

Denote

$$C_1 = \frac{\frac{k_0^*}{N} + N\gamma^2\sum_{i>k_0^*}(\lambda_i + \epsilon_d)^2}{\frac{k_0^*}{N} + N\gamma^2 \cdot \sum_{i>k_0^*}\lambda_i^2}, \quad C_2 = \frac{\epsilon_d\|\mathbf{w}^*\|_{\mathbf{I}_{k_0^*:\infty}}^2}{\|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}}^2 + \frac{1}{N^2\gamma^2}\|\mathbf{w}^*\|_{\mathbf{H}_{0:k_0^*}^{-1}}^2}.$$

Under the definition of $R_0$,

$$\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] \lesssim \|\mathbf{w}^*\|^2\epsilon_d + \epsilon_l + \frac{\sigma_A^{(q)2}}{\sigma^2/B}C_1\mathrm{EffectiveVarII}$$
$$+ (1 + C_2)C_1\mathrm{EffectiveVarI} + (1 + C_2)\mathrm{EffectiveBias}.$$

Therefore, if

$$\epsilon_l \leq O(R_0), \ \epsilon_o + \epsilon_a \leq O(\sigma^2), \ \epsilon_p \leq O\left(\frac{\sigma^2}{B\mathrm{tr}(\mathbf{H}+\epsilon_d\mathbf{I})}\right),$$

$$\epsilon_d \leq O\left(\frac{R_0}{\|\mathbf{w}^*\|^2} \wedge \sqrt{\frac{\sum_{i>k_0^*}\lambda_i^2}{d - k_0^*}} \wedge \frac{\|\mathbf{w}^*\|_{\mathbf{H}_{k_0^*:\infty}}^2 + \frac{1}{N^2\gamma^2}\|\mathbf{w}^*\|_{\mathbf{H}_{0:k_0^*}^{-1}}^2}{\|\mathbf{w}^*\|_{\mathbf{I}_{k_0^*:\infty}}^2}\right),$$

then

$$\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] \leq O(R_0).$$

$\square$

## F.6 Proof for the Multiplicative Statement in Corollary 4.3

*Proof.* We prove by applying Theorem F.1. Note that by the power-law assumption, we can estimate $k^*$ by

$$(1 + \epsilon_d)k^{*-a} \approx \frac{1}{N\gamma},$$

that is

$$k^* \approx [N\gamma(1+\epsilon_d)]^{\frac{1}{a}}. \tag{F.11}$$

Further, the power-law assumption also implies that for any positive $k$,

$$\sum_{i>k} i^{-a} \approx k^{1-a}. \tag{F.12}$$

We first compute $\frac{k^*}{N} + N\gamma^2(1+\epsilon_d)^2 \cdot \sum_{i>k^*} \lambda_i^2$. For one thing, for $i > k^*$, $\lambda_i \leq \frac{1}{N\gamma(1+\epsilon_d)}$, hence

$$\frac{k^*}{N} + N\gamma^2(1+\epsilon_d)^2 \cdot \sum_{i>k^*} \lambda_i^2 \leq \frac{d}{N}. \tag{F.13}$$

For another, applying (F.11) and (F.12), we have

$$\frac{k^*}{N} + N\gamma^2(1+\epsilon_d)^2 \cdot \sum_{i>k^*} \lambda_i^2 \leq \frac{[N\gamma(1+\epsilon_d)]^{\frac{1}{a}} + N^2\gamma^2(1+\epsilon_d)^2[N\gamma(1+\epsilon_d)]^{\frac{1-2a}{a}}}{N}. \tag{F.14}$$

Jointly, by (F.13) and (F.14), we have

$$\frac{k^*}{N} + N\gamma^2(1+\epsilon_d)^2 \cdot \sum_{i>k^*} \lambda_i^2 \leq O\left(\frac{\min\left\{d, [N\gamma(1+\epsilon_d)]^{\frac{1}{a}}\right\}}{N}\right). \tag{F.15}$$

We secondly compute $\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}}$. Assuming constant level optimal parameter (i.e., $\|\mathbf{w}_i^*\|^2 = \Theta(1), \ \forall i > 0$), we have

$$\|\mathbf{w}^*\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}^*\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} = O\left(k^* + N\gamma(1+\epsilon_d)\sum_{i>k^*}\lambda_i\right). \tag{F.16}$$

Utilizing the same technique in (F.15), we have

$$k^* + N\gamma(1+\epsilon_d)\sum_{i>k^*}\lambda_i \leq O\left(\min\left\{d, [N\gamma(1+\epsilon_d)]^{\frac{1}{a}}\right\}\right). \tag{F.17}$$

Hence, with (F.16) and (F.17), we have

$$\|\mathbf{w}^*\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}^*\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} \leq O\left(\min\left\{d, [N\gamma(1+\epsilon_d)]^{\frac{1}{a}}\right\}\right). \tag{F.18}$$

Therefore,

$$\begin{aligned}
\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} &= \frac{1}{(1+\epsilon_d)^2}\left(\|\mathbf{w}^*\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}^*\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}}\right) \\
&\leq O\left(\min\left\{d, [N\gamma(1+\epsilon_d)]^{\frac{1}{a}}\right\}\right),
\end{aligned} \tag{F.19}$$

where the equality holds by $\mathbf{w}^{(q)^*} = \mathbf{H}^{(q)^{-1}}\mathbf{H}\mathbf{w}^*$ and the inequality holds by (F.18).

We thirdly compute $\frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^{(q)^*}\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-1}} + \|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}}$. Assuming constant level optimal parameter (i.e., $\|\mathbf{w}_i^*\|^2 = \Theta(1), \ \forall i > 0$), we have

$$\frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^*\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-1}} + \|\mathbf{w}^*\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} = O\left(\frac{1}{N^2\gamma^2}\sum_{i=1}^{k^*}(1+\epsilon_d)^{-1}\lambda_i^{-1} + \sum_{i=k^*+1}^{d}(1+\epsilon_d)\lambda_i\right). \tag{F.20}$$

51

Utilizing the same technique in (F.15), we have

$$\frac{1}{N^2\gamma^2}\sum_{i=1}^{k^*}(1+\epsilon_d)^{-1}\lambda_i^{-1} + \sum_{i=k^*+1}^{d}(1+\epsilon_d)\lambda_i \le \frac{k^*}{N\gamma} + \sum_{i=k^*+1}^{d}(1+\epsilon_d)\lambda_i$$

$$\le O\left(\frac{\min\left\{d, [N\gamma(1+\epsilon_d)]^{\frac{1}{a}}\right\}}{N\gamma}\right), \tag{F.21}$$

where the first inequality uses $\lambda_i \ge \frac{1}{N\gamma(1+\epsilon_d)}, i \le k^*$. Hence,

$$\frac{1}{\gamma^2 N^2}\cdot\|\mathbf{w}^{(q)*}\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-1}} + \|\mathbf{w}^{(q)*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} \le \frac{1}{\gamma^2 N^2}\cdot\|\mathbf{w}^*\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-1}} + \|\mathbf{w}^*\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}}$$

$$\le O\left(\frac{\min\left\{d, [N\gamma(1+\epsilon_d)]^{\frac{1}{a}}\right\}}{N\gamma}\right), \tag{F.22}$$

where the first inequality holds by $\mathbf{w}^{(q)*} = \mathbf{H}^{(q)^{-1}}\mathbf{H}\mathbf{w}^*$ and the last inequality holds by (F.20) and (F.21).

Therefore, applying (F.15), (F.19) and (F.22) into Theorem F.1, we have

$$\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] \le O\left(\frac{\epsilon_d}{1+\epsilon_d} + \epsilon_l\right) + O\left(\frac{\min\left\{d, [N\gamma(1+\epsilon_d)]^{\frac{1}{a}}\right\}}{N\gamma}\right)$$

$$+O\left(\frac{\min\left\{d, [N\gamma(1+\epsilon_d)]^{\frac{1}{a}}\right\}}{N}\cdot\left(\frac{\sigma^2}{B} + (\epsilon_p + \epsilon_o + \epsilon_a) + \frac{(1+\tilde{\epsilon})\min\left\{d, [N\gamma(1+\epsilon_d)]^{\frac{1}{a}}\right\}}{N\gamma}\right)\right).$$

$\square$

## F.7 Proof for the Additive Statement in Corollary 4.3

*Proof.* We prove by applying Theorem F.2. We first consider the case $\epsilon_d + d^{-a} \ge \frac{1}{N\gamma}$. In this case,

$$k^* = \max\left\{k : \lambda_k^{(q)} \ge \frac{1}{N\gamma}\right\} = d.$$

Hence, by Theorem F.2,

VarErr

$$= \frac{\sigma_A^{(q)^2} + \frac{2\alpha_B}{N\gamma}\left(\|\mathbf{w}^{(q)*}\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}^{(q)*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}}\right)}{1 - \gamma\alpha_B\mathrm{tr}\left(\mathbf{H} + \epsilon_d\mathbf{I}\right)}\cdot\left(\frac{k^*}{N} + N\gamma^2\sum_{i>k^*}(\lambda_i + \epsilon_d)^2\right)$$

$$= \frac{\sigma_A^{(q)^2} + \frac{2\alpha_B}{N\gamma}\|\mathbf{w}^{(q)*}\|^2_{\mathbf{I}^{(q)}_{0:d}}}{1 - \gamma\alpha_B\mathrm{tr}\left(\mathbf{H} + \epsilon_d\mathbf{I}\right)}\frac{d}{N}.$$

BiasErr

$$= \frac{1}{\gamma^2 N^2}\cdot\|\mathbf{w}^{(q)*}\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-1}} + \|\mathbf{w}^{(q)*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}}$$

$$= \frac{1}{\gamma^2 N^2}\cdot\|\mathbf{w}^{(q)*}\|^2_{(\mathbf{H}^{(q)})^{-1}}.$$

Note that by $\mathbf{w}^{(q)^*} = \mathbf{H}^{(q)^{-1}}\mathbf{H}\mathbf{w}^*$, it holds

$$\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{I}^{(q)}_{0:d}} \le \|\mathbf{w}^*\|^2_{\mathbf{I}^{(q)}_{0:d}} = O(d),$$

$$\|\mathbf{w}^{(q)^*}\|^2_{(\mathbf{H}^{(q)})^{-1}} \le \|\mathbf{w}^*\|^2_{(\mathbf{H}^{(q)})^{-1}} \le O\left(\frac{d}{\epsilon_d + d^{-a}}\right).$$

Therefore,

$$
\begin{aligned}
\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] \le& O\left(\frac{\epsilon_d}{1+\epsilon_d}d + \epsilon_l\right) + O\left(\frac{d}{\gamma^2 N^2(\epsilon_d + d^{-a})}\right) \\
&+ O\left(\frac{d}{N}\left(\frac{\epsilon_o + \epsilon_a}{B} + (1 + d\epsilon_d)\epsilon_p + \frac{\sigma^2}{B} + \frac{d}{N\gamma}\right)\right).
\end{aligned}
\tag{F.23}
$$

Next, we consider $\epsilon_d + d^{-a} \le \frac{1}{N\gamma}$. We first compute $\frac{k^*}{N} + N\gamma^2\sum_{i>k^*}(\lambda_i + \epsilon_d)^2$. Note that by the power-law assumption, we can estimate $k^*$ by

$$k^{*-a} + \epsilon_d \approx \frac{1}{N\gamma}.$$

That is,

$$k^* \approx \left[\frac{1}{N\gamma} - \epsilon_d\right]^{-\frac{1}{a}}. \tag{F.24}$$

Hence, by (F.24) together with (F.12), we have

$$
\begin{aligned}
\frac{k^*}{N} + N\gamma^2\sum_{i>k^*}(\lambda_i + \epsilon_d)^2 \le& \frac{\left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}} + N^2\gamma^2\left[2\left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1-2a}{a}} + 2\epsilon_d^2\left(d - \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}}\right)\right]}{N} \\
\le& O\left(\frac{\epsilon_d^2 N^2\gamma^2\left[d - \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}}\right] + \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}}}{N}\right).
\end{aligned}
\tag{F.25}
$$

We second compute $\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}}$.

$$
\begin{aligned}
&\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}^{(q)^*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} \\
\le& \|\mathbf{w}^*\|^2_{\mathbf{I}^{(q)}_{0:k^*}} + N\gamma\|\mathbf{w}^*\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} \\
=& O\left(k^* + N\gamma\sum_{i>k^*}(\lambda_i + \epsilon_d)\right) \\
\le& O\left(\left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}} + N\gamma\left[\left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1-a}{a}} + \left(d - \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}}\right)\epsilon_d\right]\right) \\
\le& O\left(\epsilon_d N\gamma\left(d - \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}}\right) + \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}}\right),
\end{aligned}
\tag{F.26}
$$

where the first inequality holds by $\mathbf{w}^{(q)^*} = \mathbf{H}^{(q)^{-1}}\mathbf{H}\mathbf{w}^*$, the first equality holds by the constant level optimal solution assumption and the second inequality holds by (F.24) and (F.12).

We thirdly bound $\frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^{(q)*}\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-1}} + \|\mathbf{w}^{(q)*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}}$. Similar to (F.26),

$$
\begin{aligned}
&\frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^{(q)*}\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-1}} + \|\mathbf{w}^{(q)*}\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} \\
\leq& \frac{1}{\gamma^2 N^2} \cdot \|\mathbf{w}^*\|^2_{(\mathbf{H}^{(q)}_{0:k^*})^{-1}} + \|\mathbf{w}^*\|^2_{\mathbf{H}^{(q)}_{k^*:\infty}} \\
=& O\left( \frac{1}{\gamma^2 N^2} \sum_{i=1}^{k^*} \frac{1}{(\lambda_i + \epsilon_d)} + \sum_{i>k^*} (\lambda_i + \epsilon_d) \right) \\
\leq& O\left( \frac{\left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}}}{N\gamma} + \left( d - \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}} \right) \epsilon_d + \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1-a}{a}} \right) \\
\leq& O\left( \frac{\left( d - \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}} \right) \epsilon_d N\gamma + \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}}}{N\gamma} \right).
\end{aligned}
\tag{F.27}
$$

Finally, applying (F.25), (F.26) and (F.27) into Theorem F.2, we have

$$
\begin{aligned}
\mathbb{E}[\mathcal{E}(\overline{\mathbf{w}}_N)] \leq& O\left( \frac{\epsilon_d}{1+\epsilon_d} d + \epsilon_l \right) + O\left( \frac{\left( d - \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}} \right) \epsilon_d N\gamma + \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}}}{N\gamma} \right) \\
+& O\left( \frac{\epsilon_d^2 N^2 \gamma^2 \left[ d - \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}} \right] + \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}}}{N} \right) \left( \frac{\epsilon_o + \epsilon_a}{B} + (1 + d\epsilon_d)\alpha_B \epsilon_p + \frac{\sigma^2}{B} \right) \\
+& O\left( \frac{\epsilon_d^2 N^2 \gamma^2 \left[ d - \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}} \right] + \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}}}{N} \right) \frac{O\left( \epsilon_d N\gamma \left( d - \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}} \right) + \left(\frac{1}{N\gamma} - \epsilon_d\right)^{-\frac{1}{a}} \right)}{N\gamma}.
\end{aligned}
\tag{F.28}
$$

Consider (F.23) and (F.28) the proof is immediately completed.

$\square$

# G  Discussion of Assumptions

## G.1  Discussion of Assumption 3.3

Consider the standard fourth moment assumption on the full-precision data (Zou et al., 2023):

**Assumption G.1.** *Assume there exists a positive constant $\alpha_0 > 0$, such that for any PSD matrix $\mathbf{A}$, it holds that*

$$
\mathbb{E}\left[ \mathbf{x}\mathbf{x}^\top \mathbf{A} \mathbf{x}\mathbf{x}^\top \right] \preceq \alpha_0 \operatorname{tr}(\mathbf{H}\mathbf{A})\mathbf{H}.
$$

Under Assumption G.1, we are ready to verify if Assumption 3.3 can be satisfied. We begin by:

$$
\begin{aligned}
\mathbb{E}\left[\mathbf{x}^{(q)}\mathbf{x}^{(q)^{\top}}\mathbf{A}\mathbf{x}^{(q)}\mathbf{x}^{(q)^{\top}}\right] =&\mathbb{E}\left[\left(\mathbf{x}^{(q)^{\top}}\mathbf{A}\mathbf{x}^{(q)}\right)\mathbf{x}^{(q)}\mathbf{x}^{(q)^{\top}}\right]\\
\preceq&2\mathbb{E}\left[\left(\mathbf{x}^{(q)^{\top}}\mathbf{A}\mathbf{x}^{(q)}\right)\left(\mathbf{x}\mathbf{x}^{\top}+\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^{\top}}\right)\right]\\
\preceq&4\mathbb{E}\left[\left(\mathbf{x}^{\top}\mathbf{A}\mathbf{x}+\boldsymbol{\epsilon}^{(d)^{\top}}\mathbf{A}\boldsymbol{\epsilon}^{(d)}\right)\left(\mathbf{x}\mathbf{x}^{\top}+\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^{\top}}\right)\right]\\
=&4\mathbb{E}\left[\mathbf{x}\mathbf{x}^{\top}\mathbf{A}\mathbf{x}\mathbf{x}^{\top}\right]+4\mathbb{E}\left[\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^{\top}}\mathbf{A}\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^{\top}}\right]\\
&+4\mathbb{E}\left[\left(\mathbf{x}^{\top}\mathbf{A}\mathbf{x}\right)\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^{\top}}\right]+4\mathbb{E}\left[\left(\boldsymbol{\epsilon}^{(d)^{\top}}\mathbf{A}\boldsymbol{\epsilon}^{(d)}\right)\mathbf{x}\mathbf{x}^{\top}\right].
\end{aligned}
\tag{G.1}
$$

From Assumption G.1,

$$
\mathbb{E}\left[\mathbf{x}\mathbf{x}^{\top}\mathbf{A}\mathbf{x}\mathbf{x}^{\top}\right]\preceq\alpha_0\operatorname{tr}(\mathbf{H}\mathbf{A})\mathbf{H}.
\tag{G.2}
$$

Regarding $\mathbb{E}\left[\left(\boldsymbol{\epsilon}^{(d)^{\top}}\mathbf{A}\boldsymbol{\epsilon}^{(d)}\right)\mathbf{x}\mathbf{x}^{\top}\right]$,

$$
\mathbb{E}\left[\left(\boldsymbol{\epsilon}^{(d)^{\top}}\mathbf{A}\boldsymbol{\epsilon}^{(d)}\right)\mathbf{x}\mathbf{x}^{\top}\right]=\mathbb{E}\left[\mathbb{E}\left[\left(\boldsymbol{\epsilon}^{(d)^{\top}}\mathbf{A}\boldsymbol{\epsilon}^{(d)}\right)|\mathbf{x}\right]\mathbf{x}\mathbf{x}^{\top}\right]=\mathbb{E}\left[\operatorname{tr}\left(\mathbf{A}\mathbb{E}\left[\left(\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^{\top}}\right)|\mathbf{x}\right]\right)\mathbf{x}\mathbf{x}^{\top}\right].
\tag{G.3}
$$

Regarding $\mathbb{E}\left[\left(\mathbf{x}^{\top}\mathbf{A}\mathbf{x}\right)\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^{\top}}\right]$,

$$
\mathbb{E}\left[\left(\mathbf{x}^{\top}\mathbf{A}\mathbf{x}\right)\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^{\top}}\right]=\mathbb{E}\left[\operatorname{tr}\left(\mathbf{A}\mathbf{x}\mathbf{x}^{\top}\right)\mathbb{E}\left[\left(\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^{\top}}\right)|\mathbf{x}\right]\right].
\tag{G.4}
$$

Next, we provide some examples on specific quantization mechanism to exemplify the satisfaction of Assumption 3.3.

**Example G.1. (Strong Multiplicative Quantization)** *We consider a strong multiplicative quantization. In this case, there exists a constant $C'$ such that*

$$
\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^{\top}}\preceq C'\mathbf{x}\mathbf{x}^{\top}.
$$

Under Example G.1,

$$
\mathbb{E}\left[\left(\boldsymbol{\epsilon}^{(d)^{\top}}\mathbf{A}\boldsymbol{\epsilon}^{(d)}\right)\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^{\top}}\right]\preceq C'^2\mathbb{E}\left[\left(\mathbf{x}^{\top}\mathbf{A}\mathbf{x}\right)\mathbf{x}\mathbf{x}^{\top}\right],
\tag{G.5}
$$

and

$$
\mathbb{E}\left[\left(\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^{\top}}\right)|\mathbf{x}\right]\preceq C'\mathbf{x}\mathbf{x}^{\top}.
\tag{G.6}
$$

Therefore, together with (G.1), (G.2), (G.3), (G.4), (G.5) and (G.6), we have

$$
\mathbb{E}\left[\mathbf{x}^{(q)}\mathbf{x}^{(q)^{\top}}\mathbf{A}\mathbf{x}^{(q)}\mathbf{x}^{(q)^{\top}}\right]\leq4\alpha_0(1+2C'+C'^2)\operatorname{tr}(\mathbf{H}\mathbf{A})\mathbf{H}\leq4\alpha_0(1+2C'+C'^2)\operatorname{tr}(\mathbf{H}^{(q)}\mathbf{A})\mathbf{H}^{(q)}.
$$

That is, under strong multiplicative quantization Example G.1 and fourth moment Assumption G.1 on full-precision data, Assumption 3.3 is verified.

**Example G.2. (Strong Additive Quantization)** *We consider a strong additive quantization. In this case, there exist constants $C, C'$ and constant matrix $\mathbf{M}$ such that*

$$
\mathbf{D}:=\mathbb{E}\left[\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^{\top}}\right]=C\mathbf{M},\quad\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^{\top}}\preceq C'\mathbf{M}.
\tag{G.7}
$$

Under Example G.2,

$$\mathbb{E}\left[\left(\boldsymbol{\epsilon}^{(d)\top}\mathbf{A}\boldsymbol{\epsilon}^{(d)}\right)\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)\top}\right] \preceq C'\operatorname{tr}\left(\mathbf{A}\mathbf{M}\right)\mathbb{E}\left[\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)\top}\right] = \frac{C'}{C}\operatorname{tr}(\mathbf{A}\mathbf{D})\mathbf{D}. \tag{G.8}$$

Therefore, together with (G.1), (G.2), (G.3), (G.4), (G.7) and (G.8), we have

$$\mathbb{E}\left[\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\mathbf{A}\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right] \leq 4\alpha_0\operatorname{tr}(\mathbf{A}\mathbf{H})\mathbf{H} + 4\frac{C'}{C}\operatorname{tr}(\mathbf{A}\mathbf{D})\mathbf{H} + 4\frac{C'}{C}\operatorname{tr}(\mathbf{A}\mathbf{H})\mathbf{D} + 4\frac{C'}{C}\operatorname{tr}(\mathbf{A}\mathbf{D})\mathbf{D}$$
$$\leq 4\left[\alpha_0 + 3\frac{C'}{C}\right]\operatorname{tr}(\mathbf{H}^{(q)}\mathbf{A})\mathbf{H}^{(q)}.$$

That is, under strong additive quantization Example G.2 and fourth moment Assumption G.1 on full-precision data, Assumption 3.3 is verified.

## G.2  Discussion of Assumption 3.4

Consider the standard noise assumption on the full-precision data (Zou et al., 2023):

**Assumption G.2.** *There exists a constant $\sigma_0^2$ such that*

$$\mathbb{E}\left[(y - \langle \mathbf{w}^*, \mathbf{x}\rangle)^2\mathbf{x}\mathbf{x}^\top\right] \preceq \sigma_0^2\mathbf{H}.$$

Under Assumption G.2, we are ready to verify if Assumption 3.4 can be satisfied. We begin by:

$$
\begin{aligned}
&\mathbb{E}\left[(y^{(q)} - \langle \mathbf{w}^{(q)*}, \mathbf{x}^{(q)}\rangle)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right]\\
=&\mathbb{E}\left[(y^{(q)} - y + y - \langle \mathbf{w}^*, \mathbf{x}\rangle + \langle \mathbf{w}^*, \mathbf{x}\rangle - \langle \mathbf{w}^{(q)*}, \mathbf{x}^{(q)}\rangle)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right]\\
\preceq&3\mathbb{E}\left[(y^{(q)} - y)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right] + 3\mathbb{E}\left[(y - \langle \mathbf{w}^*, \mathbf{x}\rangle)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right]\\
&+3\mathbb{E}\left[(\langle \mathbf{w}^*, \mathbf{x}\rangle - \langle \mathbf{w}^{(q)*}, \mathbf{x}^{(q)}\rangle)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right]\\
\preceq&3\mathbb{E}\left[(y^{(q)} - y)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right] + 3\mathbb{E}\left[(y - \langle \mathbf{w}^*, \mathbf{x}\rangle)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right]\\
&+6\mathbb{E}\left[\langle \mathbf{w}^{(q)*} - \mathbf{w}^*, \mathbf{x}\rangle^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right] + 6\mathbb{E}\left[\langle \mathbf{w}^{(q)*}, \boldsymbol{\epsilon}^{(d)}\rangle^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right].
\end{aligned}
\tag{G.9}
$$

Next, we provide some examples on specific quantization mechanism to exemplify the satisfaction of Assumption 3.4.

**Example G.3. (Strong Multiplicative Quantization)** *In this case, we consider there exist constants $C', C''$ such that*

$$\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)\top} \preceq C'\mathbf{x}\mathbf{x}^\top,$$
$$\mathbb{E}[(y^{(q)} - y)^2|y] = C''y^2.$$

Regarding $\mathbb{E}\left[(y - \langle \mathbf{w}^*, \mathbf{x}\rangle)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right]$,

$$
\begin{aligned}
\mathbb{E}\left[(y - \langle \mathbf{w}^*, \mathbf{x}\rangle)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right] &\preceq 2\mathbb{E}\left[(y - \langle \mathbf{w}^*, \mathbf{x}\rangle)^2\mathbf{x}\mathbf{x}^\top\right] + 2\mathbb{E}\left[(y - \langle \mathbf{w}^*, \mathbf{x}\rangle)^2\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)\top}\right]\\
&\preceq 2(1 + C')\mathbb{E}\left[(y - \langle \mathbf{w}^*, \mathbf{x}\rangle)^2\mathbf{x}\mathbf{x}^\top\right]\\
&\preceq 2(1 + C')\sigma_0^2\mathbf{H},
\end{aligned}
\tag{G.10}
$$

where the second inequality holds by the definition of Example G.3 and the last inequality holds by Assumption G.2.

Regarding $\mathbb{E}\left[\langle\mathbf{w}^{(q)^*},\boldsymbol{\epsilon}^{(d)}\rangle^2\mathbf{x}^{(q)}\mathbf{x}^{(q)^\top}\right]$,

$$
\begin{aligned}
&\mathbb{E}\left[\langle\mathbf{w}^{(q)^*},\boldsymbol{\epsilon}^{(d)}\rangle^2\mathbf{x}^{(q)}\mathbf{x}^{(q)^\top}\right]\\
=&\mathbb{E}\left[\boldsymbol{\epsilon}^{(d)^\top}\mathbf{w}^{(q)^*}\mathbf{w}^{(q)^{*\top}}\boldsymbol{\epsilon}^{(d)}\mathbf{x}^{(q)}\mathbf{x}^{(q)^\top}\right]\\
\preceq&2\mathbb{E}\left[\boldsymbol{\epsilon}^{(d)^\top}\mathbf{w}^{(q)^*}\mathbf{w}^{(q)^{*\top}}\boldsymbol{\epsilon}^{(d)}\mathbf{x}\mathbf{x}^\top\right]+2\mathbb{E}\left[\boldsymbol{\epsilon}^{(d)^\top}\mathbf{w}^{(q)^*}\mathbf{w}^{(q)^{*\top}}\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^\top}\right]\\
\preceq&2C'\alpha_0\mathrm{tr}(\mathbf{w}^{(q)^*}\mathbf{w}^{(q)^{*\top}}\mathbf{H})\mathbf{H}+2C'^2\alpha_0\mathrm{tr}(\mathbf{w}^{(q)^*}\mathbf{w}^{(q)^{*\top}}\mathbf{H})\mathbf{H},
\end{aligned}
\tag{G.11}
$$

where the last inequality holds by the definition of Example G.3 and Assumption G.1.

Regarding $\mathbb{E}\left[\langle\mathbf{w}^{(q)^*}-\mathbf{w}^*,\mathbf{x}\rangle^2\mathbf{x}^{(q)}\mathbf{x}^{(q)^\top}\right]$,

$$
\begin{aligned}
\mathbb{E}\left[\langle\mathbf{w}^{(q)^*}-\mathbf{w}^*,\mathbf{x}\rangle^2\mathbf{x}^{(q)}\mathbf{x}^{(q)^\top}\right]\preceq&2\mathbb{E}\left[\langle\mathbf{w}^{(q)^*}-\mathbf{w}^*,\mathbf{x}\rangle^2\mathbf{x}\mathbf{x}^\top\right]+2\mathbb{E}\left[\langle\mathbf{w}^{(q)^*}-\mathbf{w}^*,\mathbf{x}\rangle^2\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^\top}\right]\\
\preceq&2(1+C')\alpha_0\mathrm{tr}\left((\mathbf{w}^{(q)^*}-\mathbf{w}^*)(\mathbf{w}^{(q)^*}-\mathbf{w}^*)^\top\mathbf{H}\right)\mathbf{H},
\end{aligned}
\tag{G.12}
$$

where the last inequality holds by the definition of Example G.3 and Assumption G.1.

Regarding $\mathbb{E}\left[(y^{(q)}-y)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)^\top}\right]$, if we further assume that there exists a constant $C'''$ such that $\mathbb{E}\left[y^2\mathbf{x}\mathbf{x}^\top\right]\preceq C'''\mathbf{H}$, then

$$
\begin{aligned}
\mathbb{E}\left[(y^{(q)}-y)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)^\top}\right]\preceq&2\mathbb{E}\left[(y^{(q)}-y)^2\mathbf{x}\mathbf{x}^\top\right]+2\mathbb{E}\left[(y^{(q)}-y)^2\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)^\top}\right]\\
\preceq&2(1+C')\mathbb{E}\left[(y^{(q)}-y)^2\mathbf{x}\mathbf{x}^\top\right]\\
\preceq&2(1+C')C''\mathbb{E}[y^2\mathbf{x}\mathbf{x}^\top]\\
\preceq&2(1+C')C''C'''\mathbf{H},
\end{aligned}
\tag{G.13}
$$

where the second and third inequality hold by the definition of Example G.3.

Therefore, together with (G.9), (G.10), (G.11), (G.12) and (G.13), we have

$$
\begin{aligned}
&\mathbb{E}\left[(y^{(q)}-\langle\mathbf{w}^{(q)^*},\mathbf{x}^{(q)}\rangle)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)^\top}\right]\\
\preceq&3\mathbb{E}\left[(y^{(q)}-y)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)^\top}\right]+3\mathbb{E}\left[(y-\langle\mathbf{w}^*,\mathbf{x}\rangle)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)^\top}\right]\\
&+6\mathbb{E}\left[\langle\mathbf{w}^{(q)^*}-\mathbf{w}^*,\mathbf{x}\rangle^2\mathbf{x}^{(q)}\mathbf{x}^{(q)^\top}\right]+6\mathbb{E}\left[\langle\mathbf{w}^{(q)^*},\boldsymbol{\epsilon}^{(d)}\rangle^2\mathbf{x}^{(q)}\mathbf{x}^{(q)^\top}\right]\\
\preceq&6(1+C')C''C'''\mathbf{H}+6(1+C')\sigma_0^2\mathbf{H}+12C'\alpha_0(1+C')\mathrm{tr}(\mathbf{w}^{(q)^*}\mathbf{w}^{(q)^{*\top}}\mathbf{H})\mathbf{H}\\
&+12(1+C')\alpha_0\mathrm{tr}\left((\mathbf{w}^{(q)^*}-\mathbf{w}^*)(\mathbf{w}^{(q)^*}-\mathbf{w}^*)^\top\mathbf{H}\right)\mathbf{H}\\
\preceq&\sigma^2\mathbf{H}^{(q)},
\end{aligned}
$$

where

$$
\sigma^2=6(1+C')(C''C'''+\sigma_0^2)+12C'\alpha_0(1+C')\|\mathbf{w}^{(q)^*}\|_{\mathbf{H}}^2+12(1+C')\alpha_0\|\mathbf{w}^{(q)^*}-\mathbf{w}^*\|_{\mathbf{H}}^2.
$$

That is, under strong multiplicative quantization Example G.3 and fourth moment Assumption G.2 on full-precision data, Assumption 3.4 is verified.

**Example G.4. (Strong Additive Quantization)** *In this case, we consider there exist constants $C, C', C''$ and constant matrix $\mathbf{M}$ such that*

$$\mathbf{D} := \mathbb{E}\left[\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)\top}\right] = C\mathbf{M}, \ \boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)\top} \preceq C'\mathbf{M},$$
$$\mathbb{E}[(y^{(q)} - y)^2|y] \leq C''.$$

Regarding $\mathbb{E}\left[(y - \langle\mathbf{w}^*, \mathbf{x}\rangle)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right]$, if we further assume that $\mathbb{E}\left[(y - \langle\mathbf{w}^*, \mathbf{x}\rangle)^2\right] \leq \sigma_0^2$, then

$$
\begin{aligned}
\mathbb{E}\left[(y - \langle\mathbf{w}^*, \mathbf{x}\rangle)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right] &\preceq 2\mathbb{E}\left[(y - \langle\mathbf{w}^*, \mathbf{x}\rangle)^2\mathbf{x}\mathbf{x}^\top\right] + 2\mathbb{E}\left[(y - \langle\mathbf{w}^*, \mathbf{x}\rangle)^2\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)\top}\right] \\
&\preceq 2\sigma_0^2\mathbf{H} + 2C'\sigma_0^2\mathbf{M} \\
&= 2\sigma_0^2\mathbf{H} + 2\frac{C'}{C}\sigma_0^2\mathbf{D},
\end{aligned}
\tag{G.14}
$$

where the last inequality and equality hold by the definition of Example G.4 and Assumption G.2.

Regarding $\mathbb{E}\left[(y^{(q)} - y)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right]$, by the definition of Example G.4,

$$\mathbb{E}\left[(y^{(q)} - y)^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right] \leq C''\mathbf{H}^{(q)}. \tag{G.15}$$

Regarding $\mathbb{E}\left[\langle\mathbf{w}^{(q)*}, \boldsymbol{\epsilon}^{(d)}\rangle^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right]$,

$$
\begin{aligned}
&\mathbb{E}\left[\langle\mathbf{w}^{(q)*}, \boldsymbol{\epsilon}^{(d)}\rangle^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right] \\
={}&\mathbb{E}\left[\boldsymbol{\epsilon}^{(d)\top}\mathbf{w}^{(q)*}\mathbf{w}^{(q)*\top}\boldsymbol{\epsilon}^{(d)}\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right] \\
\preceq{}&2\mathbb{E}\left[\boldsymbol{\epsilon}^{(d)\top}\mathbf{w}^{(q)*}\mathbf{w}^{(q)*\top}\boldsymbol{\epsilon}^{(d)}\mathbf{x}\mathbf{x}^\top\right] + 2\mathbb{E}\left[\boldsymbol{\epsilon}^{(d)\top}\mathbf{w}^{(q)*}\mathbf{w}^{(q)*\top}\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)\top}\right] \\
\preceq{}&2C'\mathrm{tr}(\mathbf{w}^{(q)*}\mathbf{w}^{(q)*\top}\mathbf{M})\mathbf{H} + 2C'\mathrm{tr}(\mathbf{w}^{(q)*}\mathbf{w}^{(q)*\top}\mathbf{M})\mathbf{D},
\end{aligned}
\tag{G.16}
$$

where the last inequality holds by the definition of Example G.4.

Regarding $\mathbb{E}\left[\langle\mathbf{w}^{(q)*} - \mathbf{w}^*, \mathbf{x}\rangle^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right]$,

$$
\begin{aligned}
&\mathbb{E}\left[\langle\mathbf{w}^{(q)*} - \mathbf{w}^*, \mathbf{x}\rangle^2\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right] \\
={}&\mathbb{E}\left[\mathbf{x}^\top(\mathbf{w}^{(q)*} - \mathbf{w}^*)(\mathbf{w}^{(q)*} - \mathbf{w}^*)^\top\mathbf{x}\mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right] \\
\preceq{}&2\mathbb{E}\left[\mathbf{x}^\top(\mathbf{w}^{(q)*} - \mathbf{w}^*)(\mathbf{w}^{(q)*} - \mathbf{w}^*)^\top\mathbf{x}\mathbf{x}\mathbf{x}^\top\right] + 2\mathbb{E}\left[\mathbf{x}^\top(\mathbf{w}^{(q)*} - \mathbf{w}^*)(\mathbf{w}^{(q)*} - \mathbf{w}^*)^\top\mathbf{x}\boldsymbol{\epsilon}^{(d)}\boldsymbol{\epsilon}^{(d)\top}\right] \\
\preceq{}&2\alpha_0\mathrm{tr}\left((\mathbf{w}^{(q)*} - \mathbf{w}^*)(\mathbf{w}^{(q)*} - \mathbf{w}^*)^\top\mathbf{H}\right)\mathbf{H} + 2\frac{C'}{C}\mathrm{tr}\left((\mathbf{w}^{(q)*} - \mathbf{w}^*)(\mathbf{w}^{(q)*} - \mathbf{w}^*)^\top\mathbf{H}\right)\mathbf{D},
\end{aligned}
\tag{G.17}
$$

where the last inequality holds by the definition of Example G.4 and Assumption G.1.

Therefore, together with (G.9), (G.14), (G.15), (G.16) and (G.17), we have

$$\mathbb{E}\left[(y^{(q)} - \langle \mathbf{w}^{(q)*}, \mathbf{x}^{(q)}\rangle)^2 \mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right]$$

$$\preceq 3\mathbb{E}\left[(y^{(q)} - y)^2 \mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right] + 3\mathbb{E}\left[(y - \langle \mathbf{w}^*, \mathbf{x}\rangle)^2 \mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right]$$

$$+ 6\mathbb{E}\left[\langle \mathbf{w}^{(q)*} - \mathbf{w}^*, \mathbf{x}\rangle^2 \mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right] + 6\mathbb{E}\left[\langle \mathbf{w}^{(q)*}, \boldsymbol{\epsilon}^{(d)}\rangle^2 \mathbf{x}^{(q)}\mathbf{x}^{(q)\top}\right]$$

$$\preceq 3C''\mathbf{H}^{(q)} + 6\sigma_0^2 \mathbf{H} + 6\frac{C'}{C}\sigma_0^2 \mathbf{D} + 12C'\mathrm{tr}(\mathbf{w}^{(q)*}\mathbf{w}^{(q)*\top}\mathbf{M})\mathbf{H}^{(q)}$$

$$+ (2\alpha_0 + 2\frac{C'}{C})\mathrm{tr}\left((\mathbf{w}^{(q)*} - \mathbf{w}^*)(\mathbf{w}^{(q)*} - \mathbf{w}^*)^\top \mathbf{H}\right)\mathbf{H}^{(q)}$$

$$\preceq \sigma^2 \mathbf{H}^{(q)},$$

where

$$\sigma^2 = 3C'' + 6\left(\frac{C + C'}{C} + 12C'\mathrm{tr}(\mathbf{w}^{(q)*}\mathbf{w}^{(q)*\top}\mathbf{M})\right)\sigma_0^2 + 2(\alpha_0 + \frac{C'}{C})\mathrm{tr}\left((\mathbf{w}^{(q)*} - \mathbf{w}^*)(\mathbf{w}^{(q)*} - \mathbf{w}^*)^\top \mathbf{H}\right)$$

$$= 3C'' + 6\left(\frac{C + C'}{C} + 12C'\|\mathbf{w}^{(q)*}\|_\mathbf{M}^2\right)\sigma_0^2 + 2(\alpha_0 + \frac{C'}{C})\|\mathbf{w}^{(q)*} - \mathbf{w}^*\|_\mathbf{H}^2.$$

That is, under strong additive quantization Example G.4 and noise Assumption G.2 on full-precision data, Assumption 3.4 is verified.