# THE IMPACT OF THE MINI-BATCH SIZE ON THE DYNAMICS OF SGD: VARIANCE AND BEYOND*

ANONYMOUS AUTHORS

**Abstract.** We study mini-batch stochastic gradient descent (SGD) dynamics under linear regression and deep polynomially-activated networks by focusing on the variance of the gradients only given the initial weights and mini-batch size, which is the first study of this nature. In both cases, we provide recursive relationships of the norm of the gradients and weight matrices between consecutive time steps. We further show that, in each iteration, the norm of the gradient is a polynomial in the reciprocal of the mini-batch size and a decreasing function of the mini-batch size. The results theoretically back the important intuition that smaller batch sizes yield larger variance of the stochastic gradients and lower loss function values which is a common believe among the researchers. The proof techniques exhibit explicit relationships between a variety of general functions of stochastic gradient estimators and initial weights, which is useful for further research on the dynamics of SGD. We empirically provide insights to our results on various datasets and commonly used deep network structures. We further discuss possible extensions of the approaches we build in studying the generalization ability of the deep learning models.

**Key words.** Stochastic Gradient Descent, Polynomially-activated Neural Networks

**AMS subject classifications.** 68Q25, 68R10, 68U05

**1. Introduction.** Deep learning models have achieved great success in a variety of tasks including natural language processing, computer vision, and reinforcement learning [9]. Despite their practical success, there are only limited studies of the theoretical properties of deep learning; see survey papers [39, 8] and references therein. The general problem underlying deep learning models is to optimize (minimize) a loss function, defined by the deviation of model predictions on data samples from the corresponding true labels. The prevailing method to train deep learning models is the mini-batch stochastic gradient descent algorithm and its variants [4, 5]. SGD updates model parameters by calculating a stochastic approximation of the full gradient of the loss function, based on a random selected subset of the training samples called a mini-batch.

Although SGD can converge to the minimum of a convex function [6], deep neural networks are strongly non-convex. Thus, the success of SGD in neural network training, especially the dynamics of SGD, becomes an interesting question. Some researchers approximate the dynamics of SGD by a continuous-time dynamic system [26, 25, 28, 17]. Another line of research [27, 7, 2] show that the dynamics of SGD in training over-parameterized neural networks are similar to training a linear model. However, these statements are approximate in nature and do not provide explicit formulas for calculating any specific quantities during SGD training. The mini-batch size is also a key factor deciding the dynamics of SGD. Some research focuses on how to choose an optimal mini-batch size based on different criteria [38, 11]. However, these works make strong assumptions on the loss function properties (strong or point or quasi convexity, or constant variance near stationary points) or about the formulation of the SGD algorithm (continuous time interpretation by means of differential equations). The theoretical results regarding the relationship between the mini-batch size and the variance (and other performances, like loss and generalization ability) of the SGD algorithm applied to general machine learning models are still missing.

Besides, it is well-accepted that selecting a large mini-batch size reduces the

---

*Submitted to the editors on DATE.

training time of deep learning models, as computation on large mini-batches can be better parallelized on processing units. For example, Goyal et al. [12] scale ResNet-50 [13] from a mini-batch size of 256 images and training time of 29 hours, to a larger mini-batch size of 8,192 images. Their training achieves the same level of accuracy while reducing the training time to one hour. However, noted by many researchers, larger mini-batch sizes suffer from a worse generalization ability [22, 19]. Therefore, many efforts have been made to develop specialized training procedures that achieve good generalization using large mini-batch sizes [16, 12]. Smaller batch sizes have the advantage of allegedly offering better generalization (at the expense of a higher training time). We hypothesize that, given the same initial point, smaller sizes lead to lower training loss and, unfortunately, decrease stability of the algorithm on average. The latter follows from the fact that the smaller is the batch size, more stochasticity and volatility is introduced. After all, if the batch size equals to the number of samples, there is no stochasticity in the algorithm. To this end, we conjecture that the variance of the gradient in each iteration is a decreasing function of the mini-batch size. We partially prove this conjecture in this work.

In this paper, we study the dynamics of SGD by representing related quantities only using the mini-batch size, initial points and learning rates, which are available before training. This is different from previous literature which analyzes SGD by focusing on one-step properties. In fact, the dynamics of SGD are not comparable if we merely consider the one-step behavior, as the model parameters change iteration by iteration. We are able to build general frameworks in the convex linear regression case and in a deep polynomially-activated neural network setting. The frameworks provide explicit and recursive relationships of general forms, which cover many interesting quantities regarding the dynamics of SGD.

As an application of our frameworks, we are able to prove the hypothesis about variance in the convex linear regression case and to show significant progress in a deep polynomially-activated neural network setting. We show that the variance is a polynomial in the reciprocal of the mini-batch size and that it is decreasing if the mini-batch size is larger than a threshold (further experiments reveal that this threshold can be as small as 1). The increased variance as the mini-batch size decreases should also intuitively imply convergence to lower training loss values and in turn better prediction and generalization ability (these relationships are yet to be confirmed analytically; but we provide empirical evidence to their validity).

The major contributions of this paper are as follows.

(i) For linear regression, we build a framework to recursively calculate the norm of any linear combination of sample-wise gradients between consecutive iterations (Theorem 3.2). This recursive relationship can be used to calculate any quantity related to the full or stochastic gradient or loss at any iteration with respect to the initial weights. As an application of this framework, we show that in each iteration the norm of any linear combination of sample-wise gradients can be computed by a polynomial in the reciprocal of the mini-batch size $b$ and is a decreasing function of $b$ (Theorem 3.3). As a special case, the variance of the stochastic gradient estimator and the full gradient at the iterate in step $t$ are also decreasing functions of $b$ at any iteration step $t$ (Theorem 3.4 and Corollary 3.5).

(ii) For a deep polynomially-activated neural network under a teacher-student network setting, we provide a framework for recursively calculating the trace of any product of the stochastic gradient estimators, weight matrices and other constant matrices at time step $t$ by using the variables at time step $t-1$ (Theorems 3.6 and 3.7). This explicit relationship can be used to derive the expected value of the product

of the weight matrices and stochastic gradient estimators as a polynomial in $1/b$ with coefficients a sum of products of the initial weights (Theorem 3.8). As a special case, the variance of the stochastic gradient estimator is a polynomial in $1/b$ without the constant term (Theorem 3.9) and therefore it is a decreasing function of $b$ when $b$ is large enough (Theorem 3.10). The results and proof techniques can be extended in an approximate sense to deep networks with general non-linear activation functions (Section 3.3). As a comparison, other papers that study theoretical properties of two-layer networks either fix one layer of the network, or assume the over-parameterized property of the model and they study convergence, while our paper makes no such assumptions on the model capacity. The proof also reveals the structure of the coefficients of the polynomial, and thus it serves as a tool for future work on proving other properties of the stochastic gradient estimators and weight matrices.

(iii) The proofs are involved and require several key ideas. The main one is to show a more general result than it is necessary in order to carry out the induction on time step $t$. New concepts and definitions are introduced in order to handle the more general case. Along the way we show a result of general interest establishing expectation of the product of quadratic terms of samples with general distribution intertwined with constant matrices.

(iv) We verify the theoretical results regarding the decreasing property of variance on various datasets and provide a further understanding. We also empirically show that the results extend to other widely used network structures and hold for all choices of the mini-batch sizes. We also empirically verify that, on average, in each iteration the loss function value and the generalization ability (measured by the gap between accuracy on the training and test sets) are all decreasing functions of the mini-batch size.

In conclusion, we study the dynamics of SGD under linear regression and a multi-layer polynomially-activated network setting by building frameworks that can recursively and explicitly calculate general products and sums of the stochastic gradient estimators and weights matrices between consecutive iterations. As an application of the frameworks, we focus on representing the variance of the stochastic gradient estimators by the mini-batch size, initial weights and other constant variables, and therefore prove the decreasing property of the variance of the stochastic gradient estimators. The proof techniques can also be used to derive other properties of the SGD dynamics in regard to the mini-batch size and initial weights. To the best of authors' knowledge, the work is the first one to theoretically and explicitly study the important quantities of SGD at iteration $t$ only using the initial weights and mini-batch size, under mild assumptions on the network and the loss function. We support our theoretical results by experiments. We further experiment on other state-of-the-art deep learning models and datasets to empirically show the validity of the conjectures about the impact of mini-batch size on average loss, average accuracy and the generalization ability of a model.

The rest of the manuscript is structured as follows. In Section 2 we review the literature while in Section 3 we present a general framework on how to recursively represent some functions of the stochastic gradient estimators by initial weights, under different models including linear regression, deep polynomially-activated networks, and general neural networks. We also provide applications of the presented framework in Section 3. Section 4 presents the experiments that verify our theorems and provide further insights into the impact of the mini-batch sizes on SGD dynamics. The proofs of the theorems and other technical details are available in Appendix A.

**2. Literature Review.** Stochastic gradient descent type methods are broadly used in machine learning [3, 21, 5]. The performance of SGD highly relies on the choice of the mini-batch size. It has been widely observed that choosing a large mini-batch size to train deep neural networks appears to deteriorate generalization [22]. This phenomenon exists even if the models are trained without any budget or limits, until the loss function value ceases to improve [19]. One explanation for this phenomenon is that large mini-batch SGD produces "sharp" minima that generalize worse [15, 19]. Specialized training procedures to achieve good performance with large mini-batch sizes have also been proposed [16, 12].

It is well-known that SGD has a slow asymptotic rate of convergence due to its inherent variance [18]. Variants of SGD that can reduce the variance of the stochastic gradient estimator, which yield faster convergence, have also been suggested. The use of the information of full gradients to provide variance control for stochastic gradients is addressed in [18, 34, 36]. The works in [23, 24, 35] further improve the efficiency and complexity of the algorithm by carefully controling the variance.

There is prior work focusing on studying the dynamics of SGD. Neelakantan et al. propose to add isotropic white noise to the full gradient to study the "structured" variance [31]. The works in [25, 28, 17] connect SGD with stochastic differential equations to explain the property of converged minima and generalization ability of the model. Smith et al. propose an "optimal" mini-batch size which maximizes the test set accuracy by a Bayesian approach [38]. The Stochastic Gradient Langevin Dynamics (SGLD, a variant of SGD) algorithm for non-convex optimization is studied in [43, 30].

In most of the prior work about the convergence of SGD, it is assumed that the variance of stochastic gradient estimators is upper-bounded by a linear function of the norm of the full gradient, e.g. Assumption 4.3 in [5]. Gower et al. [11] give more precise bounds of the variance under different sampling methods and Khaled et al. [20] extend them to smooth non-convex regime. These bounds are still dependent on the model parameters at the corresponding iteration. To the best of the authors' knowledge, there is no existing result which represents stochastic gradient estimators only using the initial weights and the mini-batch size. This paper partially solves this problem.

**3. Analysis.** Mini-batch SGD is a lighter-weight version of gradient descent. Suppose that we are given a loss function $L(w)$ where $w$ is the collection (vector, matrix, or tensor) of all model parameters. At each iteration $t$, instead of computing the full gradient $\nabla_w L(w_t)$, SGD randomly samples a mini-batch set $\mathcal{B}_t$ that consists of $b = |\mathcal{B}_t|$ training instances and sets $w_{t+1} \leftarrow w_t - \alpha_t \nabla_w L_{\mathcal{B}_t}(w_t)$, where the positive scalar $\alpha_t$ is the learning rate (or step size) and $\nabla_w L_{\mathcal{B}_t}(w_t)$ denotes the stochastic gradient estimator based on mini-batch $\mathcal{B}_t$.

An important property of the stochastic gradient estimator $\nabla_w L_{\mathcal{B}_t}(w_t)$ is that it is an unbiased estimator, i.e. $\mathbb{E}\nabla_w L_{\mathcal{B}_t}(w_t) = \nabla_w L(w_t)$, where the expectation is taken over all possible choices of mini-batch $\mathcal{B}_t$. However, it is unclear what is the value of[1]

$$\text{var}\left(\nabla_w L_{\mathcal{B}_t}(w_t)\right) := \mathbb{E}\left\|\nabla_w L_{\mathcal{B}_t}(w_t)\right\|^2 - \left\|\mathbb{E}\nabla_w L_{\mathcal{B}_t}(w_t)\right\|^2.$$

Intuitively, we should have $\text{var}\left(\nabla_w L_{\mathcal{B}_t}(w_t)\right) \propto \frac{n^2}{b}\text{var}\left(\nabla_w L(w_t)\right)$, where $n$ is the number of training samples and stochasticity on the right-hand side comes from

---

[1]Note that this definition is different from the variance of a vector, i.e., the covariance matrix. This "scalar" variance is a common practice in the field of optimization (e.g. equation (4.6) in [5]).

mini-batch samples behind $w_t$ [38, 11]. However, even the quantities $\nabla_w L(w_t)$ and $\mathsf{var}\left(\nabla_w L(w_t)\right)$ are still challenging to compute as we do not have direct formulas of their precise values. Besides, as we choose different $b$'s, their values are not comparable as we end up with different $w_t$'s.

A plausible idea to address these issues is to represent $\mathbb{E}\nabla_w L_{\mathcal{B}_t}(w_t)$ and $\mathsf{var}(\nabla_w L_{\mathcal{B}_t}(w_t))$ only using the fixed and known quantities $w_0, b, t$, and $\alpha_t$. In this way, we can further discover the properties, like decreasing with respect to $b$, of $\mathbb{E}\nabla_w L_{\mathcal{B}_t}(w_t)$ and $\mathsf{var}\left(\nabla_w L_{\mathcal{B}_t}(w_t)\right)$. The biggest challenge is how to connect the quantities in iteration $t$ with those of iteration 0. This is similar to discovering the properties of a stochastic differential equation at time $t$ given only the dynamics of the stochastic differential equation and the initial point.

In this section, we address these questions by recursively representing some general forms of stochastic gradient estimators under two settings: linear regression and a deep polynomially-activated network. In Section 3.1 in a linear regression setting, we provide explicit formulas for calculating any norm of the linear combination of sample-wise gradients at time step $t$. As an application of the presented recursive relationships, we therefore show that the $\mathsf{var}\left(\nabla_w L_{\mathcal{B}_t}(w_t)\right)$ is a decreasing function of the mini-batch size $b$. In Section 3.2, under a deep polynomially-activated network with teacher-student setting, we provide explicit formulas for calculating any trace of the mixed product of weight matrices and stochastic gradient estimators. With this tool, we further show that these traces are polynomials in $1/b$ with finite degree and that $\mathsf{var}\left(\nabla_w L_{\mathcal{B}_t}(w_t)\right)$ is a decreasing function of the mini-batch size $b > b_0$ for some constant $b_0$. In Section 3.3, we extend the results to general deep neural networks with mild assumptions on the activation functions in an approximate sense.

For a random matrix $M$, we define $\mathsf{var}\left(M\right) := \mathbb{E}\left\|\mathrm{vec}(M)\right\|^2 - \left\|\mathbb{E}\mathrm{vec}(M)\right\|^2$ where $\mathrm{vec}(M)$ denotes the vectorization of matrix $M$. We denote $[m:n] := \{m, m+1, \ldots, n\}$ if $m \leqslant n$, and $\varnothing$ otherwise. We use $[n] := [1:n]$ as an abbreviation. For clarity, we use the superscript $b$ to distinguish the variables with different choices of the mini-batch size $b$. In each iteration $t$, we use $\mathcal{B}_t^b$ to denote the batch of samples (or sample indices) to calculate the stochastic gradient. We denote by $\mathcal{F}_t^b$ the filtration of information before calculating the stochastic gradient in the $t$-th iteration, i.e. $\mathcal{F}_t^b := \left\{w_0, w_1^b, \ldots, w_t^b, \mathcal{B}_0^b, \ldots, \mathcal{B}_{t-1}^b\right\}$. We use $\bigotimes_{i\in[n]} A_i$ to denote the Kronecker product of matrices $A_1, \ldots, A_n$.-

**3.1. Linear Regression.** In this subsection, we discuss the dynamics of SGD applied in linear regression. Given data points $(x_1, y_1), \cdots, (x_n, y_n)$, where $x_i \in \mathbb{R}^p$ and $y_i \in \mathbb{R}$, we define the loss function to be

$$L(w) = \frac{1}{n}\sum_{i=1}^n L_i(w) = \frac{1}{n}\sum_{i=1}^n \frac{1}{2}\left(w^T x_i - y_i\right)^2,$$

where $w \in \mathbb{R}^p$ are the model parameters. We consider minimizing $L(w)$ by mini-batch SGD. Note that the bias term in the general linear regression models is omitted, however, adding the bias term does not change the result of this section. Formally, we first choose a mini-batch size $b$ and initial weights $w_0$. In each iteration $t$, we sample $\mathcal{B}_t^b$, a subset of $[n]$ with cardinality $b$, and update the parameters by $w_{t+1}^b = w_t^b - \alpha_t g_t^b$, where $g_t^b = \frac{1}{b}\sum_{i\in\mathcal{B}_t^b} \nabla L_i\left(w_t^b\right)$.

We first show the relationship between the variance of stochastic gradient $g_t^b$ and the full gradient $\nabla L\left(w_t^b\right)$ and sample-wise gradient $\nabla L_i\left(w_t^b\right), i \in [n]$, derived by considering all possible choices of the mini-batch $\mathcal{B}_t^b$. Readers should note that Lemma 3.1 actually holds for all models with $L_2$-loss, not merely linear regression (since in

the proof we do not need to know the explicit form of $L_i(w)$).

LEMMA 3.1. *Let $c_b := \frac{n-b}{b(n-1)} \geqslant 0$. For any matrix $A \in \mathbb{R}^{p \times p}$ we have*

$$\mathsf{var}\left(Ag_t^b \middle| \mathcal{F}_t^b\right) = \mathbb{E}\left[\left\|Ag_t^b\right\|^2 \middle| \mathcal{F}_t^b\right] - \left\|A\nabla L\left(w_t^b\right)\right\|^2 = c_b\left(\frac{1}{n}\sum_{i=1}^n \left\|A\nabla L_i\left(w_t^b\right)\right\|^2 - \left\|A\nabla L\left(w_t^b\right)\right\|^2\right).$$

Lemma 3.1 provides a bridge to connect the norm and variance of $g_t^b$ with sample-wise gradients $\nabla L_i\left(w_t^b\right), i \in [n]$. Therefore, if we can further discover the properties of $\nabla L_i\left(w_t^b\right), i \in [n]$, we are able to calculate the variance of $g_t^b$. Theorem 3.2 addresses this problem by showing the relationship between any linear combination of $\nabla L_i\left(w_t^b\right)$'s and $\nabla L_i\left(w_{t-1}^b\right)$'s.

THEOREM 3.2. *For any set of square matrices $\{A_1, \cdots, A_n\} \in \mathbb{R}^{p \times p}$, if we denote $A = \sum_{i=1}^n A_i x_i x_i^T$, then we have*

$$\mathbb{E}\left[\left\|\sum_{i=1}^n A_i \nabla L_i\left(w_{t+1}^b\right)\right\|^2 \middle| \mathcal{F}_0\right] = \mathbb{E}\left[\left\|\sum_{i=1}^n B_i \nabla L_i\left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right] + \frac{\alpha_t^2 c_b}{n^2}\sum_{k=1}^n\sum_{l=1}^n \mathbb{E}\left[\left\|\sum_{i=1}^n B_i^{kl} \nabla L_i\left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right],$$

*where $B_i = A_i - \frac{\alpha_t}{n}A$; $B_i^{kl} = A$ if $i = k, i \neq l$, $B_i^{kl} = A$ if $i = l, i \neq k$, and $B_i^{kl}$ equals the zero matrix, otherwise.*

Theorem 3.2 provides an explicit relationship between the norm of any linear combinations of the sample-wise gradients at time steps $t+1$ and $t$. Therefore, we can easily use it to recursively calculate this norm for all iterations $t$. As an application of this theorem, note that $c_b$ is a decreasing function of $b$, and thus we are able to show Theorem 3.3.

THEOREM 3.3. *For any non-negative integer $t$ and any matrices $A_i \in \mathbb{R}^{p \times p}, i \in [n]$, $\mathbb{E}\left[\left\|\sum_{i=1}^n A_i \nabla L_i\left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right]$ is a decreasing function of $b$ for $b \in [n]$.*

Theorem 3.3 states that the norm of any linear combinations of the sample-wise gradients is a decreasing function of $b$. Combining Lemma 3.1 which connects the variance of $g_t^b$ with the linear combination of $\nabla L_i\left(w_t^b\right)$'s, and the fact that $\nabla L\left(w_t^b\right) = \frac{1}{n}\sum_{i=1}^n \nabla L_i\left(w_t^b\right)$, we have Theorem 3.4.

THEOREM 3.4. *Fixing initial weights $w_0$, the two quantities $\mathsf{var}\left(Bg_t^b \middle| \mathcal{F}_0\right)$ and $\mathsf{var}\left(B\nabla L\left(w_t^b\right) \middle| \mathcal{F}_0\right)$ are both decreasing functions of mini-batch size $b$ for all $b \in [n]$, $t \in \mathbb{N}$, and all square matrices $B \in \mathbb{R}^{p \times p}$.*

As a special case, Corollary 3.5 guarantees that the variance of the stochastic gradient estimator is a decreasing function of $b$.

COROLLARY 3.5. *Fixing initial weights $w_0$, the two quantities $\mathsf{var}\left(g_t^b \middle| \mathcal{F}_0\right)$ and $\mathsf{var}\left(\nabla L\left(w_t^b\right) \middle| \mathcal{F}_0\right)$ are both decreasing functions of mini-batch size $b$ for all $b \in [n]$ and $t \in \mathbb{N}$.*

In conclusion, we provide a framework for calculating the explicit value of variance of the stochastic gradient estimators and the norm of any linear combination of sample-wise gradients. In fact, the presented theorems can be applied to a variety of terms, like the total loss $L(w_t^b)$, as long as it is a polynomial of degree of 2 with respect to $w_t^b$. Theorem 3.2 can be further modified to hold for higher orders of $w_t^b$ in a similar manner.

As an application of the framework, we show that the variance of the full gradient and the stochastic gradient estimators are both decreasing functions of $b$. Readers should note that the framework here is not limited to showing the decreasing property

of the variance, but can also be used in many other circumstance. For example, we can use Theorem 3.2 to induct on $t$ and easily show that $\mathbb{E}\left[\left\|\sum_{i=1}^n A_i \nabla L_i\left(w_t^b\right)\right\|^2 \Big| \mathcal{F}_0\right]$ is a polynomial of $\frac{1}{b}$ with degree at most $t$ and calculate the coefficients therein.

**3.2. Deep Networks with Polynomial Activation Functions.** In this section, we investigate the dynamics of SGD on deep networks utilizing a polynomial activation function. We present the informal theorems in this section and reserve the complete versions for the Appendix. Additionally, we provide a comprehensive proof of the two-layer linear network (which corresponds to a polynomial activation of degree one) in the Appendix, along with the necessary additions to extend the proof to the multi-layer polynomial case.

Given a distribution $\mathcal{D}$ in $\mathbb{R}^p$, we consider the population loss

(3.1)
$$\mathcal{L}(w) = \mathbb{E}_{x \sim \mathcal{D}}\left[\frac{1}{2}\left\|W_H \sigma\left(W_{H-1}\sigma\left(\cdots\sigma\left(W_1 x\right)\right)\right) - W_H^*\sigma\left(W_{H-1}^*\sigma\left(\cdots\sigma\left(W_1^* x\right)\right)\right)\right\|^2\right]$$

under the teacher-student learning framework [14] with $w = (W_1, W_2, \cdots, W_H)$ a set of weight matrices. Here $W_k \in \mathbb{R}^{p_k \times p_{k-1}}, k \in [H], p_0 = p$ are parameter matrices of the student network, $W_k^*, k \in [H]$ are the fixed ground-truth parameters of the teacher network, and $\sigma(\cdot)$ is a polynomial with degree $D$. We use online SGD to minimize the population loss $\mathcal{L}(w)$. Formally, we first choose a mini-batch size $b$ and initial weight matrices $\{W_{0,k}, k \in [H]\}$; in each iteration $t$, we independently draw a mini-batch $\mathcal{B}_t^b := \left\{x_{t,i}^b : i \in [b]\right\}$ of $b$ samples from $\mathcal{D}$ and update the weight matrices by $W_{t+1,k}^b = W_{t,k}^b - \alpha_t g_{t,k}^b$, where

$$g_{t,k}^b := \frac{1}{b}\sum_{i=1}^b \nabla_{W_{t,k}^b}\left(\frac{1}{2}\left\|W_{t,H}^b\sigma\left(W_{t,H-1}^b\sigma\left(\cdots\sigma\left(W_{t,1}^b x_{t,i}^b\right)\right)\right) - W_H^*\sigma\left(W_{H-1}^*\sigma\left(\cdots\sigma\left(W_1^* x_{t,i}^b\right)\right)\right)\right\|^2\right).$$

For a multi-set of matrices $\mathcal{M} = \{M_1, \ldots, M_n\}$, we use $\deg(A; \mathcal{M})$ to denote the number of appearances of matrix $A$ and its transpose $A^T$ in $\mathcal{M}$. Mathematically, we have $\deg(A; \mathcal{M}) := \sum_{i \in [n]}\left(\mathbb{I}\{A = M_i\} + \mathbb{I}\{A^T = M_i\}\right)$. We further denote $\deg(\mathcal{A}; \mathcal{M}) := \sum_{A \in \mathcal{A}} \deg(A; \mathcal{M})$ for any set of matrices $\mathcal{A}$. We denote $W_t^b := \left\{W_{t,k}^b, k \in [H]\right\}$, $W_{:t}^b = \bigcup_{s \in [0:t]} W_s^b$, $G_t^b := \left\{g_{t,k}^b, k \in [H]\right\}$, $G_{:t}^b = \bigcup_{s \in [0:t]} G_s^b$, and $W^* := \{W_k^*, k \in [H]\}$. We use $\mathcal{C}$ to denote the infinite set of all non-random matrices given $\mathcal{F}_0$.[2]

**3.2.1. Dynamics: Connecting Generalized Products Step by Step.** As pointed out in the Section 1, the difficulty of studying the dynamics of SGD is how to connect the quantities in iteration $t$ with fixed variables, like the initial weights $W_{0,k}^b$ and mini-batch size $b$. We overcome this challenge by carefully building the connection between (i) $g_{t,k}^b$ and $W_{t,k}^b, k \in [H]$; (ii) $W_{t,k}^b$ and $g_{t-1,k}^b, k \in [H]$. The following two theorems address these two questions by considering a term of mixed product of $W_{t,k}^b$ and $g_{t,k}^b$, respectively.

THEOREM 3.6. *Let $\mathcal{M} := \{M_{i,j} : i \in [0 : I], j \in [J]\}$ be a multi-set of matrices such that each $M_{i,j}$ or its transpose only takes value in $W_{:t}^b \bigcup G_{:t}^b \bigcup \mathcal{C}$ and $\deg\left(G_t^b; \mathcal{M}\right) = d$. Then there exist constants $I', J', L_s$ independent of $b$ and a multi-set of matrices*

---

[2] The definition of $\mathcal{C}$ here is loose to keep the main body of the paper concise. We give a more detailed definition of $\mathcal{C}$ in Appendix A.2.

310  $\mathcal{Q} = \{Q_{l,s,i,j}, l \in [L_s], i \in [0 : I'], j \in [J'], s \in [0 : d]\}$ *such that*

311  (3.2)
312

$$\mathbb{E}\left[\text{tr}\left(C\left(\bigotimes_{i\in[I]}\prod_{j\in[J]}M_{i,j}\right)\right)\prod_{j\in[J]}M_{0,j}\bigg|\mathcal{F}_t^b\right] = \sum_{s=0}^d Q_s \frac{1}{b^s}$$

*where*

$$Q_s = \sum_{l\in[L_s]} c_{l,s}\text{tr}\left(C_{l,s}\left(\bigotimes_{i\in[I']}\prod_{j\in[J']}Q_{l,s,i,j}\right)\right)\prod_{j\in[J']}Q_{l,s,0,j}, s \in [0:d], c_{l,s}$$

313  *is a constant,* $C, C_{l,s} \in \mathcal{C}$ *are constant matrices, and* $Q_{l,s,i,j} \in W_{:t}^b \bigcup G_{:t-1}^b \bigcup \mathcal{C}$.

314     Note that the randomness of $\text{tr}\left(C\left(\bigotimes_{i\in[I]}\prod_{j\in[J]}M_{i,j}\right)\right)\prod_{j\in[J]}M_{0,j}$ in (3.2) only
315  comes from $G_t^b = \left\{g_{t,k}^b, k \in [H]\right\}$ while conditioning on $\mathcal{F}_t^b$. Together with the fact
316  that each $Q_{l,s,i,j}$ involves only $W_{:t}^b \bigcup G_{:t-1}^b \bigcup \mathcal{C}$, Theorem 3.6 enables the induction
317  step from $g_{t,k}^b$ to $W_{t,k}^b$.

318     THEOREM 3.7. *Let* $\mathcal{M} := \{M_{i,j} : i \in [0 : I], j \in [J]\}$ *be a multi-set of matrices such*
319  *that each* $M_{i,j}$ *or its transpose only takes value in* $W_{:t}^b \bigcup G_{:t-1}^b \bigcup \mathcal{C}$ *and* $\deg\left(G_t^b; \mathcal{M}\right) =$
320  $d$. *Then there exist constants* $\mu_1, \ldots, \mu_S \in \mathbb{N}^+, S < \infty$ *independent of* $b$ *and a multi-set*
321  *of matrices* $\mathcal{Q} = \{Q_{s,i,j}, s \in [S], i \in [0 : I], j \in [J]\}$ *such that*

322
323
$$\text{tr}\left(C\left(\bigotimes_{i\in[I]}\prod_{j\in[J]}M_{i,j}\right)\right)\prod_{j\in[J]}M_{0,j} = \sum_{s\in[S]}\mu_s\text{tr}\left(C\left(\bigotimes_{i\in[0:I]}\prod_{j\in[J]}Q_{s,i,j}\right)\right)\prod_{j\in[J]}Q_{s,0,j},$$

324  *where* $C \in \mathcal{C}$ *is a constant matrix, and* $M_{s,i,j} \in W_{:t-1}^b \bigcup G_{:t-1}^b \bigcup \mathcal{C}$.

325     We present the complete version of these theorems and their proofs in Appendix
326  A.2. The exact values of $I', J', c_{l,s}, C_{l,s}, L_s, \alpha_s, S, Q_{l,s,i,j}$ and $Q_{l,s,i}$ are also provided
327  in the corresponding proofs.
328     In fact, these two theorems provide a recursive relationship for explicitly repre-
329  senting any quantity of the form

330  (3.3)
331
$$\text{tr}\left(C\left(\bigotimes_{i\in[I]}\prod_{j\in[J]}M_{i,j}\right)\right)\prod_{j\in[J]}M_{0,j}, \quad M_{i,j} \in W_{:t}^b \bigcup G_{:t}^b \bigcup \mathcal{C}$$

as the sum of many other terms of the same form

$$\text{tr}\left(C\left(\bigotimes_{i\in[I]}\prod_{j\in[J]}M_{i,j}\right)\right)\prod_{j\in[J]}M_{0,j} = \sum_s \mu_s'\text{tr}\left(C\left(\bigotimes_{i\in[0:I']}\prod_{j\in[J']}Q_{s,i,j}\right)\right)\prod_{j\in[J']}Q_{s,0,j},$$

332  where $Q_{s,i,j} \in W_{:t-1}^b \bigcup G_{:t-1}^b \bigcup \mathcal{C}$ and $\mu_s'$s' are some constants independent of $b$. Since
333  $Q_{s,i,j}$ no longer takes value in $W_t^b \bigcup G_t^b$, we are able to reduce the time step by one.
334  As a direct result, by recursively applying these two theorems, we are able to represent
335  the expected value (conditioning on $\mathcal{F}_0$) of the term in (3.3) using learning rates, initial
336  weights, ground-truth weights, and other constants matrices.

337     THEOREM 3.8. *Let* $\mathcal{M} := \{M_{i,j} : i \in [0 : I], j \in [J]\}$ *be a multi-set of matrices*
338  *such that each* $M_{i,j}$ *or its transpose only takes value in* $W_{:t}^b \bigcup G_{:t}^b \bigcup \mathcal{C}$. *Then there*

*exist constants $I', J', S, \overline{L}_s$ independent of $b$, $s \in [0:S]$ and a multi-set of matrices* $\mathcal{Q} = \left\{ Q_{l,s,i,j}, l \in [\overline{L}_s], s \in [S], i \in [0:I'], j \in [J'], \right\}$ *such that*

(3.4)
$$\mathbb{E}\left[ \mathrm{tr}\left( C \left( \bigotimes_{i \in [I]} \prod_{j \in [J]} M_{i,j} \right) \right) \prod_{j \in [J]} M_{0,j} \middle| \mathcal{F}_0 \right] = \sum_{s \in [S]} Q_s \frac{1}{b^s},$$

*where*

$$Q_s = \sum_{l \in [\overline{L}_s]} c_{l,s} \mathrm{tr}\left( C_{l,s} \left( \bigotimes_{i \in [I']} \prod_{j \in [J']} Q_{l,s,i,j} \right) \right) \prod_{j \in [J']} Q_{l,s,0,j}, s \in [0:S],$$

$c_{l,s}$ *is a constant,* $C, C_{l,s} \in \mathcal{C}$ *are constant matrices, and* $Q_{l,s,i,j} \in W_0^b \bigcup \mathcal{C}$.

Again, the complete version of Theorem 3.8 and the exact values of these constants and matrices are presented in Appendix A.2.

**3.2.2. Applications: Decreasing Property of the Variance of Stochastic Gradient Estimators.** In this section, we use the theorems presented in Section 3.2.1 to show some applications of this framework. It is easy to verify that $\mathrm{var}\left( g_{t,k}^b \right), \mathbb{E}\left[ \mathcal{L}(w_t^b) \right]$ and $\mathrm{var}\left( \mathcal{L}(w_t^b) \right)$ can be written as the sum of several terms in the form of the left hand side of (3.4) by further taking expectation over the random initialization of weight matrices[3]. As a special case of Theorem 3.8, Theorem 3.9 shows that the variance of the stochastic gradient estimators is a polynomial of $\frac{1}{b}$ without a constant term. This backs the important intuition that the variance is approximately inversely proportional to the mini-batch size $b$ and provide much more precise relationship between the variance and the mini-batch size $b$.

THEOREM 3.9. *Given $t \in \mathbb{N}$, value $\mathrm{var}\left( g_{t,k}^b \right), k \in [H]$ can be written as a polynomial of $\frac{1}{b}$ with degree at most $(D+1)^{(t+1)D} - 1$ with no constant term. Formally, we have $\mathrm{var}\left( g_{t,k}^b \right) = \beta_1 \frac{1}{b} + \cdots + \beta_r \frac{1}{b^r}$, where $r \leqslant 2(D+1)^{(t+1)D} - 1$ and each $\beta_i$ is a constant independent of $b$.*

One should note that the polynomial representation of $\mathrm{var}\left( g_{t,k}^b \right)$ does not have the constant term. This is intuitively correct since $\mathrm{var}\left( g_{t,k}^b \right) \to 0$ as $b \to \infty$. Therefore, to show that the variance is a decreasing function of $b$, we only need to show that the leading coefficient $\beta_1$ is non-negative. This is guaranteed by the fact that variance is always non-negative. We therefore have the next theorem.

THEOREM 3.10. *Given $t \in \mathbb{N}$, there exists a constant $b_0$ such that for all $b \geqslant b_0$, function $\mathrm{var}\left( g_{t,k}^b \right), k \in [H]$ is a decreasing function of $b$.*

The constant $b_0$ is the largest root of the equation $\beta_1 b^{r-1} + \beta_2 b^{r-2} + \cdots + \beta_r = 0$. See the proof of Theorem 3.10 in Appendix A.2 for more details. Although we cannot provide an explicit form of $b_0$, we can calculate it by the recursive relationship as provided in Theorems 3.6 and 3.7. We further numerically verify that $b_0$ is 1 in many

---

[3]For example, for $i \in [H]$, we have

$$\mathrm{var}\left( g_{t,i}^b \right) = \mathbb{E}\left[ \left\| g_{t,i}^b \right\|^2 \right] - \left\| \mathbb{E} g_{t,i}^b \right\|^2 = \mathbb{E}_{w_0}\left[ \mathbb{E}\left[ \mathrm{tr}\left( g_{t,i}^b \left( g_{t,i}^b \right)^T \right) \middle| \mathcal{F}_0 \right] \right] - \left\| \mathbb{E}_{w_0}\left[ \mathbb{E}\left[ g_{t,i}^b \middle| \mathcal{F}_0 \right] \right] \right\|^2.$$

setups (see Section 4 for more details). From the proofs we conclude that the scale of each $\beta_i$ is of the order $\mathcal{O}\left(\|M\|\right)$, where $M$ is a product of $W_{0,k}, W_k^*, k \in [H]$ and other constant matrices.

In conclusion, we provide a framework for recursively calculating the expected value of a general form that consists of stochastic gradient estimators and weight matrices at time step $t$. As an application, we use our framework to represent the variance of stochastic gradient estimators by a polynomial in $1/b$ and prove that the variance is a decreasing function of $b$ when $b$ is large. Readers should note that the framework here can handle $g_{t,k}^b$ and $W_{t,k}^b$ with any finite degree, and thus it provides much larger capability than just calculating the variance. As a result, similar to Theorems 3.9 and 3.10, we can show that the population loss $\mathcal{L}(w_t^b)$ at iteration $t$ is also a polynomial in $1/b$ and is a decreasing function of $b$ when $b$ is large.

**3.3. General Feed-forward Neural Networks.** In this section, we discuss the extensions of our framework to feed-forward networks with general (non-polynomial) activation functions.

Note that for any smooth activation function $\sigma^S$ (e.g., Sigmoid and Leaky ReLU), it's always possible to find a corresponding polynomial function, $\sigma^P$ such that it approximates $\sigma^S$ as closely as desired within a specified compact domain. This means that, regardless of the specific smooth activation function used, there exists a polynomially-activated function that can mimic its behavior within a certain range. This intuition leads to the following theorem.

THEOREM 3.11. *For any smooth activation function $\sigma^S$, $\epsilon > 0$ and time step $T \in \mathbb{N}^+$, there exists a polynomial $\sigma^P$ (depending on $\epsilon, \sigma^S$, and $T$) such that $\left\| g_{T,k}^S - g_{T,k}^P \right\| \leqslant \epsilon, k \in [H]$, where $g_{t,k}^S$ and $g_{t,k}^P$ are the stochastic gradient of the corresponding network's weight matrix on $k$-th layer at time step $t$.*

The proof of the above theorem is deferred to Appendix A.2.4. Theorem 3.11 states that the SG estimators of a general neural network can be approximated arbitrarily well by the counterpart of a polynomially-activated function at any given time step $T$. This is a significant finding as it allows us to approximate the behavior of complex neural networks using simpler polynomial functions. Furthermore, when we combine this with the theorems presented in Section 3.2, which provide an exact representation of the SG estimators of any polynomially-activated function using only information available before training, we gain the ability to approximate the SG estimators of general networks arbitrarily well using only the known information at the initial time step $t = 0$.

This approximation has profound implications for our understanding of neural network behavior and offers potential avenues for designing more advanced optimization methods. See the discussions in Section 5 for more details.

**4. Experiments.** In this section, we present numerical results to support the theorems in Section 3, to backup the hypotheses discussed in the introduction, and provide further insights into the impact of the mini-batch size on the dynamics of SGD. The experiments are conducted on four datasets and models that are relatively small due to the computational cost of using large models and datasets.

**4.1. Datasets and Settings.** For all experiments, we perform mini-batch SGD multiple times starting from the same initial weights and following the same choice of the learning rates and other hyper-parameters, if applicable. This enables us to calculate the variance of the gradient estimators and other statistics in each iteration,

418  where the randomness comes only from different samples of SGD. The learning rate
419  $\alpha_t$ is selected to be inversely proportional to iteration $t$, or fixed, depending on the
420  task at hand.
421      All models are implemented using PyTorch version 1.4 [32] and trained on NVIDIA
422  2080Ti/1080 GPUs. We have also tested several other random initial weights and
423  ground-truth weights, and learning rates, and the results and conclusions are similar
424  and not presented.

425      **4.1.1. Graduate Admission Dataset.** The Graduate Admission dataset[4] [1]
426  is to predict the chance of a graduate admission using linear regression. The dataset
427  contains 500 samples with 6 features and is normalized by mean and variance of
428  each feature. This is a popular regression dataset with clean data. We build a linear
429  regression model to predict the chance of acceptance (we include the intercept term in
430  the model) and minimize the empirical $L_2$ loss using mini-batch SGD, as stated in
431  Section 3.1.
432      For the experiment in Figure 1(a), we randomly select an initial weight vectors
433  $w_0$ and run SGD for 2,000 iterations where it appears to converge. We record all
434  statistics at every iteration. There are in total 1,000 runs behind each observation
435  which yields a p-value lower than 0.05. As for Figure 1(b), we select 20 different
436  $b$'s and run SGD from the same initial point for 40 iterations. There are in total
437  of 200,000 runs to make sure the p-value of all statistics are lower than 0.05. In all
438  experiments, the learning rate is chosen to be $\alpha_t = \frac{1}{2t}, t \in [2000]$ because this rate
439  yields a theoretical convergence guaranteed (factor $1/2$ has been fine tuned). The
440  purpose of this experiment is to empirically study the rate of decrease of the variance.
441  The theoretical study exhibited in Section 3.1 establishes the non-increasing property
442  but it does not state anything about the rate of decrease.

443      **4.1.2. Synthetic Dataset.** We build a synthetic dataset of standard normal
444  samples to study the setting in Section 3.1. We fix the teacher network with 64 input
445  neurons, 256 hidden neurons and 128 output neurons. We optimize the population $L_2$
446  loss by updating the two parameter matrices of the student network using online SGD,
447  as stated in Section 3.1. In this case we have proved the functional form of the variance
448  as a function of $b$ and show the decreasing property of the variance of the stochastic
449  gradient estimators for large mini-batch sizes. However, we do not show the decreasing
450  property for every $b$. With this experiment we confirm that the conjecture likely holds.
451  In the experiment, we randomly select two initial weight matrices $W_{0,1}, W_{0,2}$ and
452  the ground-truth weight matrices $W_1^*, W_2^*$. We run SGD for 1,000 iterations which
453  appears to be a good number for convergence while there are 1,000 runs of SGD in
454  total to again give a p-value below 0.05. We record all statistics at every iteration.
455  The learning rate is chosen to be $\alpha_t = \frac{1}{10t}, t \in [1000]$ for the same reason as in the
456  regression experiment.

457      **4.1.3. MNIST Dataset.** The MNIST dataset is to recognize digits in handwrit-
458  ten images of digits. We use all 60,000 training samples and 10,000 validation samples
459  of MNIST. The images are normalized by mapping each entry to $[-1, 1]$. We build
460  a three-layer fully connected neural network with 1024, 512 and 10 neurons in each
461  layer. For the two hidden layers, we use the ReLU activation function. The last layer
462  is the softmax layer which gives the prediction probabilities for the 10 digits. We use
463  mini-batch SGD to optimize the cross-entropy loss of the model. The model deviates
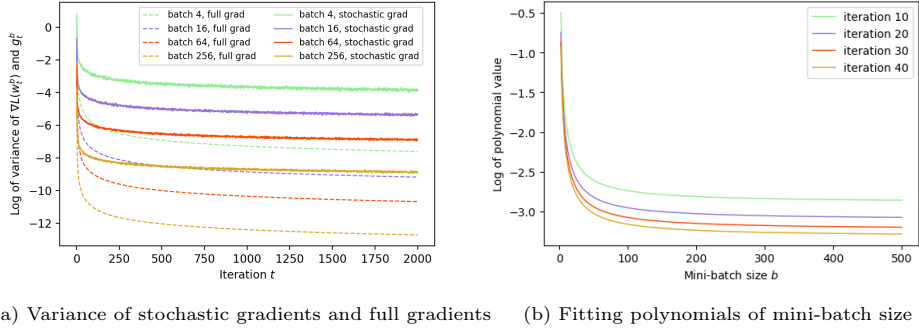464  from our analytical setting since it has non-linear activations, it has the cross-entropy

---

[4]https://www.kaggle.com/mohansacharya/graduate-admissions

(a) Variance of stochastic gradients and full gradients    (b) Fitting polynomials of mini-batch size $b$

Fig. 1: Experimental results for the Graduate Admission dataset. **Left:** $\log\left(\mathrm{var}\left(g_t^b|\mathcal{F}_0\right)\right)$ and $\log\left(\mathrm{var}\left(\nabla L(w_t^b)|\mathcal{F}_0\right)\right)$ vs iteration $t$ for 4 different mini-batch sizes. **Right:** The log of polynomial values when fitting polynomials on selected mini-batch sizes at certain iterations.

loss function (instead of $L_2$), and empirical loss (as opposed to population). MNIST is selected due to its fast training and popularity in deep learning experiments. The goal is to verify the results in this different setting and to back up our hypotheses.

We run SGD for 1,000 epochs on the training set which is enough for convergence. The learning rate is a constant set to $3 \cdot 10^{-3}$ (which has been tuned). For the experiment in Figure 4, there are in total 100 runs to give us the p-value below 0.05. For the experiment in Figure 3(a), we randomly select five different initial points and we have 50 runs for each initial point. For the experiment corresponding to Figure 3(b), we choose $\alpha = 8$ and $\sigma = 2$ as in [37]. The initial weights and other hyper-parameters are chosen to be the same as in Figure 4.

**4.1.4. Yelp Review Dataset.** The Yelp Review dataset from the Yelp Dataset Challenge [42] contains 1,569,264 samples of customer reviews with positive/negative sentiment labels. We use 10,000 samples as our training set and 1,000 samples as the validation set. We use XLNet [41] to perform sentiment classification on this dataset. Our XLNet has 6 layers, the hidden size of 384, and 12 attention heads. There are in total 35,493,122 parameters. We intentionally reduce the number of layers and hidden size of XLNet and select a relatively small size of the training and validation sets since training of XLNet is very time-consuming ([41] train on 512 TPU v3 chips for 5.5 days) and we need to train the model for multiple runs. This setting allows us to train our model in several hours on a single GPU card. We train the model using the Adam weight decay optimizer, and some other techniques, as suggested in Table 8 of [41]. This dataset represents sequential data where we further consider the hypotheses.

We randomly select a set of initial parameters and run Adam with two different mini-batch sizes of 32 and 64. For computational tractability reasons, for each mini-batch size there are in total of 100 runs and each run corresponds to 20 epochs. We record the variance of the stochastic gradient, loss and accuracy in every step of Adam. The statistics reported in Figure 5 are averaged through each epoch. In all experiments, the learning rate is set to be $4 \cdot 10^{-5}$ and the $\epsilon$ parameter of Adam is set to be $10^{-8}$ (these two have been tuned). The stochastic gradients of all parameter matrices are clipped with threshold 1 in each iteration. We use the same setup for the learning rate warm-up strategy as suggested in [41]. The maximum sequence length is set to be 128 and we pad the sequences with length smaller than 128 with zeros.
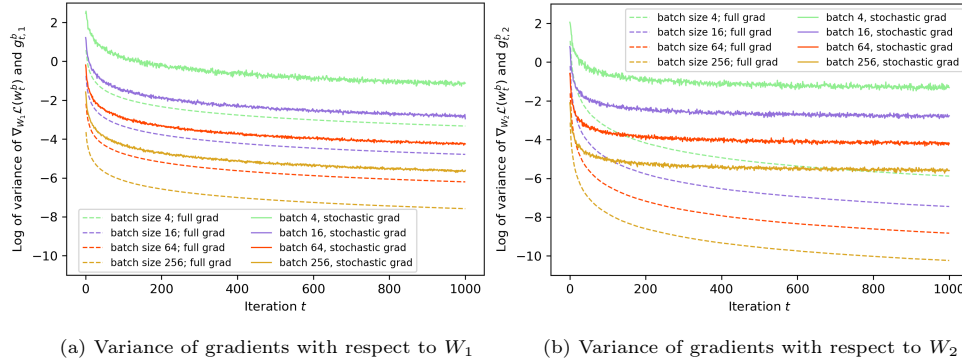
(a) Variance of gradients with respect to $W_1$          (b) Variance of gradients with respect to $W_2$

Fig. 2: Experimental results for the Synthetic dataset. **Left:** $\log\left(\mathsf{var}\left(g_{t,1}^b\big|\mathcal{F}_0\right)\right)$ and $\log\left(\mathsf{var}\left(\nabla_{W_1}\mathcal{L}(W_{t,1}^b, W_{t,2}^b)\big|\mathcal{F}_0\right)\right)$ vs iteration $t$. **Right:** $\log\left(\mathsf{var}\left(g_{t,2}^b\big|\mathcal{F}_0\right)\right)$ and $\log\left(\mathsf{var}\left(\nabla_{W_2}\mathcal{L}(W_{t,1}^b, W_{t,2}^b)\big|\mathcal{F}_0\right)\right)$ vs iteration $t$.

**4.2. Discussion.** As observed in Figure 1(a), under the linear regression setting with the Graduate Admission dataset, the variance of the stochastic gradient estimators and full gradients are all strictly decreasing functions of $b$ for all iterations. This result verifies the theorems in Section 3.1. Figure 1(b) further studies the rate of decrease of the variance. From the proofs in Section 3.1 we see that $\mathsf{var}\left(g_t^b\big|\mathcal{F}_0\right)$ is a polynomial of $\frac{1}{b}$ with degree $t+1$. Therefore, for every $t$, we can approximate this polynomial by sampling many different $b$'s and calculate the corresponding variances. We pick $b$ to cover all numbers that are either a power of 2 or multiple of 40 in $[2, 500]$ (there are a total of 21 such values) and fit a polynomial with degree 6 (an estimate from the analyses) at $t = 10, 20, 30, 40$. Figure 1(b) shows the fitted polynomials. As we observe, the value $\mathsf{var}\left(g_t^b\big|\mathcal{F}_0\right)$ (approximated by the value of the polynomial) is both decreasing with respect to the mini-batch size $b$ and iteration $t$. Further, the rate of decrease in $b$ is slower as the $b$ increasing. This provides a further insight into the dynamics of training a linear regression problem with SGD.

Under the two-layer linear network setting with the synthetic dataset, Figure 2 verifies that the variance of the stochastic gradient estimators and full gradients are all strictly decreasing functions of $b$ for all iterations. This figure also empirically shows that the constant $b_0$ in Theorem 3.10 could be as small as $b_0 = 4$. In fact, we also experiment with the mini-batch size of 1 and 2, and the decreasing property remains to hold. We also test this on multiple choices of initial weights and learning rates and this pattern remains clear.

In aforementioned two experiments we use SGD in its original form by randomly sampling mini-batches. In deep learning with large-scale training data such a strategy is computationally prohibitive and thus samples are scanned in a cyclic order which implies fixed mini-batches are processed many times. Therefore, in the next two datasets we perform standard "epoch" based training to empirically study the remaining two hypotheses discussed in the introduction (decreasing loss and error as a function of $b$) and sensitivity with respect to the initial weights. Note that we are using cross-entropy loss in the MNIST dataset and the Adam optimizer in the Yelp dataset and thus these experiments do not meet all of the assumptions of the analysis in Section 3.

As shown in Figure 3(a), we run SGD with two batch sizes 64 and 128 on five different initial weights. This plot shows that, even the smallest value of the variance

(a) Different initial weights

(b) Gap of accuracy (zoomed-in)

Fig. 3: Experimental results for the MNIST dataset. **Left:** The median, min, and max of the log of variance of the stochastic gradient estimators for two different mini-batch sizes (distinguished by colors) and five different initial weights. The solid lines show the median of all five initial weights while the highlighted regions show the min and max of the log of variance. **Right:** The gap of accuracy on training and test sets vs epochs starting from epoch 100.



(a) Log of loss for training and validation sets

(b) Log of error for training and validation sets

Fig. 4: Experimental results for the MNIST dataset. **Left:** The log of the training and validation loss vs epochs. **Right:** The log of training and validation error vs epochs. Here error is defined as one minus predicting accuracy. The plot does not show the epochs if error equals to zero.

529  among the five different initial weights with a mini-batch size of 64, is still larger than
530  the largest variance of mini-batch size 128. We observe that the sensitivity to the
531  initial weights is not large. This plot also empirically verifies our conjecture in the
532  introduction that the variance of the stochastic gradient estimators is a decreasing
533  function of the mini-batch size, for all iterations of SGD in a general deep learning
534  model.
535        In addition, we also conjecture that there exists the decreasing property for the
536  expected loss, error and the generalization ability with respect to the mini-batch size.
537  Figure 4(a) shows that the expected loss (again, randomness comes from different runs
538  of SGD through the different mini-batches with the same initial weights and learning
539  rates) on the training set is a decreasing function of $b$. However, this decreasing

(a) Variance of SG (b) Training/validation loss (c) Training - validation error

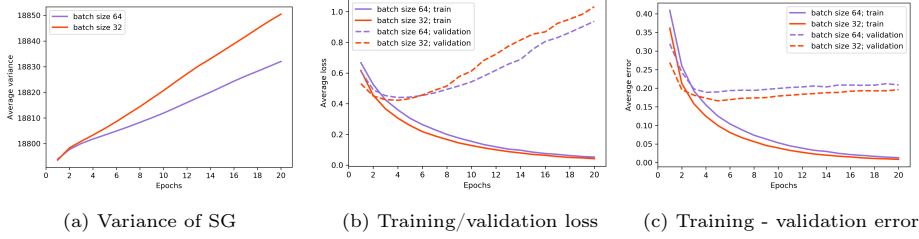Fig. 5: Experimental results for the XLNet model on the Yelp dataset. **Left:** The variance of stochastic gradient estimators vs epochs. **Middle:** The training and validation loss vs epochs. **Right:** The training and validation error vs epochs.

540 property does not hold on the validation set when the loss tends to be stable or
541 increasing, in other words, the model starts to be over-fitting. We hypothesize that
542 this is because the learned weights start to bounce around a local minimum when
543 the model is over-fitting. As the larger mini-batch size brings smaller variance, the
544 weights are closer to the local minimum found by SGD, and therefore yield a smaller
545 loss function value. Figure 4(b) shows that both the expected error on training and
546 validation sets are decreasing functions of $b$.

547 Figure 3(b) exhibits a relationship between the model's generalization ability and
548 the mini-batch size. As suggested by [37], we build a test set by distorting the 10,000
549 images of the validation set. The prediction accuracy is obtained on both training
550 and test sets and we calculate the gap between these two accuracies every 100 epochs.
551 We use this gap to measure the model generalization ability (the smaller the better).
552 Figure 3(b) shows that the gap is an increasing function of $b$ starting at epoch 500,
553 which partially aligns with our conjecture regarding the relationship between the
554 generalization ability and the mini-batch size. We also test this on multiple choices of
555 the hyper-parameters which control the degree of distortion in the test set and this
556 pattern remains clear.

557 Figure 5 shows the similar phenomenon that the variance of stochastic estimators
558 and the expected loss and error on both training and validation sets are decreasing
559 functions of $b$ even if we train XLNet using Adam. This example gives us confidence
560 that the decreasing properties are not merely restricted on shallow neural networks
561 or vanilla SGD algorithms. They actually appear in many advanced models and
562 optimization methods.

563 **5. Discussion and Future Work.** We study the dynamics of SGD by explicitly
564 representing the important quantities of SGD using the mini-batch size and initial
565 weights. For linear regression and a multi-layer polynomially-activated network, we
566 are able to build frameworks that recursively calculate general forms of the product of
567 the weight matrices and stochastic gradient estimators between consecutive iterations.
568 We further theoretically prove that the variance conjecture holds. Experiments are
569 performed on multiple models and datasets to verify our claims and their applicability
570 to practical settings. Besides, we also empirically address the conjectures about the
571 expected loss and the generalization ability.

572 We provide mathematical tools to calculate and represent the product of the
573 stochastic gradients estimators and weight matrices in the $t$-th step (and not a single
574 step), which is non-trivial and requires a sophisticated mathematical proof. These

tools can be extended to calculate any form that has a polynomial relationship to the model parameters $w_t^b$, e.g. expectation/variance of the loss function, norm of the SG estimator to any finite degree. We can also derive other properties of the dynamics of SGD by using these tools.

One possible application of the results is to help tighten the convergence rates of SGD and develop better variance reduction methods. Current analyses of SGD convergence rely on two constants $M$ and $M_V$ such that $\mathsf{var}\left(g_t^b\right) \leqslant M + M_V \left\|\nabla L(w_t^b)\right\|^2$. But it is unclear what are the exact values of $M$ and $M_V$ (see Assumption 4.3 of [5] and the context therein). It is a common practice to take relatively large $M$ and $M_V$ to make sure the above bound holds. However, this leads to a relatively poor convergence rate of the SGD algorithm. Our frameworks are able to explicitly calculate $\mathsf{var}\left(g_t^b\right)$ and $\left\|\nabla L(w_t^b)\right\|^2$ by recursive formulas and thus to provide optimal values for $M$ and $M_V$.

Another challenging research direction is to theoretically and explicitly investigate the generalization ability during training of SGD. There are existing works studying the relationship between the variance of the stochastic gradients and the generalization ability [10, 29]. Together with the frameworks developed herein, it would be possible to tighten the generalization bounds of a neural network by explicit variance and other quantities. We can further choose an optimal mini-batch size which minimizes the generalization ability by solving a polynomial equation if we have a more precise relationship between the variance and the generalization ability.

Further interesting work is to extend our techniques to more complicated and sophisticated networks as we discuss in Section 3.3. Although the underlying model of this paper corresponds to deep polynomially-activated networks in a strict manner and to general neural networks in an approximate sense, we are able to show a deeper relationship between the variance and the mini-batch size, the polynomial in $1/b$, while the common knowledge is simply that the variance is proportional to $1/b$. The extension to other optimization algorithms, like Adam and Gradient Boosting Machines, are also very attractive. We hope our theoretical framework can serve as a tool for future research of this kind.

## Appendix A. Lemmas and Proofs.

**A.1. Lemmas and Proofs of Results in Section 3.1.** For two matrices $A, B$ with the same dimension, we define the inner product $\langle A, B \rangle := \operatorname{tr}\left(A^T B\right)$.

LEMMA A.1. *Suppose that $f(x)$ and $g(x)$ are both smooth, non-negative and decreasing functions of $x \in \mathbb{R}$. Then $h(x) = f(x)g(x)$ is also a non-negative and decreasing function of $x$.*

*Proof.* It is obvious that $h(x)$ is non-negative for all $x$. The first-order derivative of $h$ is
$$h'(x) = f'(x)g(x) + f(x)g'(x) \leqslant 0,$$
and thus $h(x)$ is also a decreasing function of $x$. □

*Proof of Lemma 3.1.* Throughout the paper, We use $C_n^k = \frac{n!}{k!(n-k)!}$ to denote the combinatorial number. Note that

$$\mathbb{E}\left[g_t^b \left(g_t^b\right)^T \middle| \mathcal{F}_t^b\right] = \frac{1}{b^2}\mathbb{E}\left[\sum_{i \in \mathcal{B}_t^b} \nabla L_i\left(w_t^b\right) \sum_{i \in \mathcal{B}_t^b} \nabla L_i\left(w_t^b\right)^T \middle| \mathcal{F}_t^b\right]$$

$$= \frac{1}{b^2}\left(\frac{C_{n-1}^{b-1}}{C_n^b}\sum_{i=1}^n \nabla L_i\left(w_t^b\right)\nabla L_i\left(w_t^b\right)^T + \frac{C_{n-2}^{b-2}}{C_n^b}\sum_{i \neq j} \nabla L_i\left(w_t^b\right)\nabla L_j\left(w_t^b\right)^T\right)$$

$$= \frac{1}{b^2}\left(\frac{b}{n}\sum_{i=1}^n \nabla L_i\left(w_t^b\right)\nabla L_i\left(w_t^b\right)^T + \frac{b(b-1)}{n(n-1)}\sum_{i \neq j} \nabla L_i\left(w_t^b\right)\nabla L_j\left(w_t^b\right)^T\right)$$

$$= \frac{1}{b^2}\left(\frac{b(n-b)}{n(n-1)}\sum_{i=1}^n \nabla L_i\left(w_t^b\right)\nabla L_i\left(w_t^b\right)^T + \frac{b(b-1)}{n(n-1)}\sum_{i=1}^n \nabla L_i\left(w_t^b\right)\sum_{i=1}^n \nabla L_i\left(w_t^b\right)^T\right)$$

$$= \frac{n-b}{bn(n-1)}\sum_{i=1}^n \nabla L_i\left(w_t^b\right)\nabla L_i\left(w_t^b\right)^T + \frac{(b-1)n}{b(n-1)}\nabla L\left(w_t^b\right)\nabla L\left(w_t^b\right)^T.$$

For any $A \in \mathbb{R}^{p \times p}$, we have

$$\mathbb{E}\left[\left\|Ag_t^b\right\|^2 \middle| \mathcal{F}_t^b\right] = \mathbb{E}\left[\left(g_t^b\right)^T A^T Ag_t^b \middle| \mathcal{F}_t^b\right] = \mathbb{E}\left[\operatorname{tr}\left(\left(g_t^b\right)^T A^T Ag_t^b\right) \middle| \mathcal{F}_t^b\right] = \mathbb{E}\left[\operatorname{tr}\left(A^T Ag_t^b\left(g_t^b\right)^T\right) \middle| \mathcal{F}_t^b\right]$$

$$= \operatorname{tr}\left(A^T A\mathbb{E}\left[g_t^b\left(g_t^b\right)^T \middle| \mathcal{F}_t^b\right]\right)$$

$$= \operatorname{tr}\left(\frac{n-b}{bn(n-1)}\sum_{i=1}^n A^T A\nabla L_i\left(w_t^b\right)\nabla L_i\left(w_t^b\right)^T + \frac{(b-1)n}{b(n-1)}A^T A\nabla L\left(w_t^b\right)\nabla L\left(w_t^b\right)^T\right)$$

$$= \frac{n-b}{bn(n-1)}\sum_{i=1}^n \left\|A\nabla L_i\left(w_t^b\right)\right\|^2 + \frac{(b-1)n}{b(n-1)}\left\|A\nabla L\left(w_t^b\right)\right\|^2$$

$$= c_b\left(\frac{1}{n}\sum_{i=1}^n \left\|A\nabla L_i\left(w_t^b\right)\right\|^2 - \left\|A\nabla L\left(w_t^b\right)\right\|^2\right) + \left\|A\nabla L\left(w_t^b\right)\right\|^2.$$

Therefore, we have

$$\operatorname{var}\left(Ag_t^b \middle| \mathcal{F}_t^b\right) = \mathbb{E}\left[\left\|Ag_t^b\right\|^2 \middle| \mathcal{F}_t^b\right] - \left\|\mathbb{E}\left[Ag_t^b \middle| \mathcal{F}_t^b\right]\right\|^2 = \mathbb{E}\left[\left\|Ag_t^b\right\|^2 \middle| \mathcal{F}_t^b\right] - \left\|A\nabla L\left(w_t^b\right)\right\|^2$$

$$= c_b\left(\frac{1}{n}\sum_{i=1}^n \left\|A\nabla L_i\left(w_t^b\right)\right\|^2 - \left\|A\nabla L\left(w_t^b\right)\right\|^2\right). \qquad \square$$

LEMMA A.2. *For any set of square matrices $\{A_1, \cdots, A_n\} \in \mathbb{R}^{p \times p}$, if we denote $A = \sum_{i=1}^n A_i x_i x_i^T$, then we have*

$$\mathbb{E}\left[\left\|\sum_{i=1}^n A_i \nabla L_i\left(w_{t+1}^b\right)\right\|^2 \middle| \mathcal{F}_0\right] = \mathbb{E}\left[\left\|\sum_{i=1}^n B_i \nabla L_i\left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right] + \frac{\alpha_t^2 c_b}{n^2}\sum_{k=1}^n\sum_{l=1}^n \mathbb{E}\left[\left\|\sum_{i=1}^n B_i^{kl} \nabla L_i\left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right].$$

Here $B_i = A_i - \frac{\alpha_t}{n} A$; $B_i^{kl} = A$ if $i = k, i \neq l$, $B_i^{kl} = A$ if $i = l, i \neq k$, and $B_i^{kl}$ equals the zero matrix, otherwise.

*Proof of Lemma A.2.* Let $C_i = x_i x_i^T$ and $C = \frac{1}{n} \sum_{i=1}^n C_i$ and thus $A = \sum_{i=1}^n A_i C_i$. Then

$$\mathbb{E}\left[\left\|\sum_{i=1}^n A_i \nabla L_i \left(w_{t+1}^b\right)\right\|^2 \middle| \mathcal{F}_0\right] = \mathbb{E}\left[\mathbb{E}\left[\left\|\sum_{i=1}^n A_i \nabla L_i \left(w_{t+1}^b\right)\right\|^2 \middle| \mathcal{F}_t^b\right] \middle| \mathcal{F}_0\right]$$

$$= \mathbb{E}\left[\mathbb{E}\left[\left\|\sum_{i=1}^n A_i \left(x_i^T w_{t+1}^b - y_i\right) x_i\right\|^2 \middle| \mathcal{F}_t^b\right] \middle| \mathcal{F}_0\right] = \mathbb{E}\left[\mathbb{E}\left[\left\|\sum_{i=1}^n A_i \left(x_i^T \left(w_t^b - \alpha_t g_t^b\right) - y_i\right) x_i\right\|^2 \middle| \mathcal{F}_t^b\right] \middle| \mathcal{F}_0\right]$$

$$= \mathbb{E}\left[\mathbb{E}\left[\left\|\sum_{i=1}^n A_i \nabla L_i \left(w_t^b\right) - \alpha_t A g_t^b\right\|^2 \middle| \mathcal{F}_t^b\right] \middle| \mathcal{F}_0\right]$$

$$= \mathbb{E}\left[\left\|\sum_{i=1}^n A_i \nabla L_i \left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right] - 2\alpha_t \mathbb{E}\left[\mathbb{E}\left[\left\langle \sum_{i=1}^n A_i \nabla L_i \left(w_t^b\right), A g_t^b \right\rangle \middle| \mathcal{F}_t^b\right] \middle| \mathcal{F}_0\right] + \alpha_t^2 \mathbb{E}\left[\mathbb{E}\left[\left\|A g_t^b\right\|^2 \middle| \mathcal{F}_t^b\right] \middle| \mathcal{F}_0\right]$$

$$= \mathbb{E}\left[\left\|\sum_{i=1}^n A_i \nabla L_i \left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right] - 2\alpha_t \mathbb{E}\left[\left\langle \sum_{i=1}^n A_i \nabla L_i \left(w_t^b\right), A \nabla L \left(w_t^b\right) \right\rangle \middle| \mathcal{F}_0\right]$$

$$+ \alpha_t^2 \mathbb{E}\left[c_b \left(\frac{1}{n} \sum_{i=1}^n \left\|A \nabla L_i(w_t^b)\right\|^2 - \left\|A \nabla L(w_t^b)\right\|^2\right) + \left\|A \nabla L(w_t^b)\right\|^2 \middle| \mathcal{F}_0\right]$$

$$= \mathbb{E}\left[\left\|\sum_{i=1}^n A_i \nabla L_i \left(w_t^b\right) - \alpha_t A \nabla L(w_t^b)\right\|^2 \middle| \mathcal{F}_0\right] + \alpha_t^2 c_b \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \left\|A \nabla L_i(w_t^b)\right\|^2 - \left\|A \nabla L(w_t^b)\right\|^2 \middle| \mathcal{F}_0\right]$$

$$= \mathbb{E}\left[\left\|\sum_{i=1}^n A_i \nabla L_i \left(w_t^b\right) - \alpha_t A \nabla L(w_t^b)\right\|^2 \middle| \mathcal{F}_0\right] + \frac{\alpha_t^2 c_b}{n^2} \sum_{i \neq j} \mathbb{E}\left[\left\|A \nabla L_i \left(w_t^b\right) - A \nabla L_j \left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right]$$

$$= \mathbb{E}\left[\left\|\sum_{i=1}^n \left(A_i - \frac{\alpha_t}{n} A\right) \nabla L_i \left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right] + \frac{\alpha_t^2 c_b}{n^2} \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}\left[\left\|A \nabla L_i \left(w_t^b\right) - A \nabla L_j \left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right].$$

Therefore, if we set $B_i = A_i - \frac{\alpha_t}{n} A$ and $B_i^{kl} = \begin{cases} A & i = k, i \neq l, \\ -A & i = l, i \neq k,, \\ 0 & \text{otherwise,} \end{cases}$ we have

$$\mathbb{E}\left[\left\|\sum_{i=1}^n A_i \nabla L_i \left(w_{t+1}^b\right)\right\|^2 \middle| \mathcal{F}_0\right] = \mathbb{E}\left[\left\|\sum_{i=1}^n B_i \nabla L_i \left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right] + \frac{\alpha_t^2 c_b}{n^2} \sum_{k=1}^n \sum_{l=1}^n \mathbb{E}\left[\left\|\sum_{i=1}^n B_i^{kl} \nabla L_i \left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right]. \square$$

*Proof of Theorem 3.3.* We use induction to show this statement.

When $t = 0$, $\mathbb{E}\left[\left\|\sum_{i=1}^n A_i \nabla L_i \left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right] = \left\|\sum_{i=1}^n A_i \nabla L_i (w_0)\right\|^2$ which is invariant of $b$. Therefore, it is a decreasing function of $b$.

Suppose the statement holds for $t$. For any set of matrices $\{A_1, \ldots, A_n\}$ in $\mathbb{R}^{p \times p}$, by Theorem 3.2 we know that there exist matrices $\{B_1, \cdots, B_n\}$ and $\{B_i^{kl} : i, k, l \in [n]\}$ such that

$$\mathbb{E}\left[\left\|\sum_{i=1}^n A_i \nabla L_i \left(w_{t+1}^b\right)\right\|^2 \middle| \mathcal{F}_0\right] = \mathbb{E}\left[\left\|\sum_{i=1}^n B_i \nabla L_i \left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right] + \frac{\alpha_t^2 c_b}{n^2} \sum_{k=1}^n \sum_{l=1}^n \mathbb{E}\left[\left\|\sum_{i=1}^n B_i^{kl} \nabla L_i \left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right]. \square$$

By induction, $\mathbb{E}\left[\left\|\sum_{i=1}^n B_i \nabla L_i \left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right]$ and all $\mathbb{E}\left[\left\|\sum_{i=1}^n B_i^{kl} \nabla L_i \left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right]$ are non-negative and decreasing functions of $b$. Besides, clearly $\frac{\alpha_t^2 c_b}{n^2} = \frac{\alpha_t^2 (n-b)}{bn^3(n-1)}$ and $\frac{\alpha_t^2 c_b}{n^2} \mathbb{E}\left[\left\|\sum_{i=1}^n B_i^{kl} \nabla L_i \left(w_t^b\right)\right\|^2 \middle| \mathcal{F}_0\right]$ (via Lemma A.1) are a non-negative and decreasing function of $b$. Finally, $\mathbb{E}\left[\left\|\sum_{i=1}^n A_i \nabla L_i \left(w_{t+1}^b\right)\right\|^2 \middle| \mathcal{F}_0\right]$, as the sum of non-negative and decreasing functions in $b$, is a non-negative and decreasing function of $b$.

In order to prove Theorem 3.4, we split the task to two separate theorems about the full gradient and the stochastic gradient and prove them one by one.

THEOREM A.3. *Fixing initial weights $w_0$, $\mathsf{var}\left(B\nabla L\left(w_t^b\right)\big|\mathcal{F}_0\right)$ is a decreasing function of mini-batch size $b$ for all $b \in [n]$, $t \in \mathbb{N}$, and all square matrices $B \in \mathbb{R}^{p\times p}$.*

THEOREM A.4. *Fixing initial weights $w_0$, $\mathsf{var}\left(Bg_t^b\big|\mathcal{F}_0\right)$ is a decreasing function of mini-batch size $b$ for all $b \in [n]$, $t \in \mathbb{N}$, and all square matrices $B \in \mathbb{R}^{p\times p}$.*

*Proof of Theorem A.3.* We induct on $t$ to show that the statement holds. For $t = 0$, we have $\mathsf{var}\left(B\nabla L\left(w_t^b\right)\big|\mathcal{F}_0\right) = 0$ for any matrix $B$. Suppose the statement holds for $t - 1 \geqslant 0$. Note that from

$$\nabla L\left(w_t^b\right) = \frac{1}{n}\sum_{i=1}^n x_i\left(x_i^T w_t^b - y_i\right) = \frac{1}{n}\sum_{i=1}^n x_i\left(x_i^T\left(w_{t-1}^b - \alpha_t g_{t-1}^b\right) - y_i\right)$$

$$= \frac{1}{n}\sum_{i=1}^n x_i\left(x_i^T w_{t-1}^b - y_i\right) - \frac{\alpha_t}{n}\sum_{i=1}^n x_i x_i^T g_{t-1}^b = \nabla L\left(w_{t-1}^b\right) - \alpha_t C g_{t-1}^b,$$

we have

$$\mathsf{var}\left(B\nabla L\left(w_t^b\right)\big|\mathcal{F}_0\right) = \mathsf{var}\left(B\nabla L\left(w_{t-1}^b\right) - \alpha_t BC g_{t-1}^b\big|\mathcal{F}_0\right)$$

$$= \mathbb{E}\left[\left\|B\nabla L\left(w_{t-1}^b\right) - \alpha_t BC g_{t-1}^b\right\|^2\Big|\mathcal{F}_0^b\right] - \left\|\mathbb{E}\left[B\nabla L\left(w_{t-1}^b\right) - \alpha_t BC g_{t-1}^b\big|\mathcal{F}_0^b\right]\right\|^2$$

$$= \mathbb{E}\left[\left\|B\nabla L\left(w_{t-1}^b\right)\right\|^2 - 2\alpha_t\left\langle B\nabla L\left(w_{t-1}^b\right), BC g_{t-1}^b\right\rangle + \alpha_t^2\left\|BC g_{t-1}^b\right\|^2\Big|\mathcal{F}_0^b\right] -$$

$$\quad - \left\|\mathbb{E}\left[B\nabla L\left(w_{t-1}^b\right) - \alpha_t BC g_{t-1}^b\big|\mathcal{F}_0^b\right]\right\|^2$$

$$= \mathbb{E}\left[\left\|B\nabla L\left(w_{t-1}^b\right)\right\|^2\Big|\mathcal{F}_0\right] + \alpha_t^2\mathbb{E}\left[\mathbb{E}\left[\left\|BC g_{t-1}^b\right\|^2\Big|\mathcal{F}_{t-1}^b\right]\Big|\mathcal{F}_0^b\right] -$$

$$\quad - 2\alpha_t\mathbb{E}\left[\mathbb{E}\left[\left\langle B\nabla L\left(w_{t-1}^b\right), BC g_{t-1}^b\right\rangle\big|\mathcal{F}_{t-1}^b\right]\Big|\mathcal{F}_0^b\right] - \left\|\mathbb{E}\left[\mathbb{E}\left[B\nabla L\left(w_{t-1}^b\right) - \alpha_t BC g_{t-1}^b\big|\mathcal{F}_{t-1}^b\right]\Big|\mathcal{F}_0^b\right]\right\|^2$$

$$= \mathbb{E}\left[\left\|B\nabla L\left(w_{t-1}^b\right)\right\|^2\Big|\mathcal{F}_0\right] + \alpha_t^2\mathbb{E}\left[c_b\left(\frac{1}{n}\sum_{i=1}^n\left\|BC\nabla L_i\left(w_{t-1}^b\right)\right\|^2 - \left\|BC\nabla L\left(w_{t-1}^b\right)\right\|^2\right) + \left\|BC\nabla L\left(w_{t-1}^b\right)\right\|^2\Big|\mathcal{F}_0\right]$$

(A.1) $$\quad - 2\alpha_t\mathbb{E}\left[\left\langle B\nabla L\left(w_{t-1}^b\right), BC\nabla L\left(w_{t-1}^b\right)\right\rangle\big|\mathcal{F}_0\right] - \left\|\mathbb{E}\left[B\nabla L\left(w_{t-1}^b\right) - \alpha_t BC\nabla L\left(w_{t-1}^b\right)\big|\mathcal{F}_0^b\right]\right\|^2$$

$$= \mathbb{E}\left[\left\|B\left(I - \alpha_t C\right)\nabla L\left(w_{t-1}^b\right)\right\|^2\Big|\mathcal{F}_0^b\right] + \alpha_t^2 c_b\mathbb{E}\left[\left(\frac{1}{n}\sum_{i=1}^n\left\|BC\nabla L_i\left(w_{t-1}^b\right)\right\|^2 - \left\|BC\nabla L\left(w_{t-1}^b\right)\right\|^2\right)\Big|\mathcal{F}_0\right]$$

$$\quad - \left\|\mathbb{E}\left[B\left(I - \alpha_t C\right)\nabla L\left(w_{t-1}^b\right)\big|\mathcal{F}_0^b\right]\right\|^2$$

$$= \mathsf{var}\left(B\left(I - \alpha_t C\right)\nabla L\left(w_{t-1}^b\right)\big|\mathcal{F}_0\right) + \alpha_t^2 c_b\left(\frac{1}{n}\sum_{i=1}^n\mathbb{E}\left[\left\|BC\nabla L_i\left(w_{t-1}^b\right)\right\|^2\Big|\mathcal{F}_0\right] - \mathbb{E}\left[\left\|BC\nabla L\left(w_{t-1}^b\right)\right\|^2\Big|\mathcal{F}_0\right]\right)$$

(A.2) $$= \mathsf{var}\left(B\left(I - \alpha_t C\right)\nabla L\left(w_{t-1}^b\right)\big|\mathcal{F}_0\right) + \frac{\alpha_t^2 c_b}{n^2}\sum_{i\neq j}\mathbb{E}\left[\left\|BC\nabla L_i\left(w_{t-1}^b\right) - BC\nabla L_j\left(w_{t-1}^b\right)\right\|^2\Big|\mathcal{F}_0\right],\qquad \square$$

where (A.1) is by Lemma 3.1. By induction, we know that the first term of (A.2) is a decreasing function of $b$. Taking $A_i = BC, A_j = -BC, A_k = 0, k \in [n]\backslash\{i,j\}$ in Theorem 3.3, we know that $\mathbb{E}\left[\left\|BC\nabla L_i\left(w_{t-1}^b\right) - BC\nabla L_j\left(w_{t-1}^b\right)\right\|^2\Big|\mathcal{F}_0\right]$ is also a decreasing function of $b$. Note that $\frac{\alpha_t^2 c_b}{n^2}$ decreases as $b$ increases. By Lemma A.1 we learn that (A.2) is a decreasing function of $b$ and hence we have completed the induction.

*Proof of Theorem A.4.* We have

$$\mathsf{var}\left(Bg_t^b\,|\mathcal{F}_0\right) = \mathbb{E}\left[\left\|Bg_t^b\right\|^2\Big|\mathcal{F}_0\right] - \left\|\mathbb{E}\left[Bg_t^b|\mathcal{F}_0\right]\right\|^2$$

$$= \mathbb{E}\left[\mathbb{E}\left[\left\|Bg_t^b\right\|^2\Big|\mathcal{F}_t^b\right]\Big|\mathcal{F}_0\right] - \left\|\mathbb{E}\left[\mathbb{E}\left[Bg_t^b|\mathcal{F}_t^b\right]|\mathcal{F}_0\right]\right\|^2$$

$$= c_b\left(\frac{1}{n}\sum_{i=1}^{n}\mathbb{E}\left[\left\|B\nabla L_i\left(w_t^b\right)\right\|^2\Big|\mathcal{F}_0\right] - \mathbb{E}\left[\left\|B\nabla L\left(w_t^b\right)\right\|^2\Big|\mathcal{F}_0\right]\right)$$

$$+ \mathbb{E}\left[\left\|B\nabla L\left(w_t^b\right)\right\|^2\Big|\mathcal{F}_0\right] - \left\|\mathbb{E}\left[B\nabla L\left(w_t^b\right)|\mathcal{F}_0\right]\right\|^2$$

$$= \frac{c_b}{n^2}\sum_{i\neq j}\mathbb{E}\left[\left\|B\nabla L_i\left(w_t^b\right) - B\nabla L_j\left(w_t^b\right)\right\|^2\Big|\mathcal{F}_0\right] + \mathsf{var}\left(B\nabla L\left(w_t^b\right)|\mathcal{F}_0\right).$$

Taking $A_i = B, A_j = -B, A_k = 0, k \in [n]\backslash\{i,j\}$ in Theorem 3.3, we know that $\mathbb{E}\left[\left\|B\nabla L_i\left(w_t^b\right) - B\nabla L_j\left(w_t^b\right)\right\|^2\Big|\mathcal{F}_0\right]$ is a decreasing and non-negative function of $b$ for all $i, j \in [n]$. By Theorem A.3, we know that $\mathsf{var}\left(B\nabla L\left(w_t^b\right)|\mathcal{F}_0\right)$ is also a decreasing function of $b$. Therefore, $\mathsf{var}\left(Bg_t^b\,|\mathcal{F}_0\right)$, as the sum of two decreasing functions of $b$, is also a decreasing function of $b$. □

*Proof of Corollary 3.5.* Simply taking $B = I_p$ in Theorem 3.3 yields the proof. □

**A.2. Proofs for Results in 3.2.** We provide a comprehensive proof of the two-layer linear network in Appendix A.2.1. We defer the extension from linear networks to polynomially-activated networks in Appendix A.2.2.

**A.2.1. Two-layer Linear Networks.** Given a distribution $\mathcal{D}$ in $\mathbb{R}^p$, we consider the population loss $\mathcal{L}(w) = \mathbb{E}_{x\sim\mathcal{D}}\left[\frac{1}{2}\left\|W_2W_1x - W_2^*W_1^*x\right\|^2\right]$ under the teacher-student learning framework [14] with $w = (W_1, W_2)$ a tuple of two matrices. Here $W_1 \in \mathbb{R}^{p_1\times p}$ and $W_2 \in \mathbb{R}^{p_2\times p_1}$ are parameter matrices of the student network and $W_1^*$ and $W_2^*$ are the fixed ground-truth parameters of the teacher network. We use online SGD to minimize the population loss $\mathcal{L}(w)$. Formally, we first choose a mini-batch size $b$ and initial weight matrices $\{W_{0,1}, W_{0,2}\}$; in each iteration $t$, we independently draw a mini-batch $\mathcal{B}_t^b := \left\{x_{t,i}^b : i \in [b]\right\}$ of $b$ samples from $\mathcal{D}$ and update the weight matrices by $W_{t+1,1}^b = W_{t,1}^b - \alpha_t g_{t,1}^b$ and $W_{t+1,2}^b = W_{t,2}^b - \alpha_t g_{t,2}^b$, where

$$(A.3)\quad g_{t,1}^b := \frac{1}{b}\sum_{i=1}^{b}\nabla_{W_{t,1}^b}\left(\frac{1}{2}\left\|W_{t,2}^bW_{t,1}^bx_{t,i}^b - W_2^*W_1^*x_{t,i}^b\right\|^2\right) = \frac{1}{b}\sum_{i=1}^{b}\left(W_{t,2}^b\right)^T\mathcal{W}_t^bx_{t,i}^b\left(x_{t,i}^b\right)^T,$$

$$(A.4)\quad g_{t,2}^b := \frac{1}{b}\sum_{i=1}^{b}\nabla_{W_{t,2}^b}\left(\frac{1}{2}\left\|W_{t,2}^bW_{t,1}^bx_{t,i}^b - W_2^*W_1^*x_{t,i}^b\right\|^2\right) = \frac{1}{b}\sum_{i=1}^{b}\mathcal{W}_t^bx_{t,i}^b\left(x_{t,i}^b\right)^T\left(W_{t,1}^b\right)^T.$$

Here $\mathcal{W}_t^b := W_{t,2}^bW_{t,1}^b - W_2^*W_1^*$ denotes the gap between the product of model weights and ground-truth weights and the derivation follows from the formulas in [33].

To recap, we use $\deg\left(A;\mathcal{M}\right)$ to denote the number of appearances of matrix $A$ and its transpose $A^T$ in a multi-set of matrices $\mathcal{M} = \{M_1,\ldots,M_n\}$. Mathematically, we have $\deg\left(A;\mathcal{M}\right) := \sum_{i\in[n]}\left(\mathbb{I}\{A = A_i\} + \mathbb{I}\{A^T = A_i\}\right)$. We further denote $\deg\left(\mathcal{A};\mathcal{M}\right) := \sum_{A\in\mathcal{A}}\deg\left(A;\mathcal{M}\right)$ for any set of matrices $\mathcal{A}$. We denote $W_t^b := \{W_{t,1}^b, W_{t,2}^b\}$, $W^* := \{W_1^*, W_2^*\}$ and $G_t^b := \{g_{t,1}^b, g_{t,2}^b\}$.

In Section 3.2, we use $\mathcal{C}$ to denote the infinite set of all non-random matrices given $\mathcal{F}_0$. Here we provide the precise definition of $\mathcal{C}$ as follows. For $n \in \mathbb{N}^+$, we use $e_{n,i}, i \in [n]$ to denote the $i$-th unit vector of $\mathbb{R}^n$. We denote $\mathcal{I} = \{I_n : n \in \mathbb{N}^+\}$ as the

collection of identity matrices and we define a set of (infinite many) matrices

$$
\mathcal{C} := \left\{
\begin{array}{l}
\mathbb{E}_{x_{t,i}^b \sim \mathcal{D}, i \in [b]} \left[ (e_{p,u}^T z_0)(e_{p,v}^T \overline{z}_0) \left[ \left( y_1 \overline{y}_1^T \right) \otimes \cdots \otimes \left( y_m \overline{y}_m^T \right) \otimes \left( z_1 \overline{z}_1^T \right) \otimes \cdots \otimes \left( z_n \overline{z}_n^T \right) \right] \right] : \\[4pt]
\quad y_i = e_{p,j_1^i} \otimes \cdots \otimes e_{p,j_{m_i}^i} \otimes x_{t,s_i}^b \otimes e_{p,k_1^i} \otimes \cdots \otimes e_{p,k_{n_i}^i}, \\[4pt]
\quad \overline{y}_i = e_{p,\overline{j}_1^i} \otimes \cdots \otimes e_{p,\overline{j}_{m_i}^i} \otimes x_{t,\overline{s}_i}^b \otimes e_{p,\overline{k}_1^i} \otimes \cdots \otimes e_{p,\overline{k}_{n_i}^i}, \\[4pt]
\quad z_0 \in \left\{ x_{t,i}^b : i \in [b] \right\} \bigcup \{e_{p,u}\}, \overline{z}_0 \in \left\{ x_{t,i}^b : i \in [b] \right\} \bigcup \{e_{p,v}\}, u,v \in [p], \\[4pt]
\quad z_j, \overline{z}_j \in \left\{ x_{t,i}^b : i \in [b] \right\}, j \in [n], \\[4pt]
\quad j_\alpha^i, \overline{j}_\alpha^i, k_\beta^i, \overline{k}_\beta^i \in [p], \alpha \in [m_i], \beta \in [n_i], i \in [m], \\[4pt]
\quad m_i, n_i \in \mathbb{N}, s_i, \overline{s}_i \in [b], i \in [m], \\[4pt]
\quad m, n \in \mathbb{N}, t \in \mathbb{N}^+
\end{array}
\right\}
$$

where $p$ is the dimension of the samples and $x_{t,s}^b, s \in [b]$ are the random samples we use to build the stochastic gradient at step $t$ and thus every element of $\mathcal{C}$ is a constant matrix under $\mathcal{F}_0$. Note that $\mathcal{C}$ is a union over all $m, n, m_i, n_i \in \mathbb{N}$ and $t \in \mathbb{N}^+$. We also point out that when $z_0 = e_{p,u}, \overline{z}_0 = e_{p,v}$, the leading scalar terms are 1. We also denote $\mathcal{E} := \left\{ e_{p,i} e_{p,j}^T : i,j \in [p] \right\}$ and $\overline{\mathcal{C}} := \mathcal{C} \bigcup \mathcal{I} \bigcup \mathcal{E}$. Note that every element of $\overline{\mathcal{C}}$ is a non-random matrix under $\mathcal{F}_0$ and $\overline{\mathcal{C}}$ is an infinite set of matrices that we use in the following proofs as auxiliary matrices.

Let $g_{t,1,s}^b := \left( W_{t,2}^b \right)^T \cdot \mathcal{W}_t^b \cdot \left( x_{t,s}^b \left( x_{t,s}^b \right)^T \right)$ and $g_{t,2,s}^b := \mathcal{W}_t^b \cdot \left( x_{t,s}^b \left( x_{t,s}^b \right)^T \right) \cdot W_{t,1}^b, s \in [b]$ denote the stochastic gradient with respect to the sample $x_{t,s}^b$ at time step $t$. We have $g_{t,i}^b = \frac{1}{b} \sum_{s \in [d]} g_{t,i,s}^b, i = 1,2$. Recall that we denote $W_t^b = \left\{ W_{t,1}^b, W_{t,2}^b \right\}$, $W^* = \{W_1^*, W_2^*\}$ and $G_t^b = \{g_{t,1}^b, g_{t,2}^b\}$ in Section 3.2. We further denote $\overline{G}_t^b = \left\{ g_{t,i,s}^b : s \in [b], i = 1,2 \right\}$ and $X_t^b = \left\{ x_{t,s}^b \left( x_{t,s}^b \right)^T : s \in [b] \right\}$. For simplicity, we denote $G_{t_1:t_2}^b := \bigcup_{t=t_1}^{t_2} G_t^b$ and $W_{t_1:t_2}^b := \bigcup_{t=t_1}^{t_2} W_t^b$.

Throughout the discussion of this section, we define the term that a matrix $A$ "takes values in" or "belongs to" a multi-set $\mathcal{A}$ if either $A$ or $A^T$ are in $\mathcal{A}$. We also abuse the notation $A \in \mathcal{A}$ to denote $A$ is in $\mathcal{A}$ or $A^T$ is in $A$.

LEMMA A.5. *For matrices $M_{i,j}, i \in [m], j \in [n]$ with appropriate dimensions, we have $\bigotimes_{i \in [m]} \left( \prod_{j \in [n]} M_{i,j} \right) = \prod_{j \in [n]} \left( \bigotimes_{i \in [m]} M_{i,j} \right)$.*

*Proof.* It is easy to prove by induction on $m$ and $n$ and by the fact that $(A \otimes B)(C \otimes D) = (AC) \otimes (BD)$ for any matrices $A, B, C, D$. □

**Remark**. If we view the multi-set $\mathcal{M} := \{M_{i,j}, i \in [m], j \in [n]\}$ as a matrix of matrices

$$
\mathcal{M} : \begin{bmatrix} M_{1,1} & M_{1,2} & M_{1,3} & \cdots & M_{1,n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ M_{m,1} & M_{m,2} & M_{m,3} & \cdots & M_{m,n} \end{bmatrix},
$$

then $\bigotimes_{i \in [m]} \left( \prod_{j \in [n]} M_{i,j} \right)$ can be regarded as first multiplying the entries of $\mathcal{M}$ within each row and then using the Kronecker product to multiply all of the rows. Similarly, $\prod_{j \in [n]} \left( \bigotimes_{i \in [m]} M_{i,j} \right)$ can be regarded as first using the Kronecker product to multiply all the entries of a column, then multiplying all the rows. Lemma A.5 shows that these two calculations on multi-set $\mathcal{M}$ give the same resulting matrices. We frequently

use this lemma in the following proofs. We give illustrations of the multi-sets to help readers better understand and follow the proofs.

LEMMA A.6. *Given two distributions $\mathcal{D}_1$ and $\mathcal{D}_2$ in $\mathbb{R}^{p_1}$ and $\mathbb{R}^{p_2}$, respectively. Given $y_1, \ldots, y_m \sim \mathcal{D}_1$, $z_1, \ldots z_n \sim \mathcal{D}_2$ and constant matrices $D_0, \ldots, D_n, A_1, \ldots, A_m$ with appropriate dimensions, we have*

$$\mathbb{E}_{y_i \sim \mathcal{D}_1, z_j \sim \mathcal{D}_2} \left[ D_0 z_1 z_n^T D_n \left( z_1^T D_1 z_2 \right) \cdots \left( z_{n-1}^T D_{n-1} z_n \right) \left( y_m^T A_m y_1 \right) \left( y_1^T A_1 y_2 \right) \cdots \left( y_{m-1}^T A_{m-1} y_m \right) \right]$$

$$= \sum_{u \in [p_1], v \in [p_2]} \left[ D_0 e_{p_1, u} e_{p_2, v}^T D_n \operatorname{tr} \left( C_{u,v} \left( \left( \bigotimes_{i=0}^{m-1} A_i \right) \otimes \left( \bigotimes_{j=1}^{n-1} D_i \right) \right) \right) \right]$$

*for some constant matrices $C_{u,v}$ specified in the proof.*

*Proof.* Let $y_0 := y_m$ and $A_0 := A_m$. We have

$$\prod_{i=0}^{m-1} \left( y_i^T A_i y_{i+1} \right) \prod_{j=1}^{n-1} \left( z_j^T D_i z_{j+1} \right)$$

$$= \prod_{i=0}^{m-1} \operatorname{tr} \left( y_i^T A_i y_{i+1} \right) \prod_{j=1}^{n-1} \operatorname{tr} \left( z_j^T D_i z_{j+1} \right)$$

$$= \prod_{i=0}^{m-1} \operatorname{tr} \left( y_{i+1} y_i^T A_i \right) \prod_{j=1}^{n-1} \operatorname{tr} \left( z_{j+1} z_j^T D_i \right)$$

$$= \operatorname{tr} \left( \left( \bigotimes_{i=0}^{m-1} \left( y_{i+1} y_i^T A_i \right) \right) \otimes \left( \bigotimes_{j=1}^{n-1} \left( z_{j+1} z_j^T D_i \right) \right) \right)$$

$$= \operatorname{tr} \left( \left( \left( \bigotimes_{i=0}^{m-1} \left( y_{i+1} y_i^T \right) \right) \otimes \left( \bigotimes_{j=1}^{n-1} \left( z_{j+1} z_j^T \right) \right) \right) \cdot \left( \left( \bigotimes_{i=0}^{m-1} A_i \right) \otimes \left( \bigotimes_{j=1}^{n-1} D_i \right) \right) \right),$$

where we use the fact that $\operatorname{tr}(A) \operatorname{tr}(B) = \operatorname{tr}(A \otimes B)$ for any matrices $A$ and $B$ in the second-to-last equation and use Lemma A.5 in the last equation. Further, note that $z_1 z_n^T = \sum_{u \in [p_1], v \in [p_2]} e_{p_1, u} e_{p_2, v}^T \left( e_{p_1, u}^T z_1 \right) \left( e_{p_2, v}^T z_n \right)$, we have

$$\mathbb{E}_{y_i \sim \mathcal{D}_1, z_j \sim \mathcal{D}_2} \left[ D_0 z_1 z_n^T D_n \left( z_1^T D_1 z_2 \right) \cdots \left( z_{n-1}^T D_{n-1} z_n \right) \left( y_m^T A_m y_1 \right) \left( y_1^T A_1 y_2 \right) \cdots \left( y_{m-1}^T A_{m-1} y_m \right) \right]$$

$$= \mathbb{E}_{y_i \sim \mathcal{D}_1, z_j \sim \mathcal{D}_2} \left[ \sum_{u \in [p_1], v \in [p_2]} D_0 \left( e_{p_1, u} e_{p_2, v}^T \left( e_{p_1, u}^T z_1 \right) \left( e_{p_2, v}^T z_n \right) \right) D_n \cdot \right.$$

$$\left. \cdot \operatorname{tr} \left( \left( \left( \bigotimes_{i=0}^{m-1} \left( y_{i+1} y_i^T \right) \right) \otimes \left( \bigotimes_{j=1}^{n-1} \left( z_{j+1} z_j^T \right) \right) \right) \cdot \left( \left( \bigotimes_{i=0}^{m-1} A_i \right) \otimes \left( \bigotimes_{j=1}^{n-1} D_i \right) \right) \right) \right]$$

$$= \sum_{u \in [p_1], v \in [p_2]} \mathbb{E}_{y_i \sim \mathcal{D}_1, z_j \sim \mathcal{D}_2} \left[ D_0 e_{p_1, u} e_{p_2, v}^T D_n \cdot \operatorname{tr} \left( \left( \left( e_{p_1, u}^T z_1 \right) \left( e_{p_2, v}^T z_n \right) \bigotimes_{i=0}^{m-1} \left( y_{i+1} y_i^T \right) \right) \otimes \left( \bigotimes_{j=1}^{n-1} \left( z_{j+1} z_j^T \right) \right) \right) \cdot \right.$$

$$\left. \cdot \left( \left( \bigotimes_{i=0}^{m-1} A_i \right) \otimes \left( \bigotimes_{j=1}^{n-1} D_i \right) \right) \right)$$

$$= \sum_{u \in [p_1], v \in [p_2]} \left[ D_0 e_{p_1, u} e_{p_2, v}^T D_n \operatorname{tr} \left( \mathbb{E}_{y_i \sim \mathcal{D}_1, z_j \sim \mathcal{D}_2} \left[ \left( e_{p_1, u}^T z_1 \right) \left( e_{p_2, v}^T z_n \right) \bigotimes_{i=0}^{m-1} \left( y_{i+1} y_i^T \right) \right) \otimes \left( \bigotimes_{j=1}^{n-1} \left( z_{j+1} z_j^T \right) \right) \right] \cdot \right.$$

$$\left. \cdot \left( \left( \bigotimes_{i=0}^{m-1} A_i \right) \otimes \left( \bigotimes_{j=1}^{n-1} D_i \right) \right) \right)$$

$$= \sum_{u \in [p_1], v \in [p_2]} \left[ D_0 e_{p_1, u} e_{p_2, v}^T D_n \operatorname{tr} \left( C_{u,v} \left( \left( \bigotimes_{i=0}^{m-1} A_i \right) \otimes \left( \bigotimes_{j=1}^{n-1} D_i \right) \right) \right) \right], \qquad \square$$

where

$$C_{u,v} = \mathbb{E}_{y_i \sim \mathcal{D}_1, z_j \sim \mathcal{D}_2} \left[ \left( e_{p_1, u}^T z_1 \right) \left( e_{p_2, v}^T z_n \right) \left( \bigotimes_{i=0}^{m-1} \left( y_{i+1} y_i^T \right) \right) \otimes \left( \bigotimes_{j=1}^{n-1} \left( z_{j+1} z_j^T \right) \right) \right].$$

LEMMA A.7. *Let* $\mathcal{M} := \{M_{i,j} : i \in [0:m], j \in [n]\}$ *be a multi-set of matrices such that each* $M_{i,j}$ *or its transpose only takes value in* $W_{0:t}^b \bigcup \overline{G}_t^b \bigcup G_{0:(t-1)}^b \bigcup W^* \bigcup \overline{\mathcal{C}}$ *and* $\deg\left(\overline{G}_t^b; \mathcal{M}\right) = d$ *(here* $d, m, n$ *are constants independent of* $b$*). Then for*

$$m' := m + d - 2, \quad n' := 6mn(d+1), \quad L := 2^d p^{d'(m-1)+2},$$

*where* $d' = \deg\left(\overline{G}_t^b; \{M_{i,j} : i \in [m], j \in [n]\}\right)$, *there exist multi-sets of matrices*

$$\mathcal{Q}_l := \{Q_{l,u,v} : u \in [0:m'], v \in [n']\}, l \in [L]$$

*such that*

$$\mathbb{E}\left[\text{tr}\left(C\left(\bigotimes_{i\in[m]}\left(\prod_{j\in[n]}M_{i,j}\right)\right)\right)\prod_{j\in[n]}M_{0,j}\bigg|\mathcal{F}_t^b\right] = \sum_{l\in[L]}c_l\text{tr}\left(C_l\left(\bigotimes_{u\in[m']}\left(\prod_{v\in[n']}Q_{l,u,v}\right)\right)\right)\prod_{v\in[n']}Q_{l,0,v},$$

*where* $c_l \in \{-1, +1\}$, $C, C_l \in \mathcal{C}$ *and* $Q_{l,u,v}$ *only takes value in* $W_{0:t}^b \bigcup G_{0:(t-1)}^b \bigcup W^* \bigcup \overline{\mathcal{C}}$, $u \in [0:m'], v \in [n'], l \in [L]$. *Further, for each* $l \in [L]$ *we have*

$$\deg\left(\overline{G}_t^b; \mathcal{Q}_l\right) = 0,$$

$$\deg\left(W_t^b; \mathcal{Q}_l\right) \leqslant \deg\left(W_t^b; \mathcal{M}\right) + 3d,$$

$$\deg\left(W^*; \mathcal{Q}_l\right) \leqslant \deg\left(W^*; \mathcal{M}\right) + 2d,$$

$$\deg\left(W_t^b; \mathcal{Q}_l\right) + \deg\left(W^*; \mathcal{Q}_l\right) = \deg\left(W_t^b; \mathcal{M}\right) + \deg\left(W^*; \mathcal{M}\right) + 3d,$$

$$\deg\left(W_f^b; \mathcal{Q}_l\right) = \deg\left(W_f^b; \mathcal{M}\right), \quad f \in [0:t-1],$$

$$\deg\left(G_f^b; \mathcal{Q}_l\right) = \deg\left(G_f^b; \mathcal{M}\right), \quad f \in [0:t-1].$$

*Proof.* Let $\overline{\mathcal{M}} := \{\overline{M}_{i,j} : i \in [0:m], j \in [3n]\}$ be the multi-set of matrices such that $M_{i,j} = \overline{M}_{i,3j-2} \cdot \overline{M}_{i,3j-1} \cdot \overline{M}_{i,3j}$, where

- if $M_{i,j} \in \overline{G}_t^b$ and $M_{i,j} = g_{t,1,i_0}^b = \left(W_{t,2}^b\right)^T \mathcal{W}_t^b \left(x_{t,i_0}^b \left(x_{t,i_0}^b\right)^T\right)$ for some $i_0 \in [b]$, then we set $\overline{M}_{i,3j-2} = \left(W_{t,2}^b\right)^T, \overline{M}_{i,3j-1} = \mathcal{W}_t^b$ and $\overline{M}_{i,3j} = x_{t,i_0}^b \left(x_{t,i_0}^b\right)^T$; the case of $M_{i,j} = g_{t,2,i_0'}^b$ for some $i_0' \in [b]$ is similar;
- if $M_{i,j} \notin \overline{G}_t^b$, then we set $\overline{M}_{i,3j-2} = M_{i,j}$ and $\overline{M}_{i,3j-1} = \overline{M}_{i,3j} = I$, where $I$ is an identity matrix with an appropriate dimension[5].

Figure 6 shows the transformation from $\mathcal{M}$ to $\overline{\mathcal{M}}$. By this transformation, we have

(A.5) $$\prod_{j\in[n]}M_{i,j} = \prod_{j\in[3n]}\overline{M}_{i,j}, \quad i \in [0:m],$$

---

[5]In the following, we use $I$ to denote an identity matrix with an appropriate dimension, without specifying the dimension. Readers should be able to infer the dimension easily from the matrices that this identity matrix is multiplied with.

$$\mathcal{M}: \begin{bmatrix} M_{0,1} & M_{0,2} & \cdots & M_{0,n} \\ M_{1,1} & M_{1,2} & \cdots & M_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ M_{m,1} & M_{m,2} & \cdots & M_{m,n} \end{bmatrix}$$

$$\overline{\mathcal{M}}: \begin{bmatrix} \overline{M}_{0,1} & \overline{M}_{0,2} & \overline{M}_{0,3} & \overline{M}_{0,4} & \overline{M}_{0,5} & \overline{M}_{0,6} & \cdots & \overline{M}_{0,3n-2} & \overline{M}_{0,3n-1} & \overline{M}_{0,3n} \\ \overline{M}_{1,1} & \overline{M}_{1,2} & \overline{M}_{1,3} & \overline{M}_{1,4} & \overline{M}_{1,5} & \overline{M}_{1,6} & \cdots & \overline{M}_{1,3n-2} & \overline{M}_{1,3n-1} & \overline{M}_{1,3n} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ \overline{M}_{m,1} & \overline{M}_{m,2} & \overline{M}_{m,3} & \overline{M}_{m,4} & \overline{M}_{m,5} & \overline{M}_{m,6} & \cdots & \overline{M}_{m,3n-2} & \overline{M}_{m,3n-1} & \overline{M}_{m,3n} \end{bmatrix}$$

Fig. 6: The transformation from $\mathcal{M}$ to $\overline{\mathcal{M}}$.

where each $\overline{M}_{i,j} \in W^b_{0:t} \bigcup G^b_{0:(t-1)} \bigcup W^* \bigcup \{\mathcal{W}^b_t\} \bigcup X^b_t \bigcup \overline{\mathcal{C}}$ and

$$\deg\left(W^b_t; \overline{\mathcal{M}}\right) = \deg\left(W^b_t; \mathcal{M}\right) + \deg\left(\overline{G}^b_t, \mathcal{M}\right) = \deg\left(W^b_t; \mathcal{M}\right) + d,$$

$$\deg\left(W^*; \overline{\mathcal{M}}\right) = \deg\left(W^*; \mathcal{M}\right),$$

$$\deg\left(\mathcal{W}^b_t; \overline{\mathcal{M}}\right) = \deg\left(\overline{G}^b_t; \mathcal{M}\right) = d,$$

$$\deg\left(X^b_t; \overline{\mathcal{M}}\right) = \deg\left(\overline{G}^b_t; \mathcal{M}\right) = d,$$

$$\deg\left(\overline{G}^b_t; \overline{\mathcal{M}}\right) = 0,$$

$$\deg\left(W^b_f; \overline{\mathcal{M}}\right) = \deg\left(W^b_f; \mathcal{M}\right), \quad f \in [0:t-1],$$

$$\deg\left(G^b_f; \overline{\mathcal{M}}\right) = \deg\left(G^b_f; \mathcal{M}\right), \quad f \in [0:t-1].$$

Further, let $\widetilde{\mathcal{M}} := \left\{ \widetilde{M}_{i,j} : i \in [0:m], j \in [3mn] \right\}$ be a multi-set of matrices such that

(A.6)
$$\widetilde{M}_{i,j} := \begin{cases} \overline{M}_{i,j} & 1 \leqslant i \leqslant m, 3 \cdot (i-1) \cdot n + 1 \leqslant j \leqslant 3 \cdot i \cdot n, \\ \overline{M}_{i,j} & i = 0, 1 \leqslant j \leqslant 3n, \\ I & \text{otherwise,} \end{cases}$$

where $I$ denotes an identity matrix with an appropriate dimension. Figure 7 shows the transformation from $\overline{\mathcal{M}}$ to $\widetilde{\mathcal{M}}$.

Fig. 7: The transformation from $\overline{\mathcal{M}}$ to $\widetilde{\mathcal{M}}$.

Then we have

$$\deg\left(W_t^b; \widetilde{\mathcal{M}}\right) = \deg\left(W_t^b; \overline{\mathcal{M}}\right) = \deg\left(W_t^b; \mathcal{M}\right) + d,$$

$$\deg\left(W^*; \widetilde{\mathcal{M}}\right) = \deg\left(W^*; \overline{\mathcal{M}}\right) = \deg\left(W^*; \mathcal{M}\right),$$

$$\deg\left(\mathcal{W}_t^b; \widetilde{\mathcal{M}}\right) = \deg\left(\mathcal{W}_t^b; \overline{\mathcal{M}}\right) = d,$$

$$\deg\left(X_t^b; \widetilde{\mathcal{M}}\right) = \deg\left(X_t^b; \overline{\mathcal{M}}\right) = d,$$

$$\deg\left(W_f^b; \widetilde{\mathcal{M}}\right) = \deg\left(W_f^b; \overline{\mathcal{M}}\right) = \deg\left(W_f^b; \mathcal{M}\right), \quad f \in [0 : t-1],$$

$$\deg\left(G_f^b; \widetilde{\mathcal{M}}\right) = \deg\left(G_f^b; \overline{\mathcal{M}}\right) = \deg\left(W_f^b; \mathcal{M}\right), \quad f \in [0 : t-1].$$

By (A.5), (A.6) and Lemma A.5, we have

(A.7)
$$\bigotimes_{i\in[m]}\left(\prod_{j\in[n]} M_{i,j}\right) = \bigotimes_{i\in[m]}\left(\prod_{j\in[3n]} \overline{M}_{i,j}\right) = \bigotimes_{i\in[m]}\left(\prod_{j\in[3mn]} \widetilde{M}_{i,j}\right) = \prod_{j\in[3mn]}\left(\bigotimes_{i\in[m]} \widetilde{M}_{i,j}\right)$$

and

(A.8)
$$\prod_{j\in[n]} M_{0,j} = \prod_{j\in[3n]} \overline{M}_{0,j} = \prod_{j\in[3mn]} \widetilde{M}_{0,j}.$$

If we denote

$$d_0 := \deg\left(X_t^b; \left\{\widetilde{M}_{0,j} : j \in [3mn]\right\}\right) = \deg\left(\overline{G}_t^b; \{M_{0,j} : j \in [n]\}\right)$$

and

$$d' := \deg\left(X_t^b; \left\{\widetilde{M}_{i,j} : i \in [m], j \in [3mn]\right\}\right) = \deg\left(\overline{G}_t^b; \{M_{i,j} : i \in [m], j \in [n]\}\right),$$

then we have $d_0 + d' = \deg\left(\overline{G}_t^b, \mathcal{M}\right) = d$.

Without loss of generalization, we assume that $d_0 > 0$ and $d' > 0$ (the case of $d_0 = 0$ or $d' = 0$ are simpler than the general case we discuss below and can be derived directly from the following arguments).

Note that for any $j \in [3mn]$, the multi-set $\widetilde{\mathcal{M}}_j := \left\{ \widetilde{M}_{i,j} : i \in [m] \right\}$[6] contains at most one element that is not an identity matrix. Thus, there exist exactly $d'$ pairs of indices $(i_1, j_1), \ldots, (i_{d'}, j_{d'}), 1 \leqslant j_1 < \cdots < j_{d'} \leqslant 3mn, i_k \in [m], k \in [d']$ such that $\widetilde{M}_{i_k, j_k} = x_{t,s_k}^b \left( x_{t,s_k}^b \right)^T \in X_t^b$ for some $s_k \in [b], k \in [d']$. By (A.6), for any $k \in [d']$, $\widetilde{M}_{i,j_k}$ is an identity matrix with an appropriate dimension if $i \neq j_k, i \in [m]$ (it is easy to see that $\widetilde{M}_{i,j_k} = I_p, i \neq j_k$, since $\widetilde{M}_{i_k, j_k} = x_{t,s_k}^b \left( x_{t,s_k}^b \right)^T \in \mathbb{R}^{p \times p}$). Thus, we can write $\bigotimes_{i \in [m]} \widetilde{M}_{i,j_k}$ in the following way

$$
\bigotimes_{i \in [m]} \widetilde{M}_{i,j_k}
$$

$$
= \underbrace{I_p \otimes \cdots \otimes I_p}_{(i_k - 1)\ I_p\text{'s}} \otimes \left( x_{t,s_k}^b \left( x_{t,s_k}^b \right)^T \right) \otimes \underbrace{I_p \otimes \cdots \otimes I_p}_{(m - i_k)\ I_p\text{'s}}
$$

$$
= \left( \sum_{q_1 \in [p]} e_{p,q_1} e_{p,q_1}^T \right) \otimes \cdots \otimes \left( \sum_{q_{i_k-1} \in [p]} e_{p,q_{i_k-1}} e_{p,q_{i_k-1}}^T \right) \otimes \left( x_{t,s_k}^b \left( x_{t,s_k}^b \right)^T \right) \otimes
$$

$$
\otimes \left( \sum_{q_{i_k+1} \in [p]} e_{p,q_{i_k+1}} e_{p,q_{i_k+1}}^T \right) \otimes \cdots \otimes \left( \sum_{q_m \in [p]} e_{p,q_m} e_{p,q_m}^T \right)
$$

$$
= \sum_{q_1,\ldots,q_{i_k-1},q_{i_k+1},\ldots,q_m \in [p]} \left( e_{p,q_1} e_{p,q_1}^T \right) \otimes \cdots \otimes \left( e_{p,q_{i_k-1}} e_{p,q_{i_k-1}}^T \right) \otimes \left( x_{t,s_k}^b \left( x_{t,s_k}^b \right)^T \right) \otimes
$$

$$
\otimes \left( e_{p,q_{i_k+1}} e_{p,q_{i_k+1}}^T \right) \otimes \cdots \otimes \left( e_{p,q_m} e_{p,q_m}^T \right)
$$

$$
= \sum_{q_1,\ldots,q_{i_k-1},q_{i_k+1},\ldots,q_m \in [p]} \left( e_{p,q_1} \otimes \cdots \otimes e_{p,q_{i_k-1}} \otimes x_{t,s_k}^b \otimes e_{p,q_{i_k+1}} \otimes \cdots \otimes e_{p,q_m} \right) \cdot
$$

$$
\cdot \left( e_{p,q_1} \otimes \cdots \otimes e_{p,q_{i_k-1}} \otimes x_{t,s_k}^b \otimes e_{p,q_{i_k+1}} \otimes \cdots \otimes e_{p,q_m} \right)^T
$$

(A.9)

$$
:= \sum_{q \in [p^{m-1}]} y_{t,k,q}^b \left( y_{t,k,q}^b \right)^T,
$$

where the second-to-last equation follows from Lemma A.5 and $y_{t,k,q}^b = e_{p,q_1} \otimes \cdots \otimes e_{p,q_{i_k-1}} \otimes x_{t,s_k}^b \otimes e_{p,q_{i_k+1}} \otimes \cdots \otimes e_{p,q_m}$ with $q - 1 = (q_1 - 1) + (q_2 - 1)p + \cdots + (q_{i_k-1} - 1)p^{i_k-2} + (q_{i_k+1} - 1)p^{i_k-1} + \cdots + (q_m - 1)p^{m-2}$. [7]

---

[6]Note that $M_{0,j} \notin \widetilde{\mathcal{M}}_j, j \in [3mn]$.

[7]Intuitively, this equation gives a one-to-one mapping between $\{(q_1, \ldots, q_{i_k-1}, q_{i_k+1}, \ldots, q_m) : q_1, \ldots, q_{i_k-1}, q_{i_k+1}, \ldots, q_m \in [p]\}$ and $\{q : q \in [p^{m-1}]\}$. In fact, $q_1 - 1, \ldots, q_{i_k-1} - 1, q_{i_k+1} - 1, \ldots, q_m - 1$ are the digits of the base-$p$ representation of $q - 1$.

874    If we denote

$$A_0 := \prod_{1 \leqslant j < j_1} \left( \bigotimes_{i \in [m]} \widetilde{M}_{i,j} \right),$$

875

$$A_k := \prod_{j_k < j < j_{k+1}} \left( \bigotimes_{i \in [m]} \widetilde{M}_{i,j} \right), \qquad 1 \leqslant k \leqslant d' - 1,$$

876

$$A_{d'} := \prod_{j_{d'} < j \leqslant 3mn} \left( \bigotimes_{i \in [m]} \widetilde{M}_{i,j} \right).$$

877
878

879    Figure 8 gives an intuition on how we split the multi-set $\widetilde{\mathcal{M}}$ to form quantities
880    $A_0, A_1, \ldots, A_{d'}$.

$$\widetilde{\mathcal{M}} : \begin{bmatrix} \widetilde{M}_{0,1} & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \widetilde{M}_{0,3mn} \\ \widetilde{M}_{1,1} & \cdots & \widetilde{M}_{1,j_1-1} & \widetilde{M}_{1,j_1} & \widetilde{M}_{1,j_1+1} & \cdots & \widetilde{M}_{1,j_2-1} & \widetilde{M}_{1,j_2} & \cdots & \cdots & \widetilde{M}_{1,j_d} & \widetilde{M}_{1,j_d+1} & \cdots & \widetilde{M}_{1,3mn} \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \widetilde{M}_{m,1} & \cdots & \widetilde{M}_{m,j_1-1} & \widetilde{M}_{m,j_1} & \widetilde{M}_{m,j_1+1} & \cdots & \widetilde{M}_{m,j_2-1} & \widetilde{M}_{m,j_2} & \cdots & \cdots & \widetilde{M}_{m,j_d} & \widetilde{M}_{m,j_d+1} & \cdots & \widetilde{M}_{m,3mn} \end{bmatrix}$$

$$\underbrace{\qquad}_{A_0} \qquad \underbrace{\qquad}_{A_1} \qquad \underbrace{\qquad}_{A_d}$$

Fig. 8: The formation of $A_0, A_1, \ldots, A_d$.

881    Combining (A.7) and (A.9), we have

882
$$\mathrm{tr}\left( C \left( \bigotimes_{i \in [m]} \left( \prod_{j \in [n]} M_{i,j} \right) \right) \right) = \mathrm{tr}\left( C \left( \prod_{j \in [3mn]} \left( \bigotimes_{i \in [m]} \widetilde{M}_{i,j} \right) \right) \right)$$

883
$$= \mathrm{tr}\left( C A_0 \left( \bigotimes_{i \in [m]} \widetilde{M}_{i,j_1} \right) A_1 \left( \bigotimes_{i \in [m]} \widetilde{M}_{i,j_2} \right) A_2 \cdots A_{d'-1} \left( \bigotimes_{i \in [m]} \widetilde{M}_{i,j_{d'}} \right) A_{d'} \right)$$

(A.10)

884
$$= \mathrm{tr}\left( C A_0 \left( \sum_{q_1 \in [p^{m-1}]} y^b_{t,1,q_1} \left( y^b_{t,1,q_1} \right)^T \right) A_1 \left( \sum_{q_2 \in [p^{m-1}]} y^b_{t,2,q_2} \left( y^b_{t,1,q_1} \right)^T \right) \cdot \right.$$

885
$$\left. \cdot A_2 \cdots A_{d'-1} \left( \sum_{q_{d'} \in [p^{m-1}]} y^b_{t,d',q_{d'}} \left( y^b_{t,d',q_{d'}} \right)^T \right) A_{d'} \right)$$

886
$$= \sum_{q_1, \ldots, q_{d'} \in [p^{m-1}]} \mathrm{tr}\left( C A_0 y^b_{t,1,q_1} \left( y^b_{t,1,q_1} \right)^T A_1 y^b_{t,2,q_2} \left( y^b_{t,1,q_1} \right)^T A_2 \cdots A_{d'-1} y^b_{t,d',q_{d'}} \left( y^b_{t,d',q_{d'}} \right)^T A_{d'} \right)$$

887
$$= \sum_{q_1, \ldots, q_{d'} \in [p^{m-1}]} \left( y^b_{t,d',q_{d'}} \right)^T A_{d'} C A_0 y^b_{t,1,q_1} \left( y^b_{t,1,q_1} \right)^T A_1 y^b_{t,2,q_2} \left( y^b_{t,1,q_1} \right)^T A_2 \cdots A_{d'-1} y^b_{t,d',q_{d'}}$$

(A.11)

888
$$= \sum_{q_1, \ldots, q_{d'} \in [p^{m-1}]} \left( \left( y^b_{t,d',q_{d'}} \right)^T A_{d'} C A_0 y^b_{t,1,q_1} \right) \left( \left( y^b_{t,1,q_1} \right)^T A_1 y^b_{t,2,q_2} \right) \cdots \left( \left( y^b_{t,d'-1,q_{d'-1}} \right)^T A_{d'-1} y^b_{t,d',q_{d'}} \right),$$
889

890    where we use the fact that $\mathrm{tr}\,(AB) = \mathrm{tr}\,(BA)$ for any matrices $A$ and $B$ with appropriate
891    dimension in the second-to-last equation.
892         Similarly, there exist exactly $d_0$ indices $l_1, \ldots, l_{d_0}$ such that $1 \leqslant l_1 < \cdots < l_{d_0} \leqslant$

$3mn$ and $\widetilde{M}_{0,l_k} = x^b_{t,r_k}\left(x^b_{t,r_k}\right)^T \in X^b_t$ for some $r_k \in [b], k \in [d_0]$. If we denote

$$D_0 := \prod_{1 \leqslant j < l_1} \widetilde{M}_{0,j},$$

$$D_k := \prod_{l_k < j < l_{k+1}} \widetilde{M}_{0,j}, \qquad 1 \leqslant k \leqslant d_0 - 1,$$

$$D_{d_0} := \prod_{l_{d_0} < j \leqslant 3mn} \widetilde{M}_{0,j},$$

then we have

$$\prod_{j \in [3mn]} \widetilde{M}_{0,j} = D_0 x^b_{t,r_1}\left(x^b_{t,r_1}\right)^T D_1 x^b_{t,r_2}\left(x^b_{t,r_2}\right)^T \cdots D_{d_0-1} x^b_{t,r_{d_0}}\left(x^b_{t,r_{d_0}}\right)^T D_{d_0}$$

$$\text{(A.12)} \qquad = D_0 x^b_{t,r_1}\left(x^b_{t,r_{d_0}}\right)^T D_{d_0}\left(\left(x^b_{t,r_1}\right)^T D_1 x^b_{t,r_2}\right)\left(\left(x^b_{t,r_2}\right)^T D_2 x^b_{t,r_3}\right)\cdots\left(\left(x^b_{t,r_{d_0-1}}\right)^T D_{d_0-1} x^b_{t,r_{d_0}}\right).$$

Combining (A.11), (A.12) and by Lemma A.6, we have

$$\mathbb{E}\left[\mathrm{tr}\left(C\left(\bigotimes_{i \in [m]}\left(\prod_{j \in [n]} M_{i,j}\right)\right)\right)\prod_{j \in [n]} M_{0,j}\Big|\mathcal{F}^b_t\right]$$

$$= \mathbb{E}\left[\mathrm{tr}\left(C\left(\bigotimes_{i \in [m]}\left(\prod_{j \in [3mn]} \widetilde{M}_{i,j}\right)\right)\right)\prod_{j \in [3mn]} \widetilde{M}_{0,j}\Big|\mathcal{F}^b_t\right]$$

$$= \sum_{q_1,\ldots,q_{d'} \in [p^{m-1}]} \mathbb{E}\Big[D_0 x^b_{t,r_1}\left(x^b_{t,r_{d_0}}\right)^T D_{d_0}\left(\left(x^b_{t,r_1}\right)^T D_1 x^b_{t,r_2}\right)\cdots\left(\left(x^b_{t,r_{d_0-1}}\right)^T D_{d_0-1} x^b_{t,r_{d_0}}\right)\cdot$$

$$\cdot\left(\left(y^b_{t,d',q_{d'}}\right)^T A_{d'} C A_0 y^b_{t,1,q_1}\right)\left(\left(y^b_{t,1,q_1}\right)^T A_1 y^b_{t,2,q_2}\right)\cdots\left(\left(y^b_{t,d'-1,q_{d'-1}}\right)^T A_{d'-1} y^b_{t,d',q_{d'}}\right)\Big|\mathcal{F}^b_t\Big]$$

$$\text{(A.13)} \qquad = \sum_{q_i \in [p^{m-1}]}\sum_{p_1,p_2 \in [p]} D_0 e_{p,p_1} e^T_{p,p_2} D_{d_0} \mathrm{tr}\left(C_{q_1,\ldots,q_{d'},p_1,p_2}\left(\left(A_{d'} C A_0\right) \otimes A_1 \otimes \cdots A_{d'-1} \otimes D_1 \otimes \cdots D_{d_0-1}\right)\right),$$

where the exact value of $C_{q_1,\ldots,q_{d'},p_1,p_2}$ is available in Lemma A.6.

Finally, it remains to show that $(A_{d'} C A_0) \otimes A_1 \otimes \cdots \otimes A_{d'-1} \otimes D_1 \otimes \cdots D_{d_0-1}$ can be written in the form of $\bigotimes\left(\prod M'_{i',j'}\right)$. To this end, let $\{B_{i,j} : i \in [d-1], j \in [d+1]\}$ be a multi-set of matrices such that $B_{1,1} = A_{d'}, B_{1,2} = C, B_{i,i+2} = A_{i-1}, i \in [d'], B_{d'+i,d'+i+2} = D_i, i \in [d_0 - 1]$ and $B_{i,j} = I$ otherwise. Following is an illustration of the multi-set $\{B_{i,j} : i \in [d-1], j \in [d+1]\}$.

$$\{B_{i,j}\}_{(d-1)\times(d+1)} : \begin{bmatrix} A_{d'} & C & A_0 & I & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & I \\ I & I & I & A_1 & I & \cdots & \cdots & \cdots & \cdots & \cdots & I \\ I & I & I & I & A_2 & I & \cdots & \cdots & \cdots & \cdots & I \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \vdots \\ I & \cdots & \cdots & \cdots & \cdots & I & A_{d'-1} & I & \cdots & \cdots & I \\ I & \cdots & \cdots & \cdots & \cdots & \cdots & I & D_0 & I & \cdots & I \\ I & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & I & D_1 & \cdots & I \\ I & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & I & \ddots & I \\ I & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & I & D_{d_0-1} \end{bmatrix}$$

We have

(A.14)

$$(A_{d'} C A_0) \otimes A_1 \otimes \cdots \otimes A_{d'-1} \otimes D_1 \otimes \cdots D_{d_0-1} = \bigotimes_{i \in [d-1]}\left(\prod_{j \in [d+1]} B_{i,j}\right) = \prod_{j \in [d+1]}\left(\bigotimes_{i \in [d-1]} B_{i,j}\right).$$

Note that for each $j \in [d+1]$, there is at most one element of $\{B_{i,j} : i \in [d-1]\}$ that is not an identity matrix. We next show that, for each $j \in [d+1]$, $\bigotimes_{i\in[d-1]} B_{i,j}$ can be written as a product of the Kronecker product of some matrices of the form

(A.15)
$$\bigotimes_{i\in[d-1]} B_{i,j} = \prod_{j'\in[3mn]} \left( \bigotimes_{i'\in[m+d-2]} \widehat{M}_{j,i',j'} \right).$$

In fact, for $j = 1$ we have

$$\bigotimes_{i\in[d-1]} B_{i,1}$$

$$= A_{d'} \otimes \underbrace{I \otimes \cdots \otimes I}_{(d-2)\ I\text{'s}}$$

$$= \left[ \prod_{j_{d'}<j'\leqslant 3mn} \left( \bigotimes_{i'\in[m]} \widetilde{M}_{i',j'} \right) \right] \otimes I \otimes \cdots \otimes I$$

$$= \left[ \bigotimes_{i'\in[m]} \left( \prod_{j_{d'}<j'\leqslant 3mn} \widetilde{M}_{i',j'} \right) \right] \otimes I \otimes \cdots \otimes I$$

$$= \left( \prod_{j_{d'}<j'\leqslant 3mn} \widetilde{M}_{1,j'} \right) \otimes \cdots \otimes \left( \prod_{j_{d'}<j'\leqslant 3mn} \widetilde{M}_{m,j'} \right) \otimes \left[ \prod_{j_{d'}<j'\leqslant 3mn} I \right] \otimes \cdots \otimes \left[ \prod_{j_{d'}<j'\leqslant 3mn} I \right]$$

$$= \prod_{j_{d'}<j'\leqslant 3mn} \left[ \left( \bigotimes_{i'\in[m]} \widetilde{M}_{i',j'} \right) \otimes \underbrace{I \otimes \cdots \otimes I}_{(d-2)\ I\text{'s}} \right]$$

$$= \left\{ \prod_{j'\leqslant j_{d'}} \left[ \underbrace{I \otimes \cdots \otimes I}_{(m+d-2)\ I\text{'s}} \right] \right\} \cdot \left\{ \prod_{j_{d'}<j'\leqslant 3mn} \left[ \left( \bigotimes_{i'\in[m]} \widetilde{M}_{i',j'} \right) \otimes \underbrace{I \otimes \cdots \otimes I}_{(d-2)\ I\text{'s}} \right] \right\}$$

$$:= \prod_{j'\in[3mn]} \left( \bigotimes_{i'\in[m+d-2]} \widehat{M}_{1,i',j'} \right).$$

The case of $j = 3$ (and thus $\bigotimes_{i\in[d-1]} B_{i,3} = A_0 \otimes \underbrace{I \otimes \cdots \otimes I}_{(d-2)\ I\text{'s}}$) is similar to $j = 1$.

For $j = 2$, we have

$$\bigotimes_{i\in[d-1]} B_{i,2}$$

$$= C \otimes \underbrace{I \otimes \cdots \otimes I}_{(d-2)\ I\text{'s}} = C \otimes \underbrace{I \otimes \cdots \otimes I}_{(d-2)\ I\text{'s}} \otimes \underbrace{I \otimes \cdots \otimes I}_{(m-1)\ I\text{'s}}$$

$$= \left[ C \cdot \left( \prod_{j'\in[3mn-1]} I \right) \right] \otimes \left[ \bigotimes_{i'\in[d-2]} \left( \prod_{j'\in[3mn]} I \right) \right] \otimes \left[ \bigotimes_{i'\in[m-1]} \left( \prod_{j'\in[3mn]} I \right) \right]$$

$$= \left[ C \otimes \left( \bigotimes_{i'\in[d-2]} I \right) \otimes \left( \bigotimes_{i'\in[m-1]} I \right) \right] \cdot \prod_{j'\in[3mn-1]} \left[ \left( \bigotimes_{i'\in[d-1]} I \right) \otimes \left( \bigotimes_{i'\in[m-1]} I \right) \right]$$

$$:= \prod_{j'\in[3mn]} \left( \bigotimes_{i'\in[m+d-2]} \widehat{M}_{2,i',j'} \right).$$

For $4 \leqslant j \leqslant d' + 2$ (for clarity, we write $\bigotimes_{i\in[d-1]} B_{i,k}$ to replace $\bigotimes_{i\in[d-1]} B_{i,j}$ for

$4 \leqslant k \leqslant d' + 2$ so that we avoid the conflict of $j$ and $j_1, \ldots, j_{d'}$), we have

$$\bigotimes_{i \in [d-1]} B_{i,k}$$

$$= \underbrace{I \otimes \cdots \otimes I}_{(k-3) \ I\text{'s}} \otimes A_{k-3} \otimes \underbrace{I \otimes \cdots \otimes I}_{(d-k+1) \ I\text{'s}}$$

$$= \left( \bigotimes_{i' \in [k-3]} I \right) \otimes \left[ \prod_{j_{k-3} < j' < j_{k-2}} \left( \bigotimes_{i' \in [m]} \widetilde{M}_{i',j'} \right) \right] \otimes \left( \bigotimes_{i' \in [d-k+1]} I \right)$$

$$= \left[ \bigotimes_{i' \in [k-3]} \left( \prod_{j_{k-3} < j' < j_{k-2}} I \right) \right] \otimes \left[ \bigotimes_{i' \in [m]} \left( \prod_{j_{k-3} < j' < j_{k-2}} \widetilde{M}_{i',j'} \right) \right] \otimes \left[ \bigotimes_{i' \in [d-k+1]} \left( \prod_{j_{k-3} < j' < j_{k-2}} I \right) \right]$$

$$= \prod_{j_{k-3} < j' < j_{k-2}} \left[ \left( \bigotimes_{i' \in [k-3]} I \right) \otimes \left( \bigotimes_{i' \in [m]} \widetilde{M}_{i',j'} \right) \otimes \left( \bigotimes_{i' \in [d-k+1]} I \right) \right]$$

$$:= \prod_{j' \in [3mn]} \left( \bigotimes_{i' \in [m+d-2]} \widehat{M}_{k,i',j'} \right).$$

The case of $d' + 3 \leqslant j \leqslant d + 1$ is similar.

In conclusion, we build a multi-set of matrices

$$\widehat{\mathcal{M}} := \left\{ \widehat{M}_{j,i',j'} : j \in [d+1], i' \in [3mn], j' \in [m+d-2] \right\}$$

such that (A.15) holds for all $j \in [d+1]$ and each

$$\widehat{M}_{j,i',j'} \in \left\{ \widetilde{\mathcal{M}}_{i,j} : i \in [m], j \in [3mn], j \neq j_1, \ldots, j_{d'} \right\} \bigcup \left\{ \widetilde{\mathcal{M}}_{0,j} : j \neq l_1, \ldots, l_{d_0} \right\} \bigcup \mathcal{I}$$

only takes value in $W_{0:t}^b \bigcup G_{0:(t-1)}^b \bigcup W^* \bigcup \{\mathcal{W}_t^b\} \mathcal{I} \bigcup \{C\}$.

Further, if we denote multi-sets of matrices

$$\widehat{\mathcal{M}}_0^{p_1,p_2} := \left\{ \widehat{M}_{0,j}^{p_1,p_2} : j \in [3mn+1] \right\}, p_1, p_2 \in [p]$$

such that

(A.16)
$$\widehat{M}_{0,j}^{p_1,p_2} := \begin{cases} \widetilde{M}_{0,j} & 1 \leqslant j < l_1, \\ e_{p,p_1} e_{p,p_2}^T & j = l_1, \\ \widetilde{M}_{0,j-1} & l_{d_0} + 1 < j \leqslant 3mn + 1, \\ I & \text{otherwise}, \end{cases}$$

and by the representation of $\widehat{M}_{j,i',j'}$ above, we have

$$\mathbb{E} \left[ \text{tr} \left( C \left( \bigotimes_{i \in [m]} \left( \prod_{j \in [n]} M_{i,j} \right) \right) \right) \prod_{j \in [n]} M_{0,j} \middle| \mathcal{F}_t^b \right]$$

$$= \sum_{q_i \in [p^{m-1}]} \sum_{p_1, p_2 \in [p]} D_0 e_{p,p_1} e_{p,p_2}^T D_{d_0} \text{tr} \left( C_{q_1, \ldots, q_{d'}, p_1, p_2} \left( \left( A_{d'} C A_0 \right) \otimes A_1 \otimes \cdots A_{d'-1} \otimes D_1 \otimes \cdots D_{d_0-1} \right) \right)$$

(A.17)

$$= \sum_{q_i \in [p^{m-1}]} \sum_{p_1, p_2 \in [p]} \text{tr} \left( C_{q_1, \ldots, q_{d'}, p_1, p_2} \left( \prod_{j \in [d+1]} \prod_{j' \in [3mn]} \left( \bigotimes_{i' \in [m+d-2]} \widehat{M}_{j,i',j'} \right) \right) \right) \prod_{j \in [3mn+1]} \widehat{M}_{0,j}^{p_1,p_2}$$

and for each $p_1 \in [p]$ and $p_2 \in [p]$,

$$\deg\left(W_t^b; \widehat{\mathcal{M}}_0^{p_1,p_2}\right) + \deg\left(W_t^b; \widehat{\mathcal{M}}\right) = \deg\left(W_t^b; \widetilde{\mathcal{M}}\right) = \deg\left(W_t^b; \mathcal{M}\right) + d,$$

$$\deg\left(W^*; \widehat{\mathcal{M}}_0^{p_1,p_2}\right) + \deg\left(W^*; \widehat{\mathcal{M}}\right) = \deg\left(W^*; \widetilde{\mathcal{M}}\right) = \deg\left(W^*; \mathcal{M}\right),$$

$$\deg\left(\mathcal{W}_t^b; \widehat{\mathcal{M}}_0^{p_1,p_2}\right) + \deg\left(\mathcal{W}_t^b; \widehat{\mathcal{M}}\right) = \deg\left(\mathcal{W}_t^b; \widetilde{\mathcal{M}}\right) = d,$$

$$\deg\left(X_t^b; \widehat{\mathcal{M}}_0^{p_1,p_2}\right) + \deg\left(X_t^b; \widehat{\mathcal{M}}\right) = \sum_{j\in[3mn], j\neq j_1,\ldots,j_d} \deg\left(X_t^b; \widetilde{\mathcal{M}}_j\right) = 0,$$

$$\deg\left(W_f^b; \widehat{\mathcal{M}}_0^{p_1,p_2}\right) + \deg\left(W_f^b; \widehat{\mathcal{M}}\right) = \deg\left(W_f^b; \widetilde{\mathcal{M}}\right) = \deg\left(W_f^b; \mathcal{M}\right), \quad f \in [0:t-1]$$

$$\deg\left(G_f^b; \widehat{\mathcal{M}}_0^{p_1,p_2}\right) + \deg\left(G_f^b; \widehat{\mathcal{M}}\right) = \deg\left(G_f^b; \widetilde{\mathcal{M}}\right) = \deg\left(G_f^b; \mathcal{M}\right) \quad f \in [0:t-1].$$

For simplicity, let us denote

$$\prod_{j\in[d+1]} \prod_{i'\in[3mn]} \left( \bigotimes_{j'\in[m+d-2]} \widehat{M}_{j,i',j'} \right) := \prod_{v\in[3mn(d+1)]} \left( \bigotimes_{u\in[m+d-2]} N_{u,v} \right) = \bigotimes_{u\in[m+d-2]} \left( \prod_{v\in[3mn(d+1)]} N_{u,v} \right),$$

$$\prod_{j\in[3mn+1]} \widehat{M}_{0,j}^{p_1,p_2} := \prod_{v\in[3mn(d+1)]} N_{0,v}^{p_1,p_2},$$

where $N_{j',3mn(j-1)+i'} = \widehat{M}_{j,i',j'}, j' \in [m+d-2], j \in [d+1], i' \in [3mn], N_{0,j}^{p_1,p_2} = \widehat{M}_{0,j}^{p_1,p_2}, j \in [3mn+1], p_1, p_2 \in [p]$, and $N_{0,j}^{p_1,p_2} = I, 3mn+1 < j \leqslant 3mn(d+1), p_1, p_2 \in [p]$. Thus we have

$$\mathrm{tr}\left( C_{q_1,\ldots,q_{d'},p_1,p_2} \left( \prod_{j\in[d+1]} \prod_{j'\in[3mn]} \left( \bigotimes_{i'\in[m+d-2]} \widehat{M}_{j,i',j'} \right) \right) \right) \prod_{j\in[3mn+1]} \widehat{M}_{0,j}^{p_1,p_2}$$

(A.18)

$$= \mathrm{tr}\left( C_{q_1,\ldots,q_{d'},p_1,p_2} \left( \bigotimes_{u\in[m+d-2]} \left( \prod_{v\in[3mn(d+1)]} N_{u,v} \right) \right) \right) \prod_{v\in[3mn(d+1)]} N_{0,v}^{p_1,p_2}.$$

It remains to expand all appearance of $\mathcal{W}_t^b$ in the multi-sets

$$\mathcal{N} := \{N_{u,v} : u \in [m+d-2], v \in [3mn(d+1)]\}$$

and

$$\mathcal{N}_0^{p_1,p_2} := \left\{ N_{0,v}^{p_1,p_2} : v \in [3mn(d+1)] \right\}, p_1, p_2 \in [p].$$

In fact, for each $p_1 \in [p]$ and $p_2 \in [p]$, it is easy to see that

$$\deg\left(\mathcal{W}_t^b, \mathcal{N}_0^{p_1,p_2}\right) + \deg\left(\mathcal{W}_t^b, \mathcal{N}\right) = d.$$

Recall that $\mathcal{W}_t^b = W_{t,2}^b W_{t,1}^b - W_2^* W_1^*$. If we replace all appearance of $\mathcal{W}_t^b$ in (A.18) with $\left(W_{t,2}^b W_{t,1}^b - W_2^* W_1^*\right)$ and expand all parentheses, we have

$$\mathrm{tr}\left( C_{q_1,\ldots,q_{d'},p_1,p_2} \left( \bigotimes_{u\in[m+d-2]} \left( \prod_{v\in[3mn(d+1)]} N_{u,v} \right) \right) \right) \prod_{v\in[3mn(d+1)]} N_{0,v}^{p_1,p_2}$$

(A.19)

$$:= \sum_{l\in[2^d]} c_l \mathrm{tr}\left( C_{q_1,\ldots,q_{d'},p_1,p_2} \left( \bigotimes_{u\in[m+d-2]} \left( \prod_{v\in[6mn(d+1)]} \overline{N}_{u,v}^l \right) \right) \right) \prod_{v\in[6mn(d+1)]} \overline{N}_{0,v}^{l,p_1,p_2},$$

where $c_l \in \{-1, 1\}^8$ for $l \in [2^d]$. For each $u \in [m + d - 2]$ and $v \in [3mn(d+1)]$, the two consecutive matrices $\overline{N}^l_{u,2v-1}$ and $\overline{N}^l_{u,2v-1}$ equal to (i) either $W^b_{t,2}, W^b_{t,1}$ or $W^*_2, W^*_1$, respectively, if $N_{u,v} = \mathcal{W}^b_t$; (ii) $N_{u,v}$ and $I$, respectively. The same argument also holds for all $\overline{N}^{l,p_1,p_2}_{0,2v-1}$ and $\overline{N}^{l,p_1,p_2}_{0,2v}$, $v \in [3mn(d+1)]$. The summation comes from the fact that $\deg\left(\mathcal{W}^b_t, \mathcal{N}^{p_1,p_2}_0\right) + \deg\left(\mathcal{W}^b_t, \mathcal{N}\right) = d$ and thus we end up with $2^d$ terms of the Kronecker product of product of matrices.

Further, if we denote multi-sets of matrices

$$\overline{\mathcal{N}}^l := \left\{\overline{N}^l_{r,s} : r \in [m+d-1], s \in [6mn(d+1)]\right\}$$

and $\overline{\mathcal{N}}^{l,p_1,p_2}_0 := \left\{\overline{N}^{l,p_1,p_2}_{0,j} : j \in [6mn(d+1)]\right\}, p_1, p_2 \in [p], l \in [2^d]$, then the elements of $\overline{\mathcal{N}}^l$'s and $\overline{\mathcal{N}}^{l,p_1,p_2}_0$'s only take value in $W^b_{0:t} \bigcup G^b_{0:(t-1)} \bigcup W^* \bigcup \overline{\mathcal{C}}$. For each $l \in [2^d]$, $p_1 \in [p]$ and $p_2 \in [p]$, we have

$$\deg\left(W^b_t; \overline{\mathcal{N}}^l\right) + \deg\left(W^b_t; \overline{\mathcal{N}}^{l,p_1,p_2}_0\right) \leq \deg\left(W^b_t; \widehat{\mathcal{M}}\right) + 2\deg\left(\mathcal{W}^b_t; \widehat{\mathcal{M}}\right) = \deg\left(W^b_t; \mathcal{M}\right) + 3d,$$

$$\deg\left(W^*; \overline{\mathcal{N}}^l\right) + \deg\left(W^*; \overline{\mathcal{N}}^{l,p_1,p_2}_0\right) \leq \deg\left(W^*; \widehat{\mathcal{M}}\right) + 2\deg\left(\mathcal{W}^b_t; \widehat{\mathcal{M}}\right) = \deg\left(W^*; \mathcal{M}\right) + 2d,$$

$$\deg\left(\mathcal{W}^b_t; \overline{\mathcal{N}}^l\right) = 0,$$

$$\deg\left(W^b_f; \overline{\mathcal{N}}^l\right) + \deg\left(W^b_f; \overline{\mathcal{N}}^{l,p_1,p_2}_0\right) = \deg\left(W^b_f; \widehat{\mathcal{M}}\right) + \deg\left(W^b_f; \widehat{\mathcal{M}}^{p_1,p_2}_0\right) = \deg\left(W^b_f; \mathcal{M}\right), \quad f \in [0:t-1],$$

$$\deg\left(G^b_f; \overline{\mathcal{N}}^l\right) + \deg\left(G^b_f; \overline{\mathcal{N}}^{l,p_1,p_2}_0\right) = \deg\left(G^b_f; \widehat{\mathcal{M}}\right) + \deg\left(G^b_f; \widehat{\mathcal{M}}^{p_1,p_2}_0\right) = \deg\left(G^b_f; \mathcal{M}\right), \quad f \in [0:t-1].$$

and

$$\deg\left(W^b_t; \overline{\mathcal{N}}^{l,p_1,p_2}_0\right) + \deg\left(W^*; \overline{\mathcal{N}}^{l,p_1,p_2}_0\right) + \deg\left(W^b_t; \overline{\mathcal{N}}^l\right) + \deg\left(W^*; \overline{\mathcal{N}}^l\right)$$

$$= \deg\left(W^*; \widehat{\mathcal{M}}\right) + \deg\left(W^b_t; \widehat{\mathcal{M}}\right) + 2\deg\left(\mathcal{W}^b_t; \widehat{\mathcal{M}}\right)$$

$$= \deg\left(W^b_t; \mathcal{M}\right) + \deg\left(W^*; \mathcal{M}\right) + 3d.$$

Combining (A.17), (A.18) and (A.19), we have

$$\mathbb{E}\left[\mathrm{tr}\left(C\left(\bigotimes_{i\in[m]}\left(\prod_{j\in[n]} M_{i,j}\right)\right)\right)\Big| \mathcal{F}^b_t\right]$$

$$= \sum_{q_i\in[p^{m-1}]} \sum_{p_1,p_2\in[p]} \sum_{l\in[2^d]} c_l \mathrm{tr}\left(C_{q_1,\ldots,q_{d'},p_1,p_2}\left(\bigotimes_{u\in[m+d-2]}\left(\prod_{v\in[6mn(d+1)]}\overline{N}^l_{u,v}\right)\right)\right)\prod_{j\in[6mn(d+1)]}\overline{N}^{l,p_1,p_2}_{0,j}$$

where $C_{q_1,\ldots,q_{d'},p_1,p_2} \in \mathcal{C}$ by its definition. Obviously, there exists a one-to-one mapping between $\{(q_1,\ldots,q_{d'},p_1,p_2,l) : q_1,\ldots,q_{d'} \in [p^{m-1}], p_1, p_2 \in [p], l \in [2^d]\}$ and $\{l : l \in [L]\}, L = 2^d p^{d'(m-1)+2}$. By taking

$$\mathcal{Q}_l = \{Q_{l,u,v} : u \in [0:(m+d-2)], v \in [6mn(d+1)]\}$$

based on this one-to-one mapping, we have finished the proof. □

THEOREM A.8 (complete version of two-layer linear networks for Theorem 3.6). Let $\mathcal{M} := \{M_{i,j} : i \in [0:m], j \in [n]\}$ be a multi-set of matrices such that each $M_{i,j}$ or its transpose only takes value in $W^b_{0:t} \bigcup G^b_{0:t} \bigcup W^* \bigcup \overline{\mathcal{C}}$ and $\deg\left(G^b_t; \mathcal{M}\right) = d$ (here $d, m, n$ are constants independent of $b$). Then for

$$m' := m + d - 2, \quad n' := 6mn(d+1),$$

---

[8]In fact, $c_l = (-1)^s$, where $s$ equals to the number of appearance of $W^*_2 W^*_1$ that come from $\mathcal{W}^b_t$ in $\left\{\overline{N}^l_{u,v} : u \in [m+d-2], v \in [6mn(d+1)]\right\} \bigcup \left\{\overline{N}^{l,p_1,p_2}_{0,v} : j \in [6mn(d+1)]\right\}$.

*there exist a constant $L$ [9] independent of $b$ and multi-sets of matrices*

$$\mathcal{Q}_{l,s} := \left\{ Q_{l,s,u,v} : u \in [0:m'], v \in [n'] \right\}, l \in [L], s \in [0:d]$$

*such that*

$$\mathbb{E}\left[ \mathrm{tr}\left( C\left( \bigotimes_{i\in[m]} \left( \prod_{j\in[n]} M_{i,j} \right) \right) \right) \prod_{j\in[n]} M_{0,j} \,\Bigg|\, \mathcal{F}_t^b \right] = \widetilde{\alpha}_0 + \widetilde{\alpha}_1 \frac{1}{b} + \cdots + \widetilde{\alpha}_d \frac{1}{b^d},$$

*where*

$$\widetilde{\alpha}_s = \sum_{l\in[L]} c_{l,s} \mathrm{tr}\left( C_{l,s} \left( \bigotimes_{u\in[m']} \left( \prod_{v\in[n']} Q_{l,s,u,v} \right) \right) \right) \prod_{v\in[n']} Q_{l,s,0,v}, s \in [0:d],$$

$c_{l,s}$ *is a constant,* $C_{l,s} \in \mathcal{C}$ *and* $Q_{l,s,u,v}$ *only takes value in* $W_{0:t}^b \bigcup G_{0:(t-1)}^b \bigcup W^* \bigcup \overline{C}$.
*Further, we have*

$$\deg\left( G_t^b ; \mathcal{Q}_{l,s} \right) = 0,$$

$$\deg\left( W_t^b ; \mathcal{Q}_{l,s} \right) \leqslant \deg\left( W_t^b ; \mathcal{M} \right) + 3d,$$

$$\deg\left( W^* ; \mathcal{Q}_{l,s} \right) \leqslant \deg\left( W^* ; \mathcal{M} \right) + 2d,$$

$$\deg\left( W_t^b ; \mathcal{Q}_{l,s} \right) + \deg\left( W^* ; \mathcal{Q}_{l,s} \right) = \deg\left( W_t^b ; \mathcal{M} \right) + \deg\left( W^* ; \mathcal{M} \right) + 3d,$$

$$\deg\left( W_f^b ; \mathcal{Q}_{l,s} \right) = \deg\left( W_f^b ; \mathcal{M} \right), \quad f \in [0, t-1],$$

$$\deg\left( G_f^b ; \mathcal{Q}_{l,s} \right) = \deg\left( G_f^b ; \mathcal{M} \right), \quad f \in [0, t-1],$$

$$\deg\left( W^* ; \mathcal{Q}_{l,s} \right) = \deg\left( W^* ; \mathcal{M} \right).$$

*Proof.* Note that $\deg\left( G_t^b ; \mathcal{M} \right) = d$. By (A.3) and (A.4), replacing all appearance of $g_{t,i}^b$ by the sum of $b$ different terms $g_{t,i,s}^b, s \in [b], i \in \{1, 2\}$ in

$$\mathrm{tr}\left( C\left( \bigotimes_{i\in[m]} \left( \prod_{j\in[n]} M_{i,j} \right) \right) \right) \prod_{j\in[n]} M_{0,j},$$

we know there exists a multi-set of matrices $\mathcal{M}' = \left\{ M_{k,i,j} : k \in [b^d], i \in [0:m], j \in [n] \right\}$ such that

$$\alpha := \mathrm{tr}\left( C\left( \bigotimes_{i\in[m]} \left( \prod_{j\in[n]} M_{i,j} \right) \right) \right) \prod_{j\in[n]} M_{i,j} = \frac{1}{b^d} \sum_{k\in[b^d]} \mathrm{tr}\left( C\left( \bigotimes_{i\in[m]} \left( \prod_{j\in[n]} M_{k,i,j} \right) \right) \right) \prod_{j\in[n]} M_{k,0,j},$$

where every element $M_{k,i,j}$ of $\mathcal{M}'$ only takes value in $W_{0:t}^b \bigcup G_{0:(t-1)}^b \bigcup \overline{G}_t^b \bigcup W^* \bigcup \overline{C}$
and for each $k \in [b^d]$, we have

$$\deg\left( \overline{G}_t^b ; \mathcal{M}_k' \right) = \deg\left( G_t^b ; \mathcal{M} \right) = d,$$

$$\deg\left( W_t^b ; \mathcal{M}_k' \right) = \deg\left( W_t^b ; \mathcal{M} \right),$$

$$\deg\left( W^* ; \mathcal{M}_k' \right) = \deg\left( W^* ; \mathcal{M} \right),$$

$$\deg\left( W_f^b ; \mathcal{M}_k' \right) = \deg\left( W_f^b ; \mathcal{M} \right), \quad f \in [0, t-1],$$

$$\deg\left( G_f^b ; \mathcal{M}_k' \right) = \deg\left( G_f^b ; \mathcal{M} \right), \quad f \in [0, t-1],$$

$$\deg\left( W^* ; \mathcal{M}_k' \right) = \deg\left( W^* ; \mathcal{M} \right),$$

---

[9] The exact value of $L$ is specified later in the proof.

where multi-set $\mathcal{M}'_k := \{M_{k,i,j} : i \in [0:m], j \in [n]\}, k \in [b^d]$.

Let $\alpha_k := \mathrm{tr}\left(C\left(\bigotimes_{i\in[m]}\left(\prod_{j\in[n]} M_{k,i,j}\right)\right)\right)\prod_{j\in[n]} M_{k,0,j}, k \in [b^d]$. We split the set $\{\alpha_k : k \in [b^d]\}$ into disjoint and non-empty sets (equivalent classes) $S_1, \ldots, S_N$ such that

1. for every $i \in [N]$ and every $\overline{\alpha}_1, \overline{\alpha}_2 \in S_i$, we have $\mathbb{E}\left[\overline{\alpha}_1\big|\mathcal{F}_t^b\right] = \mathbb{E}\left[\overline{\alpha}_2\big|\mathcal{F}_t^b\right]$,
2. for every $i, j \in [N], i \neq j$ and every $\overline{\alpha}_1 \in S_i$ and $\overline{\alpha}_2 \in S_j$, we have $\mathbb{E}\left[\overline{\alpha}_1\big|\mathcal{F}_t^b\right] \neq \mathbb{E}\left[\overline{\alpha}_2\big|\mathcal{F}_t^b\right]$,
3. $\bigcup_{i=1}^{N} S_i = \{\alpha_k : k \in [b^d]\}$.

For every $r \in [N]$, let $k_r \in [b^d]$ be such that $\alpha_{k_r} \in S_r$ is a representative element of the equivalent class $S_r$ (in fact it can be any element of $S_r$). For each $r \in [N]$, we can always write $|S_r| = e_{r,0} + e_{r,1}b + \cdots + e_{r,d}b^d$ such that $e_{r,s} \in [0:b-1], s \in [0:d-1], e_{r,d} \in \{0,1\}$ (actually $e_{r,s}$'s are the digits of the base-$b$ representation of $|S_r|$). Then we have

$$\mathbb{E}\left[\alpha\big|\mathcal{F}_t^b\right] = \mathbb{E}\left[\frac{1}{b^d}\sum_{k=1}^{b^d}\alpha_k\bigg|\mathcal{F}_t^b\right] = \frac{1}{b^d}\mathbb{E}\left[\sum_{r=1}^{N}|S_r|\,\alpha_{k_r}\bigg|\mathcal{F}_t^b\right]$$

$$= \frac{1}{b^d}\mathbb{E}\left[\sum_{r=1}^{N}\left(e_{r,0} + e_{r,1}b + \cdots + e_{r,d}b^d\right)\alpha_{k_r}\bigg|\mathcal{F}_t^b\right]$$

$$= \frac{1}{b^d}\sum_{r=1}^{N}\left(e_{r,0} + e_{r,1}b + \cdots + e_{r,d}b^d\right)\mathbb{E}\left[\alpha_{k_r}\big|\mathcal{F}_t^b\right]$$

$$(A.20) \qquad = \sum_{r=1}^{N}\left(e_{r,d} + e_{r,d-1}\frac{1}{b} + \cdots + e_{r,0}\frac{1}{b^d}\right)\mathbb{E}\left[\alpha_{k_r}\big|\mathcal{F}_t^b\right].$$

It is important to note that $N$, the number of different equivalent classes, is independent of $b$. This follows from the fact that, by Lemma A.7, the possible values that $\mathbb{E}\left[\alpha_k\big|\mathcal{F}_t^b\right], k \in [b^d]$ can take only depend on the distribution $\mathcal{D}$. Thus the number of partition sets is independent of $b$.

By Lemma A.7, for each $k \in [b^d]$, there exist constants $m' = m + d - 2, n' = 6mn(d + 1), L' = 2^d p^{d(m-1)+2}$ that are independent of $b$ and multi-sets of matrices

$$\mathcal{Q}_l^k := \left\{Q_{l,u,v}^k : u \in [m'], v \in [n']\right\}, l \in [L']$$

such that

$$\mathbb{E}\left[\mathrm{tr}\left(C\left(\bigotimes_{i\in[m]}\left(\prod_{j\in[n]} M_{k,i,j}\right)\right)\right)\prod_{j\in[n]} M_{k,0,j}\bigg|\mathcal{F}_t^b\right] = \sum_{l\in[L']} c_l^k\,\mathrm{tr}\left(C_l^k\left(\bigotimes_{u\in[m']}\left(\prod_{v\in[n']} Q_{l,u,v}^k\right)\right)\right)\prod_{v\in[n']} Q_{l,0,v}^k,$$

where $c_l^k \in \{-1, +1\}$, $C_l^k \in \mathcal{C}$, $Q_{l,u,v}^k$ only takes value in $W_t^b \bigcup W^* \bigcup \mathcal{I} \bigcup \mathcal{C}$, $u \in [0 : m'], v \in [n'], l \in [L']$ and for all $k \in [b^d]$ and $l \in [L']$ we have

$$\deg\left(\overline{G}_t^b; \mathcal{Q}_l^k\right) = 0,$$

$$\deg\left(W_t^b; \mathcal{Q}_l^k\right) \leq \deg\left(W_t^b; \mathcal{M}'_k\right) + 3d = \deg\left(W_t^b; \mathcal{M}\right) + 3d,$$

$$\deg\left(W^*; \mathcal{Q}_l^k\right) \leq \deg\left(W^*; \mathcal{M}'_k\right) + 2d = \deg\left(W^*; \mathcal{M}\right) + 2d,$$

$$\deg\left(W_t^b; \mathcal{Q}_l^k\right) + \deg\left(W^*; \mathcal{Q}_l^k\right) = \deg\left(W_t^b; \mathcal{M}'_k\right) + \deg\left(W^*; \mathcal{M}'_k\right) + 3d$$

$$= \deg\left(W_t^b; \mathcal{M}\right) + \deg\left(W^*; \mathcal{M}\right) + 3d,$$

$$\deg\left(W_f^b; \mathcal{Q}_l^k\right) = \deg\left(W_f^b; \mathcal{M}'_k\right) = \deg\left(W_f^b; \mathcal{M}\right), \quad f \in [0, t-1],$$

$$\deg\left(G_f^b; \mathcal{Q}_l^k\right) = \deg\left(G_f^b; \mathcal{M}'_k\right) = \deg\left(G_f^b; \mathcal{M}\right), \quad f \in [0, t-1],$$

$$\deg\left(W^*; \mathcal{Q}_l^k\right) = \deg\left(W^*; \mathcal{M}'_k\right) = \deg\left(W^*; \mathcal{M}\right).$$

By (A.20) and the definition of equivalent classes $S_1, \ldots, S_N$, we have

$$\mathbb{E}\left[\alpha|\mathcal{F}_t^b\right] = \sum_{r=1}^{N}\left(e_{r,d} + e_{r,d-1}\frac{1}{b} + \cdots + e_{r,0}\frac{1}{b^d}\right)\mathbb{E}\left[\alpha_{k_r}|\mathcal{F}_t^b\right]$$

$$= \sum_{r=1}^{N}\left(e_{r,d} + e_{r,d-1}\frac{1}{b} + \cdots + e_{r,0}\frac{1}{b^d}\right)\mathbb{E}\left[\operatorname{tr}\left(C\left(\bigotimes_{i\in[m]}\left(\prod_{j\in[n]}M_{k_r,i,j}\right)\right)\right)\prod_{j\in[n]}M_{k_r,0,j}\bigg|\mathcal{F}_t^b\right]$$

$$= \sum_{r=1}^{N}\left[\left(e_{r,d} + e_{r,d-1}\frac{1}{b} + \cdots + e_{r,0}\frac{1}{b^d}\right)\sum_{l'\in[L']}c_{l'}^{k_r}\operatorname{tr}\left(C_{l'}^{k_r}\left(\bigotimes_{u\in[m']}\left(\prod_{v\in[n']}Q_{l',u,v}^{k_r}\right)\right)\right)\prod_{v\in[n']}Q_{l',0,v}^{k_r}\right]$$

$$= \tilde{\alpha}_0 + \tilde{\alpha}_1\frac{1}{b} + \cdots + \tilde{\alpha}_d\frac{1}{b^d},$$

where $\tilde{\alpha}_s = \sum_{r\in[N]}\sum_{l'\in[L']}e_{r,d-s}c_{l'}^{k_r}\operatorname{tr}\left(C_{l'}^{k_r}\left(\bigotimes_{u\in[m']}\left(\prod_{v\in[n']}Q_{l',u,v}^{k_r}\right)\right)\right)\prod_{v\in[n']}Q_{l',0,v}^{k_r}, s\in[0:d].$

Obviously, for each $s \in [0 : d]$, there exists an one-to-one mapping between $\{(r, l', s, u, v) : r \in [N], l' \in [L'], u \in [0 : m'], v \in [n']\}$ and

$$\left\{(l, s, u, v) : l \in [L], u \in [0 : m'], v \in [n']\right\},$$

where $L = N \cdot L'$. By taking the matrices $Q_{l,s,u,v}$ in the statement of this theorem based on this mapping, and note that both $N$ and $L'$ are independent of $b$, we finish the proof. $\qquad\square$

THEOREM A.9 (complete version of two-layer linear networks for Theorem 3.7). *Let $\mathcal{M} := \{M_{i,j} : i \in [0 : m], j \in [n]\}$ be a multi-set of matrices such that each $M_{i,j}$ or its transpose only takes value in $W_{0:t}^b \bigcup G_{0:(t-1)}^b \bigcup W^* \bigcup \overline{\mathcal{C}}$ and $\deg\left(W_t^b; \mathcal{M}\right) = d$ (here $d, m, n$ are constants independent of $b$) and $C \in \mathcal{C}$. Then there exist multi-sets of matrices $\mathcal{M}_k := \{M_{k,i,j} : i \in [0 : m], j \in [n]\}, k \in [2^d]$ such that*

$$\operatorname{tr}\left(C\left(\bigotimes_{i\in[m]}\left(\prod_{j\in[n]}M_{i,j}\right)\right)\right)\prod_{j\in[n]}M_{0,j} = \sum_{k\in[2^d]}\overline{\alpha}_k\operatorname{tr}\left(C\left(\bigotimes_{i\in[m]}\left(\prod_{j\in[n]}M_{k,i,j}\right)\right)\right)\prod_{j\in[n]}M_{k,0,j},$$

*where $\overline{\alpha}_k, k \in [2^d]$ are constants and each $M_{k,i,j}$ only takes value in*

$$W_{0:(t-1)}^b \bigcup G_{0:(t-1)}^b \bigcup W^* \bigcup \overline{\mathcal{C}}.$$

*Further, for each $k \in [2^d]$ we have*

$$\deg\left(G_{t-1}^b; \mathcal{M}_k\right) \leqslant \deg\left(G_{t-1}^b; \mathcal{M}\right) + d,$$

$$\deg\left(W_{t-1}^b; \mathcal{M}_k\right) \leqslant \deg\left(W_{t-1}^b; \mathcal{M}\right) + d,$$

$$\deg\left(G_{t-1}^b; \mathcal{M}_k\right) + \deg\left(W_{t-1}^b; \mathcal{M}_k\right) = \deg\left(G_{t-1}^b; \mathcal{M}\right) + \deg\left(W_{t-1}^b; \mathcal{M}\right) + d,$$

$$\deg\left(G_f^b; \mathcal{M}_k\right) = \deg\left(G_f^b; \mathcal{M}\right), \quad f \in [0 : (t-2)],$$

$$\deg\left(W_f^b; \mathcal{M}_k\right) = \deg\left(W_f^b; \mathcal{M}\right), \quad f \in [0 : (t-2)],$$

$$\deg\left(W^*; \mathcal{M}_k\right) = \deg\left(W^*; \mathcal{M}\right).$$

*Proof.* We simply use the fact that $W_{t,i}^b = W_{t-1,i}^b - \alpha_t g_{t-1,i}^b, i = 1, 2$. Note that $\deg\left(W_t^b; \mathcal{M}\right) = d$, by replacing all appearance of $W_{t,i}^b$ in

$$\operatorname{tr}\left(C\left(\bigotimes_{i\in[m]}\left(\prod_{j\in[n]}M_{i,j}\right)\right)\right)\prod_{j\in[n]}M_{0,j}$$

with $\left(W_{t-1,i}^b - \alpha_t g_{t-1,i}^b\right)$ and expand all the parentheses, we get $2^d$ terms in the form of

$$\operatorname{tr}\left(C\left(\bigotimes_{i\in[m]}\left(\prod_{j\in[n]}M_{k,i,j}\right)\right)\right)\prod_{j\in[n]}M_{0,j}.$$

The constant $\overline{\alpha}_k$ comes from the multiplication of $\alpha_t$'s. $\qquad\square$

THEOREM A.10 (complete version of two-layer linear networks for Theorem 3.8). *Let $\mathcal{M}^t := \left\{ M^t_{i,j} : i \in [0 : m_t], j \in [n_t] \right\}$ be a multi-set of matrices such that each $M^t_{i,j}$ or its transpose only takes value in $W^b_{0:t} \bigcup G^b_{0:t} \bigcup W^* \bigcup \overline{C}$ (here $m_t, n_t$ are constants independent of $b$) and $C_t \in \mathcal{C}$. Then there exist constants $q_t, m'_t, n'_t, L_{t,s}, s \in [0 : q_t]$ that are independent of $b$ and multi-sets of matrices*

$$\mathcal{M}^t_{l,s} := \left\{ M^t_{l,s,u,v} : u \in [0 : m'_t], v \in [n'_t] \right\}, s \in [q_t]$$

*such that*

$$\mathbb{E}\left[ \mathrm{tr}\left( C_t \left( \bigotimes_{i \in [m_t]} \left( \prod_{j \in [n_t]} M^t_{i,j} \right) \right) \right) \prod_{j \in [n_t]} M^t_{0,j} \middle| \mathcal{F}_0 \right] = \alpha_{t,0} + \alpha_{t,1}\frac{1}{b} + \cdots + \alpha_{t,q_t}\frac{1}{b^{q_t}},$$

*where*

$$\alpha_{t,s} = \sum_{l \in [L_{t,s}]} c_{t,l,s} \mathrm{tr}\left( C_{t,l,s} \left( \bigotimes_{u \in [m'_t]} \left( \prod_{v \in [n'_t]} M^t_{l,s,u,v} \right) \right) \right) \prod_{v \in [n'_t]} M^t_{l,s,0,v}, s \in [0 : q_t],$$

*$c_{t,l,s}$ is a constant, $C_{t,l,s} \in \mathcal{C}$ and $M^t_{l,s,u,v}$ only takes value in $W^b_0 \bigcup W^* \bigcup \overline{C}$. Further, we have*

$$q_t \leqslant \sum_{f \in [0:t]} \left( \frac{3^{f+1} - 1}{2} \deg\left( G^b_f; \mathcal{M}^t \right) + \frac{3^f - 1}{2} \deg\left( W^b_f; \mathcal{M}^t \right) \right).$$

*Proof.* We use induction on $t$ to show this theorem. The case of $t = 0$ is the same as the statement in Theorem A.8.

Suppose that the statement holds for $t \geqslant 0$ and we consider the case of $t + 1$. By Theorem A.8, there exist constants $\widetilde{m}_{t+1}, \widetilde{n}_{t+1}, \widetilde{L}_{t+1}$ that are independent of $b$ and multi-sets of matrices $\mathcal{Q}^{t+1}_{l,s} := \left\{ Q^{t+1}_{l,s,u,v} : u \in [0 : \widetilde{m}_{t+1}], v \in [\widetilde{n}_{t+1}] \right\}, l \in [\widetilde{L}_{t+1}], s \in [0 : d_{t+1}]$ such that

(A.21)
$$\mathbb{E}\left[ \mathrm{tr}\left( C \left( \bigotimes_{i \in [m_{t+1}]} \left( \prod_{j \in [n_{t+1}]} M^{t+1}_{i,j} \right) \right) \right) \prod_{j \in [n_{t+1}]} M^{t+1}_{0,j} \middle| \mathcal{F}^b_{t+1} \right] = \widetilde{\alpha}_{t+1,0} + \widetilde{\alpha}_{t+1,1}\frac{1}{b} + \cdots + \widetilde{\alpha}_{t+1,d_{t+1}}\frac{1}{b^{d_{t+1}}},$$

where

(A.22)
$$\widetilde{\alpha}_{t+1,s} = \sum_{l \in [\widetilde{L}_{t+1}]} \widetilde{c}_{t+1,l,s} \mathrm{tr}\left( \widetilde{C}_{t+1,l,s} \left( \bigotimes_{u \in [\widetilde{m}_{t+1}]} \left( \prod_{v \in [\widetilde{n}_{t+1}]} Q^{t+1}_{l,s,u,v} \right) \right) \right) \prod_{v \in [\widetilde{n}_{t+1}]} Q^{t+1}_{l,s,0,v}, s \in [0 : d_{t+1}],$$

$d_{t+1} := \deg\left( G^b_{t+1}; \mathcal{M}^{t+1} \right)$, $\widetilde{c}_{t+1,l,s}$ is a constant, $\widetilde{C}_{t+1,l,s} \in \mathcal{C}$ and $Q^{t+1}_{l,s,u,v}$ only takes value in $W^b_{0:(t+1)} \bigcup G^b_{0:t} \bigcup W^* \bigcup \overline{C}$. Further, we have

$$\deg\left( W^b_{t+1}; \mathcal{Q}^{t+1}_{l,s} \right) \leqslant \deg\left( W^b_{t+1}; \mathcal{M}^{t+1} \right) + 3\deg\left( G^b_{t+1}; \mathcal{M}^{t+1} \right),$$

$$\deg\left( W^*; \mathcal{Q}^{t+1}_{l,s} \right) \leqslant \deg\left( W^*; \mathcal{M}^{t+1} \right) + 2\deg\left( G^b_{t+1}; \mathcal{M}^{t+1} \right),$$

$$\deg\left( W^b_{t+1}; \mathcal{Q}^{t+1}_{l,s} \right) + \deg\left( W^*; \mathcal{Q}^{t+1}_{l,s} \right) = \deg\left( W^b_{t+1}; \mathcal{M}^{t+1} \right) + \deg\left( W^*; \mathcal{M}^{t+1} \right) + 3\deg\left( G^b_{t+1}; \mathcal{M}^{t+1} \right).$$

By Theorem A.9, for each $l \in [\widetilde{L}_{t+1}]$ and $s \in [0 : d_{t+1}]$, there exist multi-sets of matrices $\mathcal{M}_{l,s,k}^{t} := \left\{ M_{l,s,k,i,j}^{t} : i \in [0 : m_t], j \in [n_t] \right\}, k \in [2^{d_{t+1}}]$ such that

$$\mathrm{tr}\left( \widetilde{C}_{t+1,l,s} \left( \bigotimes_{u \in [\widetilde{m}_{t+1}]} \left( \prod_{v \in [\widetilde{n}_{t+1}]} Q_{l,s,u,v}^{t+1} \right) \right) \right) \prod_{v \in [\widetilde{n}_{t+1}]} Q_{l,s,0,v}^{t+1}$$

(A.23) $$= \sum_{k \in [2^{d_{t+1}}]} \overline{\alpha}_{t,k} \mathrm{tr}\left( \widetilde{C}_{t+1,l,s} \left( \bigotimes_{i \in [m_t]} \left( \prod_{j \in [n_t]} M_{l,s,k,i,j}^{t} \right) \right) \right) \prod_{j \in [n_t]} M_{l,s,k,0,j}^{t},$$

where $m_t = \widetilde{m}_{t+1}, n_t = \widetilde{n}_{t+1}, \overline{\alpha}_{t,k}, k \in [2^{d_{t+1}}]$ are constants, and each $M_{l,s,k,i,j}^{t}$ only takes value in $W_{0:t}^{b} \bigcup G_{0:t}^{b} \bigcup W^* \bigcup \overline{C}$. Further, for each $k \in [2^{d_{t+1}}]$ we have

$$\deg\left( W_t^b; \mathcal{M}_{l,s,k}^t \right) + \deg\left( G_t^b; \mathcal{M}_{l,s,k}^t \right)$$

$$= \deg\left( W_{t+1}^b; \mathcal{Q}_{l,s}^{t+1} \right) + \deg\left( W_t^b; \mathcal{Q}_{l,s}^{t+1} \right) + \deg\left( G_t^b; \mathcal{Q}_{l,s}^{t+1} \right)$$

$$\leqslant \deg\left( W_{t+1}^b; \mathcal{M}^{t+1} \right) + 3\deg\left( G_{t+1}^b; \mathcal{M}^{t+1} \right) + \deg\left( W_t^b; \mathcal{Q}_{l,s}^{t+1} \right) + \deg\left( G_t^b; \mathcal{Q}_{l,s}^{t+1} \right),$$

$$\deg\left( G_t^b; \mathcal{M}_{l,s,k}^t \right)$$

$$\leqslant \deg\left( W_{t+1}^b; \mathcal{Q}_{l,s}^{t+1} \right) + \deg\left( G_t^b; \mathcal{Q}_{l,s}^{t+1} \right)$$

$$\leqslant \deg\left( W_{t+1}^b; \mathcal{M}^{t+1} \right) + 3\deg\left( G_{t+1}^b; \mathcal{M}^{t+1} \right) + \deg\left( G_t^b; \mathcal{Q}_{l,s}^{t+1} \right),$$

and

$$\deg\left( W^*; \mathcal{M}_{l,s,k}^t \right) = \deg\left( W^*; \mathcal{Q}_{l,s}^{t+1} \right) \leqslant \deg\left( W^*; \mathcal{M}^{t+1} \right) + 2\deg\left( G_{t+1}^b; \mathcal{M}^{t+1} \right).$$

By (A.21) − (A.23), we have

$$\mathbb{E}\left[ \mathrm{tr}\left( C \left( \bigotimes_{i \in [m_{t+1}]} \left( \prod_{j \in [n_{t+1}]} M_{i,j}^{t+1} \right) \right) \right) \prod_{j \in [n_{t+1}]} M_{0,j}^{t+1} \middle| \mathcal{F}_0 \right]$$

$$= \mathbb{E}\left[ \mathbb{E}\left[ \mathrm{tr}\left( C \left( \bigotimes_{i \in [m_{t+1}]} \left( \prod_{j \in [n_{t+1}]} M_{i,j}^{t+1} \right) \right) \right) \prod_{j \in [n_{t+1}]} M_{0,j}^{t+1} \middle| \mathcal{F}_{t+1}^b \right] \middle| \mathcal{F}_0 \right]$$

$$= \mathbb{E}\left[ \widetilde{\alpha}_{t+1,0} \middle| \mathcal{F}_0 \right] + \mathbb{E}\left[ \widetilde{\alpha}_{t+1,1} \middle| \mathcal{F}_0 \right] \frac{1}{b} + \cdots + \mathbb{E}\left[ \widetilde{\alpha}_{t+1,d_{t+1}} \middle| \mathcal{F}_0 \right] \frac{1}{b^{d_{t+1}}}$$

(A.24)

$$= \sum_{l \in [\widetilde{L}_{t+1}], s \in [d_{t+1}], k \in [2^{d_{t+1}}]} \frac{\widetilde{c}_{t+1,l,s} \overline{\alpha}_{t,k}}{b^s} \mathbb{E}\left[ \mathrm{tr}\left( \widetilde{C}_{t+1,l,s} \left( \bigotimes_{i \in [m_t]} \left( \prod_{j \in [n_t]} M_{l,s,k,i,j}^{t} \right) \right) \right) \prod_{j \in [n_t]} M_{l,s,k,0,j}^{t} \middle| \mathcal{F}_0 \right].$$

By induction, for each $l \in [\widetilde{L}_{t+1}], s \in [d_{t+1}]$ and $k \in [2^{d_{t+1}}]$, there exist constants $m_0, n_0, Z, d'$ that are independent of $b$ and multi-sets of matrices

$$\mathcal{M}_{l,s,k,r,z}^{0} := \left\{ M_{l,s,k,r,z,u,v}^{0} : u \in [0 : m_0], v \in [n_0] \right\}, r \in [d'], z \in [Z]$$

such that

(A.25)

$$\mathbb{E}\left[ \mathrm{tr}\left( \widetilde{C}_{t+1,l,s} \left( \bigotimes_{i \in [m_t]} \left( \prod_{j \in [n_t]} M_{l,s,k,i,j}^{t} \right) \right) \right) \prod_{j \in [n_t]} M_{l,s,k,0,j}^{t} \middle| \mathcal{F}_0 \right] = \alpha_0' + \alpha_1' \frac{1}{b} + \cdots + \alpha_{d'}' \frac{1}{b^{d'}},$$

where
(A.26)
$$\alpha'_r = \sum_{z \in [Z]} c_{t,l,s,k,r,z} \operatorname{tr}\left(C_{t,l,s,k,r,z}\left(\bigotimes_{u \in [m_0]}\left(\prod_{v \in [n_0]} M^0_{l,s,k,r,z,u,v}\right)\right)\right)\prod_{v \in [n_0]} M^0_{l,s,k,r,z,0,v}, r \in [d'],$$

$c_{t,l,s,k,r,z}$ is a constant, $C_{t,l,s,k,r,z} \in \mathcal{C}$ and $M^0_{l,s,k,r,z,u,v}$ only takes value in

$$W^b_0 \bigcup W^* \bigcup \overline{\mathcal{C}}.$$

Further, we have

$$d' \leqslant \sum_{f \in [0:t]}\left(\frac{3^{f+1}-1}{2}\deg\left(G^b_f; \mathcal{M}^t_{l,s,k}\right) + \frac{3^f-1}{2}\deg\left(W^b_f; \mathcal{M}^t_{l,s,k}\right)\right).$$

Combining (A.24) – (A.26), we have

$$\mathbb{E}\left[\operatorname{tr}\left(C\left(\bigotimes_{i \in [m_{t+1}]}\left(\prod_{j \in [n_{t+1}]} M^{t+1}_{i,j}\right)\right)\right)\prod_{j \in [n_{t+1}]} M^{t+1}_{0,j}\,\middle|\,\mathcal{F}_0\right] = \alpha_0 + \alpha_1 \frac{1}{b} + \cdots + \alpha_q \frac{1}{b^q},$$

where $q = d_{t+1} + d'$ and for each $e \in [0:q]$,

$$\alpha_e = \sum_{l \in [\tilde{L}_{t+1}], s \in [d_{t+1}], k \in [2^{d_{t+1}}], r \in [d'], z \in [Z], r+s=e} c_{t+1,l,s}\overline{\alpha}_{t,k} c_{t,l,s,k,r,z} \cdot$$

$$\cdot \operatorname{tr}\left(C_{t,l,s,k,r,z}\left(\bigotimes_{u \in [m_0]}\left(\prod_{v \in [n_0]} M^0_{l,s,k,r,z,u,v}\right)\right)\right)\prod_{v \in [n_0]} M^0_{l,s,k,r,z,0,v}.$$

Note that

$$q = d_{t+1} + d'$$

$$\leqslant \deg\left(G^b_{t+1}; \mathcal{M}^{t+1}\right) + \sum_{f \in [0:t]}\left(\frac{3^{f+1}-1}{2}\deg\left(G^b_f; \mathcal{M}^t_{l,s,k}\right) + \frac{3^f-1}{2}\deg\left(W^b_f; \mathcal{M}^t_{l,s,k}\right)\right)$$

$$= \deg\left(G^b_{t+1}; \mathcal{M}^{t+1}\right) + \frac{3^{t+1}-1}{2}\deg\left(G^b_t; \mathcal{M}^t_{l,s,k}\right) + \frac{3^t-1}{2}\deg\left(W^b_t; \mathcal{M}^t_{l,s,k}\right)$$

$$+ \sum_{f \in [0:(t-1)]}\left(\frac{3^{f+1}-1}{2}\deg\left(G^b_f; \mathcal{M}^{t+1}\right) + \frac{3^f-1}{2}\deg\left(W^b_f; \mathcal{M}^{t+1}\right)\right)$$

$$\leqslant \deg\left(G^b_{t+1}; \mathcal{M}^{t+1}\right) +$$

$$+ \frac{3^t-1}{2}\left(\deg\left(W^b_{t+1}; \mathcal{M}^{t+1}\right) + 3\deg\left(G^b_{t+1}; \mathcal{M}^{t+1}\right) + \deg\left(W^b_t; \mathcal{Q}^{t+1}_{l,s}\right) + \deg\left(G^b_t; \mathcal{Q}^{t+1}_{l,s}\right)\right) +$$

$$+ \frac{3^{t+1}-3^t}{2}\left(\deg\left(W^b_{t+1}; \mathcal{M}^{t+1}\right) + 3\deg\left(G^b_{t+1}; \mathcal{M}^{t+1}\right) + \deg\left(G^b_t; \mathcal{Q}^{t+1}_{l,s}\right)\right) +$$

$$+ \sum_{f \in [0:(t-1)]}\left(\frac{3^{f+1}-1}{2}\deg\left(G^b_f; \mathcal{M}^{t+1}\right) + \frac{3^f-1}{2}\deg\left(W^b_f; \mathcal{M}^{t+1}\right)\right)$$

$$= \frac{3^{t+2}-1}{2}\deg\left(G^b_{t+1}; \mathcal{M}^{t+1}\right) + \frac{3^{t+1}-1}{2}\deg\left(W^b_{t+1}; \mathcal{M}^{t+1}\right) + \frac{3^{t+1}-1}{2}\deg\left(G^b_t; \mathcal{M}^{t+1}\right)$$

$$+ \frac{3^t-1}{2}\deg\left(W^b_t; \mathcal{M}^{t+1}\right) + \sum_{f \in [0:(t-1)]}\left(\frac{3^{f+1}-1}{2}\deg\left(G^b_f; \mathcal{M}^{t+1}\right) + \frac{3^f-1}{2}\deg\left(W^b_f; \mathcal{M}^{t+1}\right)\right)$$

$$= \sum_{f \in [0:(t+1)]}\left(\frac{3^{f+1}-1}{2}\deg\left(G^b_f; \mathcal{M}^{t+1}\right) + \frac{3^f-1}{2}\deg\left(W^b_f; \mathcal{M}^{t+1}\right)\right), \qquad \qquad \Box$$

which finishes the proof.

1170    THEOREM A.11 (Two-layer linear network version for Theorem 3.9).   *Given* $t \in \mathbb{N}$,
1171    *value* $\mathsf{var}\left(g_{t,i}^b\right), i = 1, 2$ *can be written as a polynomial of* $\frac{1}{b}$ *with degree at most*
1172    $3^{t+1} - 1$ *with no constant term. Formally, we have* $\mathsf{var}\left(g_{t,i}^b\right) = \beta_1 \frac{1}{b} + \cdots + \beta_r \frac{1}{b^r}$, *where*
1173    $r \leqslant 3^{t+1} - 1$ *and each* $\beta_i$ *is a constant independent of* $b$.

1174    *Proof.* We only show the case for $g_{t,1}^b$ since the proof for $g_{t,2}$ can be tackled
1175    similarly. Note that

$$1176 \qquad \mathsf{var}\left(g_{t,1}^b\right) = \mathbb{E}\left\|g_{t,1}^b\right\|^2 - \left\|\mathbb{E}\left[g_{t,1}^b\right]\right\|^2$$

$$1177 \qquad = \mathbb{E}\left[\mathbb{E}\left[\left\|g_{t,1}^b\right\|^2 \Big| \mathcal{F}_0\right]\right] - \left\|\mathbb{E}\left[\mathbb{E}\left[g_{t,1}^b | \mathcal{F}_0\right]\right]\right\|^2$$

1178    (A.27)
$$\qquad = \mathbb{E}\left[\mathbb{E}\left[\mathrm{tr}\left(\left(g_{t,1}^b\right)^T g_{t,1}^b\right) \Big| \mathcal{F}_0\right]\right] - \left\|\mathbb{E}\left[\mathbb{E}\left[g_{t,1}^b | \mathcal{F}_0\right]\right]\right\|^2.$$
1179

1180    By Theorem A.10, there exist constants $q_1, m_1', n_1', \overline{L}_{1,s}, s \in [0 : q_1]$ that are inde-
1181    pendent of $b$ and multi-sets of matrices $\mathcal{M}_{l,s}^1 := \left\{M_{l,s,u,v}^1 : u \in [m_1'], v \in [n_1']\right\}, s \in [q_1]$
1182    such that

1183    (A.28)
$$\mathbb{E}\left[\mathrm{tr}\left(\left(g_{t,1}^b\right)^T g_{t,1}^b\right) \Big| \mathcal{F}_0\right] = \alpha_{1,0} + \alpha_{1,1} \frac{1}{b} + \cdots + \alpha_{1,q_1} \frac{1}{b^{q_1}},$$

where

$$\alpha_{1,s} = \sum_{l \in [\overline{L}_{1,s}]} c_{1,l,s} \mathrm{tr}\left(C_{1,l,s}\left(\bigotimes_{u \in [m_1']}\left(\prod_{v \in [n_1']} M_{l,s,u,v}^1\right)\right)\right), s \in [0 : q_1],$$

1184    $c_{1,l,s}$ is a constant, $C_{1,l,s} \in \mathcal{C}$ and $M_{l,s,u,v}^1$ only takes value in $W_0^b \bigcup W^* \bigcup \overline{\mathcal{C}}$. Further,
1185    we have

$$1186 \qquad\qquad\qquad\qquad q_1 \leqslant 3^{t+1} - 1.$$
1187

1188    It is worth mentioning that we do not include matrices $M_{1,l,s,0,v}, v \in [n_1']$ in the
1189    multi-set $\mathcal{M}_{l,s}^1, l \in [\overline{L}_{1,s}], s \in [0 : q_1]$ because each $M_{1,l,s,0,v}$ is actually an identity
1190    matrix from the proof of the previous theorems.
1191    Similarly, there exist constants $q_2, m_2', n_2', \overline{L}_{2,s}, s \in [0 : q_2]$ that are independent of
1192    $b$ and multi-sets of matrices $\mathcal{M}_{l,s}^2 := \left\{M_{l,s,u,v}^2 : u \in [0 : m_2'], v \in [n_2']\right\}, s \in [q_2]$ such
1193    that

1194    (A.29)
$$\mathbb{E}\left[g_{t,1}^b | \mathcal{F}_0\right] = \alpha_{2,0} + \alpha_{2,1} \frac{1}{b} + \cdots + \alpha_{2,q_2} \frac{1}{b^{q_2}},$$

where

$$\alpha_{2,s} = \sum_{l \in [\overline{L}_{2,s}]} c_{2,l,s} \mathrm{tr}\left(C_{2,l,s}\left(\bigotimes_{u \in [m_2']}\left(\prod_{v \in [n_2']} M_{l,s,u,v}^2\right)\right)\right) \prod_{v \in [n_2']} M_{l,s,0,v}^2, s \in [0 : q_2],$$

1195    $c_{2,l,s}$ is a constant, $C_{2,l,s} \in \mathcal{C}$ and $M_{l,s,u,v}^2$ only takes value in $W_0^b \bigcup W^* \bigcup \overline{\mathcal{C}}$. Further,
1196    we have

$$1197 \qquad\qquad\qquad\qquad q_2 \leqslant \frac{1}{2}\left(3^{t+1} - 1\right).$$
1198

Combining $(A.27) – (A.29)$, we know there exist constants

$$\gamma_0, \ldots, \gamma_q, q = \max\{q_1, 2q_2\} \leqslant 3^{t+1} - 1$$

such that

$$\mathsf{var}\left(\left(W_{t,2}^b\right)^T W_{t,2}^b W_{t,1}^b x x^T\right) = \gamma_0 + \gamma_1 \frac{1}{b} + \cdots \gamma_q \frac{1}{b^q},$$

where

$$\gamma_s = \mathbb{E}_{W_0^t \sim \mathcal{D}'}\left[\alpha_{1,s}\right] + \sum_{u+v=s, u, v \in [0:q_2]} \mathbb{E}_{W_0^t \sim \mathcal{D}'}\left[\alpha_{2,u}\right] \mathbb{E}_{W_0^t \sim \mathcal{D}'}\left[\alpha_{2,v}\right], s \in [0:q]$$

and $\mathcal{D}'$ is the initialization distribution of $W_0^t$. Further, $\gamma_s$'s are independent of $b$.

*Proof of Theorem 3.10.* We first show that in

$$\mathsf{var}\left(g_{t,i}^b\right) = \beta_1 \frac{1}{b} + \cdots + \beta_r \frac{1}{b^r}$$

we have $\beta_1 \geqslant 0$. If $r = 1$, the statement obviously holds. Let us assume that the statement does not hold for $r > 1$, i.e. $\beta_1 < 0$. Taking $b$ large enough such that $\beta_1 b^{r-1} + \beta_2 b^{r-2} + \cdots + \beta_r < 0$ yields

$$\mathsf{var}\left(g_{t,i}^b\right) = \frac{1}{b^r}\left(\beta_1 b^{r-1} + \beta_2 b^{r-2} + \cdots + \beta_r\right) < 0,$$

which contradicts the fact that $\mathsf{var}\left(g_{t,i}^b\right) \geqslant 0$. Therefore, we have $\beta_1 \geqslant 0$.

Let $b_0$ be large enough such that for all $b \geqslant b_0$, we have $\beta_1 b^{r-1} + 2\beta_2 b^{r-2} + \cdots + r\beta_r \geqslant 0$. We denote $f(b) = \beta_1 \frac{1}{b} + \beta_2 \frac{1}{b^2} + \cdots + \beta_r \frac{1}{b^r} \geqslant 0$. For all $b > b_0$ we have

$$f'(b) = -\frac{1}{b^{r+1}}\left(\beta_1 b^{r-1} + 2\beta_2 b^{r-2} + \cdots + r\beta_r\right) \leqslant 0. \qquad \square$$

Therefore, for all $b > b_0$ we have $\left(\mathsf{var}\left(g_{t,i}^b\right)\right)' = -\frac{r}{b^{r+1}} f(b) + \frac{1}{b^r} f(b) \leqslant 0$, and thus $\mathsf{var}\left(g_{t,i}^b\right)$ is a decreasing function of $b$ for all $b > b_0$.

**A.2.2. Two-layer Networks with Quadratic Polynomial Activation Functions.** In this section, we expand the scope of the theorems found in Appendix A.2.1. While they originally applied to two-layer linear networks, we now extend them to networks utilizing quadratic polynomial activation functions. The main distinction between these scenarios lies in the incorporation of Hadamard products into the gradients by the quadratic activation functions, demanding additional consideration.

Specifically, we consider a special case of the general population loss (3.1). Here the population loss is defined as

$$\mathcal{L}(w) = \mathbb{E}_{x \sim \mathcal{D}}\left[\frac{1}{2}\left\|W_2 \sigma\left(W_1 x\right) - W_2^* \sigma\left(W_1^* x\right)\right\|^2\right]$$

and the SG estimators are defined as

$$g_{t,k}^b := \frac{1}{b}\sum_{i=1}^b \nabla_{W_{t,k}^b}\left(\frac{1}{2}\left\|W_{t,2}^b \sigma\left(W_{t,1}^b x_{t,i}^b\right) - W_2^* \sigma\left(W_1^* x_{t,i}^b\right)\right\|^2\right), \quad k = 1, 2,$$

where $\sigma(x) := \sigma_0 + \sigma_1 x + \sigma_2 x^2$ is a polynomial activation function of degree 2. This setup aligns to the $D = 2$ and $H = 2$ case as in (3.1).

Similar to (A.3) – (A.4), we rewrite the SG estimator as the sum of the product of weight matrices and other constant matrices. For example, we have

$$g_{t,1}^b = \frac{1}{b} \sum_{i=1}^b \nabla_{W_{t,1}^b} \left( \frac{1}{2} \left\| W_{t,2}^b \sigma \left( W_{t,1}^b x_{t,i}^b \right) - W_2^* \sigma \left( W_1^* x_{t,i}^b \right) \right\|^2 \right)$$

$$= \frac{1}{2b} \sum_{i=1}^b \nabla_{W_{t,1}^b} \left\| \sigma_2 W_{t,2}^b \left( \left( W_{t,1}^b x_{t,i}^b \right) \odot \left( W_{t,1}^b x_{t,i}^b \right) \right) + \sigma_1 W_{t,2}^b \left( W_{t,1}^b x_{t,i}^b \right) + \sigma_0 W_{t,2}^b \right.$$

(A.30)

$$\left. - \sigma_2 W_2^* \left( \left( W_1^b x_{t,i}^b \right) \odot \left( W_1^* x_{t,i}^b \right) \right) - \sigma_1 W_2^* \left( W_1^* x_{t,i}^b \right) - \sigma_0 W_2^* \right\|^2 .$$

We first show how to calculate the gradient of a mixed form with common and Hadamard products. With this approach, we can represent each summand of (A.30) as a summation of terms in the form of $\prod_k M_k$, where $M_k$ or its transpose only takes on values from $\{W_{t,1}^b, W_{t,2}^b, W_1^*, W_2^*, x_{t,i}^b\} \bigcup \mathcal{C}$.

We take two terms in the expansion of the summand in (A.30) as examples to show how to replace the Hadamard products by common products. We use the fact that, for any positive integer $n$ and vectors $v_1, \ldots, v_n \in \mathbb{R}^p$,

$$(A.31) \qquad v_1 \odot v_2 \odot \cdots \odot v_n = \sum_{j \in p} \left( e_{p,j}^T v_1 \right) \left( e_{p,j}^T v_2 \right) \cdots \left( e_{p,j}^T v_n \right) e_{p,j},$$

where $e_{p,j}, j \in [p]$ is the $j$-th unit vector in $\mathbb{R}^p$.

For example, we have[10]

$$\nabla_{W_{t,1}^b} \text{tr} \left( \sigma_1 \left( W_{t,1}^b x_{t,i}^b \right)^T \left( W_{t,2}^b \right)^T \sigma_2 W_{t,2}^b \left( \left( W_{t,1}^b x_{t,i}^b \right) \odot \left( W_{t,1}^b x_{t,i}^b \right) \right) \right)$$

$$= \sigma_1 \sigma_2 \sum_{j \in [p_1]} \nabla_{W_{t,1}^b} \text{tr} \left( \left( x_{t,i}^b \right)^T \left( W_{t,1}^b \right)^T \left( W_{t,2}^b \right)^T W_{t,2}^b \left( e_{p_1,j}^T W_{t,1}^b x_{t,i}^b \right) \left( e_{p_1,j}^T W_{t,1}^b x_{t,i}^b \right) e_{p_1,j} \right)$$

$$= \sigma_1 \sigma_2 \sum_{j \in [p_1]} \left[ \left( W_{t,2}^b \right)^T W_{t,2}^b e_{p_1,j}^T W_{t,1}^b x_{t,i}^b e_{p_1,j}^T W_{t,1}^b x_{t,i}^b e_{p_1,j} \left( x_{t,i}^b \right)^T \right.$$

$$+ e_{p_1,j} \left( W_{t,2}^b \right)^T W_{t,2}^b W_{t,1}^b x_{t,i}^b e_{p_1,j}^T \left( x_{t,i}^b \right)^T \left( W_{t,1}^b \right)^T e_{p_1,j} \left( x_{t,i}^b \right)^T$$

$$\left. + e_{p_1,j} \left( x_{t,i}^b \right)^T \left( W_{t,1}^b \right)^T e_{p_1,j} \left( W_{t,2}^b \right)^T W_{t,2}^b W_{t,1}^b x_{t,i}^b e_{p_1,j}^T \left( x_{t,i}^b \right)^T \right]$$

---

[10]We frequently use the fact, that for matrices $A, B, X$ with appropriate dimensions, $\nabla_X \text{tr} (AXB) = A^T B^T$ and $\nabla_X \text{tr} \left( AX^T B \right) = BA$.

and

$$\nabla_{W_{t,1}^b} \mathrm{tr}\left(\sigma_2\left[W_{t,2}^b\left(\left(W_{t,1}^b x_{t,i}^b\right)\odot\left(W_{t,1}^b x_{t,i}^b\right)\right)\right]^T \sigma_2 W_{t,2}^b\left(\left(W_{t,1}^b x_{t,i}^b\right)\odot\left(W_{t,1}^b x_{t,i}^b\right)\right)\right)$$

(A.32)

$$= \sigma_2^2 \sum_{j,k\in[p_1]} \nabla_{W_{t,1}^b} \mathrm{tr}(e_{p_1,k}^T \left(x_{t,i}^b\right)^T \left(W_{t,1}^b\right)^T e_{p_1,k}\left(x_{t,i}^b\right)^T\left(W_{t,1}^b\right)^T \cdot$$

$$\cdot e_{p_1,k}\left(W_{t,2}^b\right)^T W_{t,2}^b e_{p_1,j}^T W_{t,1}^b x_{t,i}^b e_{p_1,j}^T W_{t,1}^b x_{t,i}^b e_{p_1,j})$$

$$= \sigma_2^2 \sum_{j,k\in[p_1]}\Bigg[ e_{p_1,k}\left(x_{t,i}^b\right)^T\left(W_{t,1}^b\right)^T e_{p_1,k}\left(W_{t,2}^b\right)^T W_{t,2}^b e_{p_1,j}^T W_{t,1}^b x_{t,i}^b e_{p_1,j}^T W_{t,1}^b x_{t,i}^b e_{p_1,j} e_{p_1,k}^T\left(x_{t,i}^b\right)^T$$

$$+ e_{p_1,k}\left(W_{t,2}^b\right)^T W_{t,2}^b e_{p_1,j}^T W_{t,1}^b x_{t,i}^b e_{p_1,j}^T W_{t,1}^b x_{t,i}^b e_{p_1,j} e_{p_1,k}^T\left(x_{t,i}^b\right)^T\left(W_{t,1}^b\right)^T e_{p_1,k}\left(x_{t,i}^b\right)^T$$

$$+ e_{p_1,j}\left(W_{t,2}^b\right)^T W_{t,2}^b e_{p_1,k}^T W_{t,1}^b x_{t,i}^b e_{p_1,k}^T W_{t,1}^b x_{t,i}^b e_{p_1,k} e_{p_1,j}^T\left(x_{t,i}^b\right)^T\left(W_{t,1}^b\right)^T e_{p_1,j}\left(x_{t,i}^b\right)^T$$

(A.33)

$$+ e_{p_1,j}\left(x_{t,i}^b\right)^T\left(W_{t,1}^b\right)^T e_{p_1,j}\left(W_{t,2}^b\right)^T W_{t,2}^b e_{p_1,k}^T W_{t,1}^b x_{t,i}^b e_{p_1,k}^T W_{t,1}^b x_{t,i}^b e_{p_1,k} e_{p_1,j}^T\left(x_{t,i}^b\right)^T\Bigg].$$

In conclusion, there exist constants $J, K, \alpha_j, j\in[J]$ independent of $b$ and a multi-set of matrices $\{M_{s,i,j,k}, i\in[b], j\in[J], k\in[K], s=1,2\}$ such that

$$g_{t,s}^b = \frac{1}{b}\sum_{i\in[b]}\sum_{j\in[J]}\left(\alpha_{s,i,j}\prod_{k\in[K]} M_{s,i,j,k}\right), s=1,2,$$

where $M_{s,i,j,k}$ or its transpose only takes value in $\{W_{t,1}^b, W_{t,2}^b, W_1^*, W_2^*\}\bigcup\{x_{t,i}^b, i\in[b]\}\bigcup \mathcal{C}$.

It is worth mentioning that we can provide the exact values of $J$ and $K$, namely $J = 144p_1^2$ and $K = 15$. These numbers are determined by analyzing the most complicated term, i.e. the left-hand side of (A.33), among the expansion of summands in (A.30). Note that the summation on the right-hind side of (A.33) contributes $4p_1^2$ terms where each term is a product of 15 matrices and the expansion of a summand in (A.30) gives 36 terms of matrices' mixed products. Thus we have $J = 36\cdot 4p_1^2 = 144p_1^2$ and $K = 15$. We can use identity matrices and zeros to fill up the unused $M_{s,i,j,k}$ and $\alpha_{s,i,j}$ as needed.

This representation aligns with the right-hand side of (A.3) and (A.4), excepts the fact that we further expand the $\mathcal{W}_t^b = W_{t,2}^b W_{t,1}^b - W_2^* W_1^*$ to separate terms. Thus we can further analyze the dynamics of polynomially-activated networks in a similar manner as in Appendix A.2.1.

**A.2.3. Deep Networks with Polynomially-activated Functions.** In this section, we discuss the extension from two-layer network networks with quadratic polynomial activation functions to deep networks with polynomial activation functions of any degree. In other words, we consider the general setting where $D$ and $H$ can take arbitrary values as in (3.1).

The building block of above derivation is to represent the SG estimators as products of weights matrices, samples, and other constant matrices. However, given the arbitrary values of $D$ and $H$, the number of matrices required is much more than the case as in Appendix A.2.2.

LEMMA A.12. *There exist constants $J, K, \alpha_j, j\in[J]$ independent of $b$ and a multi-*

*set of matrices* $\{M_{s,i,j,k}, i \in [b], j \in [J], k \in [K], s \in [H]\}$ *such that, for any* $s \in [H]$,

$$g_{t,s}^b$$

$$:= \frac{1}{b} \sum_{i=1}^{b} \nabla_{W_{t,s}^b} \left( \frac{1}{2} \left\| W_{t,H}^b \sigma \left( W_{t,H-1}^b \sigma \left( \cdots \sigma \left( W_{t,1}^b x_{t,i}^b \right) \right) \right) - W_H^* \sigma \left( W_{H-1}^* \sigma \left( \cdots \sigma \left( W_1^* x_{t,i}^b \right) \right) \right) \right\|^2 \right)$$

(A.34)

$$= \frac{1}{b} \sum_{i \in [b]} \sum_{j \in [J]} \left( \alpha_{s,i,j} \prod_{k \in [K]} M_{s,i,j,k} \right),$$

*where* $M_{s,i,j,k}$ *or its transpose only takes value in* $\{W_{t,1}^b, W_{t,2}^b, W_1^*, W_2^*\} \bigcup \{x_{t,i}^b, i \in [b]\} \bigcup \mathcal{C}.$

To give an insight on the complexity of this representation, we provide the possible values of $J$ and $K$[11] in an induction fashion.

- $K = 6D^{H-1} + 4D^{H-2} + \cdots + 4D + 3.$
  In the expansion of $W_{t,2}^b \sigma \left( W_{t,1}^b x_{t,i}^b \right)$, the most complicated term[12] is

  $$W_{t,2}^b \left( W_{t,1}^b x_{t,i}^b \right)^{\odot D}.$$

  By applying (A.31), we can rewrite it as a sum of product of $3D + 2$ matrices, namely $\sum_{j_1 \in [p_1]} W_{t,2}^b \left( e_{p_1,j_1}^T W_{t,1}^b x_{t,i}^b \right)^D e_{p_1,j_1}$. Similarly, the most complicated term in the expansion of $W_{t,3}^b \sigma \left( W_{t,2}^b \sigma \left( W_{t,1}^b x_{t,i}^b \right) \right)$ is a sum of product of $D(3D + 2) + 2 = 3D^2 + 2D + 2$ matrices, namely

  $$\sum_{j_2} \left( e_{p_2,j_2}^T \left( \sum_{j_1} W_{t,2}^b \left( e_{p_1,j_1}^T W_{t,1}^b x_{t,i}^b e_{p_1,j_1} \right)^D \right) \right)^D e_{p_2,j_2}.$$

  We can use induction to prove that the number of matrices needed for layer $s$ should be $D$ times the number of matrices needed for layer $s - 1$ plus 2. For a general $H$-layer network, we require $\overline{K} := 3D^{H-1} + 2D^{H-2} + \cdots + 2D + 2$ matrices to represent the most complicated term in

  $$W_{t,H}^b \sigma \left( W_{t,H-1}^b \sigma \left( \cdots \sigma \left( W_{t,1}^b x_{t,i}^b \right) \right) \right).$$

  Thus we set $K = 2\overline{K} - 1 = 6D^{H-1} + 4D^{H-2} + \cdots + 4D + 3$ due to the square operator in the norm and minus one by taking the gradient with respect to $W_{t,s}^b$.

- $J = \left[ 2 \left( D^{H-1} + \cdots + D + 1 \right) p_1^{H-1} p_2^{H-2} \cdots p_{H-1} D^{H-1} \right]^2$
  From the derivation above, we can see that the, in the expansion of

  $$W_{t,H}^b \sigma \left( W_{t,H-1}^b \sigma \left( \cdots \sigma \left( W_{t,1}^b x_{t,i}^b \right) \right) \right),$$

  the most complicated term consists of $p_1^{H-1} p_2^{H-2} \cdots p_{H-1}$ terms of product of matrices and $W_{t,1}^b$ appears most frequently in each of these products ($D^{H-1}$ times). Besides, as there are in total of $D^{H-1} + \cdots + D + 1$ terms if simply replace the activation function $\sigma$ by the equivalent polynomial, we end up with $2 \left( D^{H-1} + \cdots + D + 1 \right) p_1^{H-1} p_2^{H-2} \cdots p_{H-1} D^{H-1}$ terms

---

[11]As we can always padding identity matrices to $M_{s,i,j,k}$, thus the values of $J$ and $K$ are not unique.

[12]We ignore the constant coefficient $\sigma_D$ here for convenience.

for $W_{t,H}^b \sigma \left( W_{t,H-1}^b \sigma \left( \cdots \sigma \left( W_{t,1}^b x_{t,i}^b \right) \right) \right) - W_H^* \sigma \left( W_{H-1}^* \sigma \left( \cdots \sigma \left( W_1^* x_{t,i}^b \right) \right) \right)$. By taking the square, we expect

$$J = \left[ 2 \left( D^{H-1} + \cdots + D + 1 \right) p_1^{H-1} p_2^{H-2} \cdots p_{H-1} D^{H-1} \right]^2.$$

Again, the representation in (A.34) aligns with the right-hand side of (A.3) and (A.4). Thus we can further analyze the dynamics of polynomially-activated networks in a similar manner as in Appendix A.2.1.

**A.2.4. Deep Networks with General Activation Functions.** In this section, we discuss the extension from a polynomially-activated network to a neural network with general activation functions under mild assumptions. Given a neural network

$$f^S(x) := W_H^S \sigma^S \left( W_{H-1}^S \cdots \sigma^S \left( W_1^S x \right) \right)$$

with the population loss

$$\mathcal{L}(w^S) = \mathbb{E}_{x \sim \mathcal{D}} \left[ \frac{1}{2} \left\| W_H^S \sigma^S \left( W_{H-1}^S \cdots \sigma^S \left( W_1^S x \right) \right) - W_H^* \sigma^S \left( W_{H-1}^* \cdots \sigma^S \left( W_1^* x \right) \right) \right\|^2 \right],$$

we define the gradient corresponding to each sample $x_{t,i}, i \in [b]$ as[13]

$$g_{t,k,i}^S := \nabla_{W_{t,k}^S} \left( \frac{1}{2} \left\| W_{t,H}^S \sigma^S \left( W_{t,H-1}^S \cdots \sigma^S \left( W_{t,1}^S x_{t,i} \right) \right) - \right.$$

$$\left. - W_H^* \sigma^S \left( W_{H-1}^* \cdots \sigma^S \left( W_1^* x_{t,i} \right) \right) \right\|^2 \right), \quad k \in [H].$$

Following Section 3.1 of [40], we define a set of intermediate variables

$$z_{t,0,i}^S = x_{t,i}, \qquad h_{t,1,i}^S = W_{t,1}^S z_{t,0,i}^S,$$

$$z_{t,1,i}^S = \sigma^S \left( h_{t,1,i}^S \right), \qquad h_{t,2,i}^S = W_{t,2}^S z_{t,1,i}^S,$$

$$\vdots, \qquad \vdots$$

$$z_{t,H-1,i}^S = \sigma^S \left( h_{t,H-1,i}^S \right), \qquad h_{t,H,i}^S = W_{t,H}^S z_{t,H-1,i}^S,$$

and $D_{t,k,i}^S = \mathrm{diag} \left( \sigma_S' \left( h_{t,k,i}^S \right) \right)$, where $\sigma_S'$ represents the derivative of the activation function $\sigma^S$ and $\mathrm{diag}(v)$ maps a vector $v$ to its corresponding diagonal representation. The SG estimators over weight matrix $W_{t,k}^S$ are given by

$$g_{t,k}^S := \frac{1}{b} \sum_{i \in [b]} g_{t,k,i}^S$$

$$= \frac{1}{b} \sum_{i \in [b]} W_{t,H}^S D_{t,H-1,i}^S \cdots W_{t,k+2}^S D_{t,k+1,i}^S W_{t,k+1}^S D_{t,k,i}^S \cdot$$

$$\cdot \left[ W_{t,H}^S \sigma^S \left( W_{t,H-1}^S \cdots \sigma^S \left( W_{t,1}^S x_{t,i} \right) \right) - \right.$$

$$\left. - W_H^* \sigma^S \left( W_{H-1}^* \cdots \sigma^S \left( W_1^* x_{t,i} \right) \right) \right] \left( z_{t,k-1,i}^S \right)^T.$$

We further assume that

---

[13] For simplicity, we remove the superscript $b$ in this section.

- $\sigma^S$ is smooth on $\mathbb{R}^p$,
- $\|x_{t,i}\|$ is bounded, i.e., there exists a positive constant $C_x$ such that $\|x_{t,i}\| \leqslant C_x, \forall t \in [T], i \in [b]$,
- $\left\|W_{t,k}^S\right\|$ is bounded, i.e., there exists a positive constant $C_W$ such that $\left\|W_{t,k}^S\right\| \leqslant C_W$ for all $x_{t,i} \sim \mathcal{D}$,
- $\left\|h_{t,k,i}^S\right\|$ is bounded, i.e., there exists a constant $C_h$ such that $\left\|h_{t,k,i}^S\right\| \leqslant C_h$ for all $x_{t,i} \sim \mathcal{D}$. [14]

We denote $\mathcal{R} := [-C_h, C_h]^p$. By the first assumption, there exists a constant $C_S$ such that $\left\|\sigma^S(x)\right\| \leqslant C_s, \forall x \in \mathcal{R}$. Note that $\|h_{t,k,i}\|_\infty \leqslant \|h_{t,k,i}\| \leqslant C_h$, thus $h_{t,k,i} \in \mathcal{R}$ for all $t \in [T], k \in [H], i \in [b]$.

We note that these assumptions hold in several of the neural network training regimes. For example, the Sigmoid function meets the first assumption with $C_S = 1$, $\mathcal{R} = [-C_h, C_h]^p$ for $C_h = C_h(C_W, C_x) < \infty$, and both Sigmoid function and its derivative are Lipschitz continuous.

Similarly, we define a polynomially-activated neural network

$$f^P(x) := W_H^P \sigma^P \left(W_{H-1}^P \cdots \sigma^P \left(W_1^P x\right)\right)$$

where $\sigma^P(\cdot)$ is a polynomial function. The loss function and SG estimators are defined similarly except for switching the superscript $S$ to $P$. We use SGD to optimize the loss of these two neural networks with the same initial points ($W_{0,k} := W_{0,k}^S = W_{0,k}^P, k \in [H]$), ground-truth weights ($W_1^*, \ldots, W_H^*$), samples ($x_{t,i}, i \in [b]$), and learning rate $\alpha_t$ in every iteration.

In the following, we show that, if the polynomial $\sigma^P$ is a good approximation of the activation function $\sigma^S$ over a closed domain $\overline{\mathcal{R}}$[15], then the SG estimators $g_{t,k}^S$ and $g_{t,k}^P, k \in [H]$ are also close enough. Formally, we have

THEOREM A.13. *For any $\epsilon > 0$ and time step $T \in \mathbb{N}^+$, there exists a polynomial $\sigma^P(\cdot)$ (depending on $\epsilon, \sigma^S$, and $T$) such that $\left\|g_{T,k}^S - g_{T,k}^P\right\| \leqslant \epsilon, k \in [H]$.*

*Outline of the Proof.* We choose a polynomial function $\sigma^P$ such that

$$\left\|\sigma^S(x) - \sigma^P(x)\right\| \leqslant \epsilon' \quad \text{and} \quad \left\|\sigma'_S(x) - \sigma'_P(x)\right\| \leqslant \epsilon'$$

both hold over $\overline{\mathcal{R}} := [-2C_h, 2C_h]^p$ and $\mathcal{O}(\epsilon') < C_h$. The exact value of $\epsilon' < 1$ is determined later[16]. In the following, we induct on $t$ to show that

(1) $\left\|W_{t,k}^S - W_{t,k}^P\right\| \leqslant \mathcal{O}(\epsilon'), k \in [H]$,

(2) $\left\|h_{t,k,i}^S - h_{t,k,i}^P\right\| \leqslant \mathcal{O}(\epsilon'), k \in [H], i \in [b]$,

(3) $\left\|z_{t,k,i}^S - z_{t,k,i}^P\right\| \leqslant \mathcal{O}(\epsilon'), k \in [H], i \in [b]$,

(4) $\left\|D_{t,k,i}^S - D_{t,k,i}^P\right\| \leqslant \mathcal{O}(\epsilon'), k \in [H], i \in [b]$,

(5) $h_{t,k,i}^P \in \overline{\mathcal{R}}, k \in [H], i \in [b]$,

---

[14] In fact, $C_h$ can be expressed as a function of $C_W, C_x$, and $\|\sigma^S(\cdot)\|$. For example, taking $C_{S,0} = C_x$ and we further find a constant $C_{S,k}$ such that $\|\sigma^S(x)\| \leqslant C_{S,k}$ holds for all $\|x\| \leqslant C_W C_{S,k-1}, k \in [H-1]$, then we have $h_{t,k,i}^S = W_{t,k}^S \sigma^S\left(h_{t,k-1,i}^S\right) \leqslant C_W C_{S,k}$. Taking $C_h = C_W \max_{k \in [H]}\{C_{S,k}\}$ satisfies the assumption.

[15] The rigorous definition of $\overline{\mathcal{R}}$ is provided in the proof.

[16] Note that this polynomial is guaranteed to exist since the general activation function $\sigma^S$ is continuous over the compact domain $\overline{\mathcal{R}}$.

(6) $\left\| g^S_{t,k} - g^P_{t,k} \right\| \leqslant \mathcal{O}(\epsilon'), k \in [H],$

where $\mathcal{O}(\cdot)$ is used to hide constants that relate to $L_S, L'_S, C_S, C_W, C_h, C_x, d_k, k \in [H]$ and are independent of $\epsilon'$. In the following, we use $(1)_t, \ldots, (5)_t$ to represent the statements at time step $t$, respectively. For $(2), (3), (4)$, and $(5)$, we use $(2)_{t,k}, \ldots, (5)_{t,k}$ to specify the statements for the $k$-th layer at time step $t$, respectively.

For $t = 0$, $(1)_t$ is obvious since $W^S_{0,k} = W^P_{0,k}, k \in [H]$.

For $t \geqslant 0$, $(1)_t \Rightarrow (2)_t, (3)_t, (4)_t$, we further induct on $k$ to prove them for any given $t$.

- $k = 1, (1)_t \Rightarrow (2)_{t,1}$

$$\left\| h^S_{t,1,i} - h^P_{t,1,i} \right\| = \left\| W^S_{t,1} z^S_{t,0,i} - W^P_{t,1} z^P_{t,0,i} \right\| \leqslant \left\| W^S_{t,1} - W^P_{t,1} \right\| \left\| x_{t,i} \right\|$$
$$\leqslant \mathcal{O}\left(\epsilon'\right) C_x = \mathcal{O}\left(\epsilon'\right).$$

- $k \in [H], (2)_{t,k} \Rightarrow (5)_{t,k}$

$$\left\| h^P_{t,k,i} \right\|_\infty \leqslant \left\| h^P_{t,k,i} \right\| \leqslant \left\| h^S_{t,k,i} - h^P_{t,k,i} \right\| + \left\| h^S_{t,k,i} \right\| \leqslant \mathcal{O}\left(\epsilon'\right) + C_h \leqslant 2C_h$$

- $k \in [H - 1], (2)_{t,k}, (5)_{t,k} \Rightarrow (3)_{t,k}$

$$\left\| z^S_{t,k,i} - z^P_{t,k,i} \right\| = \left\| \sigma^S\left(h^S_{t,k,i}\right) - \sigma^P\left(h^P_{t,k,i}\right) \right\|$$
$$\leqslant \left\| \sigma^S\left(h^S_{t,k,i}\right) - \sigma^P\left(h^S_{t,k,i}\right) \right\| + \left\| \sigma^P\left(h^S_{t,k,i}\right) - \sigma^P\left(h^P_{t,k,i}\right) \right\|$$
$$\leqslant \epsilon' + L_P \left\| h^S_{t,k,i} - h^P_{t,k,i} \right\|$$
$$\leqslant \epsilon' + L_P \mathcal{O}\left(\epsilon'\right) = \mathcal{O}\left(\epsilon'\right)$$

- $k \in [2 : H], (3)_{t,k-1} \Rightarrow (2)_{t,k}$

$$\left\| h^S_{t,k,i} - h^P_{t,k,i} \right\| = \left\| W^S_{t,k} z^S_{t,k-1,i} - W^P_{t,k} z^P_{t,k-1,i} \right\|$$
$$= \left\| W^S_{t,k} z^S_{t,k-1,i} - W^P_{t,k} z^S_{t,k-1,i} + W^P_{t,k} z^S_{t,k-1,i} - W^P_{t,k} z^P_{t,k-1,i} \right\|$$
$$\leqslant \left\| W^S_{t,k} - W^P_{t,k} \right\| \left\| z^S_{t,k-1,i} \right\| + \left\| W^P_{t,k} \right\| \left\| z^S_{t,k-1,i} - z^P_{t,k-1,i} \right\|$$
$$\leqslant \mathcal{O}\left(\epsilon'\right) \left\| \sigma^S\left(h^S_{t,k-1,i}\right) \right\| + \left( \left\| W^P_{t,k} - W^S_{t,k} \right\| + \left\| W^S_{t,k} \right\| \right) \mathcal{O}\left(\epsilon'\right)$$
$$\leqslant C_S \mathcal{O}\left(\epsilon'\right) + \left( \mathcal{O}\left(\epsilon'\right) + C_W \right) \mathcal{O}\left(\epsilon'\right)$$
$$\leqslant \mathcal{O}\left(\epsilon'\right)$$

1381  $\bullet$ $k \in [H], (2)_{t,k} \Rightarrow (4)_{t,k}$

$$\left\| D_{t,k,i}^S - D_{t,k,i}^P \right\| = \left\| \mathrm{diag}\left( \sigma_S'\left( h_{t,k,i}^S \right) \right) - \mathrm{diag}\left( \sigma_P'\left( h_{t,k,i}^P \right) \right) \right\|$$

$$= \left\| \sigma_S'\left( h_{t,k,i}^S \right) - \sigma_P'\left( h_{t,k,i}^P \right) \right\|_\infty \leqslant \left\| \sigma_S'\left( h_{t,k,i}^S \right) - \sigma_P'\left( h_{t,k,i}^P \right) \right\|$$

$$\leqslant \left\| \sigma_S'\left( h_{t,k,i}^S \right) - \sigma_P'\left( h_{t,k,i}^S \right) \right\| + \left\| \sigma_P'\left( h_{t,k,i}^S \right) - \sigma_P'\left( h_{t,k,i}^P \right) \right\|$$

$$\leqslant \epsilon' + L_P' \left\| h_{t,k,i}^S - h_{t,k,i}^P \right\|$$

$$\leqslant \epsilon' + L_P' \mathcal{O}\left( \epsilon' \right) = \mathcal{O}\left( \epsilon' \right)$$

For $t \geqslant 0, (1)_t + \cdots + (5)_t \Rightarrow (6)_t$, we denote

$$h_{t,i}^{S*} := W_H^* \sigma^S \left( W_{H-1}^* \cdots \sigma^S \left( W_t^* x_{t,i} \right) \right)$$

and

$$h_{t,i}^{P*} := W_H^* \sigma^P \left( W_{H-1}^* \cdots \sigma^P \left( W_t^* x_{t,i} \right) \right).$$

1388  Note that $\left\| g_{t,k}^S - g_{t,k}^P \right\| = \left\| \frac{1}{b} \sum_{i\in[b]} g_{t,k,i}^S - \frac{1}{b} \sum_{i\in[b]} g_{t,k,i}^P \right\| \leqslant \frac{1}{b} \sum_{i\in[b]} \left\| g_{t,k,i}^S - g_{t,k,i}^P \right\|$. For
1389  each $i \in [b]$, we have

$$\left\| g_{t,k,i}^S - g_{t,k,i}^P \right\|$$

$$= \left\| W_{t,H}^S D_{t,H-1,i}^S \cdots W_{t,k+1}^S D_{t,k,i}^S \left( h_{t,H,i}^S - h_{t,i}^{S*} \right) \left( z_{t,k-1,i}^S \right)^T \right.$$

$$\left. - W_{t,H}^P D_{t,H-1,i}^P \cdots W_{t,k+1}^P D_{t,k,i}^P \left( h_{t,H,i}^P - h_{t,i}^{P*} \right) \left( z_{t,k-1,i}^P \right)^T \right\|$$

$$\leqslant \left\| W_{t,H}^S D_{t,H-1,i}^S \cdots W_{t,k+1}^S D_{t,k,i}^S h_{t,H,i}^S \left( z_{t,k-1,i}^S \right)^T \right.$$

$$\left. - W_{t,H}^P D_{t,H-1,i}^P \cdots W_{t,k+1}^P D_{t,k,i}^P h_{t,H,i}^P \left( z_{t,k-1,i}^P \right)^T \right\| +$$

$$+ \left\| W_{t,H}^S D_{t,H-1,i}^S \cdots W_{t,k+1}^S D_{t,k,i}^S h_{t,i}^{S*} \left( z_{t,k-1,i}^S \right)^T \right.$$

(A.35)  $$\left. - W_{t,H}^P D_{t,H-1,i}^P \cdots W_{t,k+1}^P D_{t,k,i}^P h_{t,i}^{P*} \left( z_{t,k-1,i}^P \right)^T \right\|.$$

For the first item in (A.35), we have

$$\left\| W_{t,H}^S D_{t,H-1,i}^S \cdots W_{t,k+1}^S D_{t,k,i}^S h_{t,H,i}^S \left(z_{t,k-1,i}^S\right)^T - W_{t,H}^P D_{t,H-1,i}^P \cdots W_{t,k+1}^P D_{t,k,i}^P h_{t,H,i}^P \left(z_{t,k-1,i}^P\right)^T \right\|$$

$$= \left\| W_{t,H}^S D_{t,H-1,i}^S \cdots W_{t,k+1}^S D_{t,k,i}^S W_{t,H}^S z_{t,H-1,i}^S \left(z_{t,k-1,i}^S\right)^T - \right.$$

$$\left. - W_{t,H}^P D_{t,H-1,i}^P \cdots W_{t,k+1}^P D_{t,k,i}^P W_{t,H}^P z_{t,H-1,i}^P \left(z_{t,k-1,i}^P\right)^T \right\|$$

$$\leqslant \left\| W_{t,H}^S D_{t,H-1,i}^S \cdots W_{t,k+1}^S D_{t,k,i}^S W_{t,H}^S z_{t,H-1,i}^S - \right.$$

$$\left. - W_{t,H}^P D_{t,H-1,i}^P \cdots W_{t,k+1}^P D_{t,k,i}^P W_{t,H}^P z_{t,H-1,i}^P \right\| \left\| z_{t,k-1,i}^P \right\| +$$

$$+ \left\| W_{t,H}^S D_{t,H-1,i}^S \cdots W_{t,k+1}^S D_{t,k,i}^S W_{t,H}^S z_{t,H-1,i}^S \right\| \left\| z_{t,k-1,i}^S - z_{t,k-1,i}^P \right\|$$

$$\leqslant \left\| W_{t,H}^S D_{t,H-1,i}^S \cdots W_{t,k+1}^S D_{t,k,i}^S W_{t,H}^S z_{t,H-1,i}^S - \right.$$

$$\left. - W_{t,H}^P D_{t,H-1,i}^P \cdots W_{t,k+1}^P D_{t,k,i}^P W_{t,H}^P z_{t,H-1,i}^P \right\| \cdot \sqrt{d_{k-1}} \left\| z_{t,k-1,i}^P \right\|_\infty +$$

$$+ \left\| W_{t,H}^S \right\| \left\| D_{t,H-1,i}^S \right\| \cdots \left\| W_{t,k+1}^S \right\| \left\| D_{t,k,i}^S \right\| \left\| W_{t,H}^S \right\| \left\| z_{t,H-1,i}^S \right\| \mathcal{O}\left(\epsilon'\right)$$

$$\leqslant \left\| W_{t,H}^S D_{t,H-1,i}^S \cdots W_{t,k+1}^S D_{t,k,i}^S W_{t,H}^S z_{t,H-1,i}^S - \right.$$

$$\left. - W_{t,H}^P D_{t,H-1,i}^P \cdots W_{t,k+1}^P D_{t,k,i}^P W_{t,H}^P z_{t,H-1,i}^P \right\| \cdot \sqrt{d_{k-1}} C_S +$$

$$+ C_W^{H-k+1} C_S'^{H-k} \sqrt{d_{H-1}} C_S \mathcal{O}\left(\epsilon'\right)$$

$$\leqslant \left\| W_{t,H}^S D_{t,H-1,i}^S \cdots W_{t,k+1}^S D_{t,k,i}^S W_{t,H}^S - W_{t,H}^P D_{t,H-1,i}^P \cdots W_{t,k+1}^P D_{t,k,i}^P W_{t,H}^P \right\| \left\| z_{t,H-1,i}^P \right\| \cdot \sqrt{d_{k-1}} C_S$$

$$+ \left\| W_{t,H}^S D_{t,H-1,i}^S \cdots W_{t,k+1}^S D_{t,k,i}^S W_{t,H}^S z_{t,H-1,i}^S \right\| \left\| z_{t,H-1,i}^S - z_{t,H-1,i}^P \right\| \sqrt{d_{k-1}} C_S + \mathcal{O}\left(\epsilon'\right)$$

$$= \left\| W_{t,H}^S D_{t,H-1,i}^S \cdots W_{t,k+1}^S D_{t,k,i}^S W_{t,H}^S - \right.$$

$$\left. - W_{t,H}^P D_{t,H-1,i}^P \cdots W_{t,k+1}^P D_{t,k,i}^P W_{t,H}^P \right\| \left\| z_{t,H-1,i}^P \right\| \cdot \sqrt{d_{k-1}} C_S + \mathcal{O}\left(\epsilon'\right) + \mathcal{O}\left(\epsilon'\right)$$

$$\leqslant \left\| W_{t,H}^S D_{t,H-1,i}^S \cdots W_{t,k+1}^S D_{t,k,i}^S W_{t,H}^S - \right.$$

$$\left. - W_{t,H}^P D_{t,H-1,i}^P \cdots W_{t,k+1}^P D_{t,k,i}^P W_{t,H}^P \right\| \cdot \sqrt{d_{H-1} d_{k-1}} C_S^2 + \mathcal{O}\left(\epsilon'\right)$$

$$\leqslant \cdots\cdots$$

$$\leqslant \mathcal{O}\left(\epsilon'\right).$$

Similarly, we can show that the second term in (A.35) is also bounded by $\mathcal{O}\left(\epsilon'\right)$. Thus we have $\left\| g_{t,k}^S - g_{t,k}^P \right\| \leqslant \frac{1}{b} \sum_{i \in [b]} \left\| g_{t,k,i}^S - g_{t,k,i}^P \right\| \leqslant \mathcal{O}\left(\epsilon'\right)$.

For $t \geqslant 0, (1)_t + (5)_t \Rightarrow (1)_{t+1}$, we have

$$\left\| W_{t+1,k}^S - W_{t+1,k}^P \right\| = \left\| \left(W_{t,k}^S - \alpha_t g_{t,k}^S\right) - \left(W_{t,k}^P - \alpha_t g_{t,k}^P\right) \right\|$$

$$\leqslant \left\| W_{t,k}^S - W_{t,k}^P \right\| + \alpha_t \left\| g_{t,k}^S - g_{t,k}^P \right\|$$

$$\leqslant \mathcal{O}\left(\epsilon'\right) + \alpha_t \mathcal{O}\left(\epsilon'\right) = \mathcal{O}\left(\epsilon'\right).$$

With the above steps, we have finished the induction. The proof is achieved by taking $\epsilon'$ small enough such that $\mathcal{O}(\epsilon') < \epsilon$ at time step $T$. $\square$

While the above theorem only discuss the closeness of $g_{T,k}^S$ and $g_{T,k}^P$, it is worth mentioning that the same statement holds for all pairs of intermediate variables or even composition of them. In fact, we have the following generalized theorem.

THEOREM A.14. *For any $\epsilon > 0$ and time step $T \in \mathbb{N}^+$, there exists a polynomial $\sigma^P(\cdot)$ (depending on $\epsilon, \sigma^S$, and $T$) such that*

$$\left\| \mathrm{tr}\left( C\left( \bigotimes_i \prod_j M_{i,j}^S \right) \right) \prod_j M_{0,j}^S - \mathrm{tr}\left( C\left( \bigotimes_i \prod_j M_{i,j}^P \right) \right) \prod_j M_{0,j}^P \right\| < \epsilon,$$

where $M_{i,j}^S$ takes values in $W_{0:t}^S \bigcup G_{0:T}^S \bigcup W^* \bigcup \overline{\mathcal{C}}$ and $M_{i,j}^P$ takes the corresponding variable in the polynomially-activated network as of $M_{i,j}^S$.

Together with the closed-form representation of the expected value of

$$\text{tr}\left( C \left( \bigotimes_i \prod_j M_{i,j}^P \right) \right) \prod_j M_{0,j}^P$$

given $\mathcal{F}_0$, we are able to provide an approximation of $\text{tr}\left( C \left( \bigotimes_i \prod_j M_{i,j}^S \right) \right) \prod_j M_{0,j}^S$ at any time step $T$ with any precision. In other words, we have provided an approximation for a generalized form of mixed product at time step $t$ using solely the initial weights $W_0^b$ and other constant matrices. Similarly, Theorem 3.10, which shows the decreasing property of the SG estimators, can also be extended to general neural networks as well as other general neural networks.

## REFERENCES

[1] M. S. Acharya, A. Armaan, and A. S. Antony, *A comparison of regression models for prediction of graduate admissions*, in 2019 International Conference on Computational Intelligence in Data Science, 2019, pp. 1–5.

[2] Z. Allen-Zhu, Y. Li, and Y. Liang, *Learning and generalization in overparameterized neural networks, going beyond two layers*, arXiv preprint arXiv:1811.04918, (2018).

[3] L. Bottou, *Stochastic gradient learning in neural networks*, Proceedings of Neuro-Nimes, 91 (1991), p. 12.

[4] L. Bottou, *Online learning and stochastic approximations*, On-line Learning in Neural Networks, 17 (1998), p. 142.

[5] L. Bottou, F. E. Curtis, and J. Nocedal, *Optimization methods for large-scale machine learning*, SIAM Review, 60 (2018), pp. 223–311.

[6] S. Boyd, S. P. Boyd, and L. Vandenberghe, *Convex optimization*, Cambridge university press, 2004.

[7] S. Du, X. Zhai, B. Poczos, and A. Singh, *Gradient descent provably optimizes over-parameterized neural networks*, arXiv preprint arXiv:1810.02054, (2018).

[8] J. Fan, C. Ma, and Y. Zhong, *A selective overview of deep learning*, arXiv preprint arXiv:1904.05526, (2019).

[9] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*, MIT press, 2016.

[10] E. Gorbunov, F. Hanzely, and P. Richtárik, *A unified theory of sgd: Variance reduction, sampling, quantization and coordinate descent*, in International Conference on Artificial Intelligence and Statistics, 2020, pp. 680–690.

[11] R. M. Gower, N. Loizou, X. Qian, A. Sailanbayev, E. Shulgin, and P. Richtárik, *SGD: General analysis and improved rates*, in International Conference on Machine Learning, 2019, pp. 5200–5209.

[12] P. Goyal, P. Dollár, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He, *Accurate, large minibatch SGD: Training Imagenet in 1 hour*, arXiv preprint arXiv:1706.02677, (2017).

[13] K. He, X. Zhang, S. Ren, and J. Sun, *Deep residual learning for image recognition*, in Proceedings of the IEEE conference on Computer Vision and Pattern Recognition, 2016, pp. 770–778.

[14] G. Hinton, O. Vinyals, and J. Dean, *Distilling the knowledge in a neural network*, arXiv preprint arXiv:1503.02531, (2015).

[15] S. Hochreiter and J. Schmidhuber, *Flat minima*, Neural Computation, 9 (1997), pp. 1–42.

[16] E. Hoffer, I. Hubara, and D. Soudry, *Train longer, generalize better: closing the generalization gap in large batch training of neural networks*, in Advances in Neural Information Processing Systems, 2017, pp. 1731–1741.

[17] S. Jastrzebski, Z. Kenton, D. Arpit, N. Ballas, A. Fischer, Y. Bengio, and A. Storkey, *Three factors influencing minima in SGD*, arXiv preprint arXiv:1711.04623, (2017).

[18] R. Johnson and T. Zhang, *Accelerating stochastic gradient descent using predictive variance reduction*, in Advances in Neural Information Processing Systems, 2013, pp. 315–323.

[19] N. S. Keskar, J. Nocedal, P. T. P. Tang, D. Mudigere, and M. Smelyanskiy, *On large-batch training for deep learning: Generalization gap and sharp minima*, in 5th International Conference on Learning Representations, 2017, 2017.

[20] A. Khaled and P. Richtárik, *Better theory for sgd in the nonconvex world*, arXiv preprint arXiv:2002.03329, (2020).

[21] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, *Gradient-based learning applied to document recognition*, Proceedings of the Institute of Electrical and Electronics Engineers, 86 (1998), pp. 2278–2324.

[22] Y. LeCun, L. Bottou, G. B. Orr, and K.-R. Müller, *Efficient backprop*, in Neural networks: Tricks of the trade, Springer, 2012, pp. 9–48.

[23] L. Lei, C. Ju, J. Chen, and M. I. Jordan, *Non-convex finite-sum optimization via SCSG methods*, in Advances in Neural Information Processing Systems, 2017, pp. 2348–2358.

[24] M. Li, T. Zhang, Y. Chen, and A. J. Smola, *Efficient mini-batch training for stochastic optimization*, in Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2014, pp. 661–670.

[25] Q. Li, C. Tai, and W. E, *Stochastic modified equations and adaptive stochastic gradient algorithms*, in Proceedings of the 34th International Conference on Machine Learning, PMLR, 2017, pp. 2101–2110.

[26] Q. Li, C. Tai, and E. Weinan, *Stochastic modified equations and adaptive stochastic gradient algorithms*, in International Conference on Machine Learning, PMLR, 2017, pp. 2101–2110.

[27] Y. Li and Y. Liang, *Learning overparameterized neural networks via stochastic gradient descent on structured data*, arXiv preprint arXiv:1808.01204, (2018).

[28] S. Mandt, M. D. Hoffman, and D. M. Blei, *Stochastic gradient descent as approximate bayesian inference*, The Journal of Machine Learning Research, 18 (2017), pp. 4873–4907.

[29] Q. Meng, Y. Wang, W. Chen, T. Wang, Z.-M. Ma, and T.-Y. Liu, *Generalization error bounds for optimization algorithms via stability*, arXiv preprint arXiv:1609.08397, (2016).

[30] W. Mou, L. Wang, X. Zhai, and K. Zheng, *Generalization bounds of SGLD for non-convex learning: Two theoretical viewpoints*, in Conference On Learning Theory, 2018, pp. 605–638.

[31] A. Neelakantan, L. Vilnis, Q. V. Le, I. Sutskever, L. Kaiser, K. Kurach, and J. Martens, *Adding gradient noise improves learning for very deep networks*, arXiv preprint arXiv:1511.06807, (2015).

[32] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al., *Pytorch: An imperative style, high-performance deep learning library*, in Advances in Neural Information Processing Systems, 2019, pp. 8024–8035.

[33] K. Petersen, M. Pedersen, et al., *The matrix cookbook, vol. 7*, Technical University of Denmark, 15 (2008).

[34] N. L. Roux, M. Schmidt, and F. R. Bach, *A stochastic gradient method with an exponential convergence rate for finite training sets*, in Advances in Neural Information Processing Systems, 2012, pp. 2663–2671.

[35] M. Schmidt, N. Le Roux, and F. Bach, *Minimizing finite sums with the stochastic average gradient*, Mathematical Programming, 162 (2017), pp. 83–112.

[36] S. Shalev-Shwartz and T. Zhang, *Stochastic dual coordinate ascent methods for regularized loss minimization*, Journal of Machine Learning Research, 14 (2013), pp. 567–599.

[37] P. Y. Simard, D. Steinkraus, and J. C. Platt, *Best practices for convolutional neural networks applied to visual document analysis*, in Seventh International Conference on Document Analysis and Recognition, 2013, pp. 958–963.

[38] S. L. Smith and Q. V. Le, *A bayesian perspective on generalization and stochastic gradient descent*, arXiv preprint arXiv:1710.06451, (2017).

[39] R. Sun, *Optimization for deep learning: theory and algorithms*, arXiv preprint arXiv:1912.08957, (2019).

[40] R. Sun, *Optimization for deep learning: An overview*, Journal of the Operations Research Society of China, 8 (2020), pp. 249–294.

[41] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. R. Salakhutdinov, and Q. V. Le, *Xlnet: Generalized autoregressive pretraining for language understanding*, in Advances in Neural Information Processing Systems, 2019, pp. 5754–5764.

[42] X. Zhang, J. Zhao, and Y. LeCun, *Character-level convolutional networks for text classification*, in Advances in Neural Information Processing Systems, 2015, pp. 649–657.

[43] Y. Zhang, P. Liang, and M. Charikar, *A hitting time analysis of stochastic gradient langevin dynamics*, in Conference on Learning Theory, 2017, pp. 1980–2022.