

Unlocking the Theory Behind Scaling 1-Bit Neural Networks

Majid Daliri¹, Zhao Song², Chiwun Yang³

¹New York University, ²Simons Institute for the Theory of Computing, UC Berkeley, ³Sun Yat-sen University

daliri.majid@nyu.edu, magic.linuxkde@gmail.com, christiannyang37@gmail.com

Recently, 1-bit Large Language Models (LLMs) have emerged, showcasing an impressive combination of efficiency and performance that rivals traditional LLMs. Research by Wang et al. [1], Ma et al. [2] indicates that the performance of these 1-bit LLMs progressively improves as the number of parameters increases, hinting at the potential existence of a *Scaling Law in 1-bit Neural Networks*. This paper presents the **first theoretical result** that rigorously establishes this scaling law for 1-bit models. Our analysis starts with initializing a 1-bit two-layer linear network. We prove that, despite the constraint of weights restricted to $\{-1, +1\}$, its training dynamics inevitably align with kernel behavior as the network width grows. This theoretical breakthrough guarantees convergence of the 1-bit model to an arbitrarily small loss as width increases. Furthermore, we introduce the concept of the generalization difference, defined as the gap between the outputs of 1-bit networks and their full-precision counterparts, and demonstrate that this difference maintains a negligible level under the over-parameterization setting. Building on the work of Kaplan et al. [3], we examine how the training loss scales as a power-law function of the model size, dataset size, and computational resources utilized for training. Our findings underscore the promising potential of scaling 1-bit neural networks, suggesting that int1 could become the standard in future neural network precision.

1 Introduction

Large-scale neural networks, particularly Large Language Models (LLMs) [4, 5] and Large Multi-model Models (LMMs) [6, 7], are becoming increasingly relevant to our day-to-day lives, finding a huge variety of applications in both the workplace and at home [8, 9]. However, it is expensive to deploy and run these models due to their substantial computational requirements, large memory footprints, and energy consumption [10–12]. This is especially true for resource-constrained environments, such as mobile devices, edge computing, or companies with limited infrastructure [13–15]. To make these models more efficient and accessible, quantization techniques are used, which reduce the precision of the model’s parameters (such as weights and activations) from floating-point numbers to lower-bit representations (e.g., 8-bit or even lower) [16–20]. Quantization reduces the memory and computational costs of inference, enabling faster processing with less energy, while maintaining a comparable level of performance. This optimization allows language models to be more practical, scalable, and sustainable for widespread use across various platforms [21–23].

In particular, quantization techniques could be primarily divided into two methods: Post-Training Quantization (PTQ) [24–26] and Quantization-Aware Training (QAT) [1, 2, 27]. PTQ methods, including uniform and non-uniform quantization, conveniently convert pre-trained model weights and activations to lower-bit representations post-training. However, this leads to accuracy loss, especially in lower precision, as the model is not optimized for these quantized representations and significant shifts in weight distribution occur [28]. The alternative, Quantization-Aware Training (QAT), incorporates quantization during training, allowing the model to fine-tune and adapt its parameters to the quantized representation, compensating for quantization errors. Therefore, compared to PTQ, QAT maintains higher accuracy and robustness even in lower precision.

Recent studies [1, 2, 29, 30] have shown that 1-bit LLMs, most of which have matrix weights in the range of $\{-1, +1\}$, can be trained from scratch to deliver performance that rivals that of standard LLMs. These models exhibit remarkable efficiency, particularly in terms of scaling laws. Experimental results indicate that the performance of the 1-bit model improves as the number of parameters increases, a principle that mirrors the training approach utilized in standard LLMs [3]. Despite the demonstrated efficiency of quantization methods, our understanding of the training mechanism for quantization remains limited. Specifically, it remains unclear how and why the 1-bit QAT enhances learning capability as the number of neurons in the model is scaled up. In addition, we are also concerned about whether the quantization method damages the generalization ability compared to full precision networks.

In this study, we initially apply the Neural Tangent Kernel (NTK) framework to delve into the optimization and generalization issues associated with a two-layer linear network operating in 1-bit (int1) precision, as detailed in Section 4. We introduce a 1-bit quantization method to the hidden-layer weights $W \in \mathbb{R}^{d \times m}$ of the conventional NTK linear network, where d represents the input dimension and m indicates the model’s width. Our analysis reveals that the training dynamics of the 1-bit model approximate kernel behavior as the model width m expands. This key finding paves the way for an established relationship between the theoretically guaranteed loss and the model width, endowing the model with robust learning capabilities akin to kernel regression. Ultimately, the model achieves an insignificantly small training loss, contingent on setting a sufficiently large model width, selecting an appropriate learning rate, and allowing an adequate training duration.

Moreover, Section 5 provides a theoretical confirmation that, within the scaling trend, the disparities in predictions of the 1-bit model from those of the original linear network on identical inputs maintain a negligible value. We assess the error between our 1-bit linear and standard linear networks on both the training and test datasets. Our theorem demonstrates that for any input from these datasets, the absolute error between the two network predictions can be denoted as $\epsilon_{\text{quant}} \leq O(\kappa d \log(md/\delta))$ for scale coefficient $\kappa \leq 1$, model width m , dimension d and failure probability $\delta \in (0, 0.1)$. This indicates that the output behavior of the 1-bit linear model increasingly aligns with that of the standard linear model. The observed similarity on the test dataset validates the generalization similarity, suggesting the feasibility of approximating training neural networks with int1 precision equivalent to full precision.

Finally, in Section 6, we verify our theoretical results by implementing training models to learn complicated functions to compare the difference between 1-bit networks and full precision networks. Firstly, we choose a combination of difficult functions across the exponential function, trigonometric function, logarithmic function, the Lambert W function, the Gamma function, and their combination. Therefore, we sample random data points and split train and test datasets. We next compare how the training loss decreases as the model width m scales up. Besides, as shown in Section 6.3, in the trend of a growing number of parameters, the error of predictions both on training and test input likewise converge as the power-law in 1-bit networks optimization. In particular, we visualize some 1-dimension function to see how the differences of outputs are. We demonstrate the results complying with our theoretical guarantee with a negligible error.

2 Related Work

Efficient Training Methods for Quantized Networks Training large-scale neural networks with quantization introduces significant computational and memory savings, but it also presents challenges in optimization, particularly when dealing with extremely low precision formats like 1-bit or 8-bit. To address these challenges, several efficient training methods have been developed that aim to maintain accuracy while leveraging the benefits of quantization. One key method is Gradient Quantization, where the gradients during backpropagation are quantized to lower precision to reduce memory overhead and bandwidth during distributed training. Techniques like stochastic rounding are used to mitigate the impact of quantization noise, ensuring the training process remains stable and converges effectively.

Another important approach is Low-Rank Factorization [31, 32], which decomposes the large weight matrices in neural networks into smaller matrices, reducing the number of parameters that need to be updated during training. When combined with quantization, this method significantly reduces both the memory footprint and computational complexity, allowing for faster training on hardware with limited resources.

Quantization Techniques for Accelerating Language Models Beyond traditional weight and activation quantization, several advanced methods utilize quantization to enhance the efficiency of large language models (LLMs). One key approach is KV cache quantization [33–36], which reduces the memory footprint of transformer models during inference by quantizing the stored attention keys and values. This method is particularly beneficial for tasks involving long sequences, significantly speeding up inference and lowering memory consumption without a substantial loss in accuracy.

Another effective technique is mixed-precision quantization [37, 38], where different parts of the model are quantized at varying precision levels based on their sensitivity. For example, attention layers might use higher precision (e.g., 16-bit), while feedforward layers are quantized to 8-bit or lower. This balances computational efficiency and model performance. These strategies, combined with methods like activation pruning, showcase how targeted quantization can drastically accelerate LLMs while maintaining their effectiveness in real-world applications.

Neural Tangent Kernel. The study of Neural Tangent Kernel (NTK) [39] focuses on the gradient flow of neural networks during the training process, revealing that neural networks are equivalent to Gaussian processes at initialization in the infinite-width limit. This equivalence has been explored in numerous studies [40–54] that account for the robust performance and learning capabilities of over-parameterized neural networks. The kernel-based analysis framework provided by NTK is gaining popularity for its utility in elucidating the emerging abilities of large-scale neural networks. In a remarkable stride, Arora et al. [55] introduced the first exact algorithm for computing the Convolutional NTK (CNTK). This was followed by Alemohammad et al. [56] who proposed the Recurrent NTK, and Hron et al. [57] who presented the concept of infinite attention via NNGP and NTK for attention networks. These innovative works have showcased the enhanced performance achievable with the application of NTK to various neural network architectures. In a specific study, Malladi et al. [58] examined the training dynamics of fine-tuning Large Language Models (LLMs) using NTK, affirming the efficiency of such approaches.

3 Preliminary

In this section, we give the basic setups of this paper, which includes the introduction of the quantization method in this paper (Section 3.1), our NTK-style problem setup that we aim to solve in this paper (Section 3.2) and recalling the classical NTK setup for a two-layer linear network with ReLU activation function (Section 3.3).

3.1 Quantization

We first show how we reduce the computation of the inner product of two vectors from multiplication and addition operations to addition operations only, which is achieved by binarizing one of the vectors. This method could be extended to matrix multiplication easily since the basic matrix multiplication is to implement the inner product computation of two vectors in parallel. For a vector $w \in \mathbb{R}^d$, we define our quantization function as [1, 2]:

$$\text{Quant}(w) := \text{Sign}\left(\text{Ln}(w)\right) \in \{-1, +1\}^d,$$

where $\text{Ln}(w)$ is the normalization method that is given by: $\text{Ln}(w) := \frac{w - E(w) \cdot \mathbf{1}_d}{\sqrt{V(w)}} \in \mathbb{R}^d$. Specially, we use $E(w) := \frac{1}{d} \sum_{k=1}^d w_k \in \mathbb{R}$ to denote the computational expectation of vector w and use $V(w) := \|w - E(w) \cdot \mathbf{1}_d\|_2^2 \in \mathbb{R}$ to denote the corresponding variance.

Besides, the k^{th} entry of signal function $\text{Sign}(z) \in \mathbb{R}^d$ for $z \in \mathbb{R}^d$, $k \in [d]$ is define by: $\text{Sign}_k(z) := \begin{cases} +1, & z_k \geq 0 \\ -1, & z_k < 0 \end{cases}$. Hence, we have a binary vector $\text{Quant}(w)$ where each entry of it is limited in the range $\{-1, +1\}$, and we denote that $\tilde{w} := \text{Quant}(w)$ to simplify the notation. For any other vector $x \in \mathbb{R}^d$, addition operation $\sum_{k=1}^d \pm x_k$ is sufficient to compute $\langle \tilde{w}, x \rangle$. After that, we introduce the dequantization function to recover the original computation result by showing:

$$\text{Dequant}(\langle \tilde{w}, x \rangle) := \sqrt{V(w)} \cdot \langle \tilde{w}, x \rangle + E(w) \cdot \langle \mathbf{1}, x \rangle.$$

3.2 NTK Problem Setup

Data Points. We consider a supervised learning task with a training dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$, where each data point is under a mild assumption that $\|x_i\|_2 = 1$ and $y_i \leq 1$, $\forall i \in [n]$ [41]. Moreover, we are also concerned about the problem of the generalization of 1-bit models, we define the test dataset to compare 1-bit networks with standard networks, that is $\mathcal{D}_{\text{test}} := \{(x_{\text{test},i}, y_{\text{test},i})\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$, where $\|x_{\text{test},i}\|_2 = 1$ and $y_{\text{test},i} \leq 1$, $\forall i \in [n]$.

Model. Here, we use hidden-layer weights $W = [w_1, w_2, \dots, w_m] \in \mathbb{R}^{d \times m}$ and output-layer weights $a = [a_1, a_2, \dots, a_m]^T \in \mathbb{R}^m$. We consider a two-layer linear model f , which is defined as follows:

$$f(x, W, a) := \kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \text{ReLU}(\text{dq}(\langle \tilde{w}_r, x \rangle)),$$

where $\text{ReLU}(z) := \begin{cases} z, & z \geq 0 \\ 0, & z < 0 \end{cases}$, for all $z \in \mathbb{R}$, $\text{dq} : \mathbb{R} \rightarrow \mathbb{R}$ is a omitted version of dequantization

function $\text{Dequant} : \mathbb{R} \rightarrow \mathbb{R}$, and $\tilde{w}_r := \text{Quant}(w_r)$ as we denoted in previous section, $\kappa \in (0, 1]$ is a scale coefficient. Especially, we initialize each weight vector $w_r, \forall r \in [m]$ by sampling $w_r(0) \sim \mathcal{N}(0, \sigma \cdot I_d)$ with $\sigma = 1$. For output-layer a , we randomly sample $a_r \sim \text{Uniform}\{-1, +1\}$ independently for $r \in [m]$. Additionally, output-layer weight a is fixed during the training.

Training and Straight-Through Estimator (STE). The training loss is measured by quadratic ℓ_2 norm of the difference between model prediction $f(x_i, W, a)$ and ideal output vector y_i . Formally, we consider to train $W(t) = [w_1(t), w_2(t), \dots, w_m(t)] \in \mathbb{R}^{d \times m}$ for $t \geq 0$ utilizing the following loss:

$$\mathbb{L}(t) := \frac{1}{2} \cdot \sum_{i=1}^n \|f(x_i, W(t), a) - y_i\|_2^2. \quad (1)$$

Moreover, since the signal function Sign is not differentiable, we use Straight-Through Estimator (STE) to skip the signal function in back-propagation [1, 2, 59, 60], thus updating the trainable weights $W(t)$. For $t \geq 0$ and denote η as the learning rate, we omit $f_i(t) := f(x_i, W(t), a) \in \mathbb{R}, \forall i \in [n]$, the formulation to update r^{th} column of $W(t)$ for all $r \in [m]$ is given by:

$$w_r(t+1) := w_r(t) - \eta \sum_{i=1}^n (f_i(t) - y_i) \cdot \kappa a_r \mathbf{1}_{\text{dq}(\langle \tilde{w}_r, x_i \rangle) \geq 0} x_i.$$

3.3 Recalling Classic NTK Setup

We now recall the classic NTK setup for the two-layer ReLU linear regression [61–64]. The function is given by: $f'(x, W, a) := \kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \text{ReLU}(\langle w_r, x \rangle)$.

We define that $W'(0) := W(0) \in \mathbb{R}^{d \times m}$ to denote the trainable parameter for classic NTK setup, these two matrices are equal at initialization. For $t \geq 0$, we define the loss of training f' as follows: $\mathbb{L}'(t) := \frac{1}{2} \cdot \sum_{i=1}^n \|f'(x_i, W'(t), a) - y_i\|_2^2$. Then the update of $W'(t)$ is: $W'(t+1) := W'(t) - \eta \cdot \nabla_{W'(t)} \mathbb{L}'(t)$.

4 Kernel Behavior and Training Convergence

We give our convergence analysis for training 1-bit model within the framework of Neural Tangent Kernel (NTK) in this section. First, we state our theoretical results that define the kernel function

in training and show how it converges to NTK and maintains the PD (Positive Definite) property in Section 4.1. Then we demonstrate the arbitrary small loss convergence guarantee of training 1-bit model (Eq. (1)) in Section 4.2. Finally, we give a general version of our theoretical scaling law analysis in Section 4.3.

4.1 Neural Tangent Kernel

Here, we utilize the NTK to describe the training dynamic of the 1-bit model. Following pre-conditions in the previous section, we define a kernel function, that denotes $H(t) \in \mathbb{R}^{n \times n}$ (Gram matrix). Especially, the (i, j) -th entry of $H(t)$ is given by:

$$H_{i,j}(t) := \kappa^2 \frac{1}{m} x_i^\top x_j \sum_{r=1}^m \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_j \rangle) \geq 0}. \quad (2)$$

We define the formal NTK as $H^* := H(0) \in \mathbb{R}^{n \times n}$. Additionally, there's a commonly introduced assumption in NTK analysis: we denote the minimum value of eigenvalues of A with $\lambda_{\min}(A)$ for any $A \in \mathbb{R}^{n \times n}$. In our work's context, we presuppose that H is a Positive-definite (PD) matrix, meaning that $\lambda_{\min}(H^*) > 0$ [41].

1-Bit ReLU Pattern. The pattern of the Rectified Linear Unit (ReLU) function is determined by the indicator of function activation. As illustrated by Du et al. [41], in the settings of Section 3.3, the event $\mathbf{1}_{\langle w_r(0), x \rangle \geq 0} \neq \mathbf{1}_{\langle w_r(t), x \rangle \geq 0}$ happens infrequently for any $w, x \in \mathbb{R}^d$ that satisfies $\|w - w_r(0)\|_2 \leq R$. Notably, $R := \max_{r \in [m]} \|w_r(t) - w_r(0)\|_2 = \eta \|\sum_{\tau=1}^t \Delta w_r(\tau)\|_2$. In our analysis, for Eq. (2), the event $\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(0), x \rangle) \geq 0} \neq \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x \rangle) \geq 0}$ is also unlikely to occur during training.

The convergence of $H(t)$ towards H^* , as well as the property of $H(t)$ being a PD matrix for any $t \geq 0$, can be validated by the following lemma:

Lemma 4.1 (NTK convergence and PD property during the training, informal version of Lemma G.5). *Assume $\lambda_{\min}(H^*) > 0$. $\delta \in (0, 1)$, define $D := \max\{\sqrt{\log(md/\delta)}, 1\}$. Let $R \leq O(\lambda\delta/(\kappa^2 n^2 dD))$, then for any $t \geq 0$, with probability at least $1 - \delta$, we have: Part 1. $\|H(t) - H^*\|_F \leq O(\kappa^2 n^2 dR D/\delta)$. Part 2. $\lambda_{\min}(H(t)) \geq \lambda/2$.*

Proof of Lemma 4.1. The proof of Part 1 of this Lemma follows from the pattern $\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0}$ for $i \in [n]$ and $r \in [m]$ is rarely changed during the training, this habit is similar to the regular ReLU pattern $\mathbf{1}_{\langle w_r(t), x_i \rangle \geq 0}$ [41]. The proof of Part 2 of this Lemma can be obtained by plugging $R \leq O(\lambda\delta/(\kappa^2 n^2 dD))$. Please refer to Lemma G.5 for the detailed proof. \square

4.2 Training Convergence

Having confirmed the convergence of the kernel function of the 1-bit linear network during training in Lemma 4.1, we can transform the dynamics of the loss function $L(t)$ into the following **kernel behavior**:

$$\begin{aligned} L(t+1) - L(t) &= -(\mathbf{F}(t) - y)^\top H(t)(\mathbf{F}(t) - y) + C_2 + C_3 + C_4 \\ &\approx -(\mathbf{F}(t) - y)^\top H(t)(\mathbf{F}(t) - y), \end{aligned}$$

In this equation, $\mathbf{F}(t) = [f(x_1, W(t), a), \dots, f(x_n, W(t), a)]^\top \in \mathbb{R}^n$ and $y = [y_1, \dots, y_n]^\top \in \mathbb{R}^n$, while C_2, C_3, C_4 are negligible terms (please refer to Appendix I for a rigorous proof).

Further, by $\lambda_{\min}(H(t)) > 0$ (as per Part 2 of Lemma 4.1), for each optimization step $t \geq 0$, we find that $L(t+1) \leq (1 - \eta\lambda/2)L(t)$, thus ensuring a non-increase in loss. Given sufficient training iterations and an appropriately chosen learning rate, we can achieve training convergence, the confirmation of which is provided in the following section.

Theorem 4.2 (Training convergence guarantee, informal version of Theorem I.1). *Given an expected error $\epsilon > 0$. Assume $\lambda_{\min}(H^*) > 0$. $\delta \in (0, 0.1)$, define $D := \sqrt{\log(md/\delta)}$. Choose $m \geq \Omega(\lambda^{-8} n^{12} d^8 / (\delta\epsilon)^4)$, $\eta \leq O(\lambda\delta/(\kappa^2 n^2 dD))$. Then let $T \geq \Omega((\eta\lambda)^{-1} \log(ndD^2/\epsilon))$, with probability at least $1 - \delta$, we have: $L(T) \leq \epsilon$.*

Proof sketch of Theorem 4.2. We first combine $L(0) = O(\sqrt{nd}D^2)$ (Lemma I.3) and $L(t+1) \leq (1 - \eta\lambda/2)L(t)$ (Lemma I.2), then we choose a sufficient large $T \geq \Omega((\eta\lambda)^{-1} \log(ndD^2/\epsilon))$ to achieve $L(T) \leq \epsilon$. For the complete proof, please see Theorem I.1. \square

Scaling Law for 1-Bit Neural Networks. Theorem 4.2 primarily illustrates a fact for any dataset with n data points. After initializing the hidden-layer weights $W \in \mathbb{R}^{d \times m}$ from a normal distribution, and assuming the minimum eigenvalue of NTK $\lambda > 0$, we set m to be a large enough value to ensure the network is sufficiently over-parameterized. With an appropriate learning rate, the loss can be minimized in finite training time to an arbitrarily small error ϵ . This offers a crucial insight that confirms the existence of a *scaling law for 1-bit neural networks*, which is strictly bounded by the model width m and training steps T . Consequently, we present the following Proposition that elucidates the principle of training 1-bit linear networks from scratch. This proposition is built upon Theorem 4.2 and the principle of training loss that scales as a power-law with model size, dataset size, and the amount of compute used for training [3, 65].

Proposition 4.3 (Scaling Law for 1-Bit Neural Networks). $\delta \in (0, 0.1)$. Define $N := O(md)$ as the number of parameters, $D := O(n)$ as the size of training dataset, $C := O(NDT)$ as the total compute cost. Especially, we denote the scale coefficients as $\alpha := Dd \log(md/\delta)$, and we then choose $\eta \leq O(\lambda\delta/(m\kappa^2 n^2 dD))$ and $T \geq \Omega((\eta\lambda m)^{-1} \log(nd \log(md/\delta)/\epsilon))$. Thus, the training loss, denoted as L_{scale} , satisfies:

$$L_{\text{scale}} \approx \max\left\{\frac{D^3 \cdot d^{2.25}}{\lambda^2 N^{0.25}}, \frac{\alpha}{\exp(\eta\lambda C)}\right\}.$$

Proof of Proposition 4.3. This proof follows from the definitions of N , D , C and α . Then, by choosing $\eta \leq O(\lambda\delta/(mn^2 dD))$ and $T \geq \Omega((\eta\lambda m)^{-1} \log(nd \log(md/\delta)/\epsilon))$, we utilize Theorem 4.2 to obtain our proposition. \square

Proposition 4.3 demonstrates that the training loss of the prefix learning converges exponentially as we increase the computational cost C , which primarily depends on the number of parameters and the training time in prefix learning. This further suggests a potential relationship for formulating a scaling law for 1-bit neural networks.

4.3 Extensibility

We now bridge our theoretical framework to a real-world application involving a multi-layer 1-bit transformer trained on large-scale datasets. Let the full dataset be denoted as $\mathcal{D}_{\text{mat}} = \{(X_i, Y_i)\}_{i=1}^n \subset \mathbb{R}^{K \times d}$, where $X_i \in \mathbb{R}^{K \times d}$ represents a sequence of K tokens with d -dimensional embeddings, and Y_i denotes the corresponding target sequence. Here, K is the input context length.

The standard transformer architecture [10] interleaves multi-head self-attention and position-wise feed-forward layers. For an input sequence $X \in \mathbb{R}^{K \times d}$ (compactly representing K token embeddings), an N -layer transformer is defined recursively as:

$$\mathcal{F}(X) := \text{TF}_{(N)} \left(\text{TF}_{(N-1)} \left(\cdots \text{TF}_{(1)}(X + E) \cdots \right) \right),$$

where $E \in \mathbb{R}^{K \times d}$ is the positional embedding matrix, and $\text{TF}_{(\nu)} : \mathbb{R}^{K \times d} \rightarrow \mathbb{R}^{K \times d}$ for $\nu \in [N]$ denotes the ν -th transformer block. For brevity, we omit layer indices when describing a single transformer block TF , which consists of:

$$\begin{aligned} \text{Attn}(X) &:= X + \sum_{\xi=1}^h \text{dq} \left(\text{softmax} \left(\frac{\text{dq}(X \widetilde{W}_{\xi,Q} \widetilde{W}_{\xi,K}^\top X^\top)}{\sqrt{d}} \right) X \widetilde{W}_{\xi,V} \widetilde{W}_{\xi,O}^\top \right), \\ \text{FF}(X) &:= X + \text{dq} \left(\text{ReLU} \left(\text{dq}(X \widetilde{W}_1) + \mathbf{1}_K b_1^\top \right) \widetilde{W}_2^\top \right) + \mathbf{1}_K b_2^\top, \\ \text{TF}(X) &:= \text{FF}(\text{Attn}(X)), \end{aligned}$$

where: - h is the number of attention heads. - \widetilde{W} denotes 1-bit quantized weights, with $\text{dq}(\cdot)$ as the dequantization operator. - For each head $\xi \in [h]$, $\widetilde{W}_{\xi,Q}, \widetilde{W}_{\xi,K}, \widetilde{W}_{\xi,V} \in \mathbb{R}^{d \times d'}$ and $\widetilde{W}_{\xi,O} \in \mathbb{R}^{d' \times d}$

are query, key, value, and output projection matrices, respectively. - In the feed-forward network, $\widetilde{W}_1 \in \mathbb{R}^{d \times m}$ and $\widetilde{W}_2 \in \mathbb{R}^{m \times d}$ are projection matrices, with m as the hidden dimension, while $b_1 \in \mathbb{R}^m$ and $b_2 \in \mathbb{R}^d$ are bias terms.

The full parameters of the model is denoted as $\theta_{(d',h,m)} := \{W_{\nu,\xi,Q}, W_{\nu,\xi,K}, W_{\nu,\xi,V}, W_{\nu,\xi,O}, W_{\nu,2}, W_{\nu,2}, b_{\nu,1}, b_{\nu,2}\}_{(\nu,\xi) \in [L] \times [h]} + \{E\}$. Given a loss metric $\ell(\widehat{Y}, Y) := \frac{1}{2} \|\widehat{Y} - Y\|_F^2$, we define the training objective as follows:

$$\mathcal{L}(\theta_{(d',h,m)}) := \sum_{i=1}^n \ell(\mathcal{F}(X_i), Y_i). \quad (3)$$

Thus, we establish a general version of our theory:

Proposition 4.4. *Given an expected error $\epsilon > 0$ and denote the failure probability $\delta \in (0, 0.1)$. Given a dataset $\mathcal{D}_{\text{mat}} = \{(X_i, Y_i)\}_{i=1}^n \subset \mathbb{R}^{K \times d}$ and a model function $\mathcal{F} : \mathbb{R}^{K \times d} \rightarrow \mathbb{R}^{K \times d}$ with parameters set $\theta_{(d',h,m)}$. Assuming each NTK of \mathcal{F} is PD, denoted $H_{k,j}^*$ for $(k, j) \in [K] \times [d]$, $\lambda_{\min}(H_{k,j}^*) > 0$. Define $\lambda := \min_{(k,j) \in [K] \times [d]} \{\lambda_{\min}(H_{k,j}^*)\}$, we choose $m \geq \Omega(\lambda^{-8} n^{12} K^{12} d^{20} / (\delta \epsilon)^4)$. Then with a probability at least $1 - \delta$, there exists at least one first-order algorithm that minimizes Eq. (3) to ϵ .*

Proof. We consider a special case that only optimizes one feed-forward layer of the model, then solving $\mathcal{L}(\theta_{(d',h,m)})$ is just letting $n = Kdn'$ where n' represents the data size in Theorem 4.2. \square

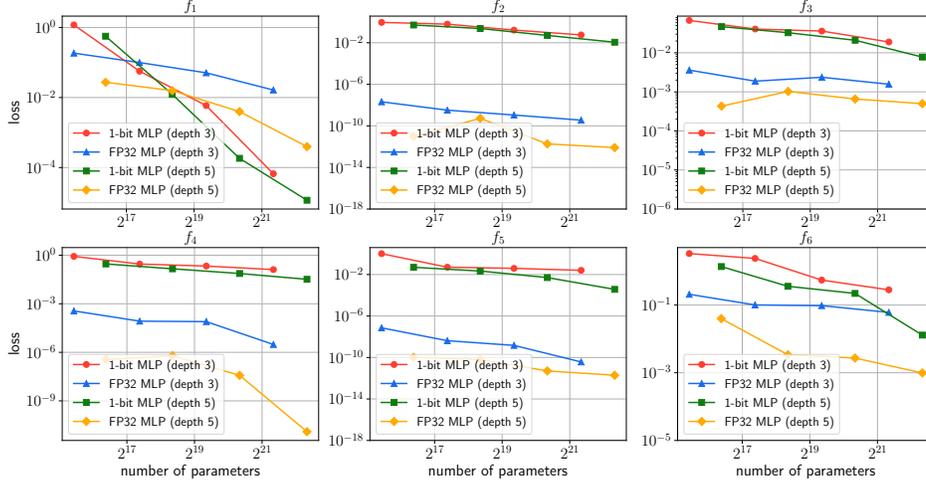


Figure 1: Verification experiment for *scaling law* for 1-bit neural networks. Minimum training loss of scaling number of parameters for MLP model to learn complicated functions f_1, f_2, f_3, f_4, f_5 and f_6 , and these function is defined in Section 6.1.

5 Generalization Similarity

In this section, we present our theoretical analysis that proves that training large-scale 1-bit neural networks is equivalent to training standard large-scale neural networks. In Section 5.1, we explain how the difference between the outputs of our 1-bit model and outputs of the standard NTK-style linear network for the same input at initialization, which is defined as function difference at initialization, will be kept in a small error while the model width (denoted as m) increase. Next, in Section 5.2, we confirm that in the trend of scaling up the model width, during the training, the predictions of 1-bit model and full precision model are also similar to a very slight error on both the training dataset and the test dataset.

5.1 Function Difference at Initialization

To begin with, at initialization, the boundary on $|f(x, W(0), a) - f'(x, W'(0), a)|$ is stated as follows:

Lemma 5.1 (Function difference at initialization, informal version of Lemma K.4). $\delta \in (0, 0.1)$. Denote $D := \sqrt{\log(md/\delta)}$. $\forall x \in \mathbb{R}^d$ that satisfies $\|x\|_2 = 1$, for any initial quantization error $\epsilon_{\text{init}} > 0$, we choose $\kappa \leq O(\epsilon_{\text{init}}/(\sqrt{d}D^2))$. Then with a probability $1 - \delta$, we have: $|f(x, W(0), a) - f'(x, W'(0), a)| \leq \epsilon_{\text{init}}$.

Proof sketch of Lemma 5.1. Due to the initialization of $W(0)$ and $W'(0)$, we then have the tail bound of the Gaussian distribution. Hence, the difference could be bounded by Hoeffding bound, we then get the result. Please refer to Lemma K.4 for the formal proof of this Lemma. \square

5.2 Generalization Similarity

We now address whether using 1-bit precision compromises the generalization ability of standard neural networks. Specifically, we use the test dataset to evaluate the **generalization similarity** - a measure of the similarity between two functions on out-of-distribution (OOD) data. This measure is designed to assess the equivalence between two functions. If, during each step of training two networks, these two training processes are deemed equivalent, then we assert that the generalization similarity is valid.

Addressing the above concern, we demonstrate that the predictions of two functions on both training and test datasets can be bounded to an arbitrarily small quantization error, provided that m is sufficiently large. Theoretically, as m scales towards infinity, the quantization error converges to 0. This finding confirms the validity of our generalization similarity measure and asserts that 1-bit precision does not compromise the generalization ability of standard neural networks.

Theorem 5.2 (Training and generalization similarity, informal version of Theorem K.1). *Let all pre-conditions in Theorem 4.2 satisfy. For any quantization error $\epsilon_{\text{quant}} > 0$, we choose $\kappa \leq O(\epsilon_{\text{quant}}/(dD^2))$. Integer $\forall t \geq 0$. For any training input $x_i \in \mathbb{R}^d$ in \mathcal{D} and any test input $x_{\text{test},i} \in \mathbb{R}^d$ in $\mathcal{D}_{\text{test}}$, with a probability at least $1 - \delta$, we have:*

- Part 1. $|f(x_i, W(t), a) - f(x_i, W(t), a)| \leq \epsilon_{\text{quant}}$.
- Part 2. $|f(x_{\text{test},i}, W(t), a) - f(x_{\text{test},i}, W(t), a)| \leq \epsilon_{\text{quant}}$.

Proof. Proof sketch of Theorem 5.2 Since we proved $|f(x, W(0), a) - f'(x, W'(0), a)| \leq \epsilon_{\text{init}}$ in Lemma 5.1, then as we choose appropriate R and learning rate η , the equations in Part 1 and Part 2 of this Theorem would be bounded by scaling m to be sufficiently large. We state the complete proof in Theorem K.1. \square

Training Equivalence. Here, we say training f and f' are equivalent since we achieve the predictions that these two functions are extremely similar by plugging an appropriate value of κ . Besides, as we proved in Theorem 4.2, this implementation would not harm the optimization of 1-bit networks. This further explains why 1-bit precision even processes better when the scales of networks are increasing, instead of turning to a training collapse. Therefore, we believe it is the theory unlocking the potential of 1-bit neural networks from the perspective of kernel-based analysis.

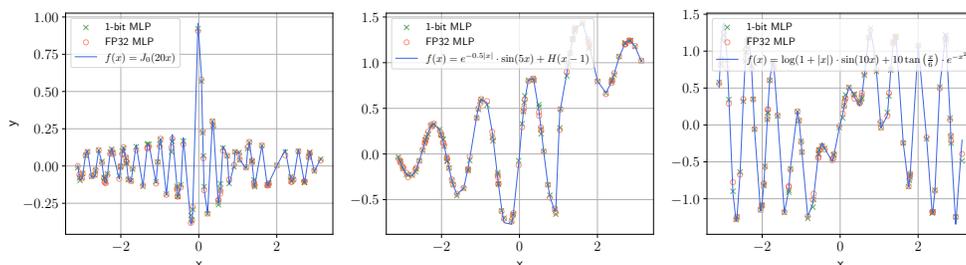


Figure 2: This plot shows the difference between the predicted and actual values of the functions on the test dataset. We tested three complex functions, as seen in the images, and the performance of the 1-bit model is nearly identical to that of the standard 32-bit floating-point model.

6 Experiments

In this section, we aim to verify our theory by evaluating how well our quantization works for learning rigorous functions and comparing it to the standard model. We designed our experiment to 1) validate the scaling law (Section 6.1), 2) visually demonstrate that the performance difference is minimal compared to the standard model, which uses full-bit precision, through visualizations of single-variable input functions (Section 6.2), and 3) show how the test and train losses decrease as the model’s parameter size increases and as the epochs progress (Section 6.3).

6.1 Verification on Scaling Law

Experiment Setup. In this experiment, we aimed to learn rigorous functions using a Multi-Layer Perceptron (MLP) with varying depths of 3 and 5 layers. The MLP models had different sizes for the hidden layers, and we measured the minimum loss achieved throughout the training process. Each model was trained for 100,000 steps. We experimented with various parameter sizes and plotted the corresponding loss functions. Additionally, we compared our method with the standard training approach using 32-bit floating-point precision.

We experimented with a variety of target functions, and for each function, the inputs x_i were randomly chosen within the range $[-1, 1]$. Specifically, each x_i was sampled from a uniform distribution over this interval to ensure that the network could handle input values across the entire domain of interest. We sampled 100 data points and trained our model over this set.

The functions we aimed to learn during the experiment are listed below:

1. $f_1(x_1, x_2, x_3, x_4, x_5) = \exp\left(\frac{1}{5} \sum_{i=1}^5 \sin^2\left(\frac{\pi x_i}{2}\right)\right)$, This function takes five inputs and applies a sinusoidal transformation followed by an exponential operation.
2. $f_2(x_1, x_2, x_3, x_4) = \ln(1 + |x_1|) + (x_2^2 - x_2) + \sin(x_3) - e^{x_4}$, the function combines logarithmic, polynomial, trigonometric, and exponential components over four input variables.
3. $f_3(x_1, x_2, x_3) = x_1 \times x_2 - x_3$, This is a simple linear function over three inputs, involving multiplication and subtraction.
4. $f_4(x_1, x_2, x_3, x_4) = x_0 \cdot \sin(x_1) + \cos(x_2) - 0.5 \cdot x_3$, A four-input function mixing trigonometric and linear terms, with coefficients applied to the terms.
5. $f_5(x_1, x_2, x_3, x_4) = \frac{x_0^2}{1 + |x_1|} - e^{x_2} + \tanh(x_3) + \sqrt{|x_0 \cdot x_2|}$, This function incorporates nonlinear operations like exponentials, hyperbolic tangents, and square roots.
6. $f_6(x_1, x_2, x_3, x_4) = \text{LambertW}(x_0 \cdot x_1) + \frac{x_2}{\log(1 + e^{x_3})} - \frac{\Gamma(x_1)}{1 + |x_0|}$, The most complex function we tested, which includes special functions like the Lambert W function and the Gamma function, alongside logarithmic and exponential components.

We compare our quantized model (using INT1, 32× smaller) to a standard non-quantized model (using 32-bit precision). For all functions (f_1 to f_6), we observe (in) that as the number of parameters increases, the loss decreases, supporting our theoretical claim that larger models lead to convergence.

Although the standard method generally performs better due to its 32-bit precision, the gap decreases as the number of parameters grows. This shows that while our method has a slightly higher loss, it remains competitive, offering significant memory and computational efficiency.

6.2 Comparison on 1-D Functions

In this experiment, we aimed to visually demonstrate the performance on highly complex functions with sharp spikes between $[-\pi, \pi]$. We sampled 100 uniformly spaced points and trained a 2-layer MLP with 20M parameters to learn the function. Additionally, we sampled 100 random points uniformly from this interval as the test dataset.

The first observation from the plot is that both the standard and 1-bit methods learn all the functions almost perfectly, with minimal difference between them. Secondly, both methods perform similarly

on these functions, which can be easily observed by comparing the scatter plots of the 1-bit and standard models. The 1-bit model requires $32\times$ less energy and computation.

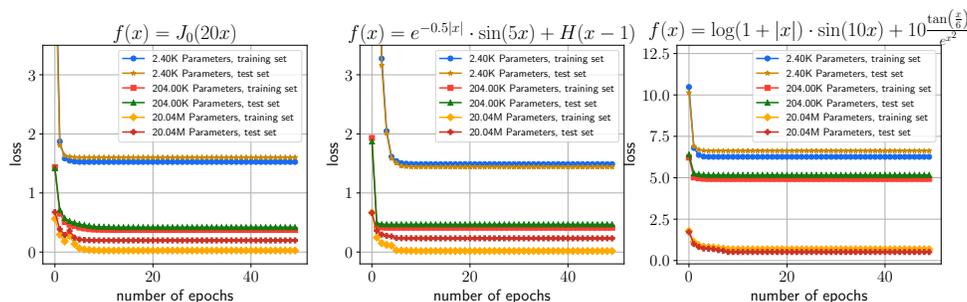


Figure 3: This plot shows the ℓ_2 difference between both the training and test points and the predicted points throughout the training phase for different model sizes and parameter counts. Each plot demonstrates how the error decreases as training progresses, highlighting the impact of model size on both training and test performance.

6.3 Evaluation on Training and Generalization Similarity

For the same set of functions, we show how the loss functions for both the train and test datasets decrease as the number of epochs increases. As the training progresses, the loss converges towards zero for models with a higher number of parameters. We experimented with models containing 2.4k, 204k, and 20M parameters, each consisting of only 2 layers.

Across all three functions, the loss decreases rapidly in the early epochs and stabilizes for both the training and test sets. Larger models with 20M parameters consistently achieve lower final losses compared to smaller models with 2.4k and 204k parameters, demonstrating the benefit of increased model size. The gap between training and test losses remains minimal, indicating strong generalization across different parameter sizes. More importantly, the key observation is that the models predict similarly on both the training and test datasets, a behavior we refer to as *generalization similarity*. This means that the models, regardless of size, behave similarly across both datasets, supporting the scaling law that increasing model size leads to better convergence and generalization, but also highlighting the consistent similarity in performance between training and testing across different functions.

7 Conclusion

In conclusion, our theoretical results confirm the scaling law for 1-bit neural networks. We demonstrated that the model achieves a small loss as the number of parameters increases. Despite the constraint of binary weights, 1-bit models show similar behavior to full-precision models as their width grows. Our experiments support this theory, showing that 1-bit networks perform nearly as well as standard models on complex functions. As the number of parameters grows, the performance gap between 1-bit and full-precision models reduces. These findings highlight that 1-bit networks are both efficient and effective, providing a strong alternative to traditional models.

Acknowledgement

The authors sincerely thank Bo Chen, Xiaoyu Li, Zhizhou Sha, Jing Xiong, Junwei Yu and Yufa Zhou for their helpful suggestions and discussion. The authors also would like to thank all anonymous reviewers for their constructive reviews that enhanced the contribution of this work.

References

- [1] Hongyu Wang, Shuming Ma, Li Dong, Shaohan Huang, Huaijie Wang, Lingxiao Ma, Fan Yang, Ruiping Wang, Yi Wu, and Furu Wei. Bitnet: Scaling 1-bit transformers for large language models. *arXiv preprint arXiv:2310.11453*, 2023.
- [2] Shuming Ma, Hongyu Wang, Lingxiao Ma, Lei Wang, Wenhui Wang, Shaohan Huang, Li Dong, Ruiping Wang, Jilong Xue, and Furu Wei. The era of 1-bit llms: All large language models are in 1.58 bits. *arXiv preprint arXiv:2402.17764*, 2024.
- [3] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [5] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.
- [6] Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. *arXiv preprint arXiv:2306.13549*, 2023.
- [7] Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng Wan, and S Yu Philip. Multimodal large language models: A survey. In *2023 IEEE International Conference on Big Data (BigData)*, pages 2247–2256. IEEE, 2023.
- [8] Jianghao Lin, Xinyi Dai, Yunjia Xi, Weiwen Liu, Bo Chen, Hao Zhang, Yong Liu, Chuhan Wu, Xiangyang Li, Chenxu Zhu, et al. How can recommender systems benefit from large language models: A survey. *arXiv preprint arXiv:2306.05817*, 2023.
- [9] Zhenjie Yang, Xiaosong Jia, Hongyang Li, and Junchi Yan. A survey of large language models for autonomous driving. *arXiv preprint arXiv:2311.01043*, 2023.
- [10] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- [11] Josh Alman and Zhao Song. Fast attention requires bounded entries. *Advances in Neural Information Processing Systems*, 36, 2023.
- [12] Zixuan Zhou, Xuefei Ning, Ke Hong, Tianyu Fu, Jiaming Xu, Shiyao Li, Yuming Lou, Luning Wang, Zhihang Yuan, Xiuhong Li, et al. A survey on efficient inference for large language models. *arXiv preprint arXiv:2404.14294*, 2024.
- [13] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [14] Yanyu Li, Geng Yuan, Yang Wen, Ju Hu, Georgios Evangelidis, Sergey Tulyakov, Yanzhi Wang, and Jian Ren. Efficientformer: Vision transformers at mobilenet speed. *Advances in Neural Information Processing Systems*, 35:12934–12949, 2022.
- [15] Jierun Chen, Shiu-hong Kao, Hao He, Weipeng Zhuo, Song Wen, Chul-Ho Lee, and S-H Gary Chan. Run, don't walk: chasing higher flops for faster neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12021–12031, 2023.

- [16] Markus Nagel, Marios Fournarakis, Rana Ali Amjad, Yelysei Bondarenko, Mart Van Baalen, and Tijmen Blankevoort. A white paper on neural network quantization. *arXiv preprint arXiv:2106.08295*, 2021.
- [17] Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022.
- [18] Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A survey of quantization methods for efficient neural network inference. In *Low-Power Computer Vision*, pages 291–326. Chapman and Hall/CRC, 2022.
- [19] Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. Awq: Activation-aware weight quantization for on-device llm compression and acceleration. *Proceedings of Machine Learning and Systems*, 6: 87–100, 2024.
- [20] Arash Ahmadian, Saurabh Dash, Hongyu Chen, Bharat Venkitesh, Zhen Stephen Gou, Phil Blunsom, Ahmet Üstün, and Sara Hooker. Intriguing properties of quantization at scale. *Advances in Neural Information Processing Systems*, 36:34278–34294, 2023.
- [21] Yelysei Bondarenko, Markus Nagel, and Tijmen Blankevoort. Understanding and overcoming the challenges of efficient transformer quantization. *arXiv preprint arXiv:2109.12948*, 2021.
- [22] Shigang Li, Kazuki Osawa, and Torsten Hoefler. Efficient quantized sparse matrix operations on tensor cores. In *SC22: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–15. IEEE, 2022.
- [23] Cong Guo, Jiaming Tang, Weiming Hu, Jingwen Leng, Chen Zhang, Fan Yang, Yunxin Liu, Minyi Guo, and Yuhao Zhu. Olive: Accelerating large language models via hardware-friendly outlier-victim pair quantization. In *Proceedings of the 50th Annual International Symposium on Computer Architecture*, pages 1–15, 2023.
- [24] Zhenhua Liu, Yunhe Wang, Kai Han, Wei Zhang, Siwei Ma, and Wen Gao. Post-training quantization for vision transformer. *Advances in Neural Information Processing Systems*, 34: 28092–28103, 2021.
- [25] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*, pages 38087–38099. PMLR, 2023.
- [26] Albert Tseng, Jerry Chee, Qingyao Sun, Volodymyr Kuleshov, and Christopher De Sa. Quip#: Even better llm quantization with hadamard incoherence and lattice codebooks, 2024. URL <https://arxiv.org/abs/2402.04396>.
- [27] Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, Raghuraman Krishnamoorthi, and Vikas Chandra. Llm-qat: Data-free quantization aware training for large language models, 2023. URL <https://arxiv.org/abs/2305.17888>.
- [28] Markus Nagel, Marios Fournarakis, Rana Ali Amjad, Yelysei Bondarenko, Mart van Baalen, and Tijmen Blankevoort. A white paper on neural network quantization, 2021. URL <https://arxiv.org/abs/2106.08295>.
- [29] Zechun Liu, Barlas Oguz, Aasish Pappu, Lin Xiao, Scott Yih, Meng Li, Raghuraman Krishnamoorthi, and Yashar Mehdad. Bit: Robustly binarized multi-distilled transformer. *Advances in neural information processing systems*, 35:14303–14316, 2022.
- [30] Rui-Jie Zhu, Yu Zhang, Ethan Sifferman, Tyler Sheaves, Yiqiao Wang, Dustin Richmond, Peng Zhou, and Jason K Eshraghian. Scalable matmul-free language modeling. *arXiv preprint arXiv:2406.02528*, 2024.

- [31] Tara N Sainath, Brian Kingsbury, Vikas Sindhwani, Ebru Arisoy, and Bhuvana Ramabhadran. Low-rank matrix factorization for deep neural network training with high-dimensional output targets. In *2013 IEEE international conference on acoustics, speech and signal processing*, pages 6655–6659. IEEE, 2013.
- [32] Yen-Chang Hsu, Ting Hua, Sungen Chang, Qian Lou, Yilin Shen, and Hongxia Jin. Language model compression with weighted low-rank factorization, 2022. URL <https://arxiv.org/abs/2207.00112>.
- [33] Coleman Hooper, Sehoon Kim, Hiva Mohammadzadeh, Michael W. Mahoney, Yakun Sophia Shao, Kurt Keutzer, and Amir Gholami. Kvquant: Towards 10 million context length llm inference with kv cache quantization, 2024. URL <https://arxiv.org/abs/2401.18079>.
- [34] Tianyi Zhang, Jonah Yi, Zhaozhuo Xu, and Anshumali Shrivastava. Kv cache is 1 bit per channel: Efficient large language model inference with coupled quantization, 2024. URL <https://arxiv.org/abs/2405.03917>.
- [35] Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong, Zhaozhuo Xu, Vladimir Braverman, Beidi Chen, and Xia Hu. Kivi: A tuning-free asymmetric 2bit quantization for kv cache. *arXiv preprint arXiv:2402.02750*, 2024.
- [36] Amir Zandieh, Majid Daliri, and Insu Han. Qjl: 1-bit quantized jl transform for kv cache quantization with zero overhead, 2024. URL <https://arxiv.org/abs/2406.03482>.
- [37] Nilesh Prasad Pandey, Markus Nagel, Mart van Baalen, Yin Huang, Chirag Patel, and Tijmen Blankevoort. A practical mixed precision algorithm for post-training quantization, 2023. URL <https://arxiv.org/abs/2302.05397>.
- [38] Chen Tang, Kai Ouyang, Zhi Wang, Yifei Zhu, Yaowei Wang, Wen Ji, and Wenwu Zhu. Mixed-precision neural network quantization via learned layer-wise importance, 2023. URL <https://arxiv.org/abs/2203.08368>.
- [39] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- [40] Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. *Advances in neural information processing systems*, 31, 2018.
- [41] Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018.
- [42] Zhao Song and Xin Yang. Quadratic suffices for over-parametrization via matrix chernoff bound. *arXiv preprint arXiv:1906.03593*, 2019.
- [43] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In *International conference on machine learning*, pages 242–252. PMLR, 2019.
- [44] Colin Wei, Jason D Lee, Qiang Liu, and Tengyu Ma. Regularization matters: Generalization and optimization of neural nets vs their induced kernel. *Advances in Neural Information Processing Systems*, 32, 2019.
- [45] Alberto Bietti and Julien Mairal. On the inductive bias of neural tangent kernels. *Advances in Neural Information Processing Systems*, 32, 2019.
- [46] Jaehoon Lee, Samuel Schoenholz, Jeffrey Pennington, Ben Adlam, Lechao Xiao, Roman Novak, and Jascha Sohl-Dickstein. Finite versus infinite neural networks: an empirical study. *Advances in Neural Information Processing Systems*, 33:15156–15172, 2020.

- [47] Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss. In *Conference on learning theory*, pages 1305–1338. PMLR, 2020.
- [48] Zhenmei Shi, Junyi Wei, and Yingyu Liang. A theoretical analysis on feature learning in neural networks: Emergence from inputs and advantage over fixed features. In *International Conference on Learning Representations*, 2021.
- [49] Mo Zhou, Rong Ge, and Chi Jin. A local convergence theory for mildly over-parameterized two-layer neural network. In *Conference on Learning Theory*, pages 4577–4632. PMLR, 2021.
- [50] Mariia Seleznova and Gitta Kutyniok. Neural tangent kernel beyond the infinite-width limit: Effects of depth and initialization. In *International Conference on Machine Learning*, pages 19522–19560. PMLR, 2022.
- [51] Yeqi Gao, Sridhar Mahadevan, and Zhao Song. An over-parameterized exponential regression. *arXiv preprint arXiv:2303.16504*, 2023.
- [52] Chenyang Li, Yingyu Liang, Zhenmei Shi, and Zhao Song. Exploring the frontiers of softmax: Provable optimization, applications in diffusion model, and beyond. *arXiv preprint arXiv:2405.03251*, 2024.
- [53] Zhenmei Shi, Junyi Wei, and Yingyu Liang. Provable guarantees for neural networks via gradient feature learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- [54] Yingyu Liang, Zhenmei Shi, Zhao Song, and Chiwun Yang. Toward infinite-long prefix in transformer. *arXiv preprint arXiv:2406.14036*, 2024.
- [55] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. *Advances in neural information processing systems*, 32, 2019.
- [56] Sina Alemohammad, Zichao Wang, Randall Balestriero, and Richard Baraniuk. The recurrent neural tangent kernel. *arXiv preprint arXiv:2006.10246*, 2020.
- [57] Jiri Hron, Yasaman Bahri, Jascha Sohl-Dickstein, and Roman Novak. Infinite attention: Nngp and ntk for deep attention networks. In *International Conference on Machine Learning*, pages 4376–4386. PMLR, 2020.
- [58] Sadhika Malladi, Alexander Wettig, Dingli Yu, Danqi Chen, and Sanjeev Arora. A kernel-based view of language model fine-tuning. In *International Conference on Machine Learning*, pages 23610–23641. PMLR, 2023.
- [59] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- [60] Penghang Yin, Jiancheng Lyu, Shuai Zhang, Stanley Osher, Yingyong Qi, and Jack Xin. Understanding straight-through estimator in training activation quantized neural nets. *arXiv preprint arXiv:1903.05662*, 2019.
- [61] Stefani Karp, Ezra Winston, Yuanzhi Li, and Aarti Singh. Local signal adaptivity: Provable feature learning in neural networks beyond kernels. *Advances in Neural Information Processing Systems*, 34:24883–24897, 2021.
- [62] Zeyuan Allen-Zhu and Yuanzhi Li. Towards understanding ensemble, knowledge distillation and self-distillation in deep learning. *arXiv preprint arXiv:2012.09816*, 2020.
- [63] Zeyuan Allen-Zhu and Yuanzhi Li. Feature purification: How adversarial training performs robust deep learning. In *2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)*, pages 977–988. IEEE, 2022.

- [64] Tianren Zhang, Chujie Zhao, Guanyu Chen, Yizhou Jiang, and Feng Chen. Feature contamination: Neural networks learn uncorrelated features and fail to generalize. *arXiv preprint arXiv:2406.03345*, 2024.
- [65] Tanishq Kumar, Zachary Ankner, Benjamin F Spector, Blake Bordelon, Niklas Muennighoff, Mansheej Paul, Cengiz Pehlevan, Christopher Ré, and Aditi Raghunathan. Scaling laws for precision. *arXiv preprint arXiv:2411.04330*, 2024.
- [66] Yichuan Deng, Zhihang Li, and Zhao Song. Attention scheme inspired softmax regression. *arXiv preprint arXiv:2304.10411*, 2023.
- [67] Yekun Ke, Xiaoyu Li, Yingyu Liang, Zhenmei Shi, and Zhao Song. Advancing the understanding of fixed point iterations in deep neural networks: A detailed analytical study. *arXiv preprint arXiv:2410.11279*, 2024.
- [68] Yaodong Yu, Sam Buchanan, Druv Pai, Tianzhe Chu, Ziyang Wu, Shengbang Tong, Benjamin Haefele, and Yi Ma. White-box transformers via sparse rate reduction. *Advances in Neural Information Processing Systems*, 36:9422–9457, 2023.
- [69] Yaodong Yu, Tianzhe Chu, Shengbang Tong, Ziyang Wu, Druv Pai, Sam Buchanan, and Yi Ma. Emergence of segmentation with minimalistic white-box transformers. *arXiv preprint arXiv:2308.16271*, 2023.
- [70] Javier Ferrando, Gabriele Sarti, Arianna Bisazza, and Marta R Costa-jussà. A primer on the inner workings of transformer-based language models. *arXiv preprint arXiv:2405.00208*, 2024.
- [71] Druv Pai, Sam Buchanan, Ziyang Wu, Yaodong Yu, and Yi Ma. Masked completion via structured diffusion with white-box transformers. In *The Twelfth International Conference on Learning Representations*, 2024.
- [72] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.
- [73] Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 1, context-free grammar. *arXiv preprint arXiv:2305.13673*, 2023.
- [74] Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.2, knowledge manipulation. *arXiv preprint arXiv:2309.14402*, 2023.
- [75] Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.1, knowledge storage and extraction. *arXiv preprint arXiv:2309.14316*, 2023.
- [76] Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 3.3, knowledge capacity scaling laws. *arXiv preprint arXiv:2404.05405*, 2024.
- [77] Dennis Wu, Jerry Yao-Chieh Hu, Weijian Li, Bo-Yu Chen, and Han Liu. Stanhop: Sparse tandem hopfield model for memory-enhanced time series prediction. *arXiv preprint arXiv:2312.17346*, 2023.
- [78] Jerry Yao-Chieh Hu, Donglin Yang, Dennis Wu, Chenwei Xu, Bo-Yu Chen, and Han Liu. On sparse modern hopfield model. *Advances in Neural Information Processing Systems*, 36, 2023.
- [79] Jerry Yao-Chieh Hu, Thomas Lin, Zhao Song, and Han Liu. On computational limits of modern hopfield models: A fine-grained complexity analysis. *arXiv preprint arXiv:2402.04520*, 2024.
- [80] Jerry Yao-Chieh Hu, Pei-Hsuan Chang, Robin Luo, Hong-Yu Chen, Weijian Li, Wei-Po Wang, and Han Liu. Outlier-efficient hopfield layers for large transformer-based models. *arXiv preprint arXiv:2404.03828*, 2024.

- [81] Chenwei Xu, Yu-Chao Huang, Jerry Yao-Chieh Hu, Weijian Li, Ammar Gilani, Hsi-Sheng Goan, and Han Liu. Bishop: Bi-directional cellular learning for tabular data with generalized sparse modern hopfield model. *arXiv preprint arXiv:2404.03830*, 2024.
- [82] Dennis Wu, Jerry Yao-Chieh Hu, Teng-Yun Hsiao, and Han Liu. Uniform memory retrieval with larger capacity for modern hopfield models. *arXiv preprint arXiv:2404.03827*, 2024.
- [83] Jerry Yao-Chieh Hu, Bo-Yu Chen, Dennis Wu, Feng Ruan, and Han Liu. Nonparametric modern hopfield models. *arXiv preprint arXiv:2404.03900*, 2024.
- [84] Josh Alman and Zhao Song. The fine-grained complexity of gradient computation for training large language models. *arXiv preprint arXiv:2402.04497*, 2024.
- [85] Josh Alman and Zhao Song. How to capture higher-order correlations? generalizing matrix softmax attention to kronecker computation. In *The Twelfth International Conference on Learning Representations*, 2024.
- [86] Insu Han, Rajesh Jayaram, Amin Karbasi, Vahab Mirrokni, David Woodruff, and Amir Zandieh. Hyperattention: Long-context attention in near-linear time. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=Eh00d2BJIM>.
- [87] Praneeth Kacham, Vahab Mirrokni, and Peilin Zhong. Polysketchformer: Fast transformers via sketches for polynomial kernels. *arXiv preprint arXiv:2310.01655*, 2023.
- [88] Raghav Addanki, Chenyang Li, Zhao Song, and Chiwun Yang. One pass streaming algorithm for super long token attention approximation in sublinear space. *arXiv preprint arXiv:2311.14652*, 2023.
- [89] Yichuan Deng, Zhao Song, and Chiwun Yang. Attention is naturally sparse with gaussian distributed input. *arXiv preprint arXiv:2404.02690*, 2024.
- [90] Zhenmei Shi, Yifei Ming, Xuan-Phi Nguyen, Yingyu Liang, and Shafiq Joty. Discovering the gems in early layers: Accelerating long-context llms with 1000x input token reduction. *arXiv preprint arXiv:2409.17422*, 2024.
- [91] Bo Chen, Xiaoyu Li, Yingyu Liang, Zhenmei Shi, and Zhao Song. Bypassing the exponential dependency: Looped transformers efficiently learn in-context by multi-step gradient descent. *arXiv preprint arXiv:2410.11268*, 2024.
- [92] Yingyu Liang, Jiangxuan Long, Zhenmei Shi, Zhao Song, and Yufa Zhou. Beyond linear approximations: A novel pruning approach for attention matrix. *arXiv preprint arXiv:2410.11261*, 2024.
- [93] Bo Chen, Yingyu Liang, Zhizhou Sha, Zhenmei Shi, and Zhao Song. Hsr-enhanced sparse attention acceleration. *arXiv preprint arXiv:2410.10165*, 2024.
- [94] Yingyu Liang, Zhizhou Sha, Zhenmei Shi, Zhao Song, and Yufa Zhou. Looped relu mlps may be all you need as practical programmable computers. *arXiv preprint arXiv:2410.09375*, 2024.
- [95] Herman Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *The Annals of Mathematical Statistics*, pages 493–507, 1952.
- [96] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *The collected works of Wassily Hoeffding*, pages 409–426, 1994.
- [97] Sergei Bernstein. On a modification of chebyshev’s inequality and of the error formula of laplace. *Ann. Sci. Inst. Sav. Ukraine, Sect. Math*, 1(4):38–49, 1924.
- [98] Aleksandr Khintchine. Über dyadische brüche. *Mathematische Zeitschrift*, 18(1):109–116, 1923.

- [99] Uffe Haagerup. The best constants in the khintchine inequality. *Studia Mathematica*, 70(3): 231–283, 1981.
- [100] David Lee Hanson and Farroll Tim Wright. A bound on tail probabilities for quadratic forms in independent random variables. *The Annals of Mathematical Statistics*, 42(3):1079–1083, 1971.
- [101] Mark Rudelson and Roman Vershynin. Hanson-wright inequality and sub-gaussian concentration. 2013.
- [102] Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of statistics*, pages 1302–1338, 2000.
- [103] Sergey Foss, Dmitry Korshunov, Stan Zachary, et al. *An introduction to heavy-tailed and subexponential distributions*, volume 6. Springer, 2011.
- [104] Joel A Tropp. Improved analysis of the subsampled randomized hadamard transform. *Advances in Adaptive Data Analysis*, 3(01n02):115–126, 2011.
- [105] Yichao Lu, Paramveer Dhillon, Dean P Foster, and Lyle Ungar. Faster ridge regression via the subsampled randomized hadamard transform. *Advances in neural information processing systems*, 26, 2013.

Appendix

Contents

1	Introduction	1
2	Related Work	2
3	Preliminary	3
3.1	Quantization	3
3.2	NTK Problem Setup	4
3.3	Recalling Classic NTK Setup	4
4	Kernel Behavior and Training Convergence	4
4.1	Neural Tangent Kernel	5
4.2	Training Convergence	5
4.3	Extensibility	6
5	Generalization Similarity	7
5.1	Function Difference at Initialization	7
5.2	Generalization Similarity	8
6	Experiments	9
6.1	Verification on Scaling Law	9
6.2	Comparison on 1-D Functions	9
6.3	Evaluation on Training and Generalization Similarity	10
7	Conclusion	10
A	More Related Work	21
B	Preliminary	21
B.1	Notations	21
B.2	Basic Facts	21
B.3	Probability Tools	22
B.4	Basic Bound	23
C	NTK Problem Setup	23
C.1	Dataset	23
C.2	Model	23
C.3	Training	25

D	Quantization	26
D.1	Quantization Functions	26
D.2	Dequantization Functions	27
D.3	Quantization Error	27
E	Patterns	28
E.1	ReLU Pattern	28
E.2	Sign Pattern	29
F	Straight-Through Estimator (STE)	29
F.1	STE Functions	29
F.2	Gradient Computation	29
G	Neural Tangent Kernel	30
G.1	Kernel Function	30
G.2	Assumption: H^* is Positive Definite	31
G.3	Kernel Convergence and PD Property	31
H	Training Dynamic	33
H.1	Decompose Loss	33
H.2	Bounding C_1	35
H.3	Bounding C_2	37
H.4	Bounding C_3	38
H.5	Bounding C_4	41
I	Inductions	43
I.1	Main Result 1: Training Convergence Guarantee	43
I.2	Induction for Loss	44
I.3	Induction for STE Gradient	47
I.4	Induction for Weights	48
J	Supplementary Setup for Classical Linear Regression	49
J.1	Model Function	49
J.2	Loss and Training	50
J.3	Induction for Weights	50
J.4	Induction for Loss	51
K	Similarities	52
K.1	Main Result 2: Training Similarity	52
K.2	Test Dataset for Generalization Evaluation	53

K.3 Function Similarity at Initialization 53

Roadmap

We initially introduce the intention of each section in the appendix here. In Appendix A, we review more prior works that relate to our work. In Appendix B, we provide the preliminary for our theoretical analysis. In Appendix C, we give the formal definition of the NTK-style problem setup we aim to solve in this paper. In Appendix D, we strictly define the quantization method we utilize for our approach. We discuss the potential pattern changing of ReLU and signal function in Appendix E. For optimizing 1-bit neural network, we state the Straight-Through Estimator method (STE) definitions in Appendix F. In Appendix G, we define NTK for our optimization problem and discuss its properties. In Appendix I, we prove the convergence guarantee of training 1-bit neural networks. In Appendix J, we review the classical setup of solving the NTK-style linear regression. We confirm the generalization similarity in Appendix K.

A More Related Work

Theoretical Approach for Understanding Modern Neural Networks. The intricate architecture of transformer-based models, coupled with the stochastic nature of their optimization processes, presents a formidable challenge in comprehending the behaviors of large language models (LLMs). However, delving into these complexities through a theoretical lens can illuminate pathways for enhancing and innovating future AI systems. This exploration encompasses various facets, including the **optimization strategies for LLMs** [52, 66, 67], the intricacies of **white-box transformers** [68–71], and the analysis of **emergent capabilities** that arise within these models [4, 72–76]. Additionally, the **modern Hopfield model** [77–83] offers a rich terrain for investigation, revealing the nuanced dynamics that govern these advanced neural networks.

Efficient Neural Networks. As the principles of scaling laws come to the forefront, contemporary neural networks are increasingly trained on expansive datasets, necessitating substantial computational resources [11, 84–94]. This demand for efficiency has spurred research into algorithms that optimize **computational complexity**, minimize **memory usage**, and enhance **alignment with GPU architectures**. Such advancements are crucial in navigating the challenges posed by the ever-growing scale of data and the intricate demands of modern AI applications, ensuring that these powerful tools remain accessible and effective in their deployment.

B Preliminary

B.1 Notations

In this paper, we use integer $m > 0$ to denote the width of neural networks, in particular, m is sufficiently large. We use integer $d > 0$ to denote the dimension of neural networks. We use integer $n > 0$ to denote the size of the training dataset.

B.2 Basic Facts

Fact B.1. For a variable $x \sim \mathcal{N}(0, \sigma^2)$, then with probability at least $1 - \delta$, we have:

$$|x| \leq C\sigma\sqrt{\log(1/\delta)}$$

Fact B.2. For an 1-Lipschitz function $f(\cdot)$, we have:

$$|f(x) - f(y)| \leq |x - y|, \forall x, y \in \mathbb{R}^d$$

Fact B.3. For a Gaussian variable $x \sim \mathcal{N}(0, \sigma^2 \cdot I_d)$ where $\sigma \in \mathbb{R}$, then for any $t > 0$, we have:

$$\Pr[x \leq t] \leq \frac{2t}{\sqrt{2\pi}\sigma}$$

Fact B.4. For a Gaussian vector $w \sim \mathcal{N}(0, \sigma^2 \cdot I_d)$ where $\sigma \in \mathbb{R}$, and a fixed vector $x \in \mathbb{R}^d$, we have:

$$w^\top x \sim \mathcal{N}(0, \sigma^2 \|x\|_2^2 \cdot I_d)$$

Fact B.5. For two matrices $H, \tilde{H} \in \mathbb{R}^{n \times n}$, we have:

$$\lambda_{\min}(\tilde{H}) \geq \lambda_{\min}(H) - \|H - \tilde{H}\|_F$$

Fact B.6. For $x \in (0, 1)$, integer $t \geq 0$, we have:

$$\sum_{\tau=1}^t (1-x)^\tau \leq -\frac{1}{\log(1-x)} \leq \frac{2}{x}$$

B.3 Probability Tools

Here, we state a probability toolkit in the following, including several helpful lemmas we'd like to use. Firstly, we provide the lemma about Chernoff bound in [95] below.

Lemma B.7 (Chernoff bound, [95]). Let $X = \sum_{i=1}^n X_i$, where $X_i = 1$ with probability p_i and $X_i = 0$ with probability $1 - p_i$, and all X_i are independent. Let $\mu = \mathbb{E}[X] = \sum_{i=1}^n p_i$. Then

- $\Pr[X \geq (1 + \delta)\mu] \leq \exp(-\delta^2\mu/3)$, $\forall \delta > 0$;
- $\Pr[X \leq (1 - \delta)\mu] \leq \exp(-\delta^2\mu/1)$, $\forall 0 < \delta < 1$.

Next, we offer the lemma about Hoeffding bound as in [96].

Lemma B.8 (Hoeffding bound, [96]). Let X_1, \dots, X_n denote n independent bounded variables in $[a_i, b_i]$ for $a_i, b_i \in \mathbb{R}$. Let $X := \sum_{i=1}^n X_i$, then we have

$$\Pr[|X - \mathbb{E}[X]| \geq t] \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

We show the lemma of Bernstein inequality as [97].

Lemma B.9 (Bernstein inequality, [97]). Let X_1, \dots, X_n denote n independent zero-mean random variables. Suppose $|X_i| \leq M$ almost surely for all i . Then, for all positive t ,

$$\Pr\left[\sum_{i=1}^n X_i \geq t\right] \leq \exp\left(-\frac{t^2/2}{\sum_{j=1}^n \mathbb{E}[X_j^2] + Mt/3}\right)$$

Then, we give the Khintchine's inequality in [98, 99] as follows:

Lemma B.10 (Khintchine's inequality, [98, 99]). Let $\sigma_1, \dots, \sigma_n$ be i.i.d sign random variables, and let z_1, \dots, z_n be real numbers. Then there are constants $C > 0$ so that for all $t > 0$

$$\Pr\left[\left|\sum_{i=1}^n z_i \sigma_i\right| \geq t \|z\|_2\right] \leq \exp(-Ct^2)$$

We give Hason-wright inequality from [100, 101] below.

Lemma B.11 (Hason-wright inequality, [100, 101]). Let $x \in \mathbb{R}^n$ denote a random vector with independent entries x_i with $\mathbb{E}[x_i] = 0$ and $|x_i| \leq K$. Let A be an $n \times n$ matrix. Then, for every $t \geq 0$

$$\Pr[|x^\top A x - \mathbb{E}[x^\top A x]| > t] \leq 2 \exp\left(-c \min\left\{t^2/(K^4 \|A\|_F^2), t/(K^2 \|A\|)\right\}\right)$$

We state Lemma 1 on page 1325 of Laurent and Massart [102].

Lemma B.12 (Lemma 1 on page 1325 of Laurent and Massart, [102]). Let $X \sim \mathcal{X}_k^2$ be a chi-squared distributed random variable with k degrees of freedom. Each one has zero mean and σ^2 variance. Then

$$\begin{aligned} \Pr[X - k\sigma^2 \geq (2\sqrt{kt} + 2t)\sigma^2] &\leq \exp(-t) \\ \Pr[X - k\sigma^2 \leq 2\sqrt{kt}\sigma^2] &\leq \exp(-t) \end{aligned}$$

Here, we provide a tail bound for sub-exponential distribution [103].

Lemma B.13 (Tail bound for sub-exponential distribution, [103]). We say $X \in \text{SE}(\sigma^2, \alpha)$ with parameters $\sigma > 0, \alpha > 0$, if

$$\mathbb{E}[e^{\lambda X}] \leq \exp(\lambda^2 \sigma^2 / 2), \forall |\lambda| < 1/\alpha.$$

Let $X \in \text{SE}(\sigma^2, \alpha)$ and $\mathbb{E}[X] = \mu$, then:

$$\Pr[|X - \mu| \geq t] \leq \exp(-0.5 \min\{t^2/\sigma^2, t/\alpha\})$$

In the following, we show the helpful lemma of matrix Chernoff bound as in [104, 105].

Lemma B.14 (Matrix Chernoff bound, [104, 105]). Let \mathcal{X} be a finite set of positive-semidefinite matrices with dimension $d \times d$, and suppose that

$$\max_{X \in \mathcal{X}} \lambda_{\max}(X) \leq B.$$

Sample $\{X_1, \dots, X_n\}$ uniformly at random from \mathcal{X} without replacement. We define μ_{\min} and μ_{\max} as follows:

$$\begin{aligned} \mu_{\min} &:= n \cdot \lambda_{\min}\left(\mathbb{E}_{X \in \mathcal{X}}(X)\right) \\ \mu_{\max} &:= n \cdot \lambda_{\max}\left(\mathbb{E}_{X \in \mathcal{X}}(X)\right). \end{aligned}$$

Then

$$\begin{aligned} \Pr\left[\lambda_{\min}\left(\sum_{i=1}^n X_i\right) \leq (1 - \delta)\mu_{\min}\right] &\leq d \cdot \exp(-\delta^2 \mu_{\min}/B) \text{ for } \delta \in (0, 1], \\ \Pr\left[\lambda_{\max}\left(\sum_{i=1}^n X_i\right) \geq (1 + \delta)\mu_{\max}\right] &\leq d \cdot \exp(-\delta^2 \mu_{\max}/(4B)) \text{ for } \delta \geq 0. \end{aligned}$$

Finally, we state Markov's inequality as below.

Lemma B.15 (Markov's inequality). If X is a non-negative random variable and $a > 0$, then the probability that X is at least a is at most the expectation of X divided by a :

$$\Pr[X \geq a] \leq \frac{\mathbb{E}[X]}{a}$$

B.4 Basic Bound

Definition B.16. For $\delta \in (0, 0.1)$ and a sufficiently large constant $C > 0$, we define:

$$D := \max\{C\sqrt{\log(md/\delta)}, 1\}$$

C NTK Problem Setup

C.1 Dataset

We consider a dataset where each data point is a tuple that includes a vector input and a scalar output. In particular, we assume that ℓ_2 norm of each input equals 1 and the absolute value of each target is not bigger than 1. We give the formal definition as follows:

Definition C.1 (Data Points). We define dataset $\mathcal{D} := \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$, where $\|x_i\|_2 = 1$ and $|y_i| \leq 1$ for any $i \in [n]$.

C.2 Model

Weights and Initialization.

Definition C.2. We give the following definitions:

- **Hidden-layer weights** $W \in \mathbb{R}^{d \times m}$. We define the hidden-layer weights $W := [w_1, w_2, \dots, w_m] \in \mathbb{R}^{d \times m}$ where $w_r \in \mathbb{R}^d, \forall r \in [m]$.

- **Output-layer weights** $a \in \mathbb{R}^m$. We define the output-layer weights $a := [a_1, a_2, \dots, a_m]^\top \in \mathbb{R}^m$, especially, vector a is fixed during the training.

Definition C.3. We give the following initializations:

- **Initialization of hidden-layer weights** $W \in \mathbb{R}^{d \times m}$. We randomly initialize $W(0) := [w_1(0), w_2(0), \dots, w_m(0)] \in \mathbb{R}^{d \times m}$, where its r -th column for $r \in [m]$ is sampled by $w_r(0) \sim \mathcal{N}(0, \sigma^2 \cdot I_d)$ with $\sigma^2 = 1$.
- **Initialization of output-layer weights** $a \in \mathbb{R}^m$. We randomly initialize $a \in \mathbb{R}^m$ where its r -th entry for $r \in [m]$ is sampled by $a_r \sim \text{Uniform}\{-1, +1\}$.

Model.

Definition C.4. For a scalar $x \in \mathbb{R}$, we define:

$$\text{ReLU}(x) = \max\{0, x\} \in \mathbb{R}$$

Definition C.5. If the following conditions hold:

- For a input vector $x \in \mathbb{R}^d$.
- For a hidden-layer weights $W \in \mathbb{R}^{d \times m}$ as Definition C.2.
- For a output-layer weights $a \in \mathbb{R}^m$ as Definition C.2.
- Let $\mathbf{q} : \mathbb{R}^d \rightarrow \{-1, +1\}^d$ be defined as Definition D.4.
- Let $\text{dq} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition D.5.
- Denote $\tilde{w}_r = \mathbf{q}(w_r) \in \{-1, +1\}^d$.
- Let $\text{ReLU} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition C.4.
- For $\kappa \in (0, 1]$.

We define:

$$f(x, W, a) := \kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \text{ReLU}\left(\text{dq}(\langle \tilde{w}_r, x \rangle)\right) \in \mathbb{R}$$

Lemma C.6. If the following conditions hold:

- For a input vector $x \in \mathbb{R}^d$.
- For a hidden-layer weights $W \in \mathbb{R}^{d \times m}$ as Definition C.2.
- For a output-layer weights $a \in \mathbb{R}^m$ as Definition C.2.
- Let $\mathbf{q} : \mathbb{R}^d \rightarrow \{-1, +1\}^d$ be defined as Definition D.4.
- Let $\text{dq} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition D.5.
- Denote $\tilde{w}_r = \mathbf{q}(w_r) \in \{-1, +1\}^d$.
- Let $\text{ReLU} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition C.4.
- Let $u : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be defined as Definition D.6.
- For $\kappa \in (0, 1]$.

Then we have:

$$f(x, W, a) := \kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \text{ReLU}\left(\langle w_r, x \rangle + \langle u(w_r), x \rangle\right)$$

Proof. We have

$$\begin{aligned}
f(x, W, a) &= \kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \text{ReLU}\left(\text{dq}(\langle \tilde{w}_r, x \rangle)\right) \\
&= \kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \text{ReLU}\left(\sqrt{V(w)} \cdot (\langle \tilde{w}, x \rangle + E(w) \cdot \langle x, \mathbf{1}_d \rangle)\right) \\
&= \kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \text{ReLU}\left(\langle w_r, x \rangle + \langle u(w_r), x \rangle\right)
\end{aligned}$$

where the first step follows from Definition C.5, the second step follows from Definition D.5, the last step follows from Definition D.6. \square

C.3 Training

Training.

Definition C.7. *If the following conditions hold:*

- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(0) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $f : \mathbb{R}^d \times \mathbb{R}^{d \times m} \times \mathbb{R}^m \rightarrow \mathbb{R}$ be defined as Definition C.5.
- For any $t \geq 0$.

We define:

$$L(W(t)) := \frac{1}{2} \cdot \sum_{i=1}^n (f(x_i, W(t), a) - y_i)^2$$

Definition C.8. *If the following conditions hold:*

- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(0) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $f : \mathbb{R}^d \times \mathbb{R}^{d \times m} \times \mathbb{R}^m \rightarrow \mathbb{R}$ be defined as Definition C.5.
- For any $t \geq 0$.
- Let $L(W(t))$ be defined as Definition C.7.
- Denote $\eta > 0$ as the learning rate.
- Let $\Delta W(t) \in \mathbb{R}^{d \times m}$ be defined as Definition F.2.

We update:

$$W(t+1) := W(t) - \eta \cdot \Delta W(t)$$

Compact Form.

Definition C.9. *If the following conditions hold:*

- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.

- Let $W(0) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $f : \mathbb{R}^d \times \mathbb{R}^{d \times m} \times \mathbb{R}^m \rightarrow \mathbb{R}$ be defined as Definition C.5.
- For any $t \geq 0$.
- Let $L(W(t))$ be defined as Definition C.7.
- Let $W(t)$ be updated by Definition C.8.

We give the following compact form of defined variables and functions:

- **Compact form of model function.** We define:

$$F(t) := [f(x_1, W(t), a), f(x_2, W(t), a), \dots, f(x_n, W(t), a)]^\top \in \mathbb{R}^n$$

- **Compact form of the input vector in the training dataset.** We define:

$$X := [x_1, x_2, \dots, x_n]^\top \in \mathbb{R}^{n \times d}$$

- **Compact form of the targets in the training dataset.** We define:

$$y := [y_1, y_2, \dots, y_n]^\top \in \mathbb{R}^n$$

- **Compact form of the training objective.** We define:

$$L(t) := \frac{1}{2} \cdot \|F(t) - y\|_2^2$$

Epecially, we have $L(t) = L(W(t))$ by simple algebras.

D Quantization

D.1 Quantization Functions

Definition D.1. For a vector $w \in \mathbb{R}^d$, we define $\text{Sign}(w) \in \{-1, +1\}^d$ where its k -th entry for $k \in [d]$ is given by:

$$\text{Sign}_k(w) := \begin{cases} -1, & \text{if } w_k < 0 \\ +1, & \text{if } w_k \geq 0 \end{cases} \in \{-1, +1\}$$

Definition D.2. For a vector $w \in \mathbb{R}^d$, we define expectation function as follows:

$$E(w) := \frac{\langle w, \mathbf{1}_d \rangle}{d} \in \mathbb{R}$$

Definition D.3. Let $E : \mathbb{R}^d \rightarrow \mathbb{R}$ be defined as Definition D.2. For a vector $w \in \mathbb{R}^d$, we define variance function as follows:

$$V(w) := \frac{1}{d} \cdot \|w - E(w) \cdot \mathbf{1}_d\|_2^2 \in \mathbb{R}$$

Definition D.4. If the following conditions hold:

- Let $\text{Sign} : \mathbb{R}^d \rightarrow \{-1, +1\}^d$ be defined as Definition D.1.
- Let $E : \mathbb{R}^d \rightarrow \mathbb{R}$ be defined as Definition D.2.
- Let $V : \mathbb{R}^d \rightarrow \mathbb{R}$ be defined as Definition D.3.
- For a weight vector $w \in \mathbb{R}^d$.

We define the quantization function as follows:

$$q(w) := \text{Sign}\left(\frac{w - E(w) \cdot \mathbf{1}_d}{\sqrt{V(w)}}\right) \in \{-1, +1\}^d$$

D.2 Dequantization Functions

Definition D.5. *If the following conditions hold:*

- Let $\mathbf{q} : \mathbb{R}^d \rightarrow \{-1, +1\}^d$ be defined as Definition D.4.
- Let $E : \mathbb{R}^d \rightarrow \mathbb{R}$ be defined as Definition D.2.
- Let $V : \mathbb{R}^d \rightarrow \mathbb{R}$ be defined as Definition D.3.
- For a weight vector $w \in \mathbb{R}^d$.
- Denote quantized vector $\tilde{w} := \mathbf{q}(w) \in \{-1, +1\}^d$.
- For a vector $x \in \mathbb{R}^d$.

We define the dequantization function as follows:

$$\text{dq}(\langle \tilde{w}, x \rangle) := \sqrt{V(w)} \cdot \langle \tilde{w}, x \rangle + E(w) \cdot \langle x, \mathbf{1}_d \rangle \in \mathbb{R}$$

D.3 Quantization Error

Definition D.6. *If the following conditions hold:*

- Let $\mathbf{q} : \mathbb{R}^d \rightarrow \{-1, +1\}^d$ be defined as Definition D.4.
- Let $E : \mathbb{R}^d \rightarrow \mathbb{R}$ be defined as Definition D.2.
- Let $V : \mathbb{R}^d \rightarrow \mathbb{R}$ be defined as Definition D.3.
- For a weight vector $w \in \mathbb{R}^d$.
- Denote quantized vector $\tilde{w} := \mathbf{q}(w) \in \{-1, +1\}^d$.
- For a vector $x \in \mathbb{R}^d$.

We define the quantization difference vector as follows:

$$u(w) := \sqrt{V(w)}\tilde{w} + E(w) \cdot \mathbf{1}_d - w \in \mathbb{R}^d$$

Lemma D.7. *If the following conditions hold:*

- Let $D > 0$ be defined as Definition B.16.
- Let $\mathbf{q} : \mathbb{R}^d \rightarrow \{-1, +1\}^d$ be defined as Definition D.4.
- Let $E : \mathbb{R}^d \rightarrow \mathbb{R}$ be defined as Definition D.2.
- Let $V : \mathbb{R}^d \rightarrow \mathbb{R}$ be defined as Definition D.3.
- For a weight vector $w \in \mathbb{R}^d$.
- Denote quantized vector $\tilde{w} := \mathbf{q}(w) \in \{-1, +1\}^d$.
- For a vector $x \in \mathbb{R}^d$ and $\|x\|_2 = 1$.
- Let $u : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be defined as Definition D.6.

Then we have:

$$\langle u(w), x \rangle \leq O\left(d(D + R)\right)$$

Proof. We define:

$$\mathbf{L}n(w) = \frac{w - E(w)\mathbf{1}_d}{\sqrt{V(w)}}$$

Then by simple algebras, we can show that:

$$\frac{1}{d} \|\text{Ln}(w)\|_2^2 = \frac{1}{d} \left\| \frac{w - E(w)\mathbf{1}_d}{\sqrt{V(w)}} \right\|_2^2 < \frac{1}{d} \frac{\|w - E(w)\mathbf{1}_d\|_2^2}{V(w)} < 1 \quad (4)$$

Thus, we obtain:

$$\begin{aligned} \|\text{Ln}(w)\|_\infty &\leq \|\text{Ln}(w)\|_2 \\ &= (\|\text{Ln}(w)\|_2^2)^{\frac{1}{2}} \\ &< \sqrt{d} \end{aligned}$$

where these steps follow from simple algebras and Eq. (4).

Finally, we can get that

$$\begin{aligned} |\langle u(w), x \rangle| &= \sqrt{V(w)} \cdot |\langle \tilde{w} - \text{Ln}(w), x \rangle| \\ &= O(D + R) \cdot |\langle \tilde{w} - \text{Ln}(w), x \rangle| \\ &\leq O(D + R) \cdot \|\tilde{w} - \text{Ln}(w)\|_2 \\ &= O(D + R) \cdot \left(\sum_{k=1}^d (\tilde{w}_k - \text{Ln}_k(w))^2 \right)^{\frac{1}{2}} \\ &\leq O(D + R) \cdot \left(\sum_{k=1}^d (\max\{\sqrt{d} - 1, 1\})^2 \right)^{\frac{1}{2}} \\ &\leq O(d(D + R)) \end{aligned}$$

where the first step follows from Definition D.6, the second step follows from Part 7 of Lemma I.6, the third step follows from Cauchy-Schwarz inequality and $\|x\|_2 = 1$, the fourth step follows from the definition of ℓ_2 norm, the fifth step follows from Definition D.1 and simple algebras, the last step follows from simple algebras. \square

E Patterns

E.1 ReLU Pattern

Definition E.1. *If the following conditions hold:*

- For any $w \in \mathbb{R}^d$.
- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(0) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3.
- Let $\text{dq} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition D.5.
- For $R > 0$.
- For $i \in [n]$ and $r \in [m]$.

We define:

$$\mathbf{A}_{i,r} := \{\exists w \in \mathbb{R}^d : \|w - w_r(0)\|_2 \leq R, \mathbf{1}_{\text{dq}(\langle w_r(0), x_i \rangle) \geq 0} \neq \mathbf{1}_{\text{dq}(\langle w, x_i \rangle) \geq 0}\}$$

Definition E.2. *Let event $\mathbf{A}_{i,r}$ for $i \in [n]$ and $r \in [m]$ be defined as Definition E.1. We define:*

$$\begin{aligned} \mathcal{S}_i &:= \{r \in [m] : \mathbb{I}\{\mathbf{A}_{i,r}\} = 0\} \\ \mathcal{S}_i^\perp &:= [m] / \mathcal{S}_i \end{aligned}$$

E.2 Sign Pattern

Definition E.3. *If the following conditions hold:*

- For any $w \in \mathbb{R}^d$.
- Let $W(0) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3.
- For $R > 0$.
- For $k \in [d]$ and $r \in [m]$.

We define:

$$\mathbb{B}_{r,k} := \{\exists w \in \mathbb{R}^d : |w_k - w_{r,k}(0)| \leq R, \mathbf{1}_{w_{r,k}(0) - E(w_r(0)) \geq 0} \neq \mathbf{1}_{w_k - E(w) \geq 0}\}$$

F Straight-Through Estimator (STE)

F.1 STE Functions

Definition F.1. *If the following conditions hold:*

- For a input vector $x \in \mathbb{R}^d$.
- For a hidden-layer weights $W \in \mathbb{R}^{d \times m}$ as Definition C.2.
- For a output-layer weights $a \in \mathbb{R}^m$ as Definition C.2.
- Let $\mathbf{q} : \mathbb{R}^d \rightarrow \{-1, +1\}^d$ be defined as Definition D.4.
- Denote $\tilde{w}_r = \mathbf{q}(w_r) \in \{-1, +1\}^d$.
- Let $\text{ReLU} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition C.4.

We define:

$$f_{\text{ste}}(x, W, a) := \kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r, x \rangle) \geq 0} \cdot \langle w_r, x \rangle \in \mathbb{R}$$

Then its compact form is given by

$$\mathbf{F}_{\text{ste}}(t) := [f_{\text{ste}}(x_1, W(t), a), f_{\text{ste}}(x_2, W(t), a), \dots, f_{\text{ste}}(x_n, W(t), a)]^\top \in \mathbb{R}^n$$

Definition F.2. *Let $W(0) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3. For any $t \geq 0$. We define:*

$$\Delta W(t) := \sum_{i=1}^n (\mathbf{F}_i(t) - y_i) \cdot \frac{d\mathbf{F}_{\text{ste},i}(t)}{dW(t)}$$

F.2 Gradient Computation

Lemma F.3. *If the following conditions hold:*

- For $i \in [n]$, $r \in [m]$ and integer $t \geq 0$.
- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $\mathbf{F}_{\text{ste}}(t)$ be defined as Definition F.1.
- Let $\mathbf{q} : \mathbb{R}^d \rightarrow \{-1, +1\}^d$ be defined as Definition D.4.

- Denote $\tilde{w}_r = \mathbf{q}(w_r) \in \{-1, +1\}^d$.
- For $\kappa \in (0, 1]$.

Then we have:

$$\frac{dF_{\text{ste},i}(t)}{dw_r(t)} = \kappa \frac{1}{\sqrt{m}} a_r \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \cdot x_i$$

Proof. This proof follows from simple calculations. □

G Neural Tangent Kernel

G.1 Kernel Function

Definition G.1. *If the following conditions hold:*

- For $i, j \in [n]$, $r \in [m]$ and integer $t \geq 0$.
- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $\mathbf{q} : \mathbb{R}^d \rightarrow \{-1, +1\}^d$ be defined as Definition D.4.
- Let $\text{dq} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition D.5.
- Denote $\tilde{w}_r = \mathbf{q}(w_r) \in \{-1, +1\}^d$.
- For $\kappa \in (0, 1]$.

We define the kernel function as $H(t) \in \mathbb{R}^{n \times n}$, where its (i, j) -th entry is given by:

$$H_{i,j}(t) := \kappa^2 \frac{1}{m} x_i^\top x_j \cdot \sum_{r=1}^m \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_j \rangle) \geq 0} \in \mathbb{R}$$

Claim G.2. *If the following conditions hold:*

- For $i, j \in [n]$, $r \in [m]$ and integer $t \geq 0$.
- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $\mathbf{q} : \mathbb{R}^d \rightarrow \{-1, +1\}^d$ be defined as Definition D.4.
- Let $\text{dq} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition D.5.
- Denote $\tilde{w}_r = \mathbf{q}(w_r) \in \{-1, +1\}^d$.
- Let $H(t) \in \mathbb{R}^{n \times n}$ be defined as Definition G.1.
- For $\kappa \in (0, 1]$.

We first define the neural tangent network as $H^* := H(0) \in \mathbb{R}^{n \times n}$, where its (i, j) -th entry is given by:

$$\begin{aligned} H_{i,j}^* &:= H_{i,j}(0) \\ &= \kappa^2 \frac{1}{m} x_i^\top x_j \cdot \sum_{r=1}^m \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(0), x_i \rangle) \geq 0} \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(0), x_j \rangle) \geq 0} \end{aligned}$$

$$\approx \kappa^2 x_i^\top x_j \cdot \mathbb{E}_{w_r \sim \mathcal{N}(0, \sigma^2 \cdot I_d)} [\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(0), x_i \rangle) \geq 0} \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(0), x_j \rangle) \geq 0}]$$

Proof. We have

$$\begin{aligned} H_{i,j}^* &= H_{i,j}(0) \\ &= \kappa^2 \frac{1}{m} x_i^\top x_j \cdot \sum_{r=1}^m \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(0), x_i \rangle) \geq 0} \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(0), x_j \rangle) \geq 0} \\ &\approx \kappa^2 x_i^\top x_j \cdot \mathbb{E}_{w_r \sim \mathcal{N}(0, \sigma^2 \cdot I_d)} [\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(0), x_i \rangle) \geq 0} \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(0), x_j \rangle) \geq 0}] \end{aligned}$$

where the first step follows from the definition of H^* , the second step follows from Definition G.1, the third step holds since $m \rightarrow +\infty$. \square

Definition G.3. *If the following conditions hold:*

- For $i, j \in [n]$, $r \in [m]$ and integer $t \geq 0$.
- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $\mathbf{q} : \mathbb{R}^d \rightarrow \{-1, +1\}^d$ be defined as Definition D.4.
- Let $\text{dq} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition D.5.
- Denote $\tilde{w}_r = \mathbf{q}(w_r) \in \{-1, +1\}^d$.
- Let \mathcal{S}_i^\perp be defined as Definition E.2.

We the pattern-changing kernel function as $H^\perp(t) \in \mathbb{R}^{n \times n}$, where its (i, j) -th entry is given by:

$$H_{i,j}^\perp(t) := \kappa^2 \frac{1}{m} x_i^\top x_j \cdot \sum_{r \in \mathcal{S}_i^\perp} \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_j \rangle) \geq 0} \in \mathbb{R}$$

G.2 Assumption: H^* is Positive Definite

Assumption G.4. *Let $H^* \in \mathbb{R}^{n \times n}$ be defined as Definition G.1. We assume that H^* is positive definite (PD), where its minimum eigenvalue is given by:*

$$\lambda := \lambda_{\min}(H^*) > 0$$

G.3 Kernel Convergence and PD Property

Lemma G.5. *If the following conditions hold:*

- Let $D > 0$ be defined as Definition B.16.
- Denote $\lambda = \lambda_{\min}(H^*) > 0$ as Assumption G.4.
- For $i, j \in [n]$, $r \in [m]$ and integer $t \geq 0$.
- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $H(t) \in \mathbb{R}^{n \times n}$ be defined as Definition G.1.

- Let $H^* \in \mathbb{R}^{n \times n}$ be defined as Claim G.2.
- $R \leq O\left(\frac{\lambda\delta}{\kappa^2 n^2 d D}\right)$.
- $\delta \in (0, 0.1)$.

Then with probability at least $1 - \delta$, we have:

- Part 1.

$$\|H(t) - H^*\|_F \leq O\left(n^2 d R \delta^{-1} D\right)$$

- Part 2.

$$\lambda_{\min}(H(t)) \geq \lambda/2$$

Proof. Proof of Part 1. Let $A_{i,r}$ be defined as Definition E.1, we first show that when $\langle w_r(0), x \rangle \geq R + O\left(d(D + R)\right)$

$$\begin{aligned} \text{dq}(\langle \tilde{w}_r(0), x_i \rangle) &= \sqrt{V(w_r(0))} \cdot \langle \tilde{w}_r(0), x_i \rangle + \langle E(w_r(0)) \cdot \mathbf{1}_d, x_i \rangle \\ &= \langle w_r(0), x_i \rangle + \langle u(w_r(0)), x_i \rangle \\ &\geq \langle w_r(0), x_i \rangle - |\langle u(w_r(0)), x_i \rangle| \\ &\geq R \end{aligned}$$

where the first step follows from Definition D.5, the second step follows from Definition D.6. the third step follows from simple algebras, the last step follows from $\langle w_r(0), x \rangle \geq R + O\left(d(D + R)\right)$ and Lemma D.7.

Thus, for any $w \in \mathbb{R}^d$ that satisfies $\|w - w_r(0)\|_2 \leq R$, we have:

$$\begin{aligned} \text{dq}(\langle \tilde{w}, x_i \rangle) &= \sqrt{V(w)} \cdot \langle \tilde{w}, x_i \rangle + \langle E(w) \cdot \mathbf{1}_d, x_i \rangle \\ &= \langle w, x_i \rangle + \langle u(w), x_i \rangle \\ &\geq \langle w, x_i \rangle - |\langle u(w), x_i \rangle| \\ &\geq \langle w_r(0), x_i \rangle - \|w - w_r(0)\|_2 - |\langle u(w), x_i \rangle| \\ &\geq 0 \end{aligned}$$

where the first step follows from Definition D.5, the second step follows from Definition D.6. the third step follows from simple algebras, the fourth step follows from Cauchy-Schwarz inequality and $\|x_i\| = 1$, the last step follows from $\|w - w_r(0)\|_2 \leq R$, $\langle w_r(0), x \rangle \geq R + O\left(d(D + R)\right)$ and Lemma D.7.

The above situation says:

$$\begin{aligned} \Pr\left[\mathbb{I}\{A_{i,r}\} = 1\right] &\leq \Pr[\langle w_r(0), x \rangle < R + O\left(d(D + R)\right)] \\ &\leq \frac{4R + O\left(d(D + R)\right)}{\sqrt{2\pi}} \\ &\leq O\left(dR(D + R)\right) \\ &\leq O\left(dRD\right) \end{aligned} \tag{5}$$

where the second step follows from anti-concentration of Gaussian (Fact B.3) and Fact B.4, the third step follows from simple algebras and the last step follows from plugging $R \leq D$.

For $i, j \in [n]$, we have

$$\mathbb{E}[|H_{i,j}(t) - H_{i,j}^*|]$$

$$\begin{aligned}
&= \mathbb{E} \left[\left| \kappa^2 \frac{1}{m} x_i^\top x_j \sum_{r=1}^m (\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_j \rangle) \geq 0} - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(0), x_i \rangle) \geq 0} \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(0), x_j \rangle) \geq 0}) \right| \right] \\
&= \kappa^2 \frac{1}{m} \sum_{r=1}^m \mathbb{E} \left[\left| \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_j \rangle) \geq 0} - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(0), x_i \rangle) \geq 0} \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(0), x_j \rangle) \geq 0} \right| \right] \\
&\leq \kappa^2 \frac{1}{m} \sum_{r=1}^m \mathbb{E} \left[\mathbb{I}\{A_{i,r} \cup A_{j,r}\} \right] \\
&\leq O(\kappa^2 dRD) \tag{6}
\end{aligned}$$

where the first step follows from Definition G.1 and Claim G.2, the second and third step follows from simple algebras, the last step follows from Eq. (5).

Then we have:

$$\begin{aligned}
\mathbb{E} \left[\sum_{i=1}^n \sum_{j=1}^n |H_{i,j}(t) - H_{i,j}^*| \right] &= \sum_{i=1}^n \sum_{j=1}^n \mathbb{E} [|H_{i,j}(t) - H_{i,j}^*|] \\
&\leq O(\kappa^2 n^2 dRD)
\end{aligned}$$

where the first step follows from simple algebras, the second step follows from Eq. (6).

Hence, by Markov's inequality (Lemma B.15), with probability at least $1 - \delta$, we have:

$$\begin{aligned}
\sum_{i=1}^n \sum_{j=1}^n |H_{i,j}(t) - H_{i,j}^*| &\leq \frac{\mathbb{E}[\sum_{i=1}^n \sum_{j=1}^n |H_{i,j}(t) - H_{i,j}^*|]}{\delta} \\
&\leq O(\kappa^2 n^2 dR\delta^{-1}(D + R))
\end{aligned}$$

We obtain:

$$\begin{aligned}
\|H(t) - H^*\|_F &\leq \|H(t) - H^*\|_1 \\
&= \sum_{i=1}^n \sum_{j=1}^n |H_{i,j}(t) - H_{i,j}^*| \\
&\leq O(\kappa^2 n^2 dR\delta^{-1}D)
\end{aligned}$$

Now following Fact B.5, we have:

$$\begin{aligned}
\lambda_{\min}(H(t)) &\geq \lambda_{\min}(H^*) - \|H(t) - H^*\|_F \\
&\geq \lambda - O(\kappa^2 n^2 dR\delta^{-1}D) \\
&\geq \lambda/2
\end{aligned}$$

where the last step follows from choosing $R \leq O(\frac{\lambda\delta}{\kappa^2 n^2 dD})$. \square

H Training Dynamic

H.1 Decompose Loss

Definition H.1. Let $W(0) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3. For any $t \geq 0$. Let $u : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be defined as Definition D.6. For $r \in [m]$. We define:

$$\mathbf{u}_r(t) := u(w_r(t))$$

Then the $F_i(t), \forall i \in [n]$ can be given by:

$$F_i(t) = \kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \cdot (\langle w_r(t), x_i \rangle + \langle \mathbf{u}_r(t), x_i \rangle)$$

Claim H.2. *If the following conditions hold:*

- For $i, j \in [n]$, $r \in [m]$ and integer $t \geq 0$.
- Let $L(t)$ be defined as Definition C.9.
- Let $F(t) \in \mathbb{R}^n$ be defined as Definition C.9.
- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $\text{dq} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition D.5.
- Denote $\tilde{w}_r = \mathbf{q}(w_r) \in \{-1, +1\}^d$.
- Let $\mathcal{S}_i, \mathcal{S}_i^\perp$ be defined as Definition E.2.
- Let $u_r(t)$ be defined as Definition H.1.
- Define

$$C_1 := -\kappa \frac{1}{\sqrt{m}} \sum_{i=1}^n \sum_{r \in \mathcal{S}_i} a_r \left(\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \langle w_r(t), x_i \rangle \right. \\ \left. - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t+1), x_i \rangle) \geq 0} \langle w_r(t+1), x_i \rangle \right) \cdot (F_i(t) - y_i)$$

- Define

$$C_2 := -\kappa \frac{1}{\sqrt{m}} \sum_{i=1}^n \sum_{r \in \mathcal{S}_i^\perp} a_r \left(\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \langle w_r(t), x_i \rangle \right. \\ \left. - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t+1), x_i \rangle) \geq 0} \langle w_r(t+1), x_i \rangle \right) \cdot (F_i(t) - y_i)$$

- Define

$$C_3 := -\kappa \frac{1}{\sqrt{m}} \sum_{i=1}^n \sum_{r=1}^m a_r \left(\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \langle u_r(t), x_i \rangle \right. \\ \left. - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t+1), x_i \rangle) \geq 0} \langle u_r(t+1), x_i \rangle \right) \cdot (F_i(t) - y_i)$$

- Define

$$C_4 := \frac{1}{2} \|F(t) - F(t+1)\|_2^2$$

- For $\kappa \in (0, 1]$.

Then we have:

$$L(t+1) = L(t) + C_1 + C_2 + C_3 + C_4$$

Proof. We have

$$\begin{aligned} L(t+1) &= \frac{1}{2} \cdot \|F(t+1) - y\|_2^2 \\ &= \frac{1}{2} \cdot \|(F(t) - y) - (F(t) - F(t+1))\|_2^2 \\ &= \frac{1}{2} \cdot (\|F(t) - y\|_2^2 - 2\langle F(t) - y, F(t) - F(t+1) \rangle + \|F(t) - F(t+1)\|_2^2) \end{aligned}$$

$$= \mathsf{L}(t) - \langle \mathsf{F}(t) - y, \mathsf{F}(t) - \mathsf{F}(t+1) \rangle + \frac{1}{2} \|\mathsf{F}(t) - \mathsf{F}(t+1)\|_2^2$$

these steps follow from simple algebras and Definition C.9.

Then for $i \in [n]$

$$\begin{aligned} & \mathsf{F}_i(t) - \mathsf{F}_i(t+1) \\ &= \kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \cdot \left(\langle w_r(t), x_i \rangle + \langle \mathbf{u}_r(t), x_i \rangle \right) \\ & \quad - \kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t+1), x_i \rangle) \geq 0} \cdot \left(\langle w_r(t+1), x_i \rangle + \langle \mathbf{u}_r(t+1), x_i \rangle \right) \\ &= \kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \left(\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \cdot \left(\langle w_r(t), x_i \rangle + \langle \mathbf{u}_r(t), x_i \rangle \right) \right. \\ & \quad \left. - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t+1), x_i \rangle) \geq 0} \cdot \left(\langle w_r(t+1), x_i \rangle + \langle \mathbf{u}_r(t+1), x_i \rangle \right) \right) \\ &= M_{1,i} + M_{2,i} + M_{3,i} \end{aligned}$$

where these steps follows from simple algebras and defining:

$$\begin{aligned} M_{1,i} &:= \kappa \frac{1}{\sqrt{m}} \sum_{r \in \mathcal{S}_i} a_r \left(\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \cdot \langle w_r(t), x_i \rangle - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t+1), x_i \rangle) \geq 0} \cdot \langle w_r(t+1), x_i \rangle \right) \\ M_{2,i} &:= \kappa \frac{1}{\sqrt{m}} \sum_{r \in \mathcal{S}_i^\perp} a_r \left(\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \cdot \langle w_r(t), x_i \rangle - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t+1), x_i \rangle) \geq 0} \cdot \langle w_r(t+1), x_i \rangle \right) \\ M_{3,i} &:= \kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \left(\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \cdot \langle \mathbf{u}_r(t), x_i \rangle - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t+1), x_i \rangle) \geq 0} \cdot \langle \mathbf{u}_r(t+1), x_i \rangle \right) \end{aligned}$$

Thus, by the definitions in Lemma conditions, we can show that

$$\mathsf{L}(t+1) = \mathsf{L}(t) + C_1 + C_2 + C_3 + C_4$$

□

H.2 Bounding C_1

Lemma H.3. *If the following conditions hold:*

- Let $D > 0$ be defined as Definition B.16.
- For $i, j \in [n]$, $r \in [m]$ and integer $t \geq 0$.
- Let $H(t) \in \mathbb{R}^{n \times n}$ be defined as Definition G.1.
- Let $H^\perp(t) \in \mathbb{R}^{n \times n}$ be defined as Definition G.3.
- Let $H^* \in \mathbb{R}^{n \times n}$ be defined as Claim G.2. Assume $\lambda_{\min}(H^*) > 0$ as Assumption G.4.
- Let $\mathsf{L}(t)$ be defined as Definition C.9.
- Let $\mathsf{F}(t) \in \mathbb{R}^n$ be defined as Definition C.9.
- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.

- Let $\text{dq} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition D.5.
- Denote $\tilde{w}_r = \mathbf{q}(w_r) \in \{-1, +1\}^d$.
- Let $\mathcal{S}_i, \mathcal{S}_i^\perp$ be defined as Definition E.2.
- Let $\mathbf{u}_r(t)$ be defined as Definition H.1.
- $\delta \in (0, 0.1)$.
- Define

$$C_1 := -\kappa \frac{1}{\sqrt{m}} \sum_{i=1}^n \sum_{r \in \mathcal{S}_i} a_r (\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \langle w_r(t), x_i \rangle - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t+1), x_i \rangle) \geq 0} \langle w_r(t+1), x_i \rangle) \cdot (\mathbf{F}_i(t) - y_i)$$

- For $\kappa \in (0, 1]$.

Then with probability at least $1 - \delta$, we have:

$$C_1 \leq \left(-\eta\kappa\lambda + O\left(\eta\kappa \frac{n^2 d R D}{\delta}\right) \right) \cdot \mathbf{L}(t)$$

Proof. We have:

$$\begin{aligned} C_1 &= -\kappa \frac{1}{\sqrt{m}} \sum_{i=1}^n \sum_{r \in \mathcal{S}_i} a_r (\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \langle w_r(t), x_i \rangle - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t+1), x_i \rangle) \geq 0} \langle w_r(t+1), x_i \rangle) \cdot (\mathbf{F}_i(t) - y_i) \\ &= -\kappa \frac{1}{\sqrt{m}} \sum_{i=1}^n \sum_{r \in \mathcal{S}_i} a_r (\langle w_r(t), x_i \rangle - \langle w_r(t+1), x_i \rangle) \cdot (\mathbf{F}_i(t) - y_i) \\ &= -\kappa^2 \eta \frac{1}{m} \sum_{i=1}^n \sum_{r \in \mathcal{S}_i} (\mathbf{F}_i(t) - y_i) \cdot \left(\sum_{j=1}^n x_i^\top x_j \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_j \rangle) \geq 0} \cdot (\mathbf{F}_j(t) - y_j) \right) \\ &= -\eta (\mathbf{F}(t) - y)^\top \cdot (H(t) - H^\perp(t)) \cdot (\mathbf{F}(t) - y) \\ &= -\eta (\mathbf{F}(t) - y)^\top \cdot H(t) \cdot (\mathbf{F}(t) - y) + \eta (\mathbf{F}(t) - y)^\top \cdot H^\perp(t) \cdot (\mathbf{F}(t) - y) \\ &\leq -\eta\lambda/2 \cdot \|\mathbf{F}(t) - y\|_2^2 + \eta \|H^\perp(t)\|_F \cdot \|\mathbf{F}(t) - y\|_2 \\ &= (-\eta\lambda + \|H^\perp(t)\|_F) \cdot \mathbf{L}(t) \end{aligned}$$

where the first step follows from definition of C_1 , the second step follows from the definition of \mathcal{S}_i (Definition E.2), the third step follows from Definition C.8 and Definition F.2, the fourth step follows from Definition G.1, Definition G.3 and simple algebras, the fifth step follows from simple algebras, the sixth step follows from Lemma G.5 and simple algebras, the last step follows from Definition C.9.

Besides, we have

$$\begin{aligned} |H_{i,j}^\perp| &= \left| \frac{1}{m} x_i^\top x_j \cdot \sum_{r \in \mathcal{S}_i^\perp} \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_j \rangle) \geq 0} \right| \\ &\leq \left| \frac{1}{m} x_i^\top x_j \cdot |\mathcal{S}_i^\perp| \right| \\ &\leq \frac{1}{m} |\mathcal{S}_i^\perp| \end{aligned} \tag{7}$$

where the first step follows from Definition G.3, the second step follows from simple algebras, the third step follows from $\|x\|_i = 1$.

We give that

$$\begin{aligned} \mathbb{E}[\sum_{i=1}^n |\mathcal{S}_i^\perp|] &= \sum_{i=1}^n \sum_{r=1}^m \Pr[\mathbb{I}\{\mathbf{A}_{i,r}\} = 1] \\ &\leq O(mndRD) \end{aligned}$$

where the first step follows from simple algebras, the second step follows from Eq. (5).

Hence, by Markov's inequality (Lemma B.15), we have

$$\sum_{i=1}^n |\mathcal{S}_i^\perp| \leq O\left(\frac{mndRD}{\delta}\right) \quad (8)$$

Thus,

$$\begin{aligned} \|H^\perp\|_F &\leq \sum_{i=1}^n \sum_{j=1}^n |H_{i,j}^\perp| \\ &\leq \frac{1}{m} \sum_{i=1}^n \sum_{j=1}^n |\mathcal{S}_i^\perp| \\ &\leq O\left(\frac{n^2 dRD}{\delta}\right) \end{aligned}$$

where the first step follows from simple algebras, the second step follows from Eq. (7), the last step follows from simple algebras and Eq. (8).

Finally, we conclude all the results, we have:

$$C_1 \leq \left(-\eta\lambda + O\left(\eta\frac{n^2 dRD}{\delta}\right)\right) \cdot \mathsf{L}(t)$$

□

H.3 Bounding C_2

Lemma H.4. *If the following conditions hold:*

- Let $D > 0$ be defined as Definition B.16.
- For $i, j \in [n]$, $r \in [m]$ and integer $t \geq 0$.
- Let $H(t) \in \mathbb{R}^{n \times n}$ be defined as Definition G.1.
- Let $H^\perp(t) \in \mathbb{R}^{n \times n}$ be defined as Definition G.3.
- Let $H^* \in \mathbb{R}^{n \times n}$ be defined as Claim G.2. Assume $\lambda_{\min}(H^*) > 0$ as Assumption G.4.
- Let $\mathsf{L}(t)$ be defined as Definition C.9.
- Let $\mathsf{F}(t) \in \mathbb{R}^n$ be defined as Definition C.9.
- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $\mathsf{dq} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition D.5.
- Denote $\tilde{w}_r = \mathsf{q}(w_r) \in \{-1, +1\}^d$.
- Let $\mathcal{S}_i, \mathcal{S}_i^\perp$ be defined as Definition E.2.

- Let $u_r(t)$ be defined as Definition H.1.
- $\delta \in (0, 0.1)$.
- Define

$$C_2 := -\kappa \frac{1}{\sqrt{m}} \sum_{i=1}^n \sum_{r \in \mathcal{S}_i^\perp} a_r \left(\mathbf{1}_{\text{dq}(\langle \bar{w}_r(t), x_i \rangle) \geq 0} \langle w_r(t), x_i \rangle \right. \\ \left. - \mathbf{1}_{\text{dq}(\langle \bar{w}_r(t+1), x_i \rangle) \geq 0} \langle w_r(t+1), x_i \rangle \right) \cdot (\mathbf{F}_i(t) - y_i)$$

- $\kappa \in (0, 1]$.

Then with probability at least $1 - \delta$, we have:

$$|C_2| \leq O\left(\eta\kappa \frac{n^{1.5}dRD}{\delta}\right) \cdot \mathbf{L}(t)$$

Proof. We have:

$$\begin{aligned} |C_2| &= \left| \kappa \frac{1}{\sqrt{m}} \sum_{i=1}^n \sum_{r \in \mathcal{S}_i^\perp} a_r \left(\mathbf{1}_{\text{dq}(\langle \bar{w}_r(t), x_i \rangle) \geq 0} \langle w_r(t), x_i \rangle \right. \right. \\ &\quad \left. \left. - \mathbf{1}_{\text{dq}(\langle \bar{w}_r(t+1), x_i \rangle) \geq 0} \langle w_r(t+1), x_i \rangle \right) \cdot (\mathbf{F}_i(t) - y_i) \right| \\ &\leq \left| \kappa \frac{1}{\sqrt{m}} \sum_{i=1}^n |\mathcal{S}_{i^\perp}| \cdot |\langle w_r(t), x_i \rangle - \langle w_r(t+1), x_i \rangle| \cdot (\mathbf{F}_i(t) - y_i) \right| \\ &\leq \left| \kappa \frac{1}{\sqrt{m}} \sum_{i=1}^n |\mathcal{S}_{i^\perp}| \cdot \|\eta \Delta w_r(t)\|_2 \cdot (\mathbf{F}_i(t) - y_i) \right| \\ &\leq \kappa \frac{1}{\sqrt{m}} \sum_{i=1}^n |\mathcal{S}_{i^\perp}| \cdot \|\eta \Delta w_r(t)\|_2 \|\mathbf{F}(t) - y\|_2 \\ &\leq \eta\kappa \frac{\sqrt{n}}{m} \sum_{i=1}^n |\mathcal{S}_{i^\perp}| \cdot \|\mathbf{F}(t) - y\|_2^2 \\ &\leq O\left(\eta\kappa \frac{n^{1.5}dRD}{\delta}\right) \cdot \mathbf{L}(t) \end{aligned}$$

where the first step follows from the definition of C_2 , the second step follows from Fact B.2 and Definition E.2 (\mathcal{S}_i^\perp), the third step follows from simple algebras and Definition C.8, the fourth step follows from simple algebras, the fifth step follows from Lemma I.4, last step follows from Eq. (8) and Definition C.9. \square

H.4 Bounding C_3

Lemma H.5. *If the following conditions hold:*

- Let $D > 0$ be defined as Definition B.16.
- For $i, j \in [n]$, $r \in [m]$ and integer $t \geq 0$.
- Let $H(t) \in \mathbb{R}^{n \times n}$ be defined as Definition G.1.
- Let $H^\perp(t) \in \mathbb{R}^{n \times n}$ be defined as Definition G.3.
- Let $H^* \in \mathbb{R}^{n \times n}$ be defined as Claim G.2. Assume $\lambda_{\min}(H^*) > 0$ as Assumption G.4.
- Let $\mathbf{L}(t)$ be defined as Definition C.9.

- Let $F(t) \in \mathbb{R}^n$ be defined as Definition C.9.
- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $\text{dq} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition D.5.
- Denote $\tilde{w}_r = \mathbf{q}(w_r) \in \{-1, +1\}^d$.
- Let $\mathcal{S}_i, \mathcal{S}_i^\perp$ be defined as Definition E.2.
- Let $\mathbf{u}_r(t)$ be defined as Definition H.1.
- $\delta \in (0, 0.1)$.
- For an error $\epsilon > 0$ and $\|F(t) - y\|_2 \geq c \cdot \epsilon$ for a sufficient small constant $c > 0$.
- Define

$$C_3 := -\kappa \frac{1}{\sqrt{m}} \sum_{i=1}^n \sum_{r=1}^m a_r \left(\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \langle \mathbf{u}_r(t), x_i \rangle \right. \\ \left. - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t+1), x_i \rangle) \geq 0} \langle \mathbf{u}_r(t+1), x_i \rangle \right) \cdot (F_i(t) - y_i)$$

- $\kappa \in (0, 1]$.

Then with probability at least $1 - \delta$, we have:

$$C_3 \leq O\left(\eta \kappa \frac{R^2 n^{1.5} \sqrt{d}}{\delta \epsilon \sqrt{m}} D\right) \cdot \mathbf{L}(t)$$

Proof. We have:

$$\begin{aligned} & |\mathbf{u}_{r,k}(t) - \mathbf{u}_{r,k}(t+1)| \\ &= \left| \sqrt{V(w_r(t))} \cdot \tilde{w}_{r,k}(t) + E(w_r(t)) - w_{r,k}(t) \right. \\ &\quad \left. - \sqrt{V(w_r(t+1))} \cdot \tilde{w}_{r,k}(t+1) - E(w_r(t+1)) + w_{r,k}(t+1) \right| \\ &\leq \left| \tilde{w}_{r,k}(t) \sqrt{V(w_r(t))} - \tilde{w}_{r,k}(t+1) \sqrt{V(w_r(t+1))} \right| \\ &\quad + |\eta E(\Delta w_r(t))| + |\eta \Delta w_{r,k}(t)| \\ &\leq \left| \tilde{w}_{r,k}(t+1) (\sqrt{V(w_r(t))} - \sqrt{V(w_r(t+1))}) \right| \\ &\quad + \left| \sqrt{V(w_r(t))} (\tilde{w}_{r,k}(t) - \tilde{w}_{r,k}(t+1)) \right| + |\eta E(\Delta w_r(t))| + |\eta \Delta w_{r,k}(t)| \\ &= Q_{1,r,k} + Q_{2,r,k} + Q_{3,r,k} + Q_{4,r,k} \end{aligned} \tag{9}$$

where the first step follows from Definition H.1, the second step follows from triangle inequality and Definition C.8, the third step follows from simple algebras, the last step follows from defining:

$$\begin{aligned} Q_{1,r,k} &:= \left| \tilde{w}_{r,k}(t+1) (\sqrt{V(w_r(t))} - \sqrt{V(w_r(t+1))}) \right| \\ Q_{2,r,k} &:= \left| \sqrt{V(w_r(t))} (\tilde{w}_{r,k}(t) - \tilde{w}_{r,k}(t+1)) \right| \\ Q_{3,r,k} &:= |\eta E(\Delta w_r(t))| \\ Q_{4,r,k} &:= |\eta \Delta w_{r,k}(t)| \end{aligned}$$

Bounding $Q_{1,r,k}$.

We have:

$$\begin{aligned}
Q_{1,r,k} &= \left| \tilde{w}_{r,k}(t+1)(\sqrt{V(w_r(t))} - \sqrt{V(w_r(t+1))}) \right| \\
&= \left| (\sqrt{V(w_r(t))} - \sqrt{V(w_r(t+1))}) \right| \\
&\leq \|w_r(t) - E(w_r(t))\mathbf{1}_d - w_r(t+1) + E(w_r(t+1))\mathbf{1}_d\|_2 \\
&\leq \|\eta\Delta w_r(t)\|_2 + \sqrt{d} \cdot |\eta E(\Delta w_r(t))| \\
&\leq \eta \frac{(1 + \sqrt{d})\sqrt{n}}{\sqrt{m}} \|F(t) - y\|_2
\end{aligned}$$

where the first step follows from the definition of $Q_{1,r,k}$, the second step follows from $\tilde{w}_{r,k}(t+1) \in \{-1, +1\}$, the third step follows from Definition D.3 and reverse triangle inequality, the fourth step follows from $\|\mathbf{1}_d\|_2 = \sqrt{d}$ and Definition C.8, the last step follows from Lemma I.4.

Bounding $Q_{2,r,k}$.

We have:

$$\begin{aligned}
Q_{2,r,k} &= \left| \sqrt{V(w_r(t))}(\tilde{w}_{r,k}(t) - \tilde{w}_{r,k}(t+1)) \right| \\
&= |\sqrt{V(w_r(t))}| \cdot |\tilde{w}_{r,k}(t) - \tilde{w}_{r,k}(t+1)| \\
&\leq \|w_r(t) - E(w_r(t))\mathbf{1}_d\| \cdot |\tilde{w}_{r,k}(t) - \tilde{w}_{r,k}(t+1)| \\
&\leq O(\sqrt{d}D + R) \cdot |\tilde{w}_{r,k}(t) - \tilde{w}_{r,k}(t+1)| \tag{10}
\end{aligned}$$

where the first step follows from the definition of $Q_{2,r,k}$, the second step follows from simple algebras, the third step follows from Definition D.3, the last step follows from Part 2 of Lemma I.6.

At the same time, we can show that

$$\begin{aligned}
&\mathbb{E}[|\tilde{w}_{r,k}(t) - \tilde{w}_{r,k}(t+1)|] \\
&\leq 2(1 - \Pr[\mathbb{I}\{\mathbf{B}_{r,k}\} = 0 \cap \mathbb{I}\{|w_{r,k}(t) - E(w_r(t))| \geq |\eta\Delta w_{r,k}(t) - \eta E(\Delta w_r(t))|\}]) \\
&\leq 2(1 - \Pr[z \geq 2R + 2\eta \frac{\sqrt{n}}{\sqrt{m}} \|F(t) - y\|_2]) \\
&= 2\Pr[z \leq 2R + 2\eta \frac{\sqrt{n}}{\sqrt{m}} \|F(t) - y\|_2] \\
&\leq O(\eta \frac{\sqrt{n}}{\sqrt{m}}) \|F(t) - y\|_2 + O(1)R \\
&\leq O(\eta \frac{R\sqrt{n}}{\epsilon\sqrt{m}}) \|F(t) - y\|_2
\end{aligned}$$

where the first step follows from Definition E.3 and simple algebras, the second step follows from defining:

$$\begin{aligned}
z &:= w_{r,k}(0) - E(w_r(0)) \\
&= \frac{d-1}{d} w_{r,k} - \frac{1}{d} \sum_{k' \in [d] \setminus \{k\}} w_{r,k'}(0) \\
&\sim \mathcal{N}\left(0, \sigma^2 \sqrt{\frac{d-1}{d}} \cdot I_d\right)
\end{aligned}$$

and the last steps follow from the anti-concentration of the Gaussian variable (Fact B.3) and $\|F(t) - y\|_2 \geq \epsilon$ by Lemma condition.

Following Markov's inequality, we get:

$$|\tilde{w}_{r,k}(t) - \tilde{w}_{r,k}(t+1)| \leq O(\eta \frac{R\sqrt{n}}{\delta\epsilon\sqrt{m}}) \|F(t) - y\|_2 \tag{11}$$

Hence,

$$Q_{2,r,k} \leq O\left(\eta \frac{R^2 \sqrt{nd}}{\delta \epsilon \sqrt{m}} D\right) \|F(t) - y\|_2$$

where this step follows from Eq. (11) and Eq. (10).

Bounding $Q_{3,r,k}$ and $Q_{4,r,k}$.

We can show that $Q_{3,r,k} \leq \eta \frac{\sqrt{n}}{\sqrt{m}} \cdot \|F(t) - y\|_2$ and $Q_{4,r,k} \leq \eta \frac{\sqrt{n}}{\sqrt{m}} \cdot \|F(t) - y\|_2$ by following Lemma I.4.

Combination. We have:

$$\mathbb{E}[C_3] = 0$$

where this step follows from the symmetry of a .

Also

$$\begin{aligned} & \left(\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \langle \mathbf{u}_r(t), x_i \rangle - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t+1), x_i \rangle) \geq 0} \langle \mathbf{u}_r(t+1), x_i \rangle \right) \\ & \leq |\langle \mathbf{u}_r(t), x_i \rangle - \langle \mathbf{u}_r(t+1), x_i \rangle| \\ & = Q_{1,r,k} + Q_{2,r,k} + Q_{3,r,k} + Q_{4,r,k} \\ & \leq O\left(\eta \frac{R^2 \sqrt{nd}}{\delta \epsilon \sqrt{m}} D\right) \|F(t) - y\|_2 \end{aligned} \quad (12)$$

where the first step follows from ReLU is a 1-Lipschitz function (Fact B.2), the last step follows from simple algebras and the combination of these terms.

By Hoeffding's inequality (Lemma B.8), with a probability at least $1 - \delta$, we have:

$$\begin{aligned} |C_3| & \leq O\left(\eta \kappa \frac{R^2 n^{1.5} \sqrt{d}}{\delta \epsilon \cdot m} \sqrt{m} D\right) \|F(t) - y\|_2^2 \\ & \leq O\left(\eta \kappa \frac{R^2 n^{1.5} \sqrt{d}}{\delta \epsilon \sqrt{m}} D\right) \cdot \mathbb{L}(t) \end{aligned}$$

□

H.5 Bounding C_4

Lemma H.6. *If the following conditions hold:*

- Let $D > 0$ be defined as Definition B.16.
- For $i, j \in [n]$, $r \in [m]$ and integer $t \geq 0$.
- Let $H(t) \in \mathbb{R}^{n \times n}$ be defined as Definition G.1.
- Let $H^\perp(t) \in \mathbb{R}^{n \times n}$ be defined as Definition G.3.
- Let $H^* \in \mathbb{R}^{n \times n}$ be defined as Claim G.2. Assume $\lambda_{\min}(H^*) > 0$ as Assumption G.4.
- Let $\mathbb{L}(t)$ be defined as Definition C.9.
- Let $F(t) \in \mathbb{R}^n$ be defined as Definition C.9.
- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $\text{dq} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition D.5.

- Denote $\tilde{w}_r = \mathbf{q}(w_r) \in \{-1, +1\}^d$.
- Let $\mathcal{S}_i, \mathcal{S}_i^\perp$ be defined as Definition E.2.
- Let $\mathbf{u}_r(t)$ be defined as Definition H.1.
- $\delta \in (0, 0.1)$.
- For an error $\epsilon > 0$ and $\|\mathbf{F}(t) - y\|_2 \geq c \cdot \epsilon$ for a sufficient small constant $c > 0$.
- Define

$$C_4 := \frac{1}{2} \|\mathbf{F}(t) - \mathbf{F}(t+1)\|_2^2$$

Then with probability at least $1 - \delta$, we have:

$$|C_4| \leq O\left(\eta^2 \kappa^2 \frac{R^4 n^2 d}{\delta^2 \epsilon^2 m} D^2\right) \mathbf{L}(t)$$

Proof. We have:

$$\begin{aligned} & |\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0}(\langle w_r(t), x_i \rangle + \langle \mathbf{u}_r(t), x_i \rangle) \\ & \quad - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t+1), x_i \rangle) \geq 0}(\langle w_r(t+1), x_i \rangle + \langle \mathbf{u}_r(t+1), x_i \rangle)| \\ & \leq |\langle \eta \Delta w_r(t), x_i \rangle + \langle \mathbf{u}_r(t), x_i \rangle - \langle \mathbf{u}_r(t+1), x_i \rangle| \\ & \leq U_{1,i,r} + U_{2,i,r} \end{aligned}$$

where the first step follows from Fact B.2, the fifth step follows from Definition C.8, and the last step follows from defining:

$$\begin{aligned} U_{1,i,r} & := \langle \eta \Delta w_r(t), x_i \rangle \\ U_{2,i,r} & := \langle \mathbf{u}_r(t), x_i \rangle - \langle \mathbf{u}_r(t+1), x_i \rangle \end{aligned}$$

For the first term $U_{1,i,r}$, we have:

$$|U_{1,i,r}| \leq \eta \frac{\sqrt{n}}{\sqrt{m}} \|\mathbf{F}(t) - y\|_2$$

this step holds since Part 2 of Lemma I.4.

For the second term $U_{2,i,r}$, we have:

$$|U_{2,i,r}| \leq O\left(\eta \frac{R^2 \sqrt{nd}}{\delta \epsilon \sqrt{m}} D\right) \|\mathbf{F}(t) - y\|_2$$

this step follows from Eq. (12) and Eq. (9).

Thus, we have:

$$\begin{aligned} C_4 & = \frac{1}{2} \|\mathbf{F}(t) - \mathbf{F}(t+1)\|_2^2 \\ & = \frac{1}{2} \sum_{i=1}^n (\mathbf{F}_i(t) - \mathbf{F}_i(t+1))^2 \\ & = \frac{1}{2} \sum_{i=1}^n \left(\kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r (U_{1,i,r} + U_{2,i,r}) \right)^2 \end{aligned}$$

Combining two terms, then by Hoeffding inequality (Lemma B.8), with a probability at least $1 - \delta$, $\mathbb{E}[\sum_{r=1}^m a_r (U_{1,i,r} + U_{2,i,r})] = 1$, we have:

$$|C_4| \leq O\left(\eta^2 \kappa^2 \frac{R^4 n^2 d}{\delta^2 \epsilon^2 m} D^2\right) \|\mathbf{F}(t) - y\|_2^2 \leq O\left(\eta^2 \kappa^2 \frac{R^4 n^2 d}{\delta^2 \epsilon^2 m} D^2\right) \mathbf{L}(t)$$

□

I Inductions

I.1 Main Result 1: Training Convergence Guarantee

Theorem I.1. *If the following conditions hold:*

- Let $D > 0$ be defined as Definition B.16.
- Given a expected error $\epsilon > 0$.
- Let $H(t) \in \mathbb{R}^{n \times n}$ be defined as Definition G.1.
- Let $H^* \in \mathbb{R}^{n \times n}$ be defined as Claim G.2. Assume $\lambda_{\min}(H^*) > 0$ as Assumption G.4.
- Let $\mathsf{L}(t)$ be defined as Definition C.9.
- Let $\mathsf{F}(t) \in \mathbb{R}^n$ be defined as Definition C.9.
- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.
- $\delta \in (0, 0.1)$, $\kappa \in (0, 1]$.
- Choose $m \geq \Omega\left(\lambda^{-8} \frac{n^{12} d^8}{\delta^4 \epsilon^4}\right)$.
- Choose $\eta \leq O\left(\lambda \frac{\delta}{\kappa^2 n^2 d D}\right)$.
- Choose $T \geq \Omega\left(\frac{1}{\eta \lambda} \log(\epsilon^{-1} n d D^2)\right)$.

Then with probability at least $1 - \delta$, we have:

$$\mathsf{L}(T) \leq \epsilon$$

Proof. **Choice of m .**

Following Lemma I.2, we have

$$m \geq \Omega\left(\lambda^{-4} \kappa^4 \frac{R^8 n^6 d^2}{\delta^4 \epsilon^4}\right)$$

Particularly, following Claim I.5, we have:

$$\begin{aligned} R &\leq \frac{4\sqrt{n}}{\lambda\sqrt{m}} \|F(0) - y\|_2 \\ &\leq \frac{4\sqrt{n}}{\lambda\sqrt{m}} \cdot O\left(\sqrt{nd}D^2\right) \\ &\leq O\left(\frac{nd}{\lambda\sqrt{m}}D^2\right) \end{aligned}$$

where the first step follows from Claim I.5, the second step follows from Lemma I.3, the third step follows from simple algebras.

Besides, by Lemma I.2, we need that

$$R \leq O\left(\frac{\lambda\delta}{\kappa^2 n^2 d D}\right)$$

where the second step follows from Definition B.16.

Thus, showing that $D^3 \leq O(m^{\frac{1}{4}})$ and $\kappa \leq 1$, we plug m as follows:

$$m \geq \Omega\left(\lambda^{-8} \frac{n^{12} d^8}{\delta^4 \epsilon^4}\right)$$

Choice of η . We have

$$\begin{aligned}\|\eta\Delta w_r(0)\|_2 &\leq \eta \frac{\sqrt{n}}{\sqrt{m}} \|F(0) - y\|_2 \\ &\leq \eta \frac{\sqrt{n}}{\sqrt{m}} O(\sqrt{nd}D^2) \\ &\leq R\end{aligned}$$

where the first step follows from Part 2 of Lemma I.4, the second step follows from Lemma I.3, the third step follows from plugging $\eta \leq O\left(\lambda \frac{\delta}{\kappa n^2 d D}\right)$ and $m \geq \Omega\left(\lambda^{-8} \frac{n^{12} d^8}{\delta^4 \epsilon^4}\right)$.

Choice of T . We have:

$$\begin{aligned}\mathsf{L}(T) \leq \epsilon &\iff (1 - \eta\lambda/2)^T \mathsf{L}(0) \leq \epsilon \\ &\iff (1 - \eta\lambda/2)^T O(\sqrt{nd}D^2) \leq \epsilon \\ &\iff (1 - \eta\lambda/2)^T \leq O\left(\frac{\epsilon}{\sqrt{nd}D^2}\right) \\ &\iff T \geq \Omega\left(\log\left(\frac{\epsilon}{\sqrt{nd}D^2}\right) / \log(1 - \eta\lambda/2)\right) \\ &\iff T \geq \Omega\left(-\frac{1}{\eta\lambda} \log\left(\frac{\epsilon}{\sqrt{nd}D^2}\right)\right) \\ &\iff T \geq \Omega\left(\frac{1}{\eta\lambda} \log(\epsilon^{-1} ndD^2)\right)\end{aligned}$$

where the first step follows from Lemma I.2, the second step follows from Lemma I.3, the third and fourth steps follow from simple algebras, the fifth step follows from Fact B.6, the sixth step follows from simple algebras. \square

I.2 Induction for Loss

Lemma I.2. *If the following conditions hold:*

- Let $D > 0$ be defined as Definition B.16.
- For $i, j \in [n]$, $r \in [m]$ and integer $t \geq 0$.
- Let $H(t) \in \mathbb{R}^{n \times n}$ be defined as Definition G.1.
- Let $H^\perp(t) \in \mathbb{R}^{n \times n}$ be defined as Definition G.3.
- Let $H^* \in \mathbb{R}^{n \times n}$ be defined as Claim G.2. Assume $\lambda_{\min}(H^*) > 0$ as Assumption G.4.
- Let $\mathsf{L}(t)$ be defined as Definition C.9.
- Let $F(t) \in \mathbb{R}^n$ be defined as Definition C.9.
- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $\text{dq} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition D.5.
- Denote $\tilde{w}_r = \mathfrak{q}(w_r) \in \{-1, +1\}^d$.
- Let $\mathcal{S}_i, \mathcal{S}_i^\perp$ be defined as Definition E.2.
- Let $u_r(t)$ be defined as Definition H.1.

- $\delta \in (0, 0.1)$.
- For an error $\epsilon > 0$ and $\|\mathbf{F}(t) - y\|_2 \geq c \cdot \epsilon$ for a sufficient small constant $c > 0$.
- $m \geq \Omega\left(\lambda^{-4} \kappa^4 \frac{R^8 n^6 d^2}{\delta^4 \epsilon^4}\right)$.
- $R \leq O\left(\frac{\lambda \delta}{\kappa^2 n^2 d D}\right)$.
- Define

$$C_1 := -\kappa \frac{1}{\sqrt{m}} \sum_{i=1}^n \sum_{r \in \mathcal{S}_i} a_r \left(\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \langle w_r(t), x_i \rangle \right. \\ \left. - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t+1), x_i \rangle) \geq 0} \langle w_r(t+1), x_i \rangle \right) \cdot (\mathbf{F}_i(t) - y_i)$$

- Define

$$C_2 := -\kappa \frac{1}{\sqrt{m}} \sum_{i=1}^n \sum_{r \in \mathcal{S}_i^\perp} a_r \left(\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \langle w_r(t), x_i \rangle \right. \\ \left. - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t+1), x_i \rangle) \geq 0} \langle w_r(t+1), x_i \rangle \right) \cdot (\mathbf{F}_i(t) - y_i)$$

- Define

$$C_3 := -\kappa \frac{1}{\sqrt{m}} \sum_{i=1}^n \sum_{r=1}^m a_r \left(\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \langle u_r(t), x_i \rangle \right. \\ \left. - \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t+1), x_i \rangle) \geq 0} \langle u_r(t+1), x_i \rangle \right) \cdot (\mathbf{F}_i(t) - y_i)$$

- Define

$$C_4 := \frac{1}{2} \|\mathbf{F}(t) - \mathbf{F}(t+1)\|_2^2$$

- $\delta \in (0, 1]$.

Then with probability at least $1 - \delta$, we have:

$$\mathbf{L}(t+1) \leq (1 - \lambda/2\eta) \cdot \mathbf{L}(t)$$

Moreover, we can show that:

$$\mathbf{L}(t) \leq (1 - \lambda/2\eta)^t \cdot \mathbf{L}(0)$$

Proof. We have:

$$\begin{aligned} \mathbf{L}(t+1) &\leq \mathbf{L}(t) + \left(-\eta\lambda + O\left(\eta \frac{n^2 d R D}{\delta}\right) + O\left(\eta \kappa \frac{n^{1.5} d R D}{\delta}\right) \right. \\ &\quad \left. + O\left(\eta \kappa \frac{R^2 n^{1.5} \sqrt{d}}{\delta \epsilon \sqrt{m}} D\right) + O\left(\eta^2 \kappa^2 \frac{R^4 n^2 d}{\delta^2 \epsilon^2 m} D^2\right) \right) \cdot \mathbf{L}(t) \\ &\leq \mathbf{L}(t) + \left(-\eta\lambda + \frac{1}{8}\eta\lambda + \frac{1}{8}\eta\lambda + \frac{1}{8}\eta\lambda + \frac{1}{8}\eta\lambda \right) \cdot \mathbf{L}(t) \\ &\leq (1 - \eta\lambda/2) \mathbf{L}(t) \end{aligned}$$

where the first step follows from Claim H.2, Lemma H.3, Lemma H.4, Lemma H.5, Lemma H.6 and $\eta\lambda \leq 1$, the second step follows from the choice of R and m , the last step follows from simple algebras.

Choice of R . We have:

$$R \leq O\left(\frac{\lambda\delta}{\kappa^2 n^2 d D}\right) \quad (13)$$

where this step is following the combination of Lemma G.5 and $O(\eta \frac{\kappa^2 n^2 d R D}{\delta} \leq \frac{1}{8} \eta \lambda)$.

Choice of m . We have:

$$\begin{aligned} \sqrt{m} &\geq \Omega\left(\lambda^{-1} \kappa \frac{R^2 n^{1.5} d^{0.5}}{\delta \epsilon} D\right) \\ \iff \sqrt{m} &\geq \Omega\left(\lambda^{-1} \kappa \frac{R^2 n^{1.5} d^{0.5}}{\delta \epsilon} m^{\frac{1}{4}}\right) \\ \iff m^{\frac{1}{4}} &\geq \Omega\left(\lambda^{-1} \kappa \frac{R^2 n^{1.5} d^{0.5}}{\delta \epsilon}\right) \\ \iff m &\geq \Omega\left(\lambda^{-4} \kappa^4 \frac{R^8 n^6 d^2}{\delta^4 \epsilon^4}\right) \end{aligned}$$

where the first step follows from plugging $O(\eta \kappa \frac{R^2 n^{1.5} \sqrt{d}}{\delta \epsilon \sqrt{m}} D) \leq \frac{1}{8} \eta \lambda$, the last three steps follow from simple algebras. \square

Lemma I.3. *If the following conditions hold:*

- Let $D > 0$ be defined as Definition B.16.
- For $i, j \in [n]$, $r \in [m]$ and integer $t \geq 0$.
- Let $\mathbf{L}(t)$ be defined as Definition C.9.
- Let $\mathbf{F}(t) \in \mathbb{R}^n$ be defined as Definition C.9.
- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $\text{dq} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition D.5.
- Denote $\tilde{w}_r = \mathbf{q}(w_r) \in \{-1, +1\}^d$.
- Let $\mathcal{S}_i, \mathcal{S}_i^\perp$ be defined as Definition E.2.
- Let $\mathbf{u}_r(t)$ be defined as Definition H.1.
- For an error $\epsilon > 0$ and $\|\mathbf{F}(t) - y\|_2 \geq c \cdot \epsilon$ for a sufficient small constant $c > 0$.

Then with probability at least $1 - \delta$, we have:

$$\|\mathbf{F}(0) - y\|_2 \leq O\left(\sqrt{n} d D^2\right)$$

Proof. We have:

$$\begin{aligned} \|\mathbf{F}(0) - y\|_2 &\leq \|\mathbf{F}(0)\|_2 + \|y\|_2 \\ &\leq \|\mathbf{F}(0)\|_2 + \sqrt{n} \\ &\leq \left(\sum_{i=1}^n |\mathbf{F}_i(0)|^2\right)^{\frac{1}{2}} + \sqrt{n} \\ &\leq \left(\sum_{i=1}^n \left|\kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \text{ReLU}\left(\text{dq}(\langle \tilde{w}_r(0), x_i \rangle)\right)\right|^2\right)^{\frac{1}{2}} + \sqrt{n} \\ &\leq O\left(\sqrt{n \log(m/\delta)} d D\right) + \sqrt{n} \end{aligned}$$

$$\leq O(\sqrt{nd}D^2)$$

where the first step follows from triangle inequality, the second step follows from $y_i \leq 1, \forall i \in [n]$ and simple algebras, the third step follows from the definition of ℓ_2 norm, the fourth step follows from Definition C.9 and Definition C.5, the last two steps follow by Hoeffding's inequality (Lemma B.8), Definition C.1 and simple algebras, and we can show that:

$$\mathbb{E}\left[\sum_{r=1}^m a_r \cdot \text{ReLU}\left(\text{dq}(\langle \tilde{w}_r(0), x_i \rangle)\right)\right] = 0$$

also,

$$\begin{aligned} \text{dq}(\langle \tilde{w}_r(0), x_i \rangle) &= \sqrt{V(w_r(0))} \cdot \langle \tilde{w}_r(0), x_i \rangle + E(w_r(0)) \langle \mathbf{1}_d, x_i \rangle \\ &\leq O(\sqrt{d}D) \cdot \sqrt{d} + O(D) \cdot \sqrt{d} \\ &\leq O(dD) \end{aligned}$$

where these steps follow from Definition D.5, Lemma I.6 and simple algebras. \square

I.3 Induction for STE Gradient

Lemma I.4. *If the following conditions hold:*

- For $i, j \in [n], r \in [m]$ and integer $t \geq 0$.
- Let $\mathbf{L}(t)$ be defined as Definition C.9.
- Let $\mathbf{F}(t) \in \mathbb{R}^n$ be defined as Definition C.9.
- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $\text{dq} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition D.5.
- Denote $\tilde{w}_r = \mathbf{q}(w_r) \in \{-1, +1\}^d$.
- Let $\mathcal{S}_i, \mathcal{S}_i^\perp$ be defined as Definition E.2.
- Let $\mathbf{u}_r(t)$ be defined as Definition H.1.
- For an error $\epsilon > 0$ and $\|\mathbf{F}(t) - y\|_2 \geq c \cdot \epsilon$ for a sufficient small constant $c > 0$.

Then with probability at least $1 - \delta$, we have:

- Part 1. $\forall k \in [d]$

$$|\Delta w_{r,k}(t)| \leq \sqrt{\frac{n}{m}} \cdot \|\mathbf{F}(t) - y\|_2$$

- Part 2.

$$\|\Delta w_r(t)\|_2 \leq \sqrt{\frac{n}{m}} \cdot \|\mathbf{F}(t) - y\|_2$$

Proof. **Proof of Part 1.** We have:

$$|\Delta w_{r,k}(t)| = \left| \kappa \frac{1}{\sqrt{m}} \sum_{i=1}^n a_r \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \cdot x_{i,k} \cdot (\mathbf{F}_i(t) - y_i) \right|$$

$$\begin{aligned}
&\leq \kappa \frac{1}{\sqrt{m}} \left(\sum_{i=1}^n (a_r \cdot \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \cdot x_{i,k})^2 \right)^{\frac{1}{2}} \cdot \|F(t) - y\|_2 \\
&\leq \sqrt{\frac{n}{m}} \cdot \|F(t) - y\|_2
\end{aligned}$$

where the first step follows from Definition F.2, the second step follows from Cauchy-Schwarz inequality, the third step follows from

$$\max_{r \in [m], i \in [n], k \in [d]} |\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_i \rangle) \geq 0} \cdot x_{i,k}| \leq 1$$

the above equation follows from simple algebras and $\|x_i\|_i = 1$.

Proof of Part 2.

By $\|x\|_i = 1, \forall i \in [n]$, this proof is trivially the same as **Proof of Part 1**. \square

I.4 Induction for Weights

Claim I.5. *If the following conditions hold:*

- For $i, j \in [n], r \in [m]$ and integer $t \geq 0$.
- Let $L(t)$ be defined as Definition C.9.
- Let $F(t) \in \mathbb{R}^n$ be defined as Definition C.9.
- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $\text{dq} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition D.5.
- Denote $\tilde{w}_r = \mathbf{q}(w_r) \in \{-1, +1\}^d$.
- Let $\mathcal{S}_i, \mathcal{S}_i^\perp$ be defined as Definition E.2.
- Let $u_r(t)$ be defined as Definition H.1.
- For an error $\epsilon > 0$ and $\|F(t) - y\|_2 \geq c \cdot \epsilon$ for a sufficient small constant $c > 0$.

Then with probability at least $1 - \delta$, we have:

$$R := \max_{t \geq 0} \max_{r \in [m]} \|w_r(0) - w_r(t)\|_2 \leq \frac{4\sqrt{n}}{\lambda\sqrt{m}} \|F(0) - y\|_2$$

Proof. We have

$$\begin{aligned}
R &= \max_{t \geq 0} \max_{r \in [m]} \|w_r(0) - w_r(t)\|_2 \\
&\leq \max_{t \geq 0} \max_{r \in [m]} \left\| \sum_{\tau=1}^t \eta \Delta w_r(\tau) \right\|_2 \\
&\leq \eta \max_{t \geq 0} \max_{r \in [m]} \sum_{\tau=1}^t \|\Delta w_r(\tau)\|_2 \\
&\leq \eta \frac{\sqrt{n}}{\sqrt{m}} \max_{t \geq 0} \sum_{\tau=1}^t \|F(\tau) - y\|_2 \\
&\leq \eta \frac{\sqrt{n}}{\sqrt{m}} \max_{t \geq 0} \sum_{\tau=1}^t (1 - \eta\lambda/2)^\tau \|F(0) - y\|_2
\end{aligned}$$

$$\leq \frac{4\sqrt{n}}{\lambda\sqrt{m}} \|F(0) - y\|_2$$

where the first step follows from the definition of R , the second step follows from Definition C.8, the third step follows from triangle inequality, the fourth step follows from Part 2 of Lemma I.4, the fifth step follows from Lemma I.2, the last step follows from Fact B.6. \square

Lemma I.6. *Let $\delta \in (0, 0.1)$. Let $D > 0$ be defined as Definition B.16. Let $E : \mathbb{R}^d \rightarrow \mathbb{R}$ be defined as Definition D.2. Let $V : \mathbb{R}^d \rightarrow \mathbb{R}$ be defined as Definition D.3. Let $W(0) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3, denote $W := [w_1, w_2, \dots, w_m] \in \mathbb{R}^{d \times m}$ satisfying $\|w_r - w_r(0)\|_2 \leq R$ where $R \geq 0$, then with a probability at least $1 - \delta$, we have*

- Part 1. $|w_{r,k}(0)| \leq O(D), \forall r \in [m], k \in [d]$.
- Part 2. $\|w_r(0)\|_2 \leq O(\sqrt{d}D), \forall r \in [m]$.
- Part 3. $\|w_r\|_2 \leq O(\sqrt{d}D + R), \forall r \in [m]$.
- Part 4. $E(w_r(0)) \leq O(D), \forall r \in [m]$.
- Part 5. $\sqrt{V(w_r(0))} \leq O(D), \forall r \in [m]$.
- Part 6. $E(w_r) \leq O(D + R), \forall r \in [m]$.
- Part 7. $\sqrt{V(w_r)} \leq O(D + R), \forall r \in [m]$.

Proof. This proof follows from the union bound of the Gaussian tail bound (Fact B.1) and some simple algebras. \square

J Supplementary Setup for Classical Linear Regression

J.1 Model Function

Definition J.1. *If the following conditions hold:*

- For a input vector $x \in \mathbb{R}^d$.
- For a hidden-layer weights $W \in \mathbb{R}^{d \times m}$ as Definition C.2.
- For a output-layer weights $a \in \mathbb{R}^m$ as Definition C.2.
- Let $\text{ReLU} : \mathbb{R} \rightarrow \mathbb{R}$ be defined as Definition C.4.
- Let $D = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- $t \geq 0$, let $W(0) \in \mathbb{R}^{d \times m}$ and $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- $W'(0) := W(0)$.
- Let $W'(t) \in \mathbb{R}^{d \times m}$ be updated as Claim J.3.
- $\kappa \in (0, 1]$.

We define:

$$f'(x, W, a) := \kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \text{ReLU}(\langle w_r, x \rangle) \in \mathbb{R}$$

Then we define the compact form of $f(x, W'(t), a)$, we define:

$$F'(t) = [f(x_1, W'(t), a), f(x_2, W'(t), a), \dots, f(x_n, W'(t), a)]^\top \in \mathbb{R}^n$$

J.2 Loss and Training

Definition J.2. *If the following conditions hold:*

- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(0) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $f' : \mathbb{R}^d \times \mathbb{R}^{d \times m} \times \mathbb{R}^m \rightarrow \mathbb{R}$ be defined as Definition J.1.
- For any $t \geq 0$.

We define:

$$\mathsf{L}'(t) := \frac{1}{2} \|\mathsf{F}'(t) - y\|_2^2$$

Claim J.3. *If the following conditions hold:*

- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(0) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3.
- Let $f' : \mathbb{R}^d \times \mathbb{R}^{d \times m} \times \mathbb{R}^m \rightarrow \mathbb{R}$ be defined as Definition J.1.
- Let $\mathsf{L}'(t)$ be defined as Definition J.2.
- For any $t \geq 0$.
- Denote $\eta > 0$ as the learning rate.

We define:

$$W'(t+1) := W'(t) - \eta \cdot \Delta W'(t)$$

Here, we also define that:

$$\begin{aligned} W'(t) &:= \frac{d}{dW'(t)} \mathsf{L}'(t) \\ &= \sum_{i=1}^n (\mathsf{F}'_i(t) - y_i) \cdot \kappa [a_1 \cdot \mathbf{1}_{\langle w'_1(t), x_i \rangle \geq 0} x_i \quad \cdots \quad a_m \cdot \mathbf{1}_{\langle w'_m(t), x_i \rangle \geq 0} x_i] \in \mathbb{R}^{d \times m} \end{aligned}$$

Proof. This proof follows from simple algebras. □

J.3 Induction for Weights

Lemma J.4 (See Corollary 4.1 and the fifth equation of page 6 in Du et al. [41]). *If the following conditions hold:*

- $t \geq 0$, let $W(0) \in \mathbb{R}^{d \times m}$ and $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- $W'(0) := W(0)$.
- Let $W'(t) \in \mathbb{R}^{d \times m}$ be updated as Claim J.3.
- $R \leq O(\frac{\lambda \delta}{\kappa^2 n^2 d D})$.

Then we have

$$\|w'_r(t) - w'_r(0)\| \leq R$$

Proof. Following Corollary 4.1 in Du et al. [41], we can show that:

$$\|w'_r(t) - w'_r(0)\| \leq \frac{4\sqrt{n}}{\sqrt{m}\lambda} \|F'(0) - y\|_2$$

Then we can complete this proof by combining the equation above with Lemma J.5 and $R \leq O(\frac{\lambda\delta}{n^2 d D})$ in Lemma conditions. \square

J.4 Induction for Loss

Lemma J.5. *If the following conditions hold:*

- Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition C.1.
- Let $W(0) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3.
- Let $a \in \mathbb{R}^m$ be initialized as Definition C.3.
- Let $f' : \mathbb{R}^d \times \mathbb{R}^{d \times m} \times \mathbb{R}^m \rightarrow \mathbb{R}$ be defined as Definition J.1.
- For any $t \geq 0$.
- $W'(0) := W(0)$.
- Let $W'(t) \in \mathbb{R}^{d \times m}$ be updated as Claim J.3.
- $\delta \in (0, 0.1)$.

Then with probability at least $1 - \delta$, we have:

$$\|F'(0) - y\|_2 \leq O(\sqrt{nd}D^2)$$

Proof. We have:

$$\begin{aligned} \|F'(0) - y\|_2 &\leq \|F'(0)\|_2 + \|y\|_2 \\ &\leq \|F'(0)\|_2 + \sqrt{n} \\ &\leq \left(\sum_{i=1}^n |F'_i(0)|^2\right)^{\frac{1}{2}} + \sqrt{n} \\ &\leq \left(\sum_{i=1}^n \left|\kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \text{ReLU}(\langle w'_r(0), x_i \rangle)\right|^2\right)^{\frac{1}{2}} + \sqrt{n} \\ &= \left(\sum_{i=1}^n \left|\kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \text{ReLU}(\langle w_r(0), x_i \rangle)\right|^2\right)^{\frac{1}{2}} + \sqrt{n} \\ &\leq O\left(\sqrt{n \log(m/\delta) d D}\right) + \sqrt{n} \\ &\leq O(\sqrt{nd}D^2) \end{aligned}$$

where the first step follows from triangle inequality, the second step follows from $y_i \leq 1, \forall i \in [n]$ and simple algebras, the third step follows from the definition of ℓ_2 norm, the fourth step follows from Definition C.9 and Definition C.5, the fifth step follows from $W'(0) = W(0)$, the last two steps follow by Hoeffding's inequality (Lemma B.8), Definition C.1, $\kappa \leq 1$ and simple algebras, and we can show that:

$$\mathbb{E}\left[\sum_{r=1}^m a_r \cdot \text{ReLU}(\langle w_r(0), x_i \rangle)\right] = 0$$

also,

$$\begin{aligned} \langle w_r(0), x_i \rangle &= \langle w_r(0), x_i \rangle \\ &\leq O(\sqrt{d}D) \leq O(dD) \end{aligned}$$

where this step follows from Lemma I.6 and simple algebras. \square

K Similarities

K.1 Main Result 2: Training Similarity

Theorem K.1. *If the following conditions hold:*

- *Let $D > 0$ be defined as Definition B.16.*
- *Given a expected error $\epsilon > 0$.*
- *Let $H^* \in \mathbb{R}^{n \times n}$ be defined as Claim G.2. Assume $\lambda_{\min}(H^*) > 0$ as Assumption G.4.*
- *Let $\mathcal{D}_{\text{test}} := \{(x_{\text{test},i}, y_{\text{test},i})\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition K.2.*
- *Let $F'(t) \in \mathbb{R}^n$ be defined as Definition J.1.*
- *Let $F(t) \in \mathbb{R}^n$ be defined as Definition C.9.*
- *Let $F'_{\text{test}}(t) \in \mathbb{R}^n$ be defined as Definition K.3.*
- *Let $F_{\text{test}}(t) \in \mathbb{R}^n$ be defined as Definition K.3.*
- *For any $t \geq 0$.*
- *Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.*
- *$W'(0) := W(0)$.*
- *Let $W'(t) \in \mathbb{R}^{d \times m}$ be updated as Claim J.3.*
- *For any error $\epsilon_{\text{quant}} > 0$.*
- *$\delta \in (0, 0.1)$.*
- *Choose $\kappa \leq O(\frac{\epsilon_{\text{quant}}}{dD^2})$.*

Then with probability at least $1 - \delta$, we have:

- *Part 1. $|F_{\text{test},i}(t) - F'_{\text{test},i}(t)| \leq \epsilon_{\text{quant}}$.*
- *Part 2. $|F_i(t) - F'_i(t)| \leq \epsilon_{\text{quant}}$.*

Proof. Proof of Part 1. We have:

$$\begin{aligned}
& |\mathbf{1}_{\langle \tilde{w}_r(t), x_{\text{test},i} \rangle \geq 0} \langle w_r(t), x_{\text{test},i} \rangle + \langle \mathbf{u}_r(t), x_{\text{test},i} \rangle \\
& \quad - \mathbf{1}_{\langle w'_r(t), x_{\text{test},i} \rangle \geq 0} \langle w'_r(t), x_{\text{test},i} \rangle| \\
& \leq |\langle w_r(t), x_{\text{test},i} \rangle + \langle \mathbf{u}_r(t), x_{\text{test},i} \rangle - \langle w'_r(t), x_{\text{test},i} \rangle| \\
& = |\langle w_r(0) - \eta \sum_{\tau=0}^{t-1} \Delta w_r(\tau), x_{\text{test},i} \rangle + \langle \mathbf{u}_r(t), x_{\text{test},i} \rangle - \langle w'_r(0) - \eta \sum_{\tau=0}^{t-1} \Delta w'_r(\tau), x_{\text{test},i} \rangle| \\
& = | - \langle \eta \sum_{\tau=0}^{t-1} \Delta w_r(\tau), x_{\text{test},i} \rangle + \langle \mathbf{u}_r(t), x_{\text{test},i} \rangle + \langle \eta \sum_{\tau=0}^{t-1} \Delta w'_r(\tau), x_{\text{test},i} \rangle | \\
& \leq |\langle \eta \sum_{\tau=0}^{t-1} \Delta w_r(\tau), x_{\text{test},i} \rangle| + |\langle \eta \sum_{\tau=0}^{t-1} \Delta w'_r(\tau), x_{\text{test},i} \rangle| + |\langle \mathbf{u}_r(t), x_{\text{test},i} \rangle| \\
& \leq R + R + |\langle \mathbf{u}_r(t), x_{\text{test},i} \rangle| \\
& \leq O(d(D + R))
\end{aligned}$$

where the first step follows from Fact B.2, the second step follows from Definition C.8 and Claim J.3, the third step follows from $w'_r(0) = w_r(0)$, the fourth step follows from triangle inequality, the fifth step follows from Claim I.5 and Lemma J.4, the last step follows from Lemma D.7 and $\delta \in (0, 0.1)$.

Then we have:

$$\begin{aligned} |F_{\text{test},i}(t) - F'_{\text{test},i}(t)| &\leq \left| \kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \left(\mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_{\text{test},i} \rangle) \geq 0} (\langle w_r(t), x_{\text{test},i} \rangle + \langle \mathbf{u}_r(t), x_{\text{test},i} \rangle) \right. \right. \\ &\quad \left. \left. - \mathbf{1}_{\langle w'_r(t), x_{\text{test},i} \rangle \geq 0} \langle w'_r(t), x_{\text{test},i} \rangle \right) \right| \\ &\leq \kappa \sqrt{\log(m/\delta)} \cdot O(d(D+R)) \\ &\leq \epsilon_{\text{quant}} \end{aligned}$$

where the first step follows from Definition K.3, the second step follows from Hoeffding's inequality (Lemma B.8), $\mathbb{E}[\sum_{r=1}^m a_r \sigma_{i,r}] = 0$, $\sigma_{i,r} \leq O\left(\frac{\sqrt{n}}{m}(D+R) + R/\delta\right)$ and defining:

$$\begin{aligned} \sigma_{i,r} &:= \left| \mathbf{1}_{\text{dq}(\langle \tilde{w}_r(t), x_{\text{test},i} \rangle) \geq 0} (\langle w_r(t), x_{\text{test},i} \rangle + \langle \mathbf{u}_r(t), x_{\text{test},i} \rangle) \right. \\ &\quad \left. - \mathbf{1}_{\langle w'_r(t), x_{\text{test},i} \rangle \geq 0} \langle w'_r(t), x_{\text{test},i} \rangle \right| \end{aligned}$$

and the last step follows from choosing

$$\kappa \leq O\left(\frac{\epsilon_{\text{quant}}}{dD^2 + dDR}\right) \leq O\left(\frac{\epsilon_{\text{quant}}}{dD^2}\right)$$

Proof of Part 2. This part can be proved in the same way as **Proof of Part 1**. □

K.2 Test Dataset for Generalization Evaluation

Definition K.2. We define test dataset $\mathcal{D}_{\text{test}} := \{(x_{\text{test},i}, y_{\text{test},i})\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$, where $\|x_{\text{test},i}\|_2 = 1$ and $y_{\text{test},i} \leq 1$ for any $i \in [n]$.

Definition K.3. If the following conditions hold:

- Let $\mathcal{D}_{\text{test}} := \{(x_{\text{test},i}, y_{\text{test},i})\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ be defined as Definition K.2.
- Let $f' : \mathbb{R}^d \times \mathbb{R}^{d \times m} \times \mathbb{R}^m \rightarrow \mathbb{R}$ be defined as Definition J.1.
- Let $f : \mathbb{R}^d \times \mathbb{R}^{d \times m} \times \mathbb{R}^m \rightarrow \mathbb{R}$ be defined as Definition C.5.
- For any $t \geq 0$.
- Let $W(t) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3 and be updated by Definition C.8.
- $W'(0) := W(0)$.
- Let $W'(t) \in \mathbb{R}^{d \times m}$ be updated as Claim J.3.

We define:

$$\begin{aligned} F'_{\text{test}}(t) &:= [f'(x_{\text{test},1}, W'(t), a), f'(x_{\text{test},2}, W'(t), a), \dots, f'(x_{\text{test},n}, W'(t), a)]^\top \\ F_{\text{test}}(t) &:= [f(x_{\text{test},1}, W(t), a), f(x_{\text{test},2}, W(t), a), \dots, f(x_{\text{test},n}, W(t), a)]^\top \end{aligned}$$

K.3 Function Similarity at Initialization

Lemma K.4. If the following conditions hold:

- Let $D > 0$ be defined as Definition B.16.
- Let $\mathbf{q} : \mathbb{R}^d \rightarrow \{-1, +1\}^d$ be defined as Definition D.4.

- Let $E : \mathbb{R}^d \rightarrow \mathbb{R}$ be defined as Definition D.2.
- Let $V : \mathbb{R}^d \rightarrow \mathbb{R}$ be defined as Definition D.3.
- For a weight vector $w \in \mathbb{R}^d$.
- Denote quantized vector $\tilde{w} := \mathbf{q}(w) \in \{-1, +1\}^d$.
- For a vector $x \in \mathbb{R}^d$ and $\|x\|_2 = 1$.
- Let $f' : \mathbb{R}^d \times \mathbb{R}^{d \times m} \times \mathbb{R}^m \rightarrow \mathbb{R}$ be defined as Definition J.1.
- Let $f : \mathbb{R}^d \times \mathbb{R}^{d \times m} \times \mathbb{R}^m \rightarrow \mathbb{R}$ be defined as Definition C.5.
- Let $W(0) \in \mathbb{R}^{d \times m}$ be initialized as Definition C.3.
- $W'(0) := W(0)$.
- $\delta \in (0, 0.1)$.
- For any error $\epsilon_{\text{init}} > 0$.
- We choose $\kappa \leq O(\epsilon_{\text{init}}/(\sqrt{d}D^2))$

Then with probability at least $1 - \delta$, we have:

$$|f(x, W(0), a) - f'(x, W'(0), a)| \leq \epsilon_{\text{init}}$$

Proof. We have:

$$\begin{aligned} & |\mathbf{1}_{\mathbf{dq}(\langle \tilde{w}_r(0), x \rangle) \geq 0} \mathbf{dq}(\langle \tilde{w}_r(0), x \rangle) \\ & \quad - \mathbf{1}_{\langle w_r(0), x \rangle \geq 0} \langle w_r(0), x \rangle| \\ & \leq |\mathbf{dq}(\langle \tilde{w}_r(0), x \rangle) - \langle w_r(0), x \rangle| \\ & \leq |\sqrt{V(w_r(0))} \langle \tilde{w}_r(0), x \rangle + E(w_r(0)) \cdot \langle \mathbf{1}_d, x \rangle - \langle w_r(0), x \rangle| \\ & \leq O(\sqrt{d}D) \end{aligned}$$

where the first step follows from Fact B.2, the second step follows from Definition D.5, the last step follows from Lemma I.6.

Then by Hoeffding inequality (Lemma B.8), with a probability at least $1 - \delta$, we have:

$$\begin{aligned} |f(x, W(0), a) - f'(x, W'(0), a)| & \leq \kappa \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \hat{\sigma}_r \\ & \leq \kappa O(\sqrt{d}D) \cdot \sqrt{\log(m/\delta)} \\ & \leq O(\kappa \sqrt{d}D^2) \end{aligned}$$

where we have:

$$\begin{aligned} \hat{\sigma}_r & := \mathbf{1}_{\mathbf{dq}(\langle \tilde{w}_r(0), x \rangle) \geq 0} \mathbf{dq}(\langle \tilde{w}_r(0), x \rangle) - \mathbf{1}_{\langle w_r(0), x \rangle \geq 0} \langle w_r(0), x \rangle \\ \mathbb{E}[\sum_{r=1}^m a_r \hat{\sigma}_r] & = 1 \\ |\hat{\sigma}_r| & \leq O(\sqrt{d}D) \end{aligned}$$

□