
Theoretical Analysis of the Selection Mechanism in Mamba: Training Dynamics and Generalization

Mugunthan Shandirasegaran

New Jersey Institute of Technology
ms3537@njit.edu

Yating Zhou

Rensselaer Polytechnic Institute
zhouy26@rpi.edu

Songyang Zhang

University of Louisiana at Lafayette
songyang.zhang@louisiana.edu

Shuai Zhang

New Jersey Institute of Technology
sz457@njit.edu

Abstract

While Transformers rely on a distinctive attention mechanism, the recent emergence of Mamba and other selective state space models (SSMs) offers a strong alternative. These models incorporate attention-like mechanisms with hardware-aware efficiency and a unique selection strategy, yet their theoretical properties remain poorly understood. In this work, we present a first-step theoretical analysis of the *selection mechanism in Mamba*. We study a simplified single-layer Mamba block trained with gradient descent on structured data containing both label-relevant and irrelevant tokens. Our results show that the gating vector dynamically aligns with label-relevant features while negating irrelevant ones, formalizing its role as an implicit feature selector. Moreover, we prove that training achieves guaranteed generalization, with explicit bounds on sample size and convergence rate. These findings offer principled insight into when and why Mamba’s selection mechanism enables efficient learning, offering a theoretical counterpoint to Transformer-centric explanations of generalization.

1 Introduction

The selective SSM [4, 3] and its variants [30, 24, 2, 12, 25], or commonly referred to as Mamba, have shown strong performance across language, vision, graphs, audio [26], medicine [27, 13, 16, 22], and genomics [15, 29, 27], revitalizing interest in non-attention architectures for sequence modeling. Unlike Transformers [23], which rely on distinctive attention weights to capture token interactions, Mamba employs a *selection mechanism* based on input-dependent gating. This design has demonstrated strong empirical performance across language and vision tasks, while offering hardware efficiency and autoregressive modeling capabilities. Yet, despite this progress, the theoretical principles underlying Mamba’s selection mechanism remain largely unexplored. In particular, it is unclear how the gating vector evolves during training, and under what conditions it enables efficient learning and generalization. Prior theoretical work has mostly focused on Transformers, analyzing how attention aligns with label-relevant features [7, 11, 14, 1, 19–21]. Comparable results for selective SSMs [5, 6, 17] are missing, leaving open fundamental questions about their learnability.

Scope of this work. We present a theoretical study of Mamba’s selection mechanism in a simplified but representative setting: a single-layer, single-head Mamba block followed by a two-layer MLP trained via gradient descent (GD). Our analysis is conducted on a simplified structured data model with both label-relevant and label-irrelevant features under token-level noise, chosen to enable a

tractable theoretical study of the gating mechanism, which constitutes a distinctive property of Mamba in contrast to Transformer-based architectures.

Connection to Transformers. Prior work has noted that Mamba can be viewed as a *gated linear attention mechanism*, where the gating term dynamically modulates how past information contributes to the current output [3]. From this perspective, our framework complements existing theoretical analyses of Transformer attention by showing how an alternative gating-based mechanism can also align with relevant features while ignoring irrelevant ones.

Contributions. Our results provide the first theoretical characterization of the training dynamics of Mamba’s gating mechanism. We prove that: (i) training with GD achieves guaranteed generalization once the sample size and number of iterations scale according to our derived sample complexity and convergence rate bounds; (ii) the gating vector aligns with class-relevant features while negating irrelevant ones, thereby formalizing its role as a feature selector (See Lemmas 5 and 6 in Appendix D); (iii) **the distance between class-relevant features**, captured by ΔL , directly impacts learning speed, highlighting the role of token order and scanning strategies in Mamba. Overall, this work establishes foundational insights into when and why Mamba’s selection mechanism enables efficient learning, offering a theoretical counterpart to Transformer-centric analyses.

2 Problem Setup and Main Theoretical Result

We study a binary classification task with training data $\{(\mathbf{X}^{(n)}, z^{(n)})\}_{n=1}^N$, where each sequence $\mathbf{X}^{(n)} = [\mathbf{x}_1^{(n)}, \dots, \mathbf{x}_L^{(n)}] \in \mathbb{R}^{d \times L}$ contains label-relevant and irrelevant tokens. Labels are determined by a majority vote over the class-relevant tokens. For example, a positive sample may contain two occurrences of the positive token o_+ and only one occurrence of the negative token o_- , whereas a negative sample would contain the opposite. We use a simplified *Mamba block* with input-dependent gating, followed by a two-layer MLP trained with SGD on hinge loss. Given a sequence \mathbf{X} , the model output is

$$F(\mathbf{X}) = \frac{1}{L} \sum_{l=1}^L \sum_{i=1}^m v_i \phi \left(\mathbf{W}_{O(i,\cdot)} \sum_{s=1}^t \left(\prod_{j=s+1}^t (1 - \sigma(\mathbf{w}_\Delta^\top \mathbf{x}_j)) \right) \cdot \sigma(\mathbf{w}_\Delta^\top \mathbf{x}_s) (\mathbf{W}_B^\top \mathbf{x}_s)^\top (\mathbf{W}_C^\top \mathbf{x}_t) \mathbf{x}_s \right), \quad (1)$$

where $\phi(\cdot)$ is the ReLU activation.

Our main result, Theorem 1, gives conditions under which the model generalizes.

Theorem 1 (Generalization of Mamba). *Suppose the model width satisfies $m \geq d^2 \log q$ for some constant $q > 0$, and the token noise level is bounded as $\tau < \mathcal{O}(\frac{1}{d})$. Then, with high probability, if the number of training samples N satisfies $N \geq \Omega\left(\frac{4L^2 d}{\eta^2 [1 + (1/2)^{\Delta L}]^2}\right)$, and the number of iterations T satisfies $T = \Theta\left(\frac{L^2}{\eta [1 + (1/2)^{\Delta L}]}\right)$, the returned model achieves guaranteed generalization.*

3 Numerical Experiments

We complement our theory with simple synthetic experiments under the majority-voting data model. The results confirm our prediction that larger ΔL slows convergence (Table 1). We also measured the cosine similarity between \mathbf{w}_Δ and feature directions, finding 0.18 for o_+ , 0.22 for o_- , and -0.08 for irrelevant features, confirming that the gating filters in relevant tokens while ignoring irrelevant ones.

Table 1: Average epochs for convergence under varying distances between class-relevant features ΔL .

N	$\Delta L=1$	2	4	6	10	25
100	1.90	4.10	18.55	22.25	32.15	47.65
200	1.50	2.50	10.60	19.50	26.65	36.50

4 Conclusion and Future Work

This paper provides a first-step theoretical analysis of the Mamba architecture by examining its gated selection mechanism under the majority-voting data model. We establish sample complexity

and convergence guarantees for training with GD, and show that the gating vector aligns with class-relevant features while ignoring irrelevant ones. To the best of our knowledge, this is the first generalization result for Mamba, highlighting its role as an implicit feature selector. Our results also reveal that token order has a significant impact on the performance, and that different scanning strategies in Mamba can lead to distinct behaviors. Future research directions include extending the analysis to more general data distributions, deeper or multi-head Mamba variants, and hybrid architectures that combine Mamba with attention.

Acknowledgments

This work was supported by the National Science Foundation (NSF) #2349879 and #2349878. We also thank all anonymous reviewers for their constructive comments.

References

- [1] Davoud Ataee Tarzanagh, Yingcong Li, Xuechen Zhang, and Samet Oymak. Max-margin token selection in attention mechanism. *Advances in neural information processing systems*, 36:48314–48362, 2023.
- [2] Ali Behrouz and Farnoosh Hashemi. Graph mamba: Towards learning on graphs with state space models. In *Proceedings of the 30th ACM SIGKDD conference on knowledge discovery and data mining*, pages 119–130, 2024.
- [3] Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality. *arXiv preprint arXiv:2405.21060*, 2024.
- [4] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- [5] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- [6] Ankit Gupta, Albert Gu, and Jonathan Berant. Diagonal state spaces are as effective as structured state spaces. *Advances in neural information processing systems*, 35:22982–22994, 2022.
- [7] Hongkang Li, Meng Wang, Sijia Liu, and Pin-Yu Chen. A theoretical understanding of shallow vision transformers: Learning, generalization, and sample complexity. *arXiv preprint arXiv:2302.06015*, 2023.
- [8] Hongkang Li, Meng Wang, Songtao Lu, Xiaodong Cui, and Pin-Yu Chen. How do nonlinear transformers learn and generalize in in-context learning? *arXiv preprint arXiv:2402.15607*, 2024.
- [9] Hongkang Li, Meng Wang, Tengfei Ma, Sijia Liu, Zaixi Zhang, and Pin-Yu Chen. What improves the generalization of graph transformers? a theoretical dive into the self-attention and positional encoding. In *International Conference on Machine Learning*, pages 28784–28829. PMLR, 2024.
- [10] Hongkang Li, Yihua Zhang, Shuai Zhang, Meng Wang, Sijia Liu, and Pin-Yu Chen. When is task vector provably effective for model editing? a generalization analysis of nonlinear transformers. *arXiv preprint arXiv:2504.10957*, 2025.
- [11] Yuchen Li, Yuanzhi Li, and Andrej Risteski. How do transformers learn topic structure: Towards a mechanistic understanding. In *International Conference on Machine Learning*, pages 19689–19729. PMLR, 2023.
- [12] Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, Jianbin Jiao, and Yunfan Liu. Vmamba: Visual state space model. *Advances in neural information processing systems*, 37:103031–103063, 2024.
- [13] Jun Ma, Feifei Li, and Bo Wang. U-mamba: Enhancing long-range dependency for biomedical image segmentation. *arXiv preprint arXiv:2401.04722*, 2024.

- [14] Samet Oymak, Ankit Singh Rawat, Mahdi Soltanolkotabi, and Christos Thrampoulidis. On the role of attention in prompt-tuning. In *International Conference on Machine Learning*, pages 26724–26768. PMLR, 2023.
- [15] Cong Qi, Hanzhang Fang, Tianxing Hu, Siqi Jiang, and Wei Zhi. Bidirectional mamba for single-cell data: Efficient context learning with biological fidelity. *arXiv preprint arXiv:2504.16956*, 2025.
- [16] Jiacheng Ruan, Jincheng Li, and Suncheng Xiang. Vm-unet: Vision mamba unet for medical image segmentation. *arXiv preprint arXiv:2402.02491*, 2024.
- [17] Jimmy TH Smith, Andrew Warrington, and Scott W Linderman. Simplified state space layers for sequence modeling. *arXiv preprint arXiv:2208.04933*, 2022.
- [18] Jiawei Sun, Shuai Zhang, Hongkang Li, and Meng Wang. Theoretical guarantees and training dynamics of contrastive learning: How misaligned data influence feature purity. In *High-dimensional Learning Dynamics 2025*, 2025. URL <https://openreview.net/forum?id=r09riCUHD9>.
- [19] Davoud Ataei Tarzanagh, Yingcong Li, Christos Thrampoulidis, and Samet Oymak. Transformers as support vector machines. *arXiv preprint arXiv:2308.16898*, 2023.
- [20] Yuandong Tian, Yiping Wang, Beidi Chen, and Simon S Du. Scan and snap: Understanding training dynamics and token composition in 1-layer transformer. *Advances in neural information processing systems*, 36:71911–71947, 2023.
- [21] Yuandong Tian, Yiping Wang, Zhenyu Zhang, Beidi Chen, and Simon Du. Joma: Demystifying multilayer transformers via joint dynamics of mlp and attention. *arXiv preprint arXiv:2310.00535*, 2023.
- [22] Ting Yu Tsai, Li Lin, Shu Hu, Ming-Ching Chang, Hongtu Zhu, and Xin Wang. Uu-mamba: Uncertainty-aware u-mamba for cardiac image segmentation. In *2024 IEEE 7th International Conference on Multimedia Information Processing and Retrieval (MIPR)*, pages 267–273, 2024. doi: 10.1109/MIPR62202.2024.00050.
- [23] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [24] Chloe Wang, Oleksii Tsepa, Jun Ma, and Bo Wang. Graph-mamba: Towards long-range graph sequence modeling with selective state spaces. *arXiv preprint arXiv:2402.00789*, 2024.
- [25] Junxiong Wang, Tushaar Gangavarapu, Jing Nathan Yan, and Alexander M Rush. Mambabyte: Token-free selective state space model. *arXiv preprint arXiv:2401.13660*, 2024.
- [26] Sarthak Yadav and Zheng-Hua Tan. Audio mamba: Selective state spaces for self-supervised audio representations. *arXiv preprint arXiv:2406.02178*, 2024.
- [27] Yubiao Yue and Zhenzhang Li. Medmamba: Vision mamba for medical image classification. *arXiv preprint arXiv:2403.03849*, 2024.
- [28] Shuai Zhang, Meng Wang, Pin-Yu Chen, Sijia Liu, Songtao Lu, and Miao Liu. Joint edge-model sparse learning is provably efficient for graph neural networks. *The Eleventh International Conference on Learning Representations*, 2023.
- [29] Qi Zhao, Ze Li, Qian Mao, Tingwei Chen, Yiran Zhang, Bingle Li, Zheng Zhao, and Xiaoya Fan. Mambacpg: an accurate model for single-cell dna methylation status imputation using mamba. *Briefings in Bioinformatics*, 26(4):bbaf360, 2025.
- [30] Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang. Vision mamba: Efficient visual representation learning with bidirectional state space model. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 62429–62442. PMLR, 2024.

A Data Model

Consider an arbitrary set of orthonormal vectors $\{\mathbf{o}_+, \mathbf{o}_-, \mathbf{o}_3, \dots, \mathbf{o}_d\}$ in \mathbb{R}^d , where \mathbf{o}_+ and \mathbf{o}_- represent discriminative features, and the remaining \mathbf{o}_j are non-discriminative (filler) features ($\mathbf{o}_+ \equiv \mathbf{o}_1$, $\mathbf{o}_- \equiv \mathbf{o}_2$). We add token-level noise to all input tokens; each token $x_l^{(n)}$ in $\mathbf{X}^{(n)}$ is a noisy version of one of the input patterns. The label is determined by a majority vote over the class-relevant features in the input sequence. In our data model, a positive sample contains two instances of \mathbf{o}_+ and one instance of \mathbf{o}_- . Let L_1^+ and L_2^+ denote the positions of the two \mathbf{o}_+ tokens (with $L_1^+ < L_2^+$), and let L^- denote the position of the \mathbf{o}_- token. We assume that L_1^+ , L_2^+ , and L^- are drawn independently from the same distribution over sequence positions.

In addition, we consider a **balanced dataset** sampled from an unknown distribution \mathcal{D} . Let $\mathcal{N}_+ = \{(X^{(n)}, z^{(n)}) : z^{(n)} = +1, n \in [N]\}$ and $\mathcal{N}_- = \{(X^{(n)}, z^{(n)}) : z^{(n)} = -1, n \in [N]\}$ denote the sets of positive and negative samples, respectively. The class imbalance satisfies $| |\mathcal{N}_+| - |\mathcal{N}_-| | = O(\sqrt{N})$.

B Related Work

State Space Models (SSMs). State space models (SSMs) have recently gained attention as a competitive alternative to attention-based architectures for sequence modeling. Early works such as S4[5, 6, 17] demonstrated that carefully parameterized state space recurrences can capture long-range dependencies with strong hardware efficiency, achieving competitive results on long-sequence language and audio benchmarks. More recent variants, including Mamba [4, 3], extend this line of research by introducing input-dependent gating, which enables dynamic selection of relevant features and further improves performance across natural language and vision tasks. Beyond 1D sequences, SSMs have also been adapted to two-dimensional data. For instance, VMamba [12] proposes SS2D, which leverages multiple scanning routes to align the ordered structure of selective 1D scans with the non-sequential nature of vision inputs, thereby enhancing contextual aggregation and boosting recognition performance. Other extensions, such as Graph Mamba [24, 2], demonstrate the adaptability of selective SSMs to non-Euclidean data domains by incorporating structured graph information. These empirical advances underscore the importance of input ordering and scanning strategies in determining the effectiveness of SSMs, a phenomenon that resonates with the structured data models considered in our theoretical analysis.

Theoretical Analysis of SSMs. While empirical studies have demonstrated the effectiveness of SSMs across language, vision, and graph domains, theoretical investigations remain relatively scarce. Existing analyses primarily focus on drawing connections between SSMs and attention-like mechanisms, such as interpreting the recurrence as a form of linear attention or highlighting similarities in feature selection[3]. Although these works provide valuable intuition, they do not address fundamental questions regarding how selective mechanisms in SSMs affect learning dynamics, sample complexity, or generalization under different data distributions. This gap motivates us to study how selective mechanisms in SSMs interact with structured data, a perspective inspired by recent theoretical works on structured models such as vision transformers and graph neural networks(GNNs) [7, 8, 28]. In particular, our work develops a theoretical framework that directly characterizes the role of the gating mechanism in Mamba under structured data models.

Feature Learning Framework. The feature learning framework provides a systematic perspective for analyzing how neural networks gradually emphasize informative features while suppressing uninformative ones during training. Several recent works have adopted this framework to study the optimization and generalization of Transformers [7, 11, 8, 10, 18], showing how attention mechanisms evolve to highlight class-relevant tokens and ignore irrelevant noise. Similar analyses have also been conducted for GNNs [28, 9], where message passing is shown to amplify discriminative structures in the graph progressively.

Despite the recent empirical success of Mamba, its feature learning properties in selective SSMs remain theoretically unexplored. To the best of our knowledge, we provide the first such characterization. In particular, we analyze how the gating mechanism in Mamba evolves under stochastic gradient descent, thereby extending the feature learning perspective to this new family of architectures.

C Main Theorems

In this section, we present the formal statements of our main theorems.

Let $\Psi = (\mathbf{v}, \mathbf{W}_O, \mathbf{w}_\Delta, \mathbf{W}_B, \mathbf{W}_C)$ denote the set of parameters to train. The generalization performance of the learned model Ψ is evaluated using the population risk $f(\Psi)$, defined as

$$f(\Psi) = f(\mathbf{v}, \mathbf{W}_O, \mathbf{w}_\Delta, \mathbf{W}_B, \mathbf{W}_C) = \mathbb{E}_{(\mathbf{X}, z) \sim \mathcal{D}} \ell(\mathbf{X}, z), \quad (2)$$

where $\ell(\mathbf{X}, z; \Psi)$ is the hinge loss function.

Theorem 1 establishes the sample complexity (3) and convergence rate (4) required to guarantee generalization when training with SGD for the majority-voting data model. In other words, the model generalizes once enough samples are available (3) and training has run for the required number of iterations (4).

Theorem 1 (Generalization of Mamba). *Let the learning rate $\eta > 0$ be a positive constant. Suppose the model width satisfies $m \geq d^2 \log q$ for some constant $q > 0$, and the token noise level is bounded as $\tau < \mathcal{O}\left(\frac{1}{d}\right)$.*

Then, with probability at least $1 - N^{-d}$, if the number of training samples N satisfies

$$N \geq \Omega\left(\frac{4L^2d}{\eta^2 \left[1 + \left(\frac{1}{2}\right)^{\Delta L}\right]^2}\right), \quad (3)$$

and the number of iterations T satisfies

$$T = \Theta\left(\frac{L^2}{\eta \left[1 + \left(\frac{1}{2}\right)^{\Delta L}\right]}\right), \quad (4)$$

the model obtained after T iterations of SGD achieves guaranteed generalization, i.e.,

$$f\left(\mathbf{v}^{(0)}, \mathbf{W}_O^{(T)}, \mathbf{w}_\Delta^{(T)}, \mathbf{W}_B^{(0)}, \mathbf{W}_C^{(0)}\right) = 0. \quad (5)$$

Theorem 2 establishes that after sufficient training, the gating vector \mathbf{w}_Δ aligns positively with the class-relevant features \mathbf{o}_+ (6) and \mathbf{o}_- (7), while its alignment with irrelevant features remains strictly negative (8). In other words, the selection mechanism implicitly acts as a feature selector, amplifying relevant tokens and suppressing irrelevant ones.

Theorem 2 (Gating Vector Alignment). *Suppose training is performed under the same conditions as in Theorem 1, with initialization where each entry of \mathbf{W}_O is drawn independently from $\mathcal{N}(0, \xi^2)$ and $\mathbf{w}_\Delta^{(0)} = 0$. If the number of training samples N satisfies $N \geq \Omega\left(\frac{4L^2d}{\eta^2 \left[1 + \left(1/2\right)^{\Delta L}\right]^2}\right)$, and the number of iterations T satisfies $T = \Theta\left(\frac{L^2}{\eta \left[1 + \left(1/2\right)^{\Delta L}\right]}\right)$,*

$$\langle \mathbf{w}_\Delta^{(T)}, \mathbf{o}_+ \rangle \geq \frac{\eta T c'^3}{16L} \left[\frac{1}{c'^2} + \left(\frac{1}{2}\right)^{\Delta L^+ - 2} \right] \quad (6)$$

$$\langle \mathbf{w}_\Delta^{(T)}, \mathbf{o}_- \rangle \geq \frac{\eta T c'^3}{16L} \left[\frac{1}{c'^2} + \left(\frac{1}{2}\right)^{\Delta L^- - 2} \right] \quad (7)$$

$$\langle \mathbf{w}_\Delta^{(T)}, \mathbf{o}_j \rangle \leq -\frac{\eta T c'^3}{16L} \left[\frac{1}{c'^2} + \left(\frac{1}{2}\right)^{\Delta L^+ - 2} \right] \left[\left(\frac{1}{2}\right)^{\Delta L^+} + \left(\frac{1}{2}\right)^{\Delta L^-} \right], \quad (8)$$

for any irrelevant feature \mathbf{o}_j .

D Useful Lemmas

Lemma 1. Suppose $p_1 \leq \langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_+ \rangle \leq q_1$, $p_1 \leq \langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_- \rangle \leq q_1$, and $p_2 \leq \langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_j \rangle \leq q_2$ for $j \neq 1, 2$. Then, for any lucky neuron $i \in \mathcal{W}(t)$ at iteration t , the following bounds hold:

(L1.1) A lower bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{W}_{O(i,\cdot)}$ at iteration t , in the direction of \mathbf{o}_+ , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \geq \frac{1}{\sqrt{m}L} \cdot \sigma(p_1) \left[1 + (1 - \sigma(q_1))^2 (1 - \sigma(q_2))^{L_2^+ - L_1^+ - 2} \right] - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (9)$$

(L1.2) An upper bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{W}_{O(i,\cdot)}$ at iteration t , in the direction of \mathbf{o}_+ , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \leq \frac{1}{\sqrt{m}L} \cdot \sigma(q_1) \left[1 + (1 - \sigma(p_1))^2 (1 - \sigma(p_2))^{L_2^+ - L_1^+ - 2} \right] + \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (10)$$

(L1.3) A lower bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{W}_{O(i,\cdot)}$ at iteration t , in the direction of \mathbf{o}_- , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \geq -\frac{1}{\sqrt{m}L} \cdot \sigma(p_1) \left[2 + (1 - \sigma(q_1))^2 (1 - \sigma(q_2))^{L_2^- - L_1^- - 2} \right] - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (11)$$

(L1.4) An upper bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{W}_{O(i,\cdot)}$ at iteration t , in the direction of \mathbf{o}_- , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (12)$$

(L1.5) An upper bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{W}_{O(i,\cdot)}$ at iteration t , in the direction of \mathbf{o}_j , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right), \quad \text{for } j \neq 1, 2. \quad (13)$$

Lemma 2. For any unlucky neuron $i \in \mathcal{K}_+ \setminus \mathcal{W}(t)$ at iteration t , the following bounds hold:

(L2.1) An upper bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{W}_{O(i,\cdot)}$ at iteration t , in the direction of \mathbf{o}_+ , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (14)$$

(L2.2) An upper bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{W}_{O(i,\cdot)}$ at iteration t , in the direction of \mathbf{o}_- , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (15)$$

(L2.3) An upper bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{W}_{O(i,\cdot)}$ at iteration t , in the direction of \mathbf{o}_j , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right), \quad \text{for } j \neq 1, 2. \quad (16)$$

Lemma 3. Suppose $p_1 \leq \langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_- \rangle \leq q_1$, $p_1 \leq \langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_+ \rangle \leq q_1$, and $p_2 \leq \langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_j \rangle \leq q_2$ for $j \neq 1, 2$. Then, for any lucky neuron $i \in \mathcal{U}(t)$ at iteration t , the following bounds hold:

(L3.1) A lower bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{W}_{O(i,\cdot)}$ at iteration t , in the direction of \mathbf{o}_- , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \geq \frac{1}{\sqrt{mL}} \cdot \sigma(p_1) \left[1 + (1 - \sigma(q_1))^2 (1 - \sigma(q_2))^{L_2^- - L_1^- - 2} \right] - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (17)$$

(L3.2) An upper bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{W}_{O(i,\cdot)}$ at iteration t , in the direction of \mathbf{o}_- , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \leq \frac{1}{\sqrt{mL}} \cdot \sigma(q_1) \left[1 + (1 - \sigma(p_1))^2 (1 - \sigma(p_2))^{L_2^- - L_1^- - 2} \right] + \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (18)$$

(L3.3) A lower bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{W}_{O(i,\cdot)}$ at iteration t , in the direction of \mathbf{o}_+ , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \geq -\frac{1}{\sqrt{mL}} \cdot \sigma(p_1) \left[2 + (1 - \sigma(q_1))^2 (1 - \sigma(q_2))^{L_2^+ - L_1^+ - 2} \right] - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (19)$$

(L3.4) An upper bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{W}_{O(i,\cdot)}$ at iteration t , in the direction of \mathbf{o}_+ , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (20)$$

(L3.5) An upper bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{W}_{O(i,\cdot)}$ at iteration t , in the direction of \mathbf{o}_j , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right), \quad \text{for } j \neq 1, 2. \quad (21)$$

Lemma 4. For any unlucky neuron $i \in \mathcal{K}_- \setminus \mathcal{U}(t)$ at iteration t , the following bounds hold:

(L4.1) An upper bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{W}_{O(i,\cdot)}$ at iteration t , in the direction of \mathbf{o}_- , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (22)$$

(L4.2) An upper bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{W}_{O(i,\cdot)}$ at iteration t , in the direction of \mathbf{o}_+ , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (23)$$

(L4.3) An upper bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{W}_{O(i,\cdot)}$ at iteration t , in the direction of \mathbf{o}_j , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right), \quad \text{for } j \neq 1, 2. \quad (24)$$

Lemma 5. Suppose $p_1 \leq \langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_+ \rangle \leq q_1$ and $r_1^* \leq \langle \mathbf{W}_{O(i,\cdot)}^{(t+1)}, \mathbf{o}_+ \rangle \leq s_1^*$. Let $|\mathcal{W}(t)| = \rho_t^+$ and $|\mathcal{U}(t)| = \rho_t^-$. Then, we have:

(L5.1) A lower bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{w}_\Delta^{(t)}$ at iteration t , in the direction of \mathbf{o}_+ , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{w}_\Delta^{(t)}}, \mathbf{o}_+ \right\rangle \geq \frac{\sigma(p_1)(1 - \sigma(q_1))}{2} \left[\frac{2r_1^*\rho_t^+}{\sqrt{m}} - \frac{\sqrt{m}}{2}s_1^* \right] - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (25)$$

Suppose $p_1 \leq \langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_- \rangle \leq q_1$ and $r_1^* \leq \langle \mathbf{W}_{O(i,\cdot)}^{(t+1)\top}, \mathbf{o}_+ \rangle \leq s_1^*$. Let $|\mathcal{W}(t)| = \rho_t^+$ and $|\mathcal{U}(t)| = \rho_t^-$. Then, we have:

(L5.2) A lower bound on the gradient of $\hat{\mathcal{L}}$ with respect to $\mathbf{w}_\Delta^{(t)}$ at iteration t , in the direction of \mathbf{o}_- , is given by

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{w}_\Delta^{(t)}}, \mathbf{o}_- \right\rangle \geq \frac{\sigma(p_1)(1-\sigma(q_1))}{2} \left[\frac{2r_1^*\rho_t^-}{\sqrt{m}} - \frac{\sqrt{m}}{2}s_1^* \right] - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (26)$$

Lemma 6. Suppose $p_1 \leq \langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_+ \rangle \leq q_1$, $p_1 \leq \langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_- \rangle \leq q_1$, $\langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_j \rangle \leq q_2$ for $j \neq 1, 2$, and $r_1^* \leq \langle \mathbf{W}_{O(i,\cdot)}^{(t)\top}, \mathbf{o}_+ \rangle$. Let $\rho_t^+ = |\mathcal{W}(t)|$ and $\rho_t^- = |\mathcal{U}(t)|$. Then we have:

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{w}_\Delta^{(t)}}, \mathbf{o}_j \right\rangle \leq -\frac{r_1^*}{2\sqrt{m}} \cdot \sigma(p_1)(1-\sigma(q_1)) \left[(1-\sigma(q_2))^{L_2^+ - L_1^+} \rho_t^+ + (1-\sigma(q_2))^{L_2^- - L_1^-} \rho_t^- \right]. \quad (27)$$

E Proof of Convergence

Proof. Suppose $\mathbf{w}_\Delta^{(0)} = \mathbf{0}$. Then, we have

$$\langle \mathbf{w}_\Delta^{(0)}, \mathbf{o}_+ \rangle = 0, \quad \langle \mathbf{w}_\Delta^{(0)}, \mathbf{o}_- \rangle = 0, \quad \text{and} \quad \langle \mathbf{w}_\Delta^{(0)}, \mathbf{o}_j \rangle = 0 \quad \forall j.$$

From Lemma 1, identify $p_1 = 0$, $q_1 = 0$, $p_2 = 0$ and $q_2 = 0$. Let L_1^+ and L_2^+ denote the positions of the two \mathbf{o}_+ tokens in the positive sample, with $L_1^+ < L_2^+$. Define the gap between them as

$$\Delta L^+ := L_2^+ - L_1^+. \quad (28)$$

Similarly, for the negative sample, let L_1^- and L_2^- be the positions of the two \mathbf{o}_- tokens, with $L_1^- < L_2^-$, and define

$$\Delta L^- := L_2^- - L_1^-. \quad (29)$$

Then, for any lucky neuron $i \in \mathcal{W}(0)$, we obtain

$$\begin{aligned} \frac{1}{2\sqrt{m}L} \left[1 + \left(\frac{1}{2}\right)^{\Delta L^+} \right] - \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right) &\leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(0)}}, \mathbf{o}_+ \right\rangle \\ &\leq \frac{1}{2\sqrt{m}L} \left[1 + \left(\frac{1}{2}\right)^{\Delta L^+} \right] + \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right). \end{aligned} \quad (30)$$

We can relax the lower bound and obtain

$$\begin{aligned} \frac{c'^2}{2\sqrt{m}L} \left[\frac{1}{c'^2} + \left(\frac{1}{2}\right)^{\Delta L^+ - 2} \right] - \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right) &\leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(0)}}, \mathbf{o}_+ \right\rangle \\ &\leq \frac{1}{2\sqrt{m}L} \left[1 + \left(\frac{1}{2}\right)^{\Delta L^+} \right] + \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right) \end{aligned} \quad (31)$$

$$\text{and} \quad \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(0)}}, \mathbf{o}_j \right\rangle \leq \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right) \quad \text{for } j \neq 1. \quad (32)$$

We assume that the number of samples in a batch $N = \text{poly}(d)$.

Suppose the initialization is

$$\mathbf{W}_{O(i,\cdot)}(0) = \delta_1 \mathbf{o}_+ + \delta_2 \mathbf{o}_- + \cdots + \delta_d \mathbf{o}_d, \quad \delta_j \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \xi^2) \quad j = 1, 2, \dots, d. \quad (33)$$

Then, after one gradient descent step, we have

$$\begin{aligned} & \delta_1 + \frac{\eta c'^2}{2\sqrt{mL}} \left[\frac{1}{c'^2} + \left(\frac{1}{2}\right)^{\Delta L^+ - 2} \right] \\ & - \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right) \leq \left\langle \mathbf{W}_{O(i,\cdot)}^\top (1), \mathbf{o}_+ \right\rangle \\ & \leq \delta_1 + \frac{\eta}{2\sqrt{mL}} \left[1 + \left(\frac{1}{2}\right)^{\Delta L^+} \right] \\ & + \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right). \end{aligned} \quad (34)$$

$$\text{and } \left\langle \mathbf{W}_{O(i,\cdot)}^\top (1), \mathbf{o}_j \right\rangle \leq \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right) \quad \text{for } j \neq 1. \quad (35)$$

By applying Lemma 3, for any lucky neuron $i \in \mathcal{U}(0)$, we obtain

$$\begin{aligned} & \delta_2 + \frac{\eta c'^2}{2\sqrt{mL}} \left[\frac{1}{c'^2} + \left(\frac{1}{2}\right)^{\Delta L^- - 2} \right] \\ & - \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right) \leq \left\langle \mathbf{W}_{O(i,\cdot)}^\top (1), \mathbf{o}_- \right\rangle \\ & \leq \delta_2 + \frac{\eta}{2\sqrt{mL}} \left[1 + \left(\frac{1}{2}\right)^{\Delta L^-} \right] \\ & + \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right). \end{aligned} \quad (36)$$

$$\text{and } \left\langle \mathbf{W}_{O(i,\cdot)}^\top (1), \mathbf{o}_j \right\rangle \leq \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right) \quad \text{for } j \neq 2. \quad (37)$$

For any unlucky neuron $i \in \mathcal{K}_- \setminus \mathcal{U}(0)$, Lemma 4 gives

$$\left\langle \mathbf{W}_{O(i,\cdot)}^\top (1), \mathbf{o}_j \right\rangle \leq \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right) \quad \text{for } \forall j. \quad (38)$$

Now consider the gradient update for \mathbf{w}_Δ . Define:

$$\begin{aligned} a &= \delta_1 + \frac{\eta c'^2}{2\sqrt{mL}} \left[\frac{1}{c'^2} + \left(\frac{1}{2}\right)^{\Delta L^+ - 2} \right] - \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right) \\ b &= \delta_1 + \frac{\eta}{2\sqrt{mL}} \left[1 + \left(\frac{1}{2}\right)^{\Delta L^+} \right] + \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right) \end{aligned}$$

Applying Lemma 5 with $p_1 = 0, q_1 = 0, r_1^* = a, s_1^* = b$, and $\rho_0^+ = |\mathcal{W}(0)|$, we get

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{w}_\Delta^{(0)}}, \mathbf{o}_+ \right\rangle \geq \frac{1}{8} \left[\frac{2a}{\sqrt{m}} \cdot \rho_0^+ - \frac{\sqrt{m}}{2} b \right] - \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right) \quad (39)$$

We can relax this lower bound and obtain

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{w}_\Delta^{(0)}}, \mathbf{o}_+ \right\rangle \geq \frac{c'}{4} \left[\frac{2a}{\sqrt{m}} \cdot \rho_0^+ - \frac{\sqrt{m}}{2} b \right] - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) =: \alpha \quad (40)$$

Since $\rho_0^+ \approx \frac{m}{2}$ under random initialization, when m is sufficiently large, the above gradient update is positive.

Let $\delta_1 = \frac{1}{\text{poly}(d)}$. Since $a \approx b$,

$$\begin{aligned} \alpha &= \frac{c'}{4} \left[\frac{2a}{\sqrt{m}} \cdot \frac{m}{2} - \frac{\sqrt{m}b}{2} \right] - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) \\ &= \frac{c'}{4} \left[\sqrt{ma} - \frac{\sqrt{ma}}{2} \right] - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) \\ &= \frac{c'}{4} \cdot \frac{\eta c'^2}{4L} \left[\frac{1}{c'^2} + \left(\frac{1}{2} \right)^{\Delta L^+ - 2} \right] - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) \\ &= \frac{\eta c'^3}{16L} \left[\frac{1}{c'^2} + \left(\frac{1}{2} \right)^{\Delta L^+ - 2} \right] - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) > 0 \end{aligned} \quad (41)$$

From Lemma 6, we also obtain

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{w}_\Delta^{(0)}}, \mathbf{o}_j \right\rangle \leq \frac{-a}{8\sqrt{m}} \left[\left(\frac{1}{2} \right)^{\Delta L^+} \cdot \rho_0^+ + \left(\frac{1}{2} \right)^{\Delta L^-} \cdot \rho_0^- \right] \quad (42)$$

where we apply the lemma with the values

$$p_1 = 0, \quad q_1 = 0, \quad q_2 = 0, \quad \text{and} \quad r_1^* = a.$$

We can relax this upper bound and obtain

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{w}_\Delta^{(0)}}, \mathbf{o}_j \right\rangle \leq \frac{-ac'}{4\sqrt{m}} \left[\left(\frac{1}{2} \right)^{\Delta L^+} \cdot \rho_0^+ + \left(\frac{1}{2} \right)^{\Delta L^-} \cdot \rho_0^- \right] =: \gamma \quad (43)$$

Taking $\rho_0^+ = \rho_0^- \approx \frac{m}{2}$, we can simplify and write

$$\begin{aligned} \gamma &= \frac{-ac'}{4\sqrt{m}} \cdot \frac{m}{2} \left[\left(\frac{1}{2} \right)^{\Delta L^+} + \left(\frac{1}{2} \right)^{\Delta L^-} \right] \\ &= -\sqrt{ma} \cdot \frac{c'}{8} \left[\left(\frac{1}{2} \right)^{\Delta L^+} + \left(\frac{1}{2} \right)^{\Delta L^-} \right] \\ &= \frac{-\eta c'^3}{16L} \left[\frac{1}{c'^2} + \left(\frac{1}{2} \right)^{\Delta L^+ - 2} \right] \left[\left(\frac{1}{2} \right)^{\Delta L^+} + \left(\frac{1}{2} \right)^{\Delta L^-} \right] - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right). \end{aligned} \quad (44)$$

Let $\langle \mathbf{w}_\Delta^{(1)}, \mathbf{o}_+ \rangle = \alpha^* \geq \alpha$, $\langle \mathbf{w}_\Delta^{(1)}, \mathbf{o}_- \rangle = \beta^* \geq \beta$, and $\langle \mathbf{w}_\Delta^{(1)}, \mathbf{o}_j \rangle = \gamma^* \leq \gamma$, where

$$\beta = \frac{c'}{4} \left[\frac{2a'}{\sqrt{m}} \cdot \rho_0^- - \frac{\sqrt{m}}{2} b' \right] - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) > 0 \quad (45)$$

$$a' = \delta_2 + \frac{\eta c'^2}{2\sqrt{m}L} \left[\frac{1}{c'^2} + \left(\frac{1}{2} \right)^{\Delta L^- - 2} \right] - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right)$$

$$b' = \delta_2 + \frac{\eta}{2\sqrt{m}L} \left[1 + \left(\frac{1}{2} \right)^{\Delta L^-} \right] + \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right)$$

Following the same approach as in (41), we can simplify and obtain

$$\beta = \frac{\eta c'^3}{16L} \left[\frac{1}{c'^2} + \left(\frac{1}{2} \right)^{\Delta L^- - 2} \right] - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) > 0 \quad (46)$$

For any lucky neuron $i \in \mathcal{W}(1)$ at iteration 2, we have

$$\begin{aligned} & \frac{c'^2}{\sqrt{mL}} \cdot \sigma(\alpha^*) \left[\frac{1}{c'^2} + (1 - \sigma(\gamma^*))^{\Delta L^+ - 2} \right] - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) \\ & \leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(1)}}, \mathbf{o}_+ \right\rangle \\ & \leq \frac{1}{\sqrt{mL}} \cdot \sigma(\alpha^*) \left[1 + (1 - \sigma(\alpha^*))^2 (1 - \sigma(\gamma^*))^{\Delta L^+ - 2} \right] + \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) \end{aligned} \quad (47)$$

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(1)}}, \mathbf{o}_j \right\rangle \leq \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) \quad \text{for } j \neq 1 \quad (48)$$

Note that, $\sigma(\alpha^*) > \frac{1}{2}$ and $\sigma(\gamma^*) < \frac{1}{2}$. This ensures $\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(0)}}, \mathbf{o}_+ \right\rangle \leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(1)}}, \mathbf{o}_+ \right\rangle$.

Thus, we obtain the following bound after the second gradient descent step.

$$\begin{aligned} & \delta_1 + \frac{\eta c'^2}{\sqrt{mL}} \left[\frac{1}{2c'^2} + \left(\frac{1}{2} \right)^{\Delta L^+ - 1} + \sigma(\alpha^*) \left(1/c'^2 + (1 - \sigma(\gamma^*))^{\Delta L^+ - 2} \right) \right] \\ & - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) =: u \\ & \leq \left\langle (\mathbf{W}_{O(i,\cdot)}^{(2)})^\top, \mathbf{o}_+ \right\rangle \\ & \leq \delta_1 + \frac{\eta}{2\sqrt{mL}} \left[1 + \left(\frac{1}{2} \right)^{\Delta L^+} + 2\sigma(\alpha^*) \left(1 + (1 - \sigma(\alpha^*))^2 (1 - \sigma(\gamma^*))^{\Delta L^+ - 2} \right) \right] \\ & + \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) =: v \end{aligned} \quad (49)$$

Similarly, applying Lemma 3 to any lucky neuron $i \in \mathcal{U}(1)$ at iteration 2, we get

$$\begin{aligned} & \frac{c'^2}{\sqrt{mL}} \cdot \sigma(\beta^*) \left[\frac{1}{c'^2} + (1 - \sigma(\gamma^*))^{\Delta L^- - 2} \right] - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) \\ & \leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(1)}}, \mathbf{o}_- \right\rangle \\ & \leq \frac{1}{\sqrt{mL}} \cdot \sigma(\beta^*) \left[1 + (1 - \sigma(\beta^*))^2 (1 - \sigma(\gamma^*))^{\Delta L^- - 2} \right] + \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) \end{aligned} \quad (50)$$

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(1)}}, \mathbf{o}_j \right\rangle \leq \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) \quad \text{for } j \neq 2 \quad (51)$$

Applying Lemma 5 with $p_1 = \alpha^*$, $q_1 = \alpha^*$, $r_1^* = u$, and $s_1^* = v$, we obtain

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{w}_\Delta^{(1)}}, \mathbf{o}_+ \right\rangle \geq \frac{\sigma(\alpha^*)c'}{2} \left[\frac{2u}{\sqrt{m}} \cdot \rho_1^+ - \frac{\sqrt{m}}{2} v \right] - \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right) =: \chi \quad (52)$$

Since $u \approx v$, we have $\chi \geq 0$.

By applying Lemma 6 with

$$p_1 = \alpha^* (= \beta^*), \quad q_1 = \alpha^* (= \beta^*), \quad q_2 = \gamma^*, \text{ and } r_1^* = u, \text{ we have}$$

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{w}_\Delta^{(t)}}, \mathbf{o}_j \right\rangle \leq -\frac{c'u}{2\sqrt{m}} \sigma(\alpha^*) \left[(1 - \sigma(\gamma^*))^{\Delta L^+} \rho_t^+ + (1 - \sigma(\gamma^*))^{\Delta L^-} \rho_t^- \right] =: \iota \quad (53)$$

Note that here we assumed the distribution of ΔL^+ is identical to ΔL^- to have $\alpha^* = \beta^*$.

Finally, for any lucky neuron $i \in \mathcal{W}(T)$, we obtain

$$\left\langle \mathbf{W}_{O(i,\cdot)}^\top \overset{(T)}{\cdot}, \mathbf{o}_+ \right\rangle \geq aT \quad (54)$$

$$\left\langle \mathbf{W}_{O(i,\cdot)}^\top \overset{(T)}{\cdot}, \mathbf{o}_j \right\rangle \leq \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right) \quad \text{for } j \neq 1 \quad (55)$$

For any lucky neuron $i \in \mathcal{U}(T)$, we obtain

$$\left\langle \mathbf{W}_{O(i,\cdot)}^\top \overset{(T)}{\cdot}, \mathbf{o}_- \right\rangle \geq aT \quad (56)$$

$$\left\langle \mathbf{W}_{O(i,\cdot)}^\top \overset{(T)}{\cdot}, \mathbf{o}_j \right\rangle \leq \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right) \quad \text{for } j \neq 2 \quad (57)$$

$$\left\langle \mathbf{w}_\Delta^{(T)}, \mathbf{o}_+ \right\rangle \geq \alpha T \quad (58)$$

$$\left\langle \mathbf{w}_\Delta^{(T)}, \mathbf{o}_- \right\rangle \geq \beta T \quad (59)$$

$$\left\langle \mathbf{w}_\Delta^{(T)}, \mathbf{o}_j \right\rangle \leq \gamma T \quad (60)$$

$$\mathbf{X}^{(n)} = \begin{bmatrix} \mathbf{x}_1^{(n)} & \mathbf{x}_2^{(n)} & \dots & \mathbf{x}_L^{(n)} \end{bmatrix}$$

Consider $z^{(n)} = +1$ as an example. The sequence $\mathbf{X}^{(n)}$ has two \mathbf{o}_+ and one \mathbf{o}_- at locations L_1^+ , L_2^+ , and L^- .

$$\begin{aligned} F(\mathbf{X}^{(n)}) &= \frac{1}{L} \sum_{l=1}^L \sum_{i=1}^m v_i \phi\left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}\right) \\ &= \frac{1}{\sqrt{m}L} \sum_{i \in \mathcal{K}^+} \sum_{l=1}^L \phi\left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}\right) - \frac{1}{\sqrt{m}L} \sum_{i \in \mathcal{K}^-} \sum_{l=1}^L \phi\left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}\right) \\ &\geq \frac{1}{\sqrt{m}L} \sum_{i \in \mathcal{W}(0)} \sum_{l=1}^L \phi\left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}\right) - \frac{1}{\sqrt{m}L} \sum_{i \in \mathcal{U}(0)} \sum_{l=1}^L \phi\left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}\right) \\ &\quad - \frac{1}{\sqrt{m}L} \sum_{i \in \mathcal{K}^- \setminus \mathcal{U}(0)} \sum_{l=1}^L \phi\left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}\right) \end{aligned} \quad (61)$$

The Mamba output $\mathbf{y}_l^{(n)}$ is defined as

$$\mathbf{y}_l^{(n)} = \sum_{\tau=0}^{l-1} \left(\prod_{r=\tau+1}^{l-1} \left(1 - \sigma(\mathbf{w}_\Delta^\top \mathbf{x}_{l-r}^{(n)}) \right) \right) \cdot \sigma(\mathbf{w}_\Delta^\top \mathbf{x}_{l-\tau}^{(n)}) \cdot (\mathbf{x}_{l-\tau}^{(n)\top} \mathbf{x}_l^{(n)}) \mathbf{x}_{l-\tau}^{(n)}. \quad (62)$$

Equivalently, letting $s = l - \tau$, we write

$$\mathbf{y}_l^{(n)} = \sum_{s=1}^l \left(\prod_{j=s+1}^l \left(1 - \sigma(\mathbf{w}_\Delta^\top \mathbf{x}_j^{(n)}) \right) \right) \cdot \sigma(\mathbf{w}_\Delta^\top \mathbf{x}_s^{(n)}) \cdot (\mathbf{x}_s^{(n)\top} \mathbf{x}_l^{(n)}) \mathbf{x}_s^{(n)}. \quad (63)$$

We now derive a lower bound for

$$\sum_{i \in \mathcal{W}(0)} \sum_{l=1}^L \phi(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l).$$

To that end, consider the aggregated projection

$$\sum_{i \in \mathcal{W}(0)} \sum_{l=1}^L \mathbf{W}_{O(i,\cdot)} \mathbf{y}_l = \sum_{i \in \mathcal{W}(0)} \sum_{l=1}^L \sum_{j=1}^d \langle \mathbf{W}_{O(i,\cdot)}^\top, \mathbf{o}_j \rangle \cdot \langle \mathbf{y}_l, \mathbf{o}_j \rangle. \quad (64)$$

For any $i \in \mathcal{W}(0)$, we know that

$$\langle \mathbf{W}_{O(i,\cdot)}^\top, \mathbf{o}_+ \rangle \geq aT. \quad (65)$$

Hence, let's obtain a lower bound for $\langle \mathbf{y}_l, \mathbf{o}_+ \rangle$

We only need to consider the cases where $\mathbf{x}_s = \mathbf{o}_+$ for some s in the range $1 \leq s \leq l$. In particular, we will focus on the following instances:

$$s = L_1^+ \text{ and } l \in \{L_1^+, L_2^+\}, \quad s = L_2^+ \text{ and } l = L_2^+.$$

After T iterations, we know

$$\langle \mathbf{w}_\Delta, \mathbf{o}_+ \rangle \geq \alpha T, \quad \langle \mathbf{w}_\Delta, \mathbf{o}_- \rangle \geq \beta T, \quad \langle \mathbf{w}_\Delta, \mathbf{o}_j \rangle \leq \gamma T \quad \text{for } j \neq 1, 2. \quad (66)$$

Therefore,

$$\langle \mathbf{y}_{L_1^+}, \mathbf{o}_+ \rangle = \sigma(\langle \mathbf{w}_\Delta, \mathbf{o}_+ \rangle) \geq \sigma(\alpha T). \quad (67)$$

We have,

$$\langle \mathbf{w}_\Delta, \mathbf{o}_+ \rangle \leq W_1 T, \quad \langle \mathbf{w}_\Delta, \mathbf{o}_- \rangle \leq W_2 T,$$

where

$$W_1 = \frac{\eta}{16L} \left[\frac{1}{2} + \left(\frac{1}{2} \right)^{\Delta L^+ + 1} \right] + \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right), \quad (68)$$

$$W_2 = \frac{\eta}{16L} \left[\frac{1}{2} + \left(\frac{1}{2} \right)^{\Delta L^- + 1} \right] + \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right). \quad (69)$$

Then we obtain the following:

$$\begin{aligned} \langle \mathbf{y}_{L_2^+}, \mathbf{o}_+ \rangle &\geq \sigma(\alpha T) + (1 - \sigma(W_1 T)) (1 - \sigma(W_2 T)) (1 - \sigma(\gamma T))^{\Delta L^+ - 2} \cdot \sigma(\alpha T) \\ &= \sigma(\alpha T) \left[1 + (1 - \sigma(W_1 T)) (1 - \sigma(W_2 T)) (1 - \sigma(\gamma T))^{\Delta L^+ - 2} \right]. \end{aligned} \quad (70)$$

We now lower bound the objective

$$\sum_{i \in \mathcal{W}(0)} \sum_{l=1}^L \phi(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l).$$

We begin with

$$\sum_{i \in \mathcal{W}(0)} \sum_{l=1}^L \phi(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l) \geq \sum_{i \in \mathcal{W}(0)} \left[\phi(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_{L_1^+}) + \phi(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_{L_2^+}) \right].$$

Note that

$$\mathbf{W}_{O(i,\cdot)} \mathbf{y}_{L_1^+} = \sum_{j=1}^d \left\langle \mathbf{W}_{O(i,\cdot)}^\top, \mathbf{o}_j \right\rangle \left\langle \mathbf{y}_{L_1^+}, \mathbf{o}_j \right\rangle,$$

and $\mathbf{y}_{L_1^+}$ has only \mathbf{o}_+ component.

Therefore,

$$\mathbf{W}_{O(i,\cdot)} \mathbf{y}_{L_1^+} = \left\langle \mathbf{W}_{O(i,\cdot)}^\top, \mathbf{o}_+ \right\rangle \left\langle \mathbf{y}_{L_1^+}, \mathbf{o}_+ \right\rangle \geq aT \cdot \sigma(\alpha T) > 0.$$

Similarly, we can write

$$\mathbf{W}_{O(i,\cdot)} \mathbf{y}_{L_2^+} \geq aT \cdot \sigma(\alpha T) \left[1 + (1 - \sigma(W_1 T)) (1 - \sigma(W_2 T)) (1 - \sigma(\gamma T))^{\Delta L^+ - 2} \right] > 0.$$

Applying $\phi(z) = z$ for positive z , we obtain

$$\begin{aligned} \phi(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_{L_1^+}) &\geq aT \cdot \sigma(\alpha T), \\ \phi(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_{L_2^+}) &\geq aT \cdot \sigma(\alpha T) \left[1 + (1 - \sigma(W_1 T)) (1 - \sigma(W_2 T)) (1 - \sigma(\gamma T))^{\Delta L^+ - 2} \right]. \end{aligned}$$

Hence,

$$\begin{aligned} \sum_{i \in \mathcal{W}(0)} \sum_{l=1}^L \phi(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l) &\geq \sum_{i \in \mathcal{W}(0)} aT \cdot \sigma(\alpha T) \\ &\quad \times \left[2 + (1 - \sigma(W_1 T)) (1 - \sigma(W_2 T)) (1 - \sigma(\gamma T))^{\Delta L^+ - 2} \right]. \end{aligned} \tag{71}$$

Next, we derive an upper bound for

$$\sum_{i \in \mathcal{U}(0)} \sum_{l=1}^L \phi(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}).$$

For any $i \in \mathcal{U}(0)$, we know that

$$0 < \langle \mathbf{W}_{O(i,\cdot)}^\top, \mathbf{o}_- \rangle \leq bT. \tag{72}$$

We now derive an upper bound for $\langle \mathbf{y}_l, \mathbf{o}_- \rangle$. We only need to consider the cases where $\mathbf{x}_s = \mathbf{o}_-$ such that $1 \leq s \leq l$.

$$s = L^- \text{ and } l = L^-.$$

$$\langle \mathbf{y}_{L^-}, \mathbf{o}_- \rangle = \sigma(\langle \mathbf{w}_\Delta, \mathbf{o}_- \rangle) \leq \sigma(W_2 T). \tag{73}$$

$$\sum_{l=1}^L \mathbf{W}_{O(i,\cdot)} \mathbf{y}_l \leq bT \cdot \sigma(W_2 T). \tag{74}$$

$$\sum_{i \in \mathcal{U}(0)} \sum_{l=1}^L \phi(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \leq \sum_{i \in \mathcal{U}(0)} bT \cdot \sigma(W_2 T). \tag{75}$$

In addition, we have

$$\sum_{i \in \mathcal{K}^- \setminus \mathcal{U}(0)} \sum_{l=1}^L \phi(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \leq \tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right). \tag{76}$$

By (61), we can write

$$F(\mathbf{X}^{(n)}) \geq \frac{1}{\sqrt{mL}} \left\{ \frac{m}{2} \cdot aT \cdot \sigma(\alpha T) \left[2 + (1 - \sigma(W_1 T))(1 - \sigma(W_2 T))(1 - \sigma(\gamma T))^{\Delta L^+ - 2} \right] \right. \\ \left. - \frac{m}{2} \cdot bT \cdot \sigma(W_2 T) - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) \right\}, \quad (77)$$

with

$$a = \frac{\eta}{\sqrt{mL}} \left[\frac{1}{2} + c'^2 \left(\frac{1}{2} \right)^{\Delta L^+ - 1} \right] - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right), \quad (78)$$

$$\text{and } b = \frac{\eta}{\sqrt{mL}} \left[\frac{1}{2} + \left(\frac{1}{2} \right)^{\Delta L^+ + 1} \right] + \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right). \quad (79)$$

$$a \approx b. \quad (80)$$

$$\begin{aligned} \alpha &= \frac{\eta c'^3}{16L} \left[\frac{1}{c'^2} + \left(\frac{1}{2} \right)^{\Delta L^+ - 2} \right] - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) \\ &= \frac{\eta}{16L} \left[c' + c'^3 \left(\frac{1}{2} \right)^{\Delta L^+ - 2} \right] - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) \\ &\approx W_1 \approx W_2. \end{aligned} \quad (81)$$

The last step assumes the distribution of ΔL^+ is identical to ΔL^- .

Therefore, we conclude that

$$F(\mathbf{X}^{(n)}) \geq \frac{1}{\sqrt{mL}} \left\{ \frac{m}{2} \cdot aT \cdot \sigma(\alpha T) \left[1 + (1 - \sigma(W_1 T))(1 - \sigma(W_2 T))(1 - \sigma(\gamma T))^{\Delta L^+ - 2} \right] \right\} \\ - \tilde{\mathcal{O}} \left(\frac{1}{\text{poly}(d)} \right) \quad (82)$$

$$F(\mathbf{X}^{(n)}) \geq C, \text{ where } C \text{ is some positive constant.} \quad (83)$$

E.1 Convergence Rate

Let's find the number of iterations T required such that $F(\mathbf{X}^{(n)}) \geq 1$, since the label is +1. We require

$$\frac{1}{\sqrt{mL}} \cdot \frac{m}{2} \cdot aT \geq 1 + \epsilon. \quad (84)$$

Substituting the value of $a \approx b = \frac{\eta}{\sqrt{mL}} \left[\frac{1}{2} + \left(\frac{1}{2} \right)^{\Delta L^+ + 1} \right]$, the condition becomes

$$\begin{aligned} \frac{\sqrt{maT}}{2L} &= \frac{\sqrt{m}}{2L} \cdot \frac{\eta}{\sqrt{mL}} \left[\frac{1}{2} + \left(\frac{1}{2} \right)^{\Delta L^+ + 1} \right] T \\ &= \frac{\eta T}{2L^2} \left[\frac{1}{2} + \left(\frac{1}{2} \right)^{\Delta L^+ + 1} \right] \geq 1 + \epsilon. \end{aligned} \quad (85)$$

Solving for T , we obtain

$$T \geq \frac{2L^2(1 + \epsilon)}{\eta \left[\frac{1}{2} + \left(\frac{1}{2} \right)^{\Delta L^+ + 1} \right]} = \frac{4L^2(1 + \epsilon)}{\eta \left[1 + \left(\frac{1}{2} \right)^{\Delta L^+} \right]} \approx \frac{4L^2}{\eta \left[1 + \left(\frac{1}{2} \right)^{\Delta L^+} \right]}. \quad (86)$$

Now, we additionally require that the sigmoid activation $\sigma(\alpha T)$ be sufficiently large, i.e.,

$$\sigma(\alpha T) \geq 1 - \epsilon. \quad (87)$$

When z is sufficiently large we can approximate

$$\sigma(z) = \frac{1}{1 + e^{-z}} \approx 1 - e^{-z}.$$

Substituting $z = \alpha T$, condition (87) becomes:

$$\begin{aligned} \sigma(\alpha T) &\approx 1 - e^{-\alpha T} \geq 1 - \epsilon, \\ e^{-\alpha T} &\leq \epsilon, \\ \alpha T &\geq -\ln(\epsilon) \\ T &\geq -\frac{\ln(\epsilon)}{\alpha}. \end{aligned} \quad (88)$$

Substituting $\alpha = \frac{\eta}{16L} \left[\frac{1}{2} + \left(\frac{1}{2} \right)^{\Delta L^+ + 1} \right]$, we get:

$$T \geq -\ln(\epsilon) \cdot \frac{16L}{\eta \left[1 + \left(\frac{1}{2} \right)^{\Delta L^+ + 1} \right]}. \quad (89)$$

Multiplying both the numerator and the denominator of the second term by 2, we obtain the final bound:

$$T \geq -\ln(\epsilon) \cdot \frac{32L}{\eta \left[1 + \left(\frac{1}{2} \right)^{\Delta L^+} \right]}. \quad (90)$$

Hence, by combining (86) and (90), we obtain

$$T \geq \max \left\{ \frac{4L^2}{\eta \left[1 + \left(\frac{1}{2} \right)^{\Delta L^+} \right]}, -\ln(\epsilon) \cdot \frac{32L}{\eta \left[1 + \left(\frac{1}{2} \right)^{\Delta L^+} \right]} \right\}. \quad (91)$$

By combining (84) and (87) with the expression for the model output $F(\mathbf{X}^{(n)})$ in (82), we obtain

$$\begin{aligned} F(\mathbf{X}^{(n)}) &\geq (1 + \epsilon) \cdot (1 - \epsilon) \cdot (1 + \epsilon) \\ &\geq (1 + \epsilon)^2 \cdot (1 - \epsilon) \\ &\geq (1 + 2\epsilon + \epsilon^2) \cdot (1 - \epsilon) \\ &\geq 1 + \epsilon - \mathcal{O}(\epsilon^2) \end{aligned} \quad (92)$$

Hence, for sufficiently small $\epsilon > 0$, the model output satisfies $F(\mathbf{X}^{(n)}) \geq 1 + \epsilon$.

E.2 Sample Complexity

Previously, we used $N = \text{poly}(d)$ to obtain the asymptotic term $\tilde{\mathcal{O}}\left(\frac{1}{\text{poly}(d)}\right)$. We now derive a sample-complexity bound that guarantees zero generalization error.

Assume that, for sufficiently small $\lambda \ll 1$,

$$\mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) \leq \lambda \cdot \frac{\eta}{2\sqrt{mL}} \left[1 + \left(\frac{1}{2} \right)^{\Delta L^+} \right]. \quad (93)$$

From this, we can derive a lower bound on the required sample size as

$$\begin{aligned} N &\geq \Omega\left(\lambda^{-2} \cdot \frac{4L^2 d}{\eta^2 \left[1 + \left(\frac{1}{2} \right)^{\Delta L^+} \right]^2}\right) \\ &\geq \Omega\left(\frac{4L^2 d}{\eta^2 \left[1 + \left(\frac{1}{2} \right)^{\Delta L^+} \right]^2}\right). \end{aligned} \quad (94)$$

□

F Proof of Lemmas

F.1 Proof of Lemma 1

Proof. The loss function for the n^{th} sample is defined as

$$\begin{aligned}\ell(\mathbf{X}^{(n)}, z^{(n)}) &= \max\{0, 1 - z^{(n)} \cdot F(\mathbf{X}^{(n)})\} \\ &= \max\left\{0, 1 - z^{(n)} \cdot \frac{1}{L} \sum_{l=1}^L \sum_{i=1}^m v_i \phi\left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}\right)\right\}.\end{aligned}\quad (95)$$

The empirical loss is denoted by $\hat{\mathcal{L}}$ and is given by

$$\hat{\mathcal{L}} = \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{X}^{(n)}, z^{(n)}). \quad (96)$$

The population loss is denoted by \mathcal{L} and is defined as

$$\mathcal{L} = \mathbb{E}_{(\mathbf{X}, z) \sim \mathcal{D}} \ell(\mathbf{X}, z). \quad (97)$$

We know that the gradient of the loss function for the n^{th} sample is

$$\begin{aligned}\frac{\partial \ell}{\partial \mathbf{W}_{O(i,\cdot)}} &= \frac{\partial \ell}{\partial F(\mathbf{X}^{(n)})} \cdot \frac{\partial F(\mathbf{X}^{(n)})}{\partial \mathbf{W}_{O(i,\cdot)}} \\ &= -\frac{z^{(n)}}{L} \sum_{l=1}^L v_i \cdot \phi'\left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}\right) \cdot \mathbf{y}_l^{(n)}.\end{aligned}\quad (98)$$

If we consider the gradient for the population loss,

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}} = -\mathbb{E} \left[\frac{z^{(n)}}{L} \sum_{l=1}^L v_i \cdot \phi'\left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}\right) \cdot \mathbf{y}_l^{(n)} \right] \quad (99)$$

$$\begin{aligned}&= -\mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi'\left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}\right) \cdot \mathbf{y}_l^{(n)} \right] \\ &\quad + \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi'\left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}\right) \cdot \mathbf{y}_l^{(n)} \right].\end{aligned}\quad (100)$$

We are given that

$$p_1 \leq \langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_+ \rangle \leq q_1, \quad p_2 \leq \langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_- \rangle \leq q_2, \quad \text{and} \quad p_3 \leq \langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_j \rangle \leq q_3 \quad \text{for } j \neq 1, 2. \quad (101)$$

The Mamba output can be written as

$$\mathbf{y}_l(t) = \sum_{s=1}^l \left(\prod_{j=s+1}^l \left(1 - \sigma(\mathbf{w}_\Delta^{(t)^\top} \mathbf{x}_j) \right) \right) \cdot \sigma(\mathbf{w}_\Delta^{(t)^\top} \mathbf{x}_s) \cdot (\mathbf{x}_s^\top \mathbf{x}_l) \mathbf{x}_s \quad (102)$$

We have to consider 4 cases.

Case I: $l = s = L_1^+$

$$\mathbf{x}_s = \mathbf{x}_l = \mathbf{o}_+ \quad (103)$$

$$\langle \mathbb{E} \mathbf{y}_{L_1^+}, \mathbf{o}_+ \rangle = \sigma(\mathbf{w}_\Delta^{(t)^\top} \mathbf{o}_+) \geq \sigma(p_1) = \frac{1}{1 + e^{-p_1}} \quad (104)$$

Case II: $l = s = L_2^+$

$$\left\langle \mathbb{E} \mathbf{y}_{L_2^+, L_2^+}, \mathbf{o}_+ \right\rangle = \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \quad (105)$$

Case III: $l = L_2^+, s = L_1^+$

$$\begin{aligned} \left\langle \mathbb{E} \mathbf{y}_{L_2^+, L_1^+}, \mathbf{o}_+ \right\rangle &= \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+)\right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-)\right) \\ &\quad \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j)\right)^{L_2^+ - L_1^+ - 2} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \end{aligned} \quad (106)$$

Combining (105) and (106), we obtain

$$\begin{aligned} \left\langle \mathbb{E} \mathbf{y}_{L_2^+}, \mathbf{o}_+ \right\rangle &= \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \\ &\quad + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+)\right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-)\right) \\ &\quad \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j)\right)^{L_2^+ - L_1^+ - 2} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+). \end{aligned} \quad (107)$$

Case IV: Others

For the other token positions, $\mathbf{x}_l \neq \mathbf{o}_+$. Since we assume orthogonality among the features, $\mathbf{y}_l = 0$.

From our initialization, for the lucky neuron $i \in \mathcal{W}(0)$, $v_i = +\frac{1}{\sqrt{m}}$. For $i \in \mathcal{W}(0)$, and $z^{(n)} = +1$, we have

$$\begin{aligned} &\left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_+ \right\rangle \\ &= \frac{1}{\sqrt{m}L} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \left[2 + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+)\right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-)\right) \right. \\ &\quad \left. \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j)\right)^{L_2^+ - L_1^+ - 2} \right]. \end{aligned} \quad (108)$$

For $z = -1$, when $l = s = L^+$, $\mathbf{x}_s = \mathbf{x}_l = \mathbf{o}_+$. Therefore, we obtain

$$\left\langle \mathbb{E} \mathbf{y}_{L^+}, \mathbf{o}_+ \right\rangle = \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \geq \sigma(p_1) \quad (109)$$

$$\left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_+ \right\rangle = \frac{1}{\sqrt{m}L} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \quad (110)$$

Therefore, combining (108) and (110),

$$\begin{aligned} \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle &= \left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_+ \right\rangle \\ &\quad - \left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_+ \right\rangle \\ &= \frac{1}{\sqrt{m}L} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \left[1 + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+)\right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-)\right) \right. \\ &\quad \left. \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j)\right)^{L_2^+ - L_1^+ - 2} \right]. \end{aligned} \quad (111)$$

We aim to bound the deviation between the gradient of the population loss and that of the empirical loss. Specifically, $\left\| \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} - \frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} \right\|_2 = \left\| \frac{1}{N} \sum_{n=1}^N \boldsymbol{\gamma}_n - \mathbb{E} \boldsymbol{\gamma}_n \right\|_2$, where

$$\boldsymbol{\gamma}_n = \frac{z^{(n)}}{L} \sum_{l=1}^L v_i \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \mathbf{y}_l^{(n)}. \quad (112)$$

Consider a fixed vector $\boldsymbol{\alpha}$ with $\|\boldsymbol{\alpha}\|_2 = 1$. We will show that $\boldsymbol{\alpha}^\top \boldsymbol{\gamma}_n$ is a sub-Gaussian random variable.

$$|\boldsymbol{\alpha}^\top \boldsymbol{\gamma}_n| \leq \|\boldsymbol{\alpha}\|_2 \cdot \|\boldsymbol{\gamma}_n\|_2 = \|\boldsymbol{\gamma}_n\|_2. \quad (113)$$

By the problem setup, we know that

$$|v_i| = \frac{1}{\sqrt{m}}, \quad |z^{(n)}| = 1, \quad \left| \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \right| \leq 1. \quad (114)$$

Recall the Mamba output,

$$\mathbf{y}_l^{(n)}(t) = \sum_{s=1}^l \left(\prod_{j=s+1}^l \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{x}_j) \right) \right) \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{x}_s) \cdot (\mathbf{x}_s^\top \mathbf{x}_l) \mathbf{x}_s. \quad (115)$$

Since $\|\mathbf{x}_s\|_2 = 1$, we get

$$\begin{aligned} \|\mathbf{y}_l^{(n)}\|_2 &\leq \sum_{s=1}^l |a^{l-s+1} \cdot (\mathbf{x}_s^\top \mathbf{x}_l)| \cdot \|\mathbf{x}_s\|_2 \\ &\leq \sum_{s=1}^l \frac{a}{1-a} \cdot 1 \cdot 1 = a' \quad (\text{where } a' \text{ denotes a constant}). \end{aligned} \quad (116)$$

Therefore, the norm of $\boldsymbol{\gamma}_n$ satisfies

$$\begin{aligned} \|\boldsymbol{\gamma}_n\|_2 &\leq \frac{1}{L} \sum_{l=1}^L |v_i| \cdot \left| \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \right| \cdot \|\mathbf{y}_l^{(n)}\|_2 \\ &\leq \frac{1}{L} \cdot \frac{1}{\sqrt{m}} \sum_{l=1}^L \|\mathbf{y}_l^{(n)}\|_2 \\ &\leq \frac{1}{L} \cdot \frac{1}{\sqrt{m}} \cdot \sum_{l=1}^L a' = \frac{a'}{\sqrt{m}}. \end{aligned} \quad (117)$$

Hence,

$$|\boldsymbol{\alpha}^\top \boldsymbol{\gamma}_n| \leq \frac{a'}{\sqrt{m}} \quad (\text{bounded}). \quad (118)$$

This implies that $\boldsymbol{\alpha}^\top \boldsymbol{\gamma}_n$ is sub-Gaussian with variance proxy

$$\sigma^2 = \mathcal{O} \left(\frac{1}{m} \right). \quad (119)$$

Now consider the independent sub-Gaussian variables $\boldsymbol{\alpha}^\top \boldsymbol{\gamma}_1, \dots, \boldsymbol{\alpha}^\top \boldsymbol{\gamma}_N$, each bounded as

$$-\frac{1}{\sqrt{m}} \leq \boldsymbol{\alpha}^\top \boldsymbol{\gamma}_n \leq \frac{1}{\sqrt{m}}. \quad (120)$$

Applying Hoeffding's inequality, for any $q > 0$,

$$\mathbb{P} \left(\left| \frac{1}{N} \sum_{n=1}^N \boldsymbol{\alpha}^\top \boldsymbol{\gamma}_n - \mathbb{E} \boldsymbol{\alpha}^\top \boldsymbol{\gamma}_n \right| \gtrsim \sqrt{\frac{q \log N}{mN}} \right) \leq N^{-q}. \quad (121)$$

Observe that this can be written as

$$\frac{1}{N} \sum_{n=1}^N \boldsymbol{\alpha}^\top \boldsymbol{\gamma}_n - \mathbb{E} \boldsymbol{\alpha}^\top \boldsymbol{\gamma}_n = \boldsymbol{\alpha}^\top \left(\frac{1}{N} \sum_{n=1}^N \boldsymbol{\gamma}_n - \mathbb{E} \boldsymbol{\gamma}_n \right) := \boldsymbol{\alpha}^\top \boldsymbol{\zeta}. \quad (122)$$

Therefore, by Hoeffding's inequality (cf. (121)),

$$\mathbb{P} \left(|\boldsymbol{\alpha}^\top \boldsymbol{\zeta}| \gtrsim \sqrt{\frac{q \log N}{mN}} \right) \leq N^{-q}. \quad (123)$$

To bound $\|\boldsymbol{\zeta}\|_2$, we use the dual norm identity

$$\|\boldsymbol{\zeta}\|_2 = \sup_{\|\boldsymbol{\alpha}\|_2=1} \boldsymbol{\alpha}^\top \boldsymbol{\zeta}. \quad (124)$$

We apply an ε -cover argument to obtain

$$\begin{aligned} \sup_{\|\boldsymbol{\alpha}\|_2=1} \boldsymbol{\alpha}^\top \boldsymbol{\zeta} &\leq \frac{1}{1-\varepsilon} \max_{\boldsymbol{\alpha} \in \mathcal{C}_\varepsilon} \boldsymbol{\alpha}^\top \boldsymbol{\zeta} \\ &\leq 2 \max_{\boldsymbol{\alpha} \in \mathcal{C}_{1/2}} \boldsymbol{\alpha}^\top \boldsymbol{\zeta}. \end{aligned} \quad (125)$$

We have shown that for any fixed $\boldsymbol{\alpha}$,

$$\mathbb{P} \left(|\boldsymbol{\alpha}^\top \boldsymbol{\zeta}| \gtrsim \sqrt{\frac{q \log N}{mN}} \right) \leq N^{-q}. \quad (126)$$

Therefore, for all fixed $\boldsymbol{\alpha} \in \mathcal{C}_{1/2}$,

$$|\boldsymbol{\alpha}^\top \boldsymbol{\zeta}| \gtrsim \sqrt{\frac{q \log N}{mN}} \quad \text{with probability at most } N^{-q}. \quad (127)$$

Then,

$$\max_{\boldsymbol{\alpha} \in \mathcal{C}_{1/2}} |\boldsymbol{\alpha}^\top \boldsymbol{\zeta}| \gtrsim \sqrt{\frac{q \log N}{mN}} \quad \text{with probability at most } |\mathcal{C}_{1/2}| N^{-q}. \quad (128)$$

Recall that the covering number satisfies

$$|\mathcal{C}_\varepsilon| \leq \left(\frac{3B}{\varepsilon} \right)^d. \quad (129)$$

For $B = 1$ and $\varepsilon = \frac{1}{2}$, we have

$$|\mathcal{C}_{1/2}| \leq 6^d. \quad (130)$$

We can therefore write

$$\mathbb{P} \left(\|\boldsymbol{\zeta}\|_2 \gtrsim \sqrt{\frac{q \log N}{mN}} \right) \leq 6^d \cdot N^{-q}. \quad (131)$$

We want this probability to be sufficiently small. Set $q = d$, so that

$$\mathbb{P} \left(\|\boldsymbol{\zeta}\|_2 \gtrsim 2 \sqrt{\frac{d \log N}{mN}} \right) \leq \left(\frac{N}{6} \right)^{-d}. \quad (132)$$

Hence, the deviation is bounded with high probability:

$$\|\zeta\|_2 > \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) \quad \text{with probability at most } \mathcal{O}(N^{-d}). \quad (133)$$

Or equivalently, with probability at most $\mathcal{O}(N^{-d})$,

$$\left\| \frac{1}{N} \sum_{n=1}^N \gamma_n - \mathbb{E}\gamma_n \right\|_2 > \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (134)$$

That is, with high probability $1 - \mathcal{O}(N^{-d})$, we have

$$\left\| \frac{1}{N} \sum_{n=1}^N \gamma_n - \mathbb{E}\gamma_n \right\|_2 \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (135)$$

Using the identities

$$-\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} = \frac{1}{N} \sum_{n=1}^N \gamma_n, \quad -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} = \mathbb{E}\gamma_n, \quad (136)$$

we conclude that, with high probability,

$$\left\| \left(-\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} \right) - \left(-\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} \right) \right\|_2 = \left\| \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} - \frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} \right\|_2 \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (137)$$

Using the Cauchy–Schwarz inequality, we have

$$\begin{aligned} \left| \left\langle \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} - \frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \right| &\leq \left\| \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} - \frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} \right\|_2 \cdot \|\mathbf{o}_+\|_2 \\ &= \left\| \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} - \frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} \right\|_2 \quad (\text{since } \|\mathbf{o}_+\|_2 = 1) \\ &\leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \end{aligned} \quad (138)$$

Therefore, we obtain

$$\begin{aligned} \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) &\leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \\ &\leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle + \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \end{aligned} \quad (139)$$

By pairing (111) with the given the conditions on \mathbf{w}_Δ in (101), we can write

$$\frac{1}{\sqrt{mL}} \cdot \sigma(p_1) \left[1 + (1 - \sigma(q_1)) \cdot (1 - \sigma(q_2)) \cdot (1 - \sigma(q_3))^{L_2^+ - L_1^+ - 2} \right] \leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \quad (140)$$

$$\text{and } \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \leq \frac{1}{\sqrt{mL}} \cdot \sigma(q_1) \left[1 + (1 - \sigma(p_1)) \cdot (1 - \sigma(p_2)) \cdot (1 - \sigma(p_3))^{L_2^+ - L_1^+ - 2} \right]. \quad (141)$$

Therefore, we can obtain the lower bound and the upper bound of $\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle$ as

$$\begin{aligned} & \frac{1}{\sqrt{mL}} \cdot \sigma(p_1) \left[1 + (1 - \sigma(q_1)) \cdot (1 - \sigma(q_2)) \cdot (1 - \sigma(q_3))^{L_2^+ - L_1^+ - 2} \right] - \mathcal{O} \left(\sqrt{\frac{d \log N}{mN}} \right) \\ & \leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \end{aligned} \quad (142)$$

$$\begin{aligned} \text{and } & \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \leq \frac{1}{\sqrt{mL}} \cdot \sigma(q_1) \left[1 + (1 - \sigma(p_1)) \cdot (1 - \sigma(p_2)) \cdot (1 - \sigma(p_3))^{L_2^+ - L_1^+ - 2} \right] \\ & + \mathcal{O} \left(\sqrt{\frac{d \log N}{mN}} \right). \end{aligned} \quad (143)$$

This concludes the proof of (9) and (10) in Lemma 1.

To obtain $\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle$, we have to consider $\mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right]$.

If $\mathbf{W}_{O(i,\cdot)}^{(t)} \mathbf{o}_- > 0$,

$$\begin{aligned} & \left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_- \right\rangle \\ & = \frac{1}{\sqrt{mL}} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \left[2 + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \right) \right. \\ & \quad \left. \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j) \right)^{L_2^- - L_1^- - 2} \right]. \end{aligned} \quad (144)$$

If $\mathbf{W}_{O(i,\cdot)}^{(t)} \mathbf{o}_- \leq 0$,

$$\left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_- \right\rangle = 0. \quad (145)$$

From (100), We know that

$$\begin{aligned} & \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle = \left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_- \right\rangle \\ & \quad - \left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_- \right\rangle. \end{aligned} \quad (146)$$

Hence, combining both cases, we conclude

$$\begin{aligned} & -\frac{1}{\sqrt{mL}} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \left[2 + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \right) \right. \\ & \quad \left. \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j) \right)^{L_2^- - L_1^- - 2} \right] \\ & \leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \leq 0. \end{aligned} \quad (147)$$

From (137), similar to (139), we can write

$$\begin{aligned}
& \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) \\
& \leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \\
& \leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle + \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \tag{148}
\end{aligned}$$

Hence,

$$\begin{aligned}
& -\frac{1}{\sqrt{mL}} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \left[2 + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \right) \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \right) \right. \\
& \quad \left. \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j) \right)^{L_2^- - L_1^- - 2} \right] - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) \\
& \leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \tag{149}
\end{aligned}$$

This concludes the proof of (11) and (12) in Lemma 1.

Now consider $\left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle$ for $j \neq 1, 2$.

$$\begin{aligned}
\left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle &= \left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_j \right\rangle \\
&\quad - \left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_j \right\rangle \\
&:= \langle I_1, \mathbf{o}_j \rangle - \langle I_2, \mathbf{o}_j \rangle. \tag{150}
\end{aligned}$$

Because \mathbf{o}_j for $j \neq 1, 2$ is identical in both I_1 and I_2 , $\langle I_1, \mathbf{o}_j \rangle - \langle I_2, \mathbf{o}_j \rangle = 0$. Hence, $\left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle = 0$. From (137), similar to (139), we can write

$$\begin{aligned}
& \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) \\
& \leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle \\
& \leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle + \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \tag{151}
\end{aligned}$$

Therefore,

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) \text{ for } j \neq 1, 2. \tag{152}$$

This concludes the proof of (13) in Lemma 1. \square

F.2 Proof of Lemma 2

Proof. By definition, for any unlucky neuron $i \in \mathcal{K}_+ \setminus \mathcal{W}(0)$, we have

$$\mathbf{W}_{O(i,\cdot)} \mathbf{o}_+ \leq 0. \quad (153)$$

We first consider the alignment with \mathbf{o}_+ . That is,

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle. \quad (154)$$

The gradient is given in (99). We only need to consider the cases where $\left\langle \mathbf{y}_l^{(n)}, \mathbf{o}_+ \right\rangle > 0$. However, since $\mathbf{W}_{O(i,\cdot)} \mathbf{o}_+ \leq 0$, we have

$$\phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) = 0. \quad (155)$$

$$\begin{aligned} \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle &= \left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_+ \right\rangle \\ &\quad - \left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_+ \right\rangle \\ &= 0. \end{aligned} \quad (156)$$

We know by (139),

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle + \mathcal{O} \left(\sqrt{\frac{d \log N}{mN}} \right). \quad (157)$$

Hence,

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \leq \mathcal{O} \left(\sqrt{\frac{d \log N}{mN}} \right). \quad (158)$$

We now analyze the alignment with \mathbf{o}_- . To obtain the bound on $\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle$, we consider the expectation $\mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right]$.

If $\mathbf{W}_{O(i,\cdot)}^{(t)} \mathbf{o}_- > 0$, the inner product satisfies

$$\begin{aligned} &\left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_- \right\rangle \\ &= \frac{1}{\sqrt{mL}} \cdot \sigma \left(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_- \right) \left[2 + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \right) \right. \\ &\quad \left. \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j) \right)^{L_2^- - L_1^- - 2} \right]. \end{aligned} \quad (159)$$

If $\mathbf{W}_{O(i,\cdot)}^{(t)} \mathbf{o}_- \leq 0$, then

$$\left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_- \right\rangle = 0. \quad (160)$$

From (100), We know that

$$\begin{aligned} \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle &= \left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_- \right\rangle \\ &\quad - \left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_- \right\rangle. \end{aligned} \quad (161)$$

Hence, combining both cases, we conclude

$$\begin{aligned} &- \frac{1}{\sqrt{mL}} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \left[2 + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \right) \right. \\ &\quad \left. \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j) \right)^{L_2^- - L_1^- - 2} \right] \\ &\leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \leq 0. \end{aligned} \quad (162)$$

From (137), similar to (139), we can write

$$\begin{aligned} &\left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle - \mathcal{O} \left(\sqrt{\frac{d \log N}{mN}} \right) \\ &\leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \\ &\leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle + \mathcal{O} \left(\sqrt{\frac{d \log N}{mN}} \right). \end{aligned} \quad (163)$$

Hence,

$$\begin{aligned} &- \frac{1}{\sqrt{mL}} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \left[2 + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \right) \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \right) \right. \\ &\quad \left. \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j) \right)^{L_2^- - L_1^- - 2} \right] - \mathcal{O} \left(\sqrt{\frac{d \log N}{mN}} \right) \\ &\leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \leq \mathcal{O} \left(\sqrt{\frac{d \log N}{mN}} \right). \end{aligned} \quad (164)$$

Now consider $\left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle$ for $j \neq 1, 2$.

$$\begin{aligned} \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle &= \left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_j \right\rangle \\ &\quad - \left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_j \right\rangle \\ &:= \langle I_1, \mathbf{o}_j \rangle - \langle I_2, \mathbf{o}_j \rangle. \end{aligned} \quad (165)$$

Because \mathbf{o}_j for $j \neq 1, 2$ is identical in both I_1 and I_2 , $\langle I_1, \mathbf{o}_j \rangle - \langle I_2, \mathbf{o}_j \rangle = 0$. Hence, $\left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle = 0$. From (137), similar to (139), we can write

$$\begin{aligned} & \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) \\ & \leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle \\ & \leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle + \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \end{aligned} \quad (166)$$

Therefore,

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) \text{ for } j \neq 1, 2. \quad (167)$$

□

E.3 Proof of Lemma 3

Proof. We know that the gradient of the loss function for the n^{th} sample is

$$\begin{aligned} \frac{\partial \ell}{\partial \mathbf{W}_{O(i,\cdot)}} &= \frac{\partial \ell}{\partial F(\mathbf{X}^{(n)})} \cdot \frac{\partial F(\mathbf{X}^{(n)})}{\partial \mathbf{W}_{O(i,\cdot)}} \\ &= -\frac{z^{(n)}}{L} \sum_{l=1}^L v_i \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \cdot \mathbf{y}_l^{(n)}. \end{aligned} \quad (168)$$

If we consider the gradient for the population loss,

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}} = -\mathbb{E}\left[\frac{z^{(n)}}{L} \sum_{l=1}^L v_i \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \cdot \mathbf{y}_l^{(n)}\right] \quad (169)$$

$$\begin{aligned} &= -\mathbb{E}_{z=+1}\left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \cdot \mathbf{y}_l^{(n)}\right] \\ &\quad + \mathbb{E}_{z=-1}\left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \cdot \mathbf{y}_l^{(n)}\right]. \end{aligned} \quad (170)$$

We are given that

$$p_1 \leq \left\langle \mathbf{w}_{\Delta}^{(t)}, \mathbf{o}_- \right\rangle \leq q_1, \quad p_2 \leq \left\langle \mathbf{w}_{\Delta}^{(t)}, \mathbf{o}_+ \right\rangle \leq q_2, \quad \text{and} \quad p_3 \leq \left\langle \mathbf{w}_{\Delta}^{(t)}, \mathbf{o}_j \right\rangle \leq q_3 \quad \text{for } j \neq 1, 2. \quad (171)$$

The Mamba output can be written as

$$\mathbf{y}_l(t) = \sum_{s=1}^l \left(\prod_{j=s+1}^l \left(1 - \sigma(\mathbf{w}_{\Delta}^{(t)\top} \mathbf{x}_j) \right) \right) \cdot \sigma(\mathbf{w}_{\Delta}^{(t)\top} \mathbf{x}_s) \cdot (\mathbf{x}_s^\top \mathbf{x}_l) \mathbf{x}_s \quad (172)$$

We have to consider 4 cases.

Case I: $l = s = L_1^-$

$$\mathbf{x}_s = \mathbf{x}_l = \mathbf{o}_- \quad (173)$$

$$\left\langle \mathbb{E} \mathbf{y}_{L_1^+}, \mathbf{o}_- \right\rangle = \sigma(\mathbf{w}_{\Delta}^{(t)\top} \mathbf{o}_-) \geq \sigma(p_1) = \frac{1}{1 + e^{-p_1}} \quad (174)$$

Case II: $l = s = L_2^-$

$$\left\langle \mathbb{E} \mathbf{y}_{L_2^-, L_2^-}, \mathbf{o}_- \right\rangle = \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \quad (175)$$

Case III: $l = L_2^-, s = L_1^-$

$$\begin{aligned} \left\langle \mathbb{E} \mathbf{y}_{L_2^-, L_1^-}, \mathbf{o}_- \right\rangle &= \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \right) \\ &\quad \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j) \right)^{L_2^- - L_1^- - 2} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \end{aligned} \quad (176)$$

Combining (175) and (176), we obtain

$$\left\langle \mathbb{E} \mathbf{y}_{L_2^-}, \mathbf{o}_- \right\rangle = \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \quad (177)$$

$$\begin{aligned} &+ \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \right) \\ &\quad \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j) \right)^{L_2^- - L_1^- - 2} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-). \end{aligned} \quad (178)$$

Case IV: Others

For the other token positions, $\mathbf{x}_l \neq \mathbf{o}_-$. Since we assume orthogonality among the features, $\mathbf{y}_l = 0$.

From our initialization, for the lucky neuron $i \in \mathcal{U}(0)$, $v_i = -\frac{1}{\sqrt{m}}$. For $i \in \mathcal{U}(0)$, and $z^{(n)} = -1$, we have

$$\begin{aligned} &\left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_- \right\rangle \\ &= \frac{1}{\sqrt{m}L} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \left[2 + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \right) \right. \\ &\quad \left. \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j) \right)^{L_2^- - L_1^- - 2} \right]. \end{aligned} \quad (179)$$

For $z = +1$, when $l = s = L^-$, $\mathbf{x}_s = \mathbf{x}_l = \mathbf{o}_-$. Therefore, we obtain

$$\left\langle \mathbb{E} \mathbf{y}_{L^-}, \mathbf{o}_- \right\rangle = \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \geq \sigma(p_1) \quad (180)$$

$$\left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_- \right\rangle = \frac{1}{\sqrt{m}L} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \quad (181)$$

Therefore, combining (179) and (181),

$$\begin{aligned}
\left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle &= \left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_- \right\rangle \\
&\quad - \left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_- \right\rangle \\
&= \frac{1}{\sqrt{mL}} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \left[1 + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \right) \right. \\
&\quad \left. \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j) \right)^{L_2^- - L_1^- - 2} \right].
\end{aligned} \tag{182}$$

We want to bound the deviation between the population loss gradient and the empirical loss gradient.

That is $\left\| \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} - \frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} \right\|_2 = \left\| \frac{1}{N} \sum_{n=1}^N \boldsymbol{\gamma}_n - \mathbb{E} \boldsymbol{\gamma}_n \right\|_2$, where

$$\boldsymbol{\gamma}_n = \frac{z^{(n)}}{L} \sum_{l=1}^L v_i \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \mathbf{y}_l^{(n)}. \tag{183}$$

That is, w.h.p. $1 - \mathcal{O}(N^{-d})$, we have $\left\| \frac{1}{N} \sum_{n=1}^N \boldsymbol{\gamma}_n - \mathbb{E} \boldsymbol{\gamma}_n \right\|_2 \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right)$. Hence, we can write

$$\begin{aligned}
&\left| \left\langle \left(-\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} \right) - \left(-\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} \right), \mathbf{o}_- \right\rangle \right| \\
&= \left| \left\langle \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}} - \frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \right| \\
&\leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right).
\end{aligned} \tag{184}$$

Therefore, we obtain

$$\begin{aligned}
\left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) &\leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \\
&\leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle + \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right).
\end{aligned} \tag{185}$$

By pairing (182) with the given conditions on \mathbf{w}_Δ in (171), we can write

$$\frac{1}{\sqrt{mL}} \cdot \sigma(p_1) \left[1 + (1 - \sigma(q_1)) \cdot (1 - \sigma(q_2)) \cdot (1 - \sigma(q_3))^{L_2^- - L_1^- - 2} \right] \leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \tag{186}$$

$$\text{and } \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \leq \frac{1}{\sqrt{mL}} \cdot \sigma(q_1) \left[1 + (1 - \sigma(p_1)) \cdot (1 - \sigma(p_2)) \cdot (1 - \sigma(p_3))^{L_2^- - L_1^- - 2} \right]. \tag{187}$$

Therefore, we can obtain the lower bound and the upper bound of $\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle$ as

$$\begin{aligned} & \frac{1}{\sqrt{mL}} \cdot \sigma(p_1) \left[1 + (1 - \sigma(q_1)) \cdot (1 - \sigma(q_2)) \cdot (1 - \sigma(q_3))^{L_2^- - L_1^- - 2} \right] - \mathcal{O} \left(\sqrt{\frac{d \log N}{mN}} \right) \\ & \leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \end{aligned} \quad (188)$$

$$\begin{aligned} \text{and } & \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \leq \frac{1}{\sqrt{mL}} \cdot \sigma(q_1) \left[1 + (1 - \sigma(p_1)) \cdot (1 - \sigma(p_2)) \cdot (1 - \sigma(p_3))^{L_2^- - L_1^- - 2} \right] \\ & + \mathcal{O} \left(\sqrt{\frac{d \log N}{mN}} \right). \end{aligned} \quad (189)$$

This concludes the proof of (17) and (18) in Lemma 3.

To obtain $\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle$, we have to consider $\mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right]$.

If $\mathbf{W}_{O(i,\cdot)}^{(t)} \mathbf{o}_+ > 0$,

$$\begin{aligned} & \left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_+ \right\rangle \\ & = \frac{1}{\sqrt{mL}} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \left[2 + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \right) \right. \\ & \quad \left. \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j) \right)^{L_2^+ - L_1^+ - 2} \right]. \end{aligned} \quad (190)$$

If $\mathbf{W}_{O(i,\cdot)}^{(t)} \mathbf{o}_+ \leq 0$,

$$\left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_+ \right\rangle = 0. \quad (191)$$

From (170), We know that

$$\begin{aligned} & \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle = \left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_+ \right\rangle \\ & \quad - \left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_+ \right\rangle. \end{aligned} \quad (192)$$

Hence, combining both cases, we conclude

$$\begin{aligned} & -\frac{1}{\sqrt{mL}} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \left[2 + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \right) \right. \\ & \quad \left. \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j) \right)^{L_2^+ - L_1^+ - 2} \right] \\ & \leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \leq 0. \end{aligned} \quad (193)$$

From (137), similar to (139), we can write

$$\begin{aligned}
& \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) \\
& \leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \\
& \leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle + \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \tag{194}
\end{aligned}$$

Hence,

$$\begin{aligned}
& -\frac{1}{\sqrt{mL}} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \left[2 + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \right) \right. \\
& \quad \left. \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j) \right)^{L_2^+ - L_1^+ - 2} \right] - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) \\
& \leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \tag{195}
\end{aligned}$$

This concludes the proof of (19) and (20) in Lemma 3.

Now consider $\left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle$ for $j \neq 1, 2$.

$$\begin{aligned}
\left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle &= \left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_j \right\rangle \\
&\quad - \left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_j \right\rangle \\
&:= \langle I_1, \mathbf{o}_j \rangle - \langle I_2, \mathbf{o}_j \rangle. \tag{196}
\end{aligned}$$

Because \mathbf{o}_j for $j \neq 1, 2$ is identical in both I_1 and I_2 , $\langle I_1, \mathbf{o}_j \rangle - \langle I_2, \mathbf{o}_j \rangle = 0$. Hence, $\left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle = 0$. From (137), similar to (139), we can write

$$\begin{aligned}
& \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) \\
& \leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle \\
& \leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle + \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \tag{197}
\end{aligned}$$

Therefore,

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) \text{ for } j \neq 1, 2. \tag{198}$$

This concludes the proof of (21) in Lemma 3. \square

F.4 Proof of Lemma 4

Proof. By definition, for any unlucky neuron $i \in \mathcal{K}_- \setminus \mathcal{U}(0)$, we have

$$\mathbf{W}_{O(i,\cdot)} \mathbf{o}_- \leq 0. \quad (199)$$

We first consider the alignment with \mathbf{o}_- . That is,

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle. \quad (200)$$

The gradient is given in (169). We only need to consider the cases where $\left\langle \mathbf{y}_l^{(n)}, \mathbf{o}_- \right\rangle > 0$. However, since $\mathbf{W}_{O(i,\cdot)} \mathbf{o}_- \leq 0$, we have

$$\phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) = 0. \quad (201)$$

$$\begin{aligned} \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle &= \left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_- \right\rangle \\ &\quad - \left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_- \right\rangle \\ &= 0. \end{aligned} \quad (202)$$

We know by (139),

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle + \mathcal{O} \left(\sqrt{\frac{d \log N}{mN}} \right). \quad (203)$$

Hence,

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_- \right\rangle \leq \mathcal{O} \left(\sqrt{\frac{d \log N}{mN}} \right). \quad (204)$$

We now analyze the alignment with \mathbf{o}_+ . To obtain the bound on $\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle$, we consider the expectation $\mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right]$.

If $\mathbf{W}_{O(i,\cdot)}^{(t)} \mathbf{o}_+ > 0$, the inner product satisfies

$$\begin{aligned} &\left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_+ \right\rangle \\ &= \frac{1}{\sqrt{mL}} \cdot \sigma \left(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+ \right) \left[2 + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \right) \right. \\ &\quad \left. \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j) \right)^{L_2^+ - L_1^+ - 2} \right]. \end{aligned} \quad (205)$$

If $\mathbf{W}_{O(i,\cdot)}^{(t)} \mathbf{o}_+ \leq 0$, then

$$\left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_+ \right\rangle = 0. \quad (206)$$

From (100), We know that

$$\begin{aligned} \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle &= \left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_+ \right\rangle \\ &\quad - \left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_+ \right\rangle. \end{aligned} \quad (207)$$

Hence, combining both cases, we conclude

$$\begin{aligned} &- \frac{1}{\sqrt{mL}} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \left[2 + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \right) \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \right) \right. \\ &\quad \left. \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j) \right)^{L_2^+ - L_1^+ - 2} \right] \\ &\leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \leq 0. \end{aligned} \quad (208)$$

From (137), similar to (139), we can write

$$\begin{aligned} &\left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) \\ &\leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \\ &\leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle + \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \end{aligned} \quad (209)$$

Hence,

$$\begin{aligned} &- \frac{1}{\sqrt{mL}} \cdot \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \left[2 + \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_+) \right) \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_-) \right) \right. \\ &\quad \left. \cdot \left(1 - \sigma(\mathbf{w}_\Delta^{(t)\top} \mathbf{o}_j) \right)^{L_2^+ - L_1^+ - 2} \right] - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) \\ &\leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_+ \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \end{aligned} \quad (210)$$

Now consider $\left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle$ for $j \neq 1, 2$.

$$\begin{aligned} \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle &= \left\langle \mathbb{E}_{z=+1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_j \right\rangle \\ &\quad - \left\langle \mathbb{E}_{z=-1} \left[\sum_{l=1}^L \frac{1}{L} v_i \cdot \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \mathbf{y}_l^{(n)} \right], \mathbf{o}_j \right\rangle \\ &:= \langle I_1, \mathbf{o}_j \rangle - \langle I_2, \mathbf{o}_j \rangle. \end{aligned} \quad (211)$$

Because \mathbf{o}_j for $j \neq 1, 2$ is identical in both I_1 and I_2 , $\langle I_1, \mathbf{o}_j \rangle - \langle I_2, \mathbf{o}_j \rangle = 0$. Hence, $\left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle = 0$. From (137), similar to (139), we can write

$$\begin{aligned} & \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) \\ & \leq \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle \\ & \leq \left\langle -\frac{\partial \mathcal{L}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle + \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \end{aligned} \quad (212)$$

Therefore,

$$\left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_{O(i,\cdot)}^{(t)}}, \mathbf{o}_j \right\rangle \leq \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) \text{ for } j \neq 1, 2. \quad (213)$$

□

F.5 Proof of Lemma 5

Proof. The gradient of the loss with respect to \mathbf{w}_Δ for the n^{th} sample is given by

$$\begin{aligned} \frac{\partial \ell}{\partial \mathbf{w}_\Delta} &= -\frac{z^{(n)}}{L} \cdot \sum_{i=1}^m \sum_{l=1}^L v_i \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \sum_{s=1}^l \left(\mathbf{W}_B^\top \mathbf{x}_s^{(n)} \right)^\top \left(\mathbf{W}_C^\top \mathbf{x}_l^{(n)} \right) \left(\mathbf{W}_{O(i,\cdot)} \mathbf{x}_s^{(n)} \right) \\ &\quad \cdot \sigma \left(\mathbf{w}_\Delta^\top \mathbf{x}_s^{(n)} \right) \cdot \prod_{r=s+1}^l \left(1 - \sigma \left(\mathbf{w}_\Delta^\top \mathbf{x}_r^{(n)} \right) \right) \\ &\quad \cdot \left[\left(1 - \sigma \left(\mathbf{w}_\Delta^\top \mathbf{x}_s^{(n)} \right) \right) \mathbf{x}_s^{(n)} - \sum_{j=s+1}^l \left(1 - \sigma \left(\mathbf{w}_\Delta^\top \mathbf{x}_j^{(n)} \right) \right) \mathbf{x}_j^{(n)} \right] \\ &:= -\frac{z^{(n)}}{L} \cdot \sum_{i=1}^m \sum_{l=1}^L v_i \phi' \left(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)} \right) \cdot \sum_{s=1}^l \mathbf{I}_{l,s}^{(n)}. \end{aligned} \quad (214)$$

We define the gradient summand $\mathbf{I}_{l,s}^{(n)}$ as

$$\mathbf{I}_{l,s}^{(n)} = \beta_{s,s} \cdot \mathbf{x}_s^{(n)} - \sum_{j=s+1}^l \beta_{s,j} \mathbf{x}_j^{(n)}, \quad (215)$$

where the coefficients $\beta_{s,s}$ and $\beta_{s,j}$ are given by

$$\begin{aligned} \beta_{s,s} &= (\mathbf{W}_B^\top \mathbf{x}_s^{(n)})^\top (\mathbf{W}_C^\top \mathbf{x}_l^{(n)}) (\mathbf{W}_{O(i,\cdot)} \mathbf{x}_s^{(n)}) \sigma(\mathbf{w}_\Delta^\top \mathbf{x}_s^{(n)}) \\ &\quad \times \left[\prod_{r=s+1}^l \left(1 - \sigma(\mathbf{w}_\Delta^\top \mathbf{x}_r^{(n)}) \right) \right] (1 - \sigma(\mathbf{w}_\Delta^\top \mathbf{x}_s^{(n)})). \end{aligned} \quad (216)$$

and

$$\begin{aligned} \beta_{s,j} &= (\mathbf{W}_B^\top \mathbf{x}_s^{(n)})^\top (\mathbf{W}_C^\top \mathbf{x}_l^{(n)}) (\mathbf{W}_{O(i,\cdot)} \mathbf{x}_s^{(n)}) \sigma(\mathbf{w}_\Delta^\top \mathbf{x}_s^{(n)}) \\ &\quad \times \left[\prod_{r=s+1}^l \left(1 - \sigma(\mathbf{w}_\Delta^\top \mathbf{x}_r^{(n)}) \right) \right] (1 - \sigma(\mathbf{w}_\Delta^\top \mathbf{x}_j^{(n)})). \end{aligned} \quad (217)$$

If we consider the gradient of the empirical loss,

$$\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{w}_\Delta} = -\frac{1}{N} \sum_{n=1}^N \frac{z^{(n)}}{L} \cdot \sum_{i=1}^m \sum_{l=1}^L v_i \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \cdot \sum_{s=1}^l \mathbf{I}_{l,s}^{(n)}. \quad (218)$$

We are given that

$$p_1 \leq \langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_+ \rangle \leq q_1, \quad \text{and} \quad r_1^* \leq \langle \mathbf{W}_{O(i,\cdot)}^{(t+1)^\top}, \mathbf{o}_+ \rangle \leq s_1^*. \quad (219)$$

From our initialization, for all $i \in \mathcal{K}^+$, we have $v_i = \frac{1}{\sqrt{m}}$. This gives

$$\left\langle -\frac{\partial \ell}{\partial \mathbf{w}_\Delta}, \mathbf{o}_+ \right\rangle = \frac{z^{(n)}}{L} \sum_{i=1}^m \sum_{l=1}^L \frac{1}{\sqrt{m}} \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \sum_{s=1}^l \left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle. \quad (220)$$

Averaging over the training samples, the inner product of the empirical gradient becomes

$$\begin{aligned} \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{w}_\Delta}, \mathbf{o}_+ \right\rangle &= \frac{1}{N} \sum_{n=1}^N \frac{z^{(n)}}{L} \cdot \sum_{i=1}^m \sum_{l=1}^L v_i \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \cdot \sum_{s=1}^l \left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle \\ &= \frac{1}{N} \sum_{n:z^{(n)}=+1} \frac{1}{L} \left[\sum_{i \in \mathcal{K}_+} \sum_{l=1}^L \frac{1}{\sqrt{m}} \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \sum_{s=1}^l \left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle \right. \\ &\quad \left. + \sum_{i \in \mathcal{K}_-} \sum_{l=1}^L \left(-\frac{1}{\sqrt{m}} \right) \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \sum_{s=1}^l \left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle \right] \\ &\quad + \frac{1}{N} \sum_{n:z^{(n)}=-1} \frac{-1}{L} \left[\sum_{i \in \mathcal{K}_+} \sum_{l=1}^L \frac{1}{\sqrt{m}} \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \sum_{s=1}^l \left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle \right. \\ &\quad \left. + \sum_{i \in \mathcal{K}_-} \sum_{l=1}^L \left(-\frac{1}{\sqrt{m}} \right) \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \sum_{s=1}^l \left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle \right]. \end{aligned} \quad (221)$$

First, we focus on the contribution from the samples where $z^{(n)} = +1$, for which we seek a lower bound. We analyze the inner terms by considering four cases.

Case I: $l = L_1^+, s = L_1^+$

Since $l = s$ and $\mathbf{x}_s = \mathbf{o}_+$, it follows from (215) that

$$\left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle = \beta_{s,s}. \quad (222)$$

Using (216), with $\mathbf{W}_B = \mathbf{W}_C = I$ and $\mathbf{x}_l = \mathbf{x}_s = \mathbf{o}_+$, we obtain

$$\left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle = \beta_{s,s} = \langle \mathbf{W}_{O(i,\cdot)}^{(t+1)^\top}, \mathbf{o}_+ \rangle \cdot \sigma(\langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_+ \rangle) \cdot (1 - \sigma(\langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_+ \rangle)). \quad (223)$$

Given the conditions in (219), we can write

$$\left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle \geq r_1^* \cdot \sigma(p_1) \cdot (1 - \sigma(q_1)). \quad (224)$$

Case II: $l = L_2^+, s = L_2^+$

This configuration yields the same result as in Case I. We again obtain

$$\left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle \geq r_1^* \cdot \sigma(p_1) \cdot (1 - \sigma(q_1)). \quad (225)$$

Case III: $l = L_2^+$, $s = L_1^+$ Comparing (216) with (217), we see that the two expressions differ only in their last term. In this setting, \mathbf{x}_j equals \mathbf{o}_+ only when $j = L_2^+$. Consequently, $\mathbf{x}_s = \mathbf{x}_j = \mathbf{o}_+$, which implies $\beta_{s,s} = \beta_{s,j}$. Hence,

$$\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \rangle = \beta_{s,s} - \beta_{s,j} = 0. \quad (226)$$

Case IV: Others

For the other token positions, $\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \rangle = 0$ due to orthogonality among the features.

Combining the above, the total contribution becomes

$$\sum_{s=1}^l \langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \rangle \geq 2r_1^* \cdot \sigma(p_1) \cdot (1 - \sigma(q_1)). \quad (227)$$

We now bound the entire sum over all tokens:

$$\frac{1}{L} \sum_{l=1}^L \frac{1}{\sqrt{m}} \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l) \sum_{s=1}^l \langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \rangle \geq \frac{1}{L} \sum_{l=1}^L \frac{1}{\sqrt{m}} \cdot 1 \cdot 2r_1^* \cdot \sigma(p_1) \cdot (1 - \sigma(q_1)). \quad (228)$$

Let $\rho_t^+ = |\mathcal{W}(t)|$ be the number of contributing neurons. Then the total contribution from the active neurons is lower bounded as

$$\frac{1}{L} \sum_{i \in \mathcal{K}_+} \sum_{l=1}^L v_i \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l) \sum_{s=1}^l \langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \rangle \geq \frac{2r_1^* \cdot \sigma(p_1) \cdot (1 - \sigma(q_1))}{\sqrt{m}} \cdot \rho_t^+. \quad (229)$$

Next, we consider $z^{(n)} = -1$ for $i \in \mathcal{K}_+$. For $z^{(n)} = -1$, the negative sample contains two \mathbf{o}_- and one \mathbf{o}_+ at positions L_1^-, L_2^-, L^+ , respectively.

Let $l = L^+$. Then $\mathbf{x}_l = \mathbf{o}_+$, and \mathbf{y}_l contains a component in the direction of \mathbf{o}_+ . Since $\mathbf{W}_{O(i,\cdot)}$ has an \mathbf{o}_+ component for $i \in \mathcal{K}_+$, this contributes to the gradient.

Let $s = l$, so that

$$\mathbf{I}_{l,s}^{(n)} = \beta_{s,s} \cdot \mathbf{x}_l. \quad (230)$$

We now seek an upper bound for this contribution. From the initial conditions in (219), we know

$$\langle \mathbf{W}_{O(i,\cdot)}^{(t+1)\top}, \mathbf{o}_+ \rangle \leq s_1^*. \quad (231)$$

Hence, we obtain

$$\langle \mathbf{I}_{l,s}, \mathbf{o}_+ \rangle \leq s_1^* \cdot \sigma(q_1) \cdot (1 - \sigma(p_1)). \quad (232)$$

The maximum number of such contributing neurons is $\frac{m}{2}$. Therefore, the total contribution is bounded above by

$$\begin{aligned} \frac{1}{L} \sum_{i \in \mathcal{K}_+} \sum_{l=1}^L v_i \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l) \sum_{s=1}^l \langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \rangle &\leq \frac{s_1^* \cdot \sigma(q_1) \cdot (1 - \sigma(p_1))}{\sqrt{m}} \cdot \frac{m}{2} \\ &= \frac{\sqrt{m} \cdot s_1^* \cdot \sigma(q_1) \cdot (1 - \sigma(p_1))}{2}. \end{aligned} \quad (233)$$

Thirdly, let us consider the contribution for $z^{(n)} = +1$ from $i \in \mathcal{K}_-$. From our initialization, for $i \in \mathcal{K}_-$, $v_i = -\frac{1}{\sqrt{m}}$. For $z^{(n)} = +1$, we seek an upper bound on the contribution from such neurons.

Let $z^{(n)} = +1$. To maximize the term $\mathbf{W}_{O(i,\cdot)} \mathbf{x}_s^{(n)}$ in (216), we consider $l = L^-$ since $\mathbf{W}_{O(i,\cdot)}$ has a large component in the \mathbf{o}_- direction. Then $\mathbf{x}_l = \mathbf{o}_- \Rightarrow \mathbf{y}_l$ contains the \mathbf{o}_- feature.

However, in this case, $\mathbf{x}_s = \mathbf{o}_- = \mathbf{x}_l$, and due to orthogonality,

$$\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \rangle = 0. \quad (234)$$

Hence, we only need to consider time steps $l = L_1^+, L_2^+$, where \mathbf{o}_+ features appear.

Recall that

$$\left\langle -\frac{\partial \ell}{\partial \mathbf{w}_\Delta}, \mathbf{o}_+ \right\rangle = \frac{1}{L} \sum_{i=1}^m \sum_{l=1}^L -\frac{1}{\sqrt{m}} \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l) \sum_{s=1}^l \left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle. \quad (235)$$

We analyze the inner contributions case by case.

Case I: $l = L_1^+, s = L_1^+$

Given that

$$\mathbf{W}_{O(i,\cdot)} \mathbf{o}_+ \leq \delta_1 + \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right) =: c, \quad (236)$$

we obtain

$$\left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle \leq c \cdot \sigma(q_1) \cdot (1 - \sigma(p_1)). \quad (237)$$

Case II: $l = L_2^+, s = L_2^+$

This configuration yields the same bound:

$$\left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle \leq c \cdot \sigma(q_1) \cdot (1 - \sigma(p_1)). \quad (238)$$

Case III: $l = L_2^+, s = L_1^+$

In this case, the contribution vanishes:

$$\left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle = 0. \quad (239)$$

Case IV: Others

For the other token positions, $\left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle = 0$ due to orthogonality among the features.

Thus, the total contribution from each $i \in \mathcal{K}^-$ satisfies

$$\sum_{s=1}^l \left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle \leq 2c \cdot \sigma(q_1) \cdot (1 - \sigma(p_1)). \quad (240)$$

The maximum number of such contributing neurons is $\frac{m}{2}$, so the full contribution is bounded by

$$\begin{aligned} \frac{1}{\sqrt{m}L} \sum_{i \in \mathcal{K}^-} \sum_{l=1}^L \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l) \sum_{s=1}^l \left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle &\leq \frac{2c \cdot \sigma(q_1) \cdot (1 - \sigma(p_1))}{\sqrt{m}} \cdot \frac{m}{2} \\ &= \sqrt{mc} \cdot \sigma(q_1) \cdot (1 - \sigma(p_1)). \end{aligned} \quad (241)$$

Therefore, the overall contribution is

$$-\frac{1}{\sqrt{m}L} \sum_{i \in \mathcal{K}^-} \sum_{l=1}^L \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l) \sum_{s=1}^l \left\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \right\rangle \geq -\sqrt{mc} \cdot \sigma(q_1) \cdot (1 - \sigma(p_1)). \quad (242)$$

Finally, we consider $z^{(n)} = -1$ for $i \in \mathcal{K}_-$. For $z^{(n)} = -1$, we want a lower bound since $v_i = -\frac{1}{\sqrt{m}}$.

We could consider $l = L^+ \Rightarrow \mathbf{x}_l = \mathbf{o}_+$, and write

$$\left\langle \mathbf{W}_{O(i,\cdot)}, \mathbf{o}_+ \right\rangle \geq \delta_1 - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \quad (243)$$

However, the minimum number of such contributing neurons is not tractable. Thus, if we consider the worst case where $\mathbf{W}_{O(i,\cdot)}$ for $i \in \mathcal{K}^-$ does not learn the \mathbf{o}_+ feature, the obvious lower bound is zero:

$$\frac{1}{L} \sum_{i \in \mathcal{K}^-} \sum_{l=1}^L \frac{1}{\sqrt{m}} \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l) \sum_{s=1}^l \langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_+ \rangle \geq 0. \quad (244)$$

We now combine the bounds for the four terms identified in equation (221), corresponding to the contributions from: (i) \mathcal{K}_+ with $z^{(n)} = +1$ (229), (ii) \mathcal{K}_+ with $z^{(n)} = -1$ (233), (iii) \mathcal{K}_- with $z^{(n)} = +1$ (242), and (iv) \mathcal{K}_- with $z^{(n)} = -1$ (244). We assume the batch is balanced, so the number of positive and negative samples is equal, with each class contributing $\frac{N}{2}$ samples. Then we have

$$\begin{aligned} \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{w}_\Delta}, \mathbf{o}_+ \right\rangle &\geq \frac{1}{2} \left[\frac{2r_1^* \cdot \sigma(p_1) (1 - \sigma(q_1))}{\sqrt{m}} \cdot \rho_t^+ - \sqrt{m} \cdot c \cdot \sigma(q_1) (1 - \sigma(p_1)) \right. \\ &\quad \left. - \frac{\sqrt{m} \cdot s_1^* \cdot \sigma(q_1) (1 - \sigma(p_1))}{2} + 0 \right] \\ &= \frac{\sigma(p_1) (1 - \sigma(q_1)) r_1^* \cdot \rho_t^+}{\sqrt{m}} - \frac{\sigma(q_1) (1 - \sigma(p_1)) s_1^* \cdot \sqrt{m}}{2} - \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right). \end{aligned} \quad (245)$$

where we have used the fact $\frac{\sqrt{m}}{2} \cdot \sigma(q_1) (1 - \sigma(p_1)) \cdot c = \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right)$ since $c = \mathcal{O}\left(\sqrt{\frac{d \log N}{mN}}\right)$.

□

F.6 Proof of Lemma 6

Proof. The gradient is given in (214).

Let's consider the alignment with \mathbf{o}_k for $k \neq 1, 2$.

$$\left\langle -\frac{\partial \ell}{\partial \mathbf{w}_\Delta}, \mathbf{o}_k \right\rangle = \frac{z^{(n)}}{L} \sum_{i=1}^m \sum_{l=1}^L v_i \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \sum_{s=1}^l \langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_k \rangle \quad (246)$$

From our initialization, for all $i \in \mathcal{K}^+$, we have $v_i = \frac{1}{\sqrt{m}}$.

We first consider the case $z^{(n)} = +1$ for $i \in \mathcal{K}^+$. Since $\mathbf{W}_{O(i,\cdot)}$, for $i \in \mathcal{K}^+$ has a large \mathbf{o}_+ component, we have to consider the token features with \mathbf{o}_+ . For $z^{(n)} = +1$, only when $l = L_2^+$, $s = L_1^+$ we have $\mathbf{x}_l = \mathbf{x}_s = \mathbf{o}_+$. Therefore, $\mathbf{W}_{O(i,\cdot)} \mathbf{x}_s$ is significant. Hence, we have

$$\begin{aligned} \langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_k \rangle &= - \sum_{j=s+1}^l \beta_{s,j} \langle \mathbf{x}_j^{(n)}, \mathbf{o}_k \rangle \\ &\leq -\beta_{s,s+1} \quad (\text{Assuming W.L.O.G. } \mathbf{x}_{s+1}^{(n)} = \mathbf{o}_k) \\ &\leq -\langle \mathbf{W}_{O(i,\cdot)}^{(t+1)\top}, \mathbf{o}_+ \rangle \cdot \sigma\left(\langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_+ \rangle\right) \cdot \left(1 - \sigma\left(\langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_+ \rangle\right)\right) \\ &\quad \cdot \left(1 - \sigma\left(\langle \mathbf{w}_\Delta^{(t)}, \mathbf{o}_k \rangle\right)\right)^{L_2^+ - L_1^+}. \end{aligned} \quad (247)$$

Using the conditions in (219), we can write

$$\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_k \rangle \leq -r_1^* \cdot \sigma(p_1) \cdot (1 - \sigma(q_1)) \cdot (1 - \sigma(q_2))^{L_2^+ - L_1^+}. \quad (248)$$

Hence, we obtain

$$\begin{aligned} \frac{1}{L} \sum_{l=1}^L \frac{1}{\sqrt{m}} \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \sum_{s=1}^l \langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_k \rangle &\leq \frac{1}{L} \sum_{l=1}^L \frac{1}{\sqrt{m}} \cdot 1 \\ &\cdot \left[-r_1^* \cdot \sigma(p_1) \cdot (1 - \sigma(q_1)) \cdot (1 - \sigma(q_2))^{L_2^+ - L_1^+} \right]. \end{aligned} \quad (249)$$

Let $\rho_t^+ = |\mathcal{W}(t)|$ be the number of contributing neurons. Then the total contribution from \mathcal{K}_+ neurons is bounded as

$$\begin{aligned} \frac{1}{L} \sum_{i \in \mathcal{K}_+} \sum_{l=1}^L v_i \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \sum_{s=1}^l \langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_k \rangle &\leq -\frac{r_1^*}{\sqrt{m}} \cdot \sigma(p_1) (1 - \sigma(q_1)) \\ &\cdot (1 - \sigma(q_2))^{L_2^+ - L_1^+} \cdot \rho_t^+. \end{aligned} \quad (250)$$

Similarly for $i \in \mathcal{K}^-$, for $z^{(n)} = -1$, when $l = L_2^-$, $s = L_1^-$ the contribution is significant.

$$\langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_k \rangle \leq -r_1^* \cdot \sigma(r_2) \cdot (1 - \sigma(s_2)) \cdot (1 - \sigma(q_2))^{L_2^- - L_1^-}. \quad (251)$$

Hence, we obtain

$$\begin{aligned} \frac{1}{L} \sum_{l=1}^L \frac{1}{\sqrt{m}} \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \sum_{s=1}^l \langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_k \rangle &\leq \frac{1}{L} \sum_{l=1}^L \frac{1}{\sqrt{m}} \cdot 1 \\ &\cdot \left[-r_1^* \cdot \sigma(r_2) \cdot (1 - \sigma(s_2)) \cdot (1 - \sigma(q_2))^{L_2^- - L_1^-} \right]. \end{aligned} \quad (252)$$

Let $\rho_t^- = |\mathcal{U}(t)|$ be the number of contributing neurons. Then the total contribution from \mathcal{K}_- neurons is bounded as

$$\begin{aligned} \frac{1}{L} \sum_{i \in \mathcal{K}_-} \sum_{l=1}^L v_i \cdot \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \sum_{s=1}^l \langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_k \rangle &\leq -\frac{r_1^*}{\sqrt{m}} \cdot \sigma(r_2) (1 - \sigma(s_2)) \\ &\cdot (1 - \sigma(q_2))^{L_2^- - L_1^-} \cdot \rho_t^-. \end{aligned} \quad (253)$$

Putting it together, We know

$$\begin{aligned} \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{w}_\Delta}, \mathbf{o}_k \right\rangle &= \frac{1}{N} \sum_{n=1}^N \frac{z^{(n)}}{L} \cdot \sum_{i=1}^m \sum_{l=1}^L v_i \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \cdot \sum_{s=1}^l \langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_k \rangle \\ &\leq \frac{1}{N} \sum_{n:z^{(n)}=+1} \frac{1}{L} \left[\sum_{i \in \mathcal{K}_+} \sum_{l=1}^L \frac{1}{\sqrt{m}} \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \sum_{s=1}^l \langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_k \rangle \right] \\ &\quad - \frac{1}{N} \sum_{n:z^{(n)}=-1} \frac{1}{L} \left[\sum_{i \in \mathcal{K}_-} \sum_{l=1}^L \frac{-1}{\sqrt{m}} \phi'(\mathbf{W}_{O(i,\cdot)} \mathbf{y}_l^{(n)}) \sum_{s=1}^l \langle \mathbf{I}_{l,s}^{(n)}, \mathbf{o}_k \rangle \right] \end{aligned} \quad (254)$$

We now combine the bounds for the two terms identified in equation (254), corresponding to the contributions from: (i) \mathcal{K}_+ with $z^{(n)} = +1$ (250), and (ii) \mathcal{K}_- with $z^{(n)} = -1$ (253). We assume the batch is balanced, so the number of positive and negative samples is equal, with each class contributing $\frac{N}{2}$ samples. Then we have

$$\begin{aligned} \left\langle -\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{w}_\Delta^{(t)}}, \mathbf{o}_k \right\rangle &\leq -\frac{r_1^*}{2\sqrt{m}} \left[\sigma(p_1) (1 - \sigma(q_1)) (1 - \sigma(q_2))^{L_2^+ - L_1^+} \rho_t^+ + \right. \\ &\quad \left. \sigma(r_2) (1 - \sigma(s_2)) (1 - \sigma(q_2))^{L_2^- - L_1^-} \rho_t^- \right] \end{aligned} \quad (255)$$

□