

---

# Transformers as Support Vector Machines

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The transformer architecture has led to revolutionary advancements in NLP. The at-  
2 tention layer within the transformer admits a sequence of input tokens  $X$  and makes  
3 them interact through pairwise similarities computed as  $\text{softmax}(XQK^T X^T)$ ,  
4 where  $(K, Q)$  are the trainable key-query parameters. In this work, we estab-  
5 lish a formal equivalence between the optimization geometry of self-attention and  
6 a hard-margin SVM problem that separates optimal input tokens from non-optimal  
7 tokens using linear constraints on the outer-products of token pairs. This formalism  
8 allows us to characterize the implicit bias of 1-layer transformers optimized with  
9 gradient descent: **(1)** Optimizing the attention layer, parameterized by  $(K, Q)$ , with  
10 vanishing regularization, converges in direction to an SVM solution minimizing the  
11 nuclear norm of the combined parameter  $W := KQ^T$ . Instead, directly parameteriz-  
12 ing by  $W$  minimizes a Frobenius norm SVM objective. **(2)** Complementing this, for  
13  $W$ -parameterization, we prove the local/global directional convergence of gradient  
14 descent under suitable geometric conditions, and propose a more general SVM  
15 equivalence that predicts the implicit bias of attention with nonlinear heads/MLPs.

## 16 1 Introduction

17 Self-attention, the central component of the transformer architecture, has revolutionized NLP  
18 [VSP<sup>+</sup>17]. This mechanism has proven highly effective in capturing long-range dependencies, which  
19 is essential for applications arising in NLP [KT19, BMR<sup>+</sup>20, RSR<sup>+</sup>20], computer vision [FXM<sup>+</sup>21,  
20 LLC<sup>+</sup>21, TCD<sup>+</sup>21, CSL<sup>+</sup>23], and reinforcement learning [JLL21, CLR<sup>+</sup>21, WWX<sup>+</sup>22]. Remarkable  
21 success of the self-attention mechanism and transformers has paved the way for the development of  
22 LLMs such as GPT4 [Ope23], Bard [Goo23], LLaMA [TLI<sup>+</sup>23], and ChatGPT [Ope22].

23 **Q:** Can we characterize the optimization landscape and implicit bias of transformers?

24 We address this question by rigorously connecting the optimization geometry of the attention layer  
25 and a hard max-margin SVM problem, namely (**Att-SVM**), that separates and selects the optimal  
26 tokens from each input sequence. This formalism follows [TLZO23], which sheds light on the  
27 intricacies of self-attention. Throughout, given input sequences  $X, Z \in \mathbb{R}^{T \times d}$  with length  $T$  and  
28 embedding dimension  $d$ , we study the core cross-attention and self-attention models:

$$f_{\text{cross}}(X, Z) := \mathbb{S}(ZQK^T X^T)XV, \quad f_{\text{self}}(X) := \mathbb{S}(XQK^T X^T)XV.$$

29 Here,  $K, Q \in \mathbb{R}^{d \times m}$ ,  $V \in \mathbb{R}^{d \times v}$  are the trainable key, query, value matrices respectively;  $\mathbb{S}(\cdot)$  denotes  
30 the softmax nonlinearity. Note that self-attention is a special instance of the cross-attention by setting  
31  $Z \leftarrow X$ . To expose our main results, suppose the first token of  $Z$ , denoted by  $z$ , is used for prediction.  
32 Concretely, given a dataset  $(Y_i, X_i, z_i)_{i=1}^n$  with labels  $Y_i \in \{-1, 1\}$  and inputs  $X_i \in \mathbb{R}^{T \times d}$ ,  $z_i \in \mathbb{R}^d$ , we  
33 consider the empirical risk minimization with a loss  $\ell(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ , defined as follows:

$$\mathcal{L}(K, Q) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i \cdot f(X_i, z_i)), \quad \text{where } f(X_i, z_i) = h\left(X_i^T \mathbb{S}\left(X_i K Q^T z_i\right)\right). \quad (1)$$

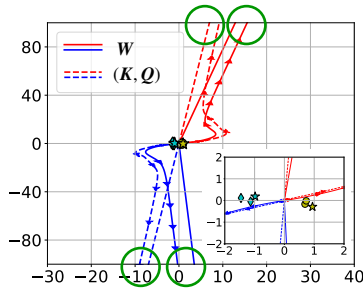


Figure 1: GD convergence of attention weights. Markers represent tokens; lines depict attention-SVM directions mapped to  $\mathbf{z}$ ; arrows illustrate GD paths converging towards these SVM directions. Green circles denote GD  $\leftrightarrow$  SVM pairings.

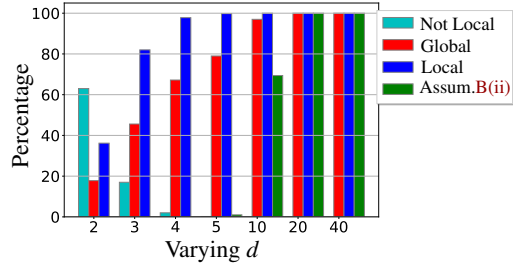


Figure 2: Percentage of different convergence types when training  $\mathbf{W}$ . Red and blue bars represent the percentages of convergence to globally and locally-optimal SVM solutions; teal are complements of the blue; green depict Assum. B(ii).

34 Here,  $h(\mathbf{x}) = \mathbf{v}^\top \mathbf{x}$  is the linear prediction head and  $f(\cdot)$  precisely represents a one-layer transformer.  
 35 The softmax operation, due to its nonlinear nature, poses a significant challenge when optimizing (1).  
 36 In this study, we focus on optimizing the attention weights ( $\mathbf{K}$ ,  $\mathbf{Q}$  or  $\mathbf{W}$ ) and overcome such challenges  
 37 to establish a fundamental SVM equivalence. The paper’s main contributions are as follows:

- 38 • **Implicit bias of the attention layer (Sec. 2).** Optimizing the attention parameters  $\mathbf{W}$  or ( $\mathbf{K}$ ,  $\mathbf{Q}$ )  
 39 with vanishing regularization converges in direction towards a solution of (Att-SVM) or (Att-SVM $_{\star}$ )  
 40 with the Frobenius norm or the nuclear norm objective, respectively. To our knowledge, this is the first  
 41 result to formally distinguish the optimization dynamics of  $\mathbf{W}$  vs ( $\mathbf{K}$ ,  $\mathbf{Q}$ ) parameterizations, revealing  
 42 the low-rank bias of the latter.
- 43 • **Convergence of gradient descent (Sec. 3).** We prove the local/global directional convergence  
 44 of gradient descent for optimizing the attention layer parameterized by  $\mathbf{W}$  under suitable geometric  
 45 conditions. Beyond these, we propose a more general SVM equivalence with nonlinear head, which  
 46 predicts the implicit bias of attention trained by gradient descent.

### 47 1.1 Preliminaries

48 **Optimization algorithms.** Given a parameter  $R > 0$ , we define the regularized path solution as  
 49 ( $\mathbf{W}$ -RP) and ( $\mathbf{KQ}$ -RP). For GD, with appropriate  $\eta > 0$ , we describe the optimization process as  
 50 ( $\mathbf{W}$ -GD) and ( $\mathbf{KQ}$ -GD). Here for ( $\mathbf{W}$ -RP) and ( $\mathbf{W}$ -GD),  $\mathcal{L}(\mathbf{Q}, \mathbf{K})$  is replaced with  $\mathcal{L}(\mathbf{W})$  with  $\mathbf{W} := \mathbf{KQ}^\top$ .

51 Given $\mathbf{W}(0) \in \mathbb{R}^{d \times d}$ , $\eta > 0$ , for $k \geq 0$ do: $\mathbf{W}(k+1) = \mathbf{W}(k) - \eta \nabla \mathcal{L}(\mathbf{W}(k)).$ (W-GD)	Given $\mathbf{Q}(0), \mathbf{K}(0) \in \mathbb{R}^{d \times m}$ , $\eta > 0$ , for $k \geq 0$ do: $\begin{bmatrix} \mathbf{K}(k+1) \\ \mathbf{Q}(k+1) \end{bmatrix} = \begin{bmatrix} \mathbf{K}(k) \\ \mathbf{Q}(k) \end{bmatrix} - \eta \begin{bmatrix} \nabla_{\mathbf{K}} \mathcal{L}(\mathbf{K}(k), \mathbf{Q}(k)) \\ \nabla_{\mathbf{Q}} \mathcal{L}(\mathbf{K}(k), \mathbf{Q}(k)) \end{bmatrix}.$ (KQ-GD)
52 Given $R > 0$ , find $d \times d$ matrix: $\bar{\mathbf{W}}_R = \arg \min_{\ \mathbf{W}\ _F \leq R} \mathcal{L}(\mathbf{W}).$ (W-RP)	Given $R > 0$ , find $d \times m$ matrices: $(\bar{\mathbf{K}}_R, \bar{\mathbf{Q}}_R) = \arg \min_{\ \mathbf{K}\ _F^2 + \ \mathbf{Q}\ _F^2 \leq 2R} \mathcal{L}(\mathbf{K}, \mathbf{Q}).$ (KQ-RP)

53 **Definition 1 (Token Score and Optimality)** Given a prediction head  $\mathbf{v} \in \mathbb{R}^d$ , the score of a token  
 54  $\mathbf{x}_{it}$  of input  $\mathbf{X}_i$  is defined as  $\gamma_{it} = Y_i \cdot \mathbf{v}^\top \mathbf{x}_{it}$ . The optimal token for each input  $\mathbf{X}_i$  is given by the index  
 55  $\text{opt}_i \in \arg \max_{t \in [T]} \gamma_{it}$  for all  $i \in [n]$ .

56 By introducing token scores and identifying optimal tokens, we can better understand the importance  
 57 of individual tokens and their impact on the overall objective. Next, we present SVM problems.

- 58 • **Hard-margin SVM for  $\mathbf{W}$ -parameterization.** Equipped with the set of optimal indices  $(\text{opt}_i)_{i=1}^n$   
 59 as per Definition 1, we introduce the following SVM formulation associated to  $\mathbf{W}$ -parameterization:

$$\mathbf{W}^{mm} = \arg \min_{\mathbf{W}} \|\mathbf{W}\|_F \quad \text{s.t.} \quad (\mathbf{x}_{i\text{opt}_i} - \mathbf{x}_{it})^\top \mathbf{W} \mathbf{z}_i \geq 1 \quad \text{for all } t \neq \text{opt}_i, i \in [n]. \quad (\text{Att-SVM})$$

61 Throughout, we assume the SVM problems are feasible. We also note that GD can provably converge  
 62 to an SVM solution over locally-optimal tokens, as detailed in Section 3.2.

63 • **SVM problem for  $(K, Q)$ -parameterization.** The objective function has an extra layer of noncon-  
 64 vexity as  $(K, Q)$  corresponds to a matrix factorization of  $W$ . Fortunately, our experiments reveal that  
 65 GD is indeed biased towards the global minima. This yields the following  $W$ -parameterized SVM  
 66 with nuclear norm objective:

$$\mathbf{W}_\star^{mm} \in \arg \min_{\text{rank}(W) \leq m} \|\mathbf{W}\|_\star \quad \text{s.t.} \quad (\mathbf{x}_{i\text{opt}_i} - \mathbf{x}_{it})^\top \mathbf{W} \mathbf{z}_i \geq 1 \quad \text{for all } t \neq \text{opt}_i, i \in [n]. \quad (\text{Att-SVM}_\star)$$

67  
 68 Above, the nonconvex rank constraint arises from the fact that the rank of  $W = KQ^\top$  is at most  $m$ .  
 69 Lemma 1, presented below, demonstrates that this guarantee holds whenever  $n \leq m$ .

70 **Lemma 1** Any optimal solution of (Att-SVM) or (Att-SVM $_\star$ ) is at most rank  $n$ . More precisely, the  
 71 row space of  $W^{mm}$  or  $W_\star^{mm}$  lies within  $\text{span}(\{\mathbf{z}_i\}_{i=1}^n)$ .

## 72 2 Understanding Implicit Bias of Self-Attention

73 We start by establishing the global convergence of regularized paths.

74 **Assumption A** Over any bounded interval  $[a, b]$ : (i)  $\ell : \mathbb{R} \rightarrow \mathbb{R}$  is strictly decreasing; (ii) The  
 75 derivative  $\ell'$  is bounded as  $|\ell'(u)| \leq M_1$ ; (iii)  $\ell'$  is  $M_0$ -Lipschitz continuous.

76 **Theorem 1** Suppose Assumption A holds, optimal indices  $(\text{opt}_i)_{i=1}^n$  are unique. Let  $W^{mm}$  be the  
 77 unique solution of (Att-SVM), and let  $W_\star^{mm}$  be the solution set of (Att-SVM $_\star$ ) with nuclear norm  
 78 achieving objective  $C_\star$ . Then, Algorithms W-RP and KQ-RP, respectively, satisfy:

- 79 •  $W$ -parameterization has Frobenius norm bias:  $\lim_{R \rightarrow \infty} \frac{\bar{W}_R}{R} = \frac{W^{mm}}{\|W^{mm}\|_F}$ .
- 80 •  $(K, Q)$ -parameterization has nuclear norm bias:  $\lim_{R \rightarrow \infty} \text{dist}\left(\frac{\bar{K}_R \bar{Q}_R^\top}{R}, \frac{W_\star^{mm}}{C_\star}\right) = 0$ .

81 Theorem 1 shows that the RP of the  $W$  and  $(K, Q)$ -parameterization converge to the max-margin  
 82 solutions of (Att-SVM) and (Att-SVM $_\star$ ) with Frobenius and nuclear norm objectives, respectively.  
 83 This result is the first to distinguish the optimization dynamics of  $W$  and  $(K, Q)$  parameterizations,  
 84 revealing the low-rank bias of the latter. To study the RP theory predictivity of the implicit bias  
 85 exhibited by GD, we examine the GD paths in Figure 1, where  $n = d = 2, T = 3$ . The teal and  
 86 yellow markers correspond to tokens from  $X_1, X_2$ , and the stars indicate the optimal tokens. We  
 87 illustrate the iterations of the attention weight in the form of  $W \mathbf{z}_i$  and  $KQ^\top \mathbf{z}_i, i = 1, 2$ . The red/blue  
 88 solid lines delineate the directions of  $W^{mm} \mathbf{z}_1 / W^{mm} \mathbf{z}_2$ ; the red/blue dashed lines show the directions  
 89 of  $W_\star^{mm} \mathbf{z}_1 / W_\star^{mm} \mathbf{z}_2$ ; the arrows denote the corresponding directions of gradient evolution. Figure 1  
 90 provides a clear depiction of the incremental alignment of  $W(k)$  and  $K(k)Q(k)^\top$  with their respective  
 91 attention SVM solutions as  $k$  increases. This strongly supports the assertions of Theorem 1.

## 92 3 Convergence and Implicit Bias of Gradient Descent

### 93 3.1 Global convergence

94 In this section, we will establish conditions that guarantee the global convergence of GD.

95 **Lemma 2** Under Assumption A,  $\nabla \mathcal{L}(W)$  is  $L_W$ -Lipschitz continuous, where  $L_W := \frac{1}{n} \sum_{i=1}^n a_i b_i$ , and  
 96  $a_i = \|\mathbf{v}\| \|\mathbf{z}_i\|^2 \|\mathbf{X}_i\|^3, b_i = M_0 \|\mathbf{v}\| \|\mathbf{X}_i\| + 3M_1$  for all  $i \in [n]$ .

97 **Assumption B** Optimal tokens' indices  $(\text{opt}_i)_{i=1}^n$  are unique and one of the following conditions on  
 98 the tokens holds: For all  $t \neq \text{opt}_i$  and  $i \in [n]$ , (i) the tokens' scores, as defined in Def. 1, satisfy  
 99  $\gamma_{it} = \gamma_{i\text{opt}_i} < \gamma_{i\text{opt}_i}$ . (ii) all tokens are support vectors, i.e.,  $(\mathbf{x}_{i\text{opt}_i} - \mathbf{x}_{it})^\top W^{mm} \mathbf{z}_i = 1$ ;

100 Here, we provide conditions for achieving global convergence towards the max-margin direction  
 101  $W^{mm}$  based on token score constraints and over-parameterization. For the former, we provide precise  
 102 theoretical guarantees. For the latter, we provide strong empirical evidence.

103 **(I) Global convergence under score constraints.** Our next result establishes the global convergence  
 104 of GD to the max-margin direction  $W^{mm}$  under Assumption B(i) that non-optimal tokens have  
 105 identical scores but lower than the score of the optimal token.

106 **Theorem 2** Suppose Assumption A on the loss  $\ell$  and Assumption B(i) on the tokens' score hold.  
 107 Then, Algorithm W-GD with  $\eta \leq 1/L_W$  and any starting point  $W(0)$  satisfies  $\lim_{k \rightarrow \infty} \frac{W(k)}{\|W(k)\|_F} = \frac{W^{mm}}{\|W^{mm}\|_F}$ .

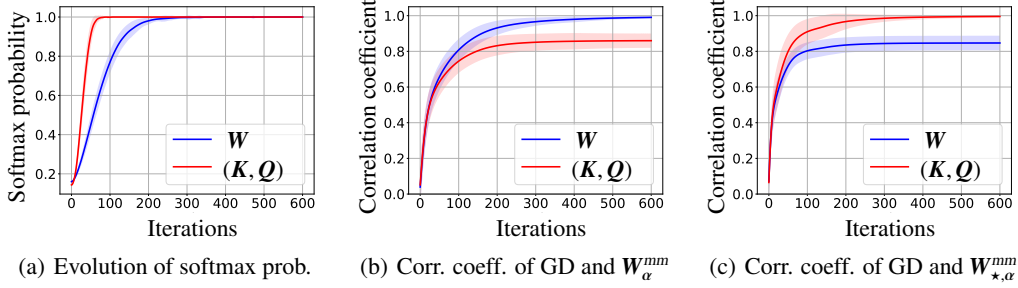


Figure 3: Local convergence behaviour of GD when training  $W$  or  $(K, Q)$  with random data.

108 **(II) Global convergence via overparameterization.** Considering that Assumption B(ii) is anticipated  
 109 to hold as the dimension  $d$  increases, the norm of the GD solution is bound to diverge to infinity. This  
 110 satisfies a prerequisite for converging towards the globally-optimal SVM direction  $W^{mm}$ . The trend  
 111 depicted in Figure 2, where the percentage of global convergence (red bars) approaches 100% and  
 112 Assumption B(ii) holds with higher probability (green bars) as  $d$  grows, reinforces this insight.

### 113 3.2 Local convergence

114 **Definition 2 (Local Optimality)** Fix token indices  $\alpha = (\alpha_i)_{i=1}^n$ . Solve (Att-SVM) with  $(opt_i)_{i=1}^n$   
 115 replaced with  $\alpha$  to obtain  $W_\alpha^{mm}$ . Consider the set  $\mathcal{T}_i \subset [T]$  such that  $(x_{i\alpha_i} - x_{it})^\top W_\alpha^{mm} z_i = 1$ . If for all  
 116  $i \in [n]$  and  $t \in \mathcal{T}_i$  scores per Def. 1 obey  $\gamma_{i\alpha_i} > \gamma_{it}$ ,  $W_\alpha^{mm}$  is called a locally-optimal direction.

117 To provide a basis for discussing local convergence of GD, we establish a cone centered around  $W_\alpha^{mm}$ :  
 118 For  $\mu \in (0, 1)$  and  $R > 0$ , we define  $C_{\mu,R}(W_\alpha^{mm}) := \{\|W\|_F \geq R \mid \langle W/\|W\|_F, W_\alpha^{mm}/\|W_\alpha^{mm}\|_F \rangle \geq 1 - \mu\}$ .

119 **Theorem 3** Suppose Assumption A holds, and let  $\alpha = (\alpha_i)_{i=1}^n$  be locally optimal tokens and  $W_\alpha^{mm}$   
 120 be a locally-optimal direction according to Def. 2. Then, Algorithm W-GD with  $\eta \leq 1/L_W$  and any  
 121  $W(0) \in C_{\mu,R}(W_\alpha^{mm})$  satisfies  $\lim_{k \rightarrow \infty} \|W(k)\|_F = \infty$  and  $\lim_{k \rightarrow \infty} \frac{W(k)}{\|W(k)\|_F} = \frac{W_\alpha^{mm}}{\|W_\alpha^{mm}\|_F}$ .

122 This theorem indicates that if GD is initiated within  $C_{\mu,R}(W_\alpha^{mm})$ , it will converge in the direction of  
 123  $W_\alpha^{mm}/\|W_\alpha^{mm}\|_F$ . Importantly, Theorem 3 does not make any assumptions on the tokens as opposed to  
 124 Theorem 2. In Figure 3 we consider setting where  $n = 6$ ,  $T = 8$ , and  $d = 10$ . In Fig. 3(a) we calculate  
 125 the softmax probabilities, which result in probability 1, indicating that attention weights succeed in  
 126 selecting one token per input. Following Def. 2 let  $\alpha = (\alpha_i)_{i=1}^n$  be the token indices selected by GD  
 127 and denote  $W_{*,\alpha}^{mm}$  as the corresponding SVM solution of (Att-SVM $_*$ ). Figs. 3(b) and 3(c) illustrate the  
 128 correlation coefficients of attention weights with respect to  $W_\alpha^{mm}$  and  $W_{*,\alpha}^{mm}$ . The results demonstrate  
 129 that  $W(KQ^\top)$  ultimately reaches a 1 correlation with  $W_\alpha^{mm}$  ( $W_{*,\alpha}^{mm}$ ), which validates Theorem 3.

### 130 3.3 Implicit bias under MLP nonlinearity

131 So far, we focus on the setting that  $h(\cdot)$  is linear and attention selects a single token per sequence.  
 132 In this section, we analyze the scenarios where  $h(\cdot)$  is nonlinear and nonconvex, and GD solution is  
 133 composed by multiple tokens. Suppose optimal solution outputs softmax probability of  $s_i^*$ ,  $i \in [n]$ .  
 134 Intuitively,  $W(k)$  should be decomposed into two components via

$$135 W(k) \approx W^{fin} + \|W(k)\|_F \cdot \bar{W}^{mm}. \quad (2)$$

136 where  $W^{fin}$  is the finite component and  $\bar{W}^{mm}$  is the directional component with  $\|\bar{W}^{mm}\|_F = 1$ . Define  
 137 the selected set  $O_i \subseteq [T]$  to be the indices  $s_{it}^* \neq 0$  and the masked set as  $\bar{O}_i = [T] - O_i$ .

138 **Finite component ( $W^{fin}$ ):** The job of  $W^{fin}$  is to assign nonzero softmax probabilities within each  $s_i^*$ .  
 139 Then,  $W^{fin}$  should satisfy the linear constraints:

$$140 (x_{it} - x_{i\tau})^\top W^{fin} z_i = \log(s_{it}^*/s_{i\tau}^*) \quad \text{for all } t, \tau \in O_i, i \in [n]. \quad (3)$$

141 **Directional component ( $\bar{W}^{mm}$ ):** While  $W^{fin}$  creates the composition by allocating the nonzero  
 142 softmax probabilities, it does not explain sparsity of attention map. This is the role of  $\bar{W}^{mm}$ , and we  
 143 obtain the following convex generalized SVM formulation

$$144 W^{mm} = \arg \min_W \|W\|_F \quad \text{subj. to} \quad \begin{cases} \forall t \in O_i, \tau \in \bar{O}_i : (x_{it} - x_{i\tau})^\top W z_i \geq 1, \\ \forall t, \tau \in O_i : (x_{it} - x_{i\tau})^\top W z_i = 0, \end{cases} \quad \forall 1 \leq i \leq n, \quad (4)$$

and  $\bar{W}^{mm} = W^{mm}/\|W^{mm}\|_F$ . It is important to note that (4) offers a substantial generalization beyond  
 the scope of the previous sections. Remarkably, in Appendix B, we empirically demonstrate that this  
 general form indeed seems to predict the implicit bias of gradient descent with MLPs.

## References

- [ACDS23] Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. *arXiv preprint arXiv:2306.00297*, 2023.
- [ACHL19] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [ALH21] Navid Azizan, Sahin Lale, and Babak Hassibi. Stochastic mirror descent on overparameterized nonlinear models. *IEEE Transactions on Neural Networks and Learning Systems*, 33(12):7717–7727, 2021.
- [ASA<sup>+</sup>22] Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. *arXiv:2211.15661*, 2022.
- [AW20a] Ehsan Amid and Manfred K Warmuth. Winnowing with gradient descent. In *Conference on Learning Theory*, pages 163–182. PMLR, 2020.
- [AW20b] Ehsan Amid and Manfred KK Warmuth. Reparameterizing mirror descent as gradient descent. *Advances in Neural Information Processing Systems*, 33:8430–8439, 2020.
- [BALA<sup>+</sup>23] Enric Boix-Adsera, Etai Littwin, Emmanuel Abbe, Samy Bengio, and Joshua Susskind. Transformers learn through gradual rank increase. *arXiv preprint arXiv:2306.07042*, 2023.
- [BCW<sup>+</sup>23] Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. *arXiv preprint arXiv:2306.04637*, 2023.
- [BGVV20] Guy Blanc, Neha Gupta, Gregory Valiant, and Paul Valiant. Implicit regularization for deep neural networks driven by an ornstein-uhlenbeck like process. In *Conference on learning theory*, pages 483–513. PMLR, 2020.
- [BMR<sup>+</sup>20] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and et al. Language models are few-shot learners. In *Advances in neural information processing systems*, volume 33, pages 1877–1901, 2020.
- [BV22] Pierre Baldi and Roman Vershynin. The quarks of attention. *arXiv preprint arXiv:2202.08371*, 2022.
- [Car21] Marcus Carlsson. von neumann’s trace inequality for hilbert–schmidt operators. *Expositiones Mathematicae*, 39(1):149–157, 2021.
- [CDL16] Jianpeng Cheng, Li Dong, and Mirella Lapata. Long short-term memory-networks for machine reading. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 551–561, Austin, Texas, November 2016. Association for Computational Linguistics.
- [CDS01] Scott Shaobing Chen, David L Donoho, and Michael A Saunders. Atomic decomposition by basis pursuit. *SIAM review*, 43(1):129–159, 2001.
- [CLR<sup>+</sup>21] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. In *Advances in Neural Information Processing Systems*, volume 34, pages 15084–15097, 2021.
- [CRT06] Emmanuel J Candès, Justin Romberg, and Terence Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on information theory*, 52(2):489–509, 2006.

- 191 [CSL<sup>+</sup>23] Yingyi Chen, Xi Shen, Yahui Liu, Qinghua Tao, and Johan AK Suykens. Jigsaw-vit:  
192 Learning jigsaw puzzles in vision transformer. *Pattern Recognition Letters*, 166:53–60,  
193 2023.
- 194 [DCL21] Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. Attention is not all you  
195 need: Pure attention loses rank doubly exponentially with depth. In *International  
196 Conference on Machine Learning*, pages 2793–2803. PMLR, 2021.
- 197 [DML21] Alex Damian, Tengyu Ma, and Jason Lee. Label noise sgd provably prefers flat global  
198 minimizers. *arXiv preprint arXiv:2106.06530*, 2021.
- 199 [Don06] David L Donoho. Compressed sensing. *IEEE Transactions on information theory*,  
200 52(4):1289–1306, 2006.
- 201 [EGKZ22] Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases  
202 and variable creation in self-attention mechanisms. In *International Conference on  
203 Machine Learning*, pages 5793–5831. PMLR, 2022.
- 204 [ENM22] Tolga Ergen, Behnam Neyshabur, and Harsh Mehta. Convexifying transformers: Im-  
205 proving optimization and understanding of transformer networks. *arXiv:2211.11052*,  
206 2022.
- 207 [FGBM23] Hengyu Fu, Tianyu Guo, Yu Bai, and Song Mei. What can a single attention layer learn?  
208 a study through the random features lens. *arXiv preprint arXiv:2307.11353*, 2023.
- 209 [FXM<sup>+</sup>21] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra  
210 Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *Proceedings of  
211 the IEEE/CVF International Conference on Computer Vision*, pages 6824–6835, 2021.
- 212 [GLSS18] Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing  
213 implicit bias in terms of optimization geometry. In *International Conference on Machine  
214 Learning*, pages 1832–1841. PMLR, 2018.
- 215 [Goo23] Google. Try bard, an ai experiment by google. <https://bard.google.com>, 2023.
- 216 [GRS<sup>+</sup>23] Angeliki Giannou, Shashank Rajput, Jy-yong Sohn, Kangwook Lee, Jason D Lee,  
217 and Dimitris Papailiopoulos. Looped transformers as programmable computers.  
218 *arXiv:2301.13196*, 2023.
- 219 [GWB<sup>+</sup>17] Suriya Gunasekar, Blake E Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur,  
220 and Nati Srebro. Implicit regularization in matrix factorization. *Advances in neural  
221 information processing systems*, 30, 2017.
- 222 [Hah20] Michael Hahn. Theoretical limitations of self-attention in neural sequence models.  
223 *Transactions of the Association for Computational Linguistics*, 8:156–171, 2020.
- 224 [HWLM20] Jeff Z HaoChen, Colin Wei, Jason D Lee, and Tengyu Ma. Shape matters: Under-  
225 standing the implicit bias of the noise covariance. *arXiv preprint arXiv:2006.08680*,  
226 2020.
- 227 [JDST20] Ziwei Ji, Miroslav Dudík, Robert E Schapire, and Matus Telgarsky. Gradient descent  
228 follows the regularization path for general losses. In *Conference on Learning Theory*,  
229 pages 2109–2136. PMLR, 2020.
- 230 [JLL21] Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one  
231 big sequence modeling problem. *Advances in neural information processing systems*,  
232 34:1273–1286, 2021.
- 233 [JSL22] Samy Jelassi, Michael Eli Sander, and Yuanzhi Li. Vision transformers provably learn  
234 spatial structure. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun  
235 Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- 236 [JST21] Ziwei Ji, Nathan Srebro, and Matus Telgarsky. Fast margin maximization via dual  
237 acceleration. In *International Conference on Machine Learning*, pages 4860–4869.  
238 PMLR, 2021.

- 239 [JT18] Ziwei Ji and Matus Telgarsky. Risk and parameter convergence of logistic regression.  
240 *arXiv preprint arXiv:1803.07300*, 2018.
- 241 [JT19] Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable  
242 data. In *Conference on Learning Theory*, pages 1772–1798. PMLR, 2019.
- 243 [JT20] Ziwei Ji and Matus Telgarsky. Directional convergence and alignment in deep learning.  
244 In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances*  
245 *in Neural Information Processing Systems*, volume 33, pages 17176–17186. Curran  
246 Associates, Inc., 2020.
- 247 [JT21] Ziwei Ji and Matus Telgarsky. Characterizing the implicit bias via a primal-dual analysis.  
248 In *Algorithmic Learning Theory*, pages 772–804. PMLR, 2021.
- 249 [KPOT21] Ganesh Ramachandra Kini, Orestis Paraskevas, Samet Oymak, and Christos Thrampoulidis. Label-imbalanced and group-sensitive classification under overparameterization. *Advances in Neural Information Processing Systems*, 34:18970–18983, 2021.
- 252 [KT19] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training  
253 of deep bidirectional transformers for language understanding. In *Proceedings of*  
254 *NAACL-HLT*, pages 4171–4186, 2019.
- 255 [LFS<sup>+</sup>17] Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen  
256 Zhou, and Yoshua Bengio. A structured self-attentive sentence embedding. In *International*  
257 *Conference on Learning Representations*, 2017.
- 258 [LIPO23] Yingcong Li, M Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers  
259 as algorithms: Generalization and stability in in-context learning. In *International*  
260 *Conference on Machine Learning*, 2023.
- 261 [LLC<sup>+</sup>21] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and  
262 Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows.  
263 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages  
264 10012–10022, 2021.
- 265 [LMZ18] Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in over-  
266 parameterized matrix sensing and neural networks with quadratic activations. In *Con-*  
267 *ference On Learning Theory*, pages 2–47. PMLR, 2018.
- 268 [LR20] TENG YUAN LIANG and ALEXANDER RAKHLIN. Just interpolate: Kernel “ridge-  
269 less” regression can generalize. *The Annals of Statistics*, 48(3):1329–1347, 2020.
- 270 [LWA22] Zhiyuan Li, Tianhao Wang, and Sanjeev Arora. What happens after SGD reaches  
271 zero loss? –a mathematical framework. In *International Conference on Learning*  
272 *Representations*, 2022.
- 273 [LWLC23] Hongkang Li, Meng Wang, Sijia Liu, and Pin-Yu Chen. A theoretical understanding of  
274 shallow vision transformers: Learning, generalization, and sample complexity. *arXiv*  
275 *preprint arXiv:2302.06015*, 2023.
- 276 [LWM19] Yuanzhi Li, Colin Wei, and Tengyu Ma. Towards explaining the regularization effect of  
277 initial large learning rate in training neural networks. *arXiv preprint arXiv:1907.04595*,  
278 2019.
- 279 [MRG<sup>+</sup>20] William Merrill, Vivek Ramanujan, Yoav Goldberg, Roy Schwartz, and Noah Smith.  
280 Effects of parameter norm growth during transformer training: Inductive bias from  
281 gradient descent. *arXiv preprint arXiv:2010.09697*, 2020.
- 282 [MWG<sup>+</sup>20] Edward Moroshko, Blake E Woodworth, Suriya Gunasekar, Jason D Lee, Nati Srebro,  
283 and Daniel Soudry. Implicit bias in deep linear classification: Initialization scale vs  
284 training accuracy. *Advances in neural information processing systems*, 33:22182–22193,  
285 2020.



- 286 [NLG<sup>+</sup>19] Mor Shpigel Nacson, Jason Lee, Suriya Gunasekar, Pedro Henrique Pamplona Savarese,  
287 Nathan Srebro, and Daniel Soudry. Convergence of gradient descent on separable data.  
288 In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages  
289 3420–3428. PMLR, 2019.
- 290 [NLL<sup>+</sup>23] Lorenzo Noci, Chuning Li, Mufan Bill Li, Bobby He, Thomas Hofmann, Chris Mad-  
291 dison, and Daniel M Roy. The shaped transformer: Attention models in the infinite  
292 depth-and-width limit. *arXiv preprint arXiv:2306.17759*, 2023.
- 293 [NNH<sup>+</sup>23] Tan Minh Nguyen, Tam Minh Nguyen, Nhat Ho, Andrea L Bertozzi, Richard Baraniuk,  
294 and Stanley Osher. A primal-dual framework for transformers and neural networks. In  
295 *The Eleventh International Conference on Learning Representations*, 2023.
- 296 [OH10] Samet Oymak and Babak Hassibi. New null space results and recovery thresholds for  
297 matrix rank minimization. *arXiv preprint arXiv:1011.6326*, 2010.
- 298 [OMFH11] Samet Oymak, Karthik Mohan, Maryam Fazel, and Babak Hassibi. A simplified  
299 approach to recovery conditions for low rank matrices. In *2011 IEEE International*  
300 *Symposium on Information Theory Proceedings*, pages 2318–2322. IEEE, 2011.
- 301 [Ope22] OpenAI. OpenAI: Introducing ChatGPT. <https://openai.com/blog/chatgpt>,  
302 2022.
- 303 [Ope23] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 304 [ORST23] Samet Oymak, Ankit Singh Rawat, Mahdi Soltanolkotabi, and Christos Thrampoulidis.  
305 On the role of attention in prompt-tuning. In *International Conference on Machine*  
306 *Learning*, 2023.
- 307 [PHD20] Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during  
308 the terminal phase of deep learning training. *Proceedings of the National Academy of*  
309 *Sciences*, 117(40):24652–24663, 2020.
- 310 [PTDU16] Ankur Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. A decomposable  
311 attention model for natural language inference. In *Proceedings of the 2016 Conference*  
312 *on Empirical Methods in Natural Language Processing*, pages 2249–2255, Austin,  
313 Texas, November 2016. Association for Computational Linguistics.
- 314 [PXS18] Romain Paulus, Caiming Xiong, and Richard Socher. A deep reinforced model for  
315 abstractive summarization. In *International Conference on Learning Representations*,  
316 2018.
- 317 [QQ19] Qian Qian and Xiaoyuan Qian. The implicit bias of adagrad on separable data. *Advances*  
318 *in Neural Information Processing Systems*, 32, 2019.
- 319 [RSR<sup>+</sup>20] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael  
320 Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learn-  
321 ing with a unified text-to-text transformer. *Journal of Machine Learning Research*,  
322 21(1):5485–5551, 2020.
- 323 [RXH11] Benjamin Recht, Weiyu Xu, and Babak Hassibi. Null space conditions and thresholds  
324 for rank minimization. *Mathematical programming*, 127:175–202, 2011.
- 325 [RZH03] Saharon Rosset, Ji Zhu, and Trevor Hastie. Margin maximizing loss functions. *Advances*  
326 *in neural information processing systems*, 16, 2003.
- 327 [SATA22] Haoyuan Sun, Kwangjun Ahn, Christos Thrampoulidis, and Navid Azizan. Mirror  
328 descent maximizes generalized margin and can be implemented efficiently. *Advances*  
329 *in Neural Information Processing Systems*, 35:31089–31101, 2022.
- 330 [SEO<sup>+</sup>22] Arda Sahiner, Tolga Ergen, Batu Ozturkler, John Pauly, Morteza Mardani, and Mert  
331 Pilanci. Unraveling attention via convex duality: Analysis and interpretations of vision  
332 transformers. In *International Conference on Machine Learning*, pages 19050–19088.  
333 PMLR, 2022.



- 334 [SHN<sup>+</sup>18] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan  
335 Srebro. The implicit bias of gradient descent on separable data. *The Journal of Machine*  
336 *Learning Research*, 19(1):2822–2878, 2018.
- 337 [SRJ04] Nathan Srebro, Jason Rennie, and Tommi Jaakkola. Maximum-margin matrix factoriza-  
338 tion. *Advances in neural information processing systems*, 17, 2004.
- 339 [SS21] Dominik Stöger and Mahdi Soltanolkotabi. Small random initialization is akin to  
340 spectral learning: Optimization and generalization guarantees for overparameterized  
341 low-rank matrix reconstruction. *Advances in Neural Information Processing Systems*,  
342 34:23831–23843, 2021.
- 343 [TBS<sup>+</sup>16] Stephen Tu, Ross Boczar, Max Simchowitz, Mahdi Soltanolkotabi, and Ben Recht.  
344 Low-rank solutions of linear matrix equations via procrustes flow. In *International*  
345 *Conference on Machine Learning*, pages 964–973. PMLR, 2016.
- 346 [TCD<sup>+</sup>21] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablay-  
347 rolles, and Hervé Jégou. Training data-efficient image transformers & distillation  
348 through attention. In *International Conference on Machine Learning*, pages 10347–  
349 10357. PMLR, 2021.
- 350 [TG07] Joel A Tropp and Anna C Gilbert. Signal recovery from random measurements via  
351 orthogonal matching pursuit. *IEEE Transactions on information theory*, 53(12):4655–  
352 4666, 2007.
- 353 [Tib96] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the*  
354 *Royal Statistical Society Series B: Statistical Methodology*, 58(1):267–288, 1996.
- 355 [TKVB22] Christos Thrampoulidis, Ganesh Ramachandra Kini, Vala Vakilian, and Tina Behnia. Im-  
356 balance trouble: Revisiting neural-collapse geometry. *Advances in Neural Information*  
357 *Processing Systems*, 35:27225–27238, 2022.
- 358 [TLI<sup>+</sup>23] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux,  
359 Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar,  
360 et al. Llama: Open and efficient foundation language models. *arXiv preprint*  
361 *arXiv:2302.13971*, 2023.
- 362 [TLZO23] Davoud Ataee Tarzanagh, Yingcong Li, Xuechen Zhang, and Samet Oymak. Margin  
363 maximization in attention mechanism. *arXiv preprint arXiv:2306.13596*, 2023.
- 364 [TVS23] Nadav Timor, Gal Vardi, and Ohad Shamir. Implicit regularization towards rank  
365 minimization in relu networks. In *International Conference on Algorithmic Learning*  
366 *Theory*, pages 1429–1459. PMLR, 2023.
- 367 [TWCD23] Yuandong Tian, Yiping Wang, Beidi Chen, and Simon Du. Scan and snap: Understand-  
368 ing training dynamics and token composition in 1-layer transformer. *arXiv:2305.16380*,  
369 2023.
- 370 [VKR19] Tomas Vaskevicius, Varun Kanade, and Patrick Rebeschini. Implicit regularization  
371 for optimal sparse recovery. *Advances in Neural Information Processing Systems*,  
372 32:2972–2983, 2019.
- 373 [VSP<sup>+</sup>17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N  
374 Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in*  
375 *neural information processing systems*, 30, 2017.
- 376 [WGL<sup>+</sup>20] Blake Woodworth, Suriya Gunasekar, Jason D Lee, Edward Moroshko, Pedro Savarese,  
377 Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in over-  
378 parametrized models. In *Conference on Learning Theory*, pages 3635–3673. PMLR,  
379 2020.
- 380 [WMCL21] Bohan Wang, Qi Meng, Wei Chen, and Tie-Yan Liu. The implicit bias for adaptive  
381 optimization algorithms on homogeneous neural networks. In *International Conference*  
382 *on Machine Learning*, pages 10849–10858. PMLR, 2021.

383 [WMZ<sup>+</sup>21] Bohan Wang, Qi Meng, Huishuai Zhang, Ruoyu Sun, Wei Chen, and Zhi-Ming Ma.  
384 Momentum doesn't change the implicit bias. *arXiv preprint arXiv:2110.03891*, 2021.

385 [WWX<sup>+</sup>22] Haixu Wu, Jialong Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Flowformer:  
386 Linearizing transformers with conservation flows. In *International Conference on*  
387 *Machine Learning*, pages 24226–24242, 2022.

388 [YKM20] Chulhee Yun, Shankar Krishnan, and Hossein Mobahi. A unifying view on implicit  
389 bias in training linear neural networks. *arXiv preprint arXiv:2010.02501*, 2020.

390 [ZFB23] Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear  
391 models in-context. *arXiv preprint arXiv:2306.09927*, 2023.

392 [ZWB<sup>+</sup>21] Difan Zou, Jingfeng Wu, Vladimir Braverman, Quanquan Gu, Dean P Foster, and Sham  
393 Kakade. The benefits of implicit regularization from sgd in least squares problems.  
394 *Advances in Neural Information Processing Systems*, 34:5456–5468, 2021.

395 [ZY05] Tong Zhang and Bin Yu. Boosting with early stopping: Convergence and consistency.  
396 *Annals of Statistics*, page 1538, 2005.

397 **Roadmap.** The appendix is organized as follows:

- 398 • Appendix **A** provides related work.
- 399 • Appendix **B** provides detailed discussion and experimental evaluation about Section 3.3.
- 400 • Appendix **C** provides auxiliary lemmas.
- 401 • Appendix **D** presents the proof for the global regularization path analysis (Section 2).
- 402 • Appendix **E** presents the proofs for the gradient descent convergences (Section 3).
- 403 • Appendix **F** provides additional experiments and their discussion.
- 404 • Appendix **G** discusses potential further directions.

## 405 Contents

406	<b>A Related work</b>	<b>11</b>
407	A.1 Implicit Regularization, Matrix Factorization, Sparsity . . . . .	11
408	A.2 Attention Mechanism and Transformers . . . . .	12
409	<b>B Understanding Multi-token Compositions: Toward A More General Max-Margin and</b>	
410	<b>Directional Convergence Theory</b>	<b>12</b>
411	B.1 When does attention select multiple tokens? . . . . .	15
412	B.2 Proof of Lemma 3 . . . . .	16
413	<b>C Auxiliary Lemmas</b>	<b>17</b>
414	C.1 Proof of Lemma 1 . . . . .	17
415	C.2 Proof of Lemma 2 . . . . .	17
416	C.3 Useful Lemmas . . . . .	19
417	<b>D Global Regularization Path</b>	<b>20</b>
418	D.1 Proof of Theorem 1 . . . . .	20

419	<b>E Convergence of Gradient Descent</b>	<b>22</b>
420	E.1 Divergence of norm of the iterates $\mathbf{W}(k)$ . . . . .	22
421	E.2 Global Convergence of Gradient Descent . . . . .	24
422	E.2.1 Proof of Theorem 2. . . . .	25
423	E.3 Local Convergence of Gradient Descent . . . . .	26
424	E.3.1 Proof of Theorem 3 . . . . .	30
425	<b>F Supporting Experiments</b>	<b>33</b>
426	<b>G Discussion, Future Directions, and Open Problems</b>	<b>36</b>

## 427 A Related work

### 428 A.1 Implicit Regularization, Matrix Factorization, Sparsity

429 Extensive research has delved into gradient descent’s implicit bias in separable classification  
430 tasks, often using logistic or exponentially-tailed losses for margin maximization [SHN<sup>+</sup>18,  
431 GLSS18, NLG<sup>+</sup>19, JT21, KPOT21, MWG<sup>+</sup>20, JT20]. The findings have also been extended  
432 to non-separable data using gradient-based techniques [JT18, JT19, JDST20]. Implicit bias  
433 in regression problems and losses has been investigated, utilizing methods like mirror descent  
434 [WGL<sup>+</sup>20, GLSS18, YKM20, VKR19, AW20a, AW20b, ALH21, SATA22]. Stochastic gradient  
435 descent has also been a subject of interest regarding its implicit bias [LWM19, BGVV20, LR20,  
436 HWLM20, LWA22, DML21, ZWB<sup>+</sup>21]. This extends to the implicit bias of adaptive and momentum-  
437 based methods [QQ19, WMZ<sup>+</sup>21, WMCL21, JST21].

438 In linear classification, GD iterations on logistic loss and separable datasets converge to the hard  
439 margin SVM solution [SHN<sup>+</sup>18, RZH03, ZY05]. The attention layer’s softmax nonlinearity behaves  
440 similarly, potentially favoring margin-maximizing solutions. Yet, the layer operates on tokens  
441 in input sequences, not for direct classification. Its bias leans toward an (**Att-SVM**), selecting  
442 relevant tokens while suppressing others. However, formalizing this intuition presents significant  
443 challenges: Firstly, our problem is nonconvex (even in terms of the  $\mathbf{W}$ -parameterization), introducing  
444 new challenges and complexities. Secondly, it requires the introduction of novel concepts such as  
445 locally-optimal tokens, demanding a tailored analysis focused on the cones surrounding them. Our  
446 findings on the implicit bias of  $(\mathbf{K}, \mathbf{Q})$ -parameterization share conceptual similarities with [SRJ04],  
447 which proposes and analyzes a max-margin matrix factorization problem. Similar problems have  
448 also been studied more recently in the context of neural-collapse phenomena [PHD20] through  
449 an analysis of the implicit bias and regularization path of the unconstrained features model with  
450 cross-entropy loss [TKVB22]. However, a fundamental distinction from these works lies in the fact  
451 that attention solves a different max-margin problem that separate tokens. Moreover, our results  
452 on  $(\mathbf{K}, \mathbf{Q})$ -parameterization are inherently connected to the rich literature on low-rank factorization  
453 [GWB<sup>+</sup>17, ACHL19, TVS23, TBS<sup>+</sup>16, SS21], stimulating further research. [TLZO23] is the first  
454 work to establish the connection between attention and SVM, which is closest to our work. Here,  
455 we augment their framework, initially developed for a simpler attention model, to transformers by  
456 providing the first guarantees for self/cross-attention layers, nonlinear prediction heads, and realistic  
457 global convergence guarantees. While our Assumption (i) and local-convergence analysis align with  
458 [TLZO23], our contributions in global convergence analysis, benefits of overparameterization, and  
459 the generalized SVM-equivalence in Section B are unique to this work.

460 It is well-known that attention map (i.e. softmax outputs) act as a feature selection mechanism and  
461 reveal the tokens that are relevant to classification. On the other hand, sparsity and lasso regression  
462 (i.e.  $\ell_1$  penalization) [Don06, Tib96, TG07, CDS01, CRT06] have been pivotal tools in the statistics  
463 literature for feature selection. Softmax and lasso regression exhibit interesting parallels: The Softmax  
464 output  $s = \mathbb{S}(\mathbf{X}\mathbf{W}\mathbf{z})$  obeys  $\|s\|_{\ell_1} = 1$  by design. Softmax is also highly receptive to being sparse  
465 because decreasing the temperature (i.e. scaling up the weights  $\mathbf{W}$ ) eventually leads to a one-hot vector  
466 unless all logits are equal. We (also, [TLZO23]) have used these intuitions to formalize attention as a  
467 *token selection mechanism*. This aspect is clearly visible in our primary SVM formulation (**Att-SVM**)

468 which selects precisely one token from each input sequence (i.e. hard attention). Section B has also  
 469 demonstrated how (Gen-SVM) can explain more general sparsity patterns by precisely selecting  
 470 desired tokens and suppressing others. We hope that this SVM-based token-selection viewpoint will  
 471 motivate future work and deeper connections to the broader feature-selection and compressed sensing  
 472 literature.

## 473 A.2 Attention Mechanism and Transformers

474 Transformers, as highlighted by [VSP<sup>+</sup>17], revolutionized the domains of NLP and machine transla-  
 475 tion. Prior work on self-attention [CDL16, PTDU16, PXS18, LFS<sup>+</sup>17] laid the foundation for this  
 476 transformative paradigm. In contrast to conventional models like MLPs and CNNs, self-attention mod-  
 477 els employ global interactions to capture feature representations, resulting in exceptional empirical  
 478 performance.

479 Despite their achievements, the mechanisms and learning processes of attention layers remain  
 480 enigmatic. Recent investigations [EGKZ22, SEO<sup>+</sup>22, ENM22, BV22, DCL21] have concentrated  
 481 on specific aspects such as sparse function representation, convex relaxations, and expressive power.  
 482 Expressivity discussions concerning hard-attention [Hah20] or attention-only architectures [DCL21]  
 483 are connected to our findings when  $h(\cdot)$  is linear. In fact, our work reveals how linear  $h$  results  
 484 in attention’s optimization dynamics to collapse on a single token whereas nonlinear  $h$  provably  
 485 requires attention to select and compose multiple tokens. This supports the benefits of the MLP layer  
 486 for expressivity of transformers. There is also a growing body of research aimed at a theoretical  
 487 comprehension of in-context learning and the role played by the attention mechanism [ASA<sup>+</sup>22,  
 488 LIPO23, ACDS23, ZFB23, BCW<sup>+</sup>23, GRS<sup>+</sup>23]. [SEO<sup>+</sup>22] investigate self-attention with linear  
 489 activation instead of softmax, while [ENM22] approximate softmax using a linear operation with  
 490 unit simplex constraints. Their primary goal is to derive convex reformulations for training problems  
 491 grounded in empirical risk minimization (ERM). In contrast, our methodologies, detailed in equations  
 492 (W-ERM) and (KQ-ERM), delve into the nonconvex domain.

493 [MRG<sup>+</sup>20, BALA<sup>+</sup>23] offer insights into the implicit bias of optimizing transformers. Specifically,  
 494 [MRG<sup>+</sup>20] provide empirical evidence that an increase in attention weights results in a sparser  
 495 softmax, which aligns with our theoretical framework. [BALA<sup>+</sup>23] study incremental learning and  
 496 furnish both theory and numerical evidence that increments of the softmax attention weights ( $KQ^T$ )  
 497 are low-rank. Our theory aligns with this concept, as the SVM formulation of ( $K, Q$ ) parameterization  
 498 inherently exhibits low-rank properties through the nuclear norm objective, rank- $m$  constraint, and  
 499 implicit constraint induced by Lemma 1.

500 Several recent works [JSL22, LWLC23, TWCD23, NLL<sup>+</sup>23, ORST23, NNH<sup>+</sup>23, FGBM23] aim to  
 501 delineate the optimization and generalization dynamics of transformers. However, their findings usu-  
 502 ally apply under strict statistical assumptions about the data, while our study offers a comprehensive  
 503 optimization-theoretic analysis of the attention model, establishing a formal linkage to max-margin  
 504 problems and SVM geometry. This allows our findings to encompass the problem geometry and apply  
 505 to diverse datasets. Overall, the max-margin equivalence provides a fundamental comprehension of  
 506 the optimization geometry of transformers, offering a framework for prospective research endeavors,  
 507 as outlined in the subsequent section.

## 508 B Understanding Multi-token Compositions: Toward A More General 509 Max-Margin and Directional Convergence Theory

510 So far, our theory has focused on the setting where the attention layer selects a single optimal token  
 511 within each sequence. As we have discussed, this is theoretically well-justified under linear head  
 512 assumption and certain nonlinear generalizations. On the other hand, for arbitrary nonconvex  $h(\cdot)$   
 513 or multilayer transformer architectures, it is expected that attention will select multiple tokens per  
 514 sequence. This motivates us to ask:

515 **Q:** What is the implicit bias and the form of  $W(k)$  when the GD solution is  
 516 composed by multiple tokens?

517 In this section, our goal is to derive and verify the generalized behavior of GD. Let  $\mathbf{o}_i = X_i^\top s_i^W$   
 518 denote the token generated by the attention layer where  $s_i^W = \mathbb{S}(X_i W z_i)$  are the softmax probabilities.

519 Suppose GD trajectory converges to achieve the risk  $\mathcal{L}_\star = \min_W \mathcal{L}(W)$ . Suppose the eventual token  
 520 composition achieving  $\mathcal{L}_\star$  is given by

$$\mathbf{o}_i^\star = \mathbf{X}_i^\top \mathbf{s}_i^\star,$$

521 where  $\mathbf{s}_i^\star$  are the eventual softmax probability vectors that dictate the token composition. Since  
 522 attention maps are sparse in practice, we are interested in the scenario where  $\mathbf{s}_i^\star$  is sparse i.e. it  
 523 contains some zero entries. This can only be accomplished by letting  $\|W\|_F \rightarrow \infty$ . However, unlike  
 524 the earlier sections, we wish to allow for arbitrary  $\mathbf{s}_i^\star$  rather than a one-hot vector which selects a  
 525 single token.

526 To proceed, we aim to understand the form of GD solution  $W(k)$  responsible for composing  $\mathbf{o}_i^\star$  via  
 527 the softmax map  $\mathbf{s}_i^\star$  as  $R \rightarrow \infty$ . Intuitively,  $W(k)$  should be decomposed into two components via

$$W(k) \approx W^{\text{fin}} + \|W(k)\|_F \cdot \bar{W}^{\text{mm}}. \quad (5)$$

528 where  $W^{\text{fin}}$  is the finite component and  $\bar{W}^{\text{mm}}$  is the directional component with  $\|\bar{W}^{\text{mm}}\|_F = 1$ . Define  
 529 the selected set  $O_i \subseteq [T]$  to be the indices  $s_{it}^\star \neq 0$  and the masked (i.e. suppressed) set as  $\bar{O}_i = [T] - O_i$   
 530 where softmax entries are zero. In the context of earlier sections, we could also call these the *optimal*  
 531 *set* and the *non-optimal set*, respectively.

532 • **Finite component:** The job of  $W^{\text{fin}}$  is to assign nonzero softmax probabilities within each  $\mathbf{s}_i^\star$ .  
 533 This is accomplished by ensuring that,  $W^{\text{fin}}$  induces the probabilities of  $\mathbf{s}_i^\star$  over  $O_i$  by satisfying the  
 534 softmax equations

$$\frac{e^{\mathbf{x}_{it}^\top W^{\text{fin}} \mathbf{z}_i}}{e^{\mathbf{x}_{i\tau}^\top W^{\text{fin}} \mathbf{z}_i}} = e^{(\mathbf{x}_{it} - \mathbf{x}_{i\tau})^\top W^{\text{fin}} \mathbf{z}_i} = s_{it}^\star / s_{i\tau}^\star$$

535 for  $t, \tau \in O_i$ . Consequently, this  $W^{\text{fin}}$  should satisfy the following linear constraints

$$(\mathbf{x}_{it} - \mathbf{x}_{i\tau})^\top W^{\text{fin}} \mathbf{z}_i = \log(s_{it}^\star / s_{i\tau}^\star) \quad \text{for all } t, \tau \in O_i, i \in [n]. \quad (6)$$

536 • **Directional component:** While  $W^{\text{fin}}$  creates the composition by allocating the nonzero softmax  
 537 probabilities, it does not explain sparsity of attention map. This is the role of  $\bar{W}^{\text{mm}}$ , which is  
 538 responsible for selecting the selected tokens  $O_i$  and suppressing the masked ones  $\bar{O}_i$  by assigning  
 539 zero softmax probability to them. To predict direction component, we build on the theory developed  
 540 in earlier sections. Concretely, there are two constraints  $\bar{W}^{\text{mm}}$  should satisfy

- 541 1. **Equal similarity over selected tokens:** For all  $t, \tau \in O_i$ , we have that  $(\mathbf{x}_{it} - \mathbf{x}_{i\tau})^\top W \mathbf{z}_i = 0$ .  
 542 This way, softmax scores assigned by  $W^{\text{fin}}$  are not disturbed by the directional component and  
 543  $W^{\text{fin}} + R\bar{W}^{\text{mm}}$  will still satisfy the softmax equations (6).
- 544 2. **Max-margin against masked tokens:** For all  $t \in O_i, \tau \in \bar{O}_i$ , enforce the margin constraint  
 545  $(\mathbf{x}_{it} - \mathbf{x}_{i\tau})^\top W \mathbf{z}_i \geq 1$  subject to minimum norm  $\|W\|_F$ .

546 Combining these, we obtain the following convex generalized SVM formulation

$$W^{\text{mm}} = \arg \min_W \|W\|_F \quad \text{subj. to} \quad \begin{cases} \forall t \in O_i, \tau \in \bar{O}_i : (\mathbf{x}_{it} - \mathbf{x}_{i\tau})^\top W \mathbf{z}_i \geq 1, & \forall 1 \leq i \leq n. \\ \forall t, \tau \in O_i : (\mathbf{x}_{it} - \mathbf{x}_{i\tau})^\top W \mathbf{z}_i = 0, & \end{cases} \quad (\text{Gen-SVM})$$

547

548 and set the normalized direction in (5) to  $\bar{W}^{\text{mm}} = W^{\text{mm}} / \|W^{\text{mm}}\|_F$ .

549 It is important to note that **(Gen-SVM)** offers a substantial generalization beyond the scope of the  
 550 previous sections, where the focus was on selecting a single token from each sequence, as described  
 551 in the main formulation **(Att-SVM)**. This broader solution class introduces a more flexible approach  
 552 to the problem.

553 We present experiments showcasing the predictive power of the **(Gen-SVM)** equivalence in nonlinear  
 554 scenarios. We conducted these experiments on random instances using an MLP denoted as  $h(\cdot)$ ,  
 555 which takes the form of  $\mathbf{1}^\top \text{ReLU}(\mathbf{x})$ . We begin by detailing the preprocessing step and our setup. For  
 556 the attention SVM equivalence analytical prediction, clear definitions of the selected and masked sets  
 557 are crucial. These sets include token indices with nonzero and zero softmax outputs, respectively.  
 558 However, practically, reaching a precisely zero output is not feasible. Hence, we define the selected  
 559 set as tokens with softmax outputs exceeding  $10^{-3}$ , and the masked set as tokens with softmax outputs

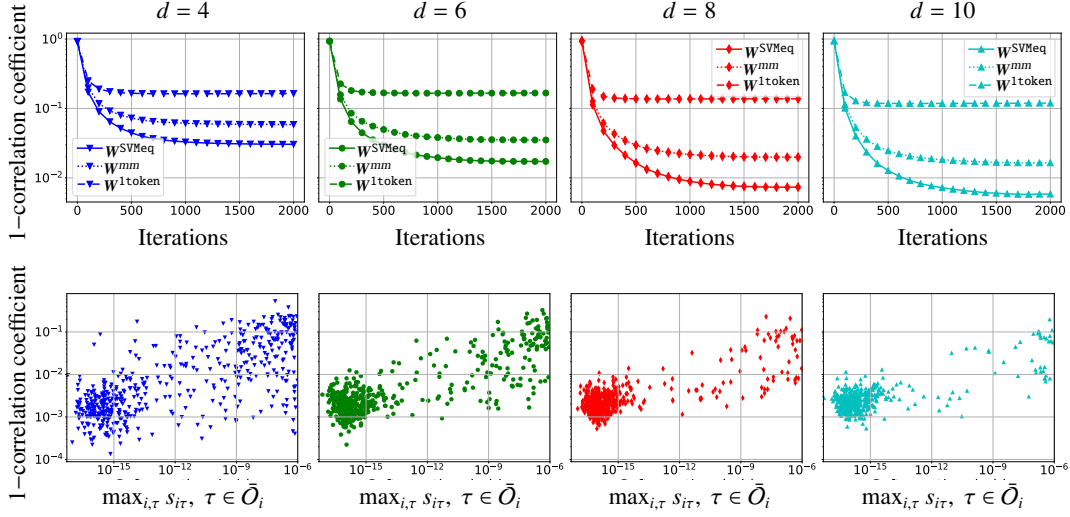


Figure 4: Behavior of GD with nonlinear nonconvex prediction head and multi-token compositions. **Upper:** The correlation between GD solution and three distinct baselines: ( $\cdots$ )  $W^{mm}$  obtained from (**Gen-SVM**); (—)  $W^{SVMeq}$  obtained by calculating  $W^{fin}$  and determining the best linear combination  $W^{fin} + \gamma W^{mm}$  that maximizes correlation with the GD solution; and (- -)  $W^{1token}$  obtained by solving (**Att-SVM**) and selecting the highest probability token from the GD solution. **Lower:** Scatterplot of the largest softmax probability over masked tokens (per our  $s_{i\tau} \leq 10^{-6}$  criteria) vs correlation coefficient.

560 below  $10^{-6}$ . We also excluded instances with softmax outputs falling between  $10^{-6}$  and  $10^{-3}$   
561 to distinctly separate the concepts of *selected* and *masked* sets, thereby enhancing the predictive accuracy  
562 of the attention SVM equivalence. In addition to the filtering process, we focus on scenarios where  
563 the label  $Y = -1$  exists to enforce *non-convexity* of prediction head  $Y_i \cdot h(\cdot)$ . It is worth mentioning  
564 that when all labels are 1, due to the convexity of  $Y_i \cdot h(\cdot)$ , GD tends to select one token per input,  
565 and Equations (**Gen-SVM**) and (**Att-SVM**) yield the same solutions. The results are displayed in  
566 Figure 4, where  $n = 3$ ,  $T = 4$ , and  $d$  varies within 4, 6, 8, 10. We conduct 500 random trials for  
567 different choices of  $d$ , each involving  $x_{it}$ ,  $z_i$ , and  $v$  randomly sampled from the unit sphere. We apply  
568 normalized GD with a step size  $\eta = 0.1$  and run 2000 iterations for each trial.

569 • Figure 4 (upper) illustrates the correlation evolution between the GD solution and three distinctive  
570 baselines: ( $\cdots$ )  $W^{mm}$  obtained from (**Gen-SVM**); (—)  $W^{SVMeq}$  obtained by calculating  $W^{fin}$   
571 and determining the best linear combination  $W^{fin} + \gamma W^{mm}$  that maximizes correlation with the GD  
572 solution; and (- -)  $W^{1token}$  obtained by solving (**Att-SVM**) and selecting the highest probability token  
573 from the GD solution. For clearer visualization, the logarithmic scale of correlation misalignment is presented  
574 in Figure 4. In essence, our findings show that  $W^{1token}$  yields unsatisfactory outcomes, whereas  
575  $W^{mm}$  attains a significant correlation coefficient in alignment with our expectations. Ultimately,  
576 our comprehensive SVM-equivalence  $W^{SVMeq}$  further enhances correlation, lending support to our  
577 analytical formulas. It’s noteworthy that SVM-equivalence displays higher predictability in a larger  $d$   
578 regime (with an average correlation exceeding 0.99). This phenomenon might be attributed to more  
579 frequent directional convergence in higher dimensions, with overparameterization contributing to a  
580 smoother loss landscape, thereby expediting optimization.

581 • Figure 4 (lower) offers a scatterplot overview of the 500 random problem instances that were  
582 solved. The  $x$ -axis represents the largest softmax probability over the masked set, denoted as  $\max_{i,\tau} s_{i\tau}$   
583 where  $\tau \in \bar{O}_i$ . Meanwhile, the  $y$ -axis indicates the predictivity of the SVM-equivalence, quantified as  
584  $1 - \text{corr\_coef}(W, W^{SVMeq})$ . From this analysis, two significant observations arise. Primarily, there  
585 exists an inverse correlation between softmax probability and SVM-predictivity. This correlation  
586 is intuitive, as higher softmax probabilities signify a stronger divergence from our desired *masked*  
587 *set* state (ideally set to 0). Secondly, as dimensionality ( $d$ ) increases, softmax probabilities over the  
588 masked set tend to converge towards the range of  $10^{-15}$  (effectively zero). Simultaneously, attention  
589 SVM-predictivity improves, creating a noteworthy correlation.



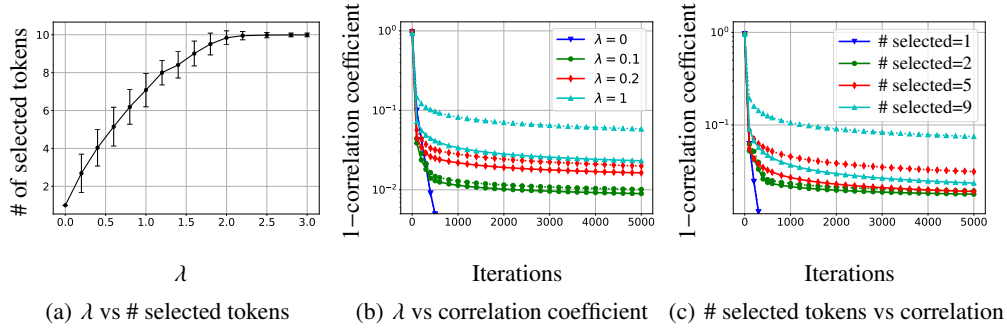


Figure 5: Behavior of GD when selecting multiple tokens. **(a)** The number of selected tokens increases with  $\lambda$ . **(b)** Predictivity of attention SVM solutions for varying  $\lambda$ ; Dotted curves depict the correlation corresponding to  $\mathbf{W}^{mm}$  calculated via (Gen-SVM) and solid curves represent the correlation to  $\mathbf{W}^{SVM^{eq}}$ , which incorporates the  $\mathbf{W}^{fin}$  correction. **(c)** Similar to (b), but evaluating correlations over different numbers of selected tokens.

### 590 B.1 When does attention select multiple tokens?

591 In this section, we provide a concrete example where the optimal solution indeed requires combining  
 592 multiple tokens in a nontrivial fashion. Here, by nontrivial we mean that, we select more than 1  
 593 tokens from an input sequence but we don't select all of its tokens. Recall that, for linear prediction  
 594 head, attention will ideally select the single token with largest score for almost all datasets. Perhaps  
 595 not surprisingly, this behavior will not persist for nonlinear prediction heads. For instance in Figure 4,  
 596 the GD output  $\mathbf{W}$  aligned better in direction with  $\mathbf{W}^{mm}$  than  $\mathbf{W}^{1token}$ . Specifically, here we prove that  
 597 if we make the function  $h_Y(\mathbf{x}) := Y \cdot h(\mathbf{x})$  concave, then optimal softmax map can select multiple  
 598 tokens in a controllable fashion.  $h_Y(\mathbf{x})$  can be viewed as generalization of the linear score function  
 599  $Y \cdot \mathbf{v}^\top \mathbf{x}$ . In the example below, we induce concavity by incorporating a small  $-\lambda \|\mathbf{x}\|^2$  term within a  
 600 linear prediction head and setting  $h(\mathbf{x}) = \mathbf{v}^\top \mathbf{x} - \lambda \|\mathbf{x}\|^2$  with  $Y = 1$ .

601 **Lemma 3** Given  $\mathbf{v} \in \mathbb{R}^d$ , recall the score vector  $\boldsymbol{\gamma} = \mathbf{X}\mathbf{v}$ . Without losing generality, assume  $\boldsymbol{\gamma}$  is  
 602 non-increasing. Define the vector of score gaps  $\boldsymbol{\gamma}^{gap} \in \mathbb{R}^{T-1}$  with entries  $\gamma_i^{gap} = \gamma_i - \gamma_{i+1}$ . Suppose  
 603 all tokens within the input sequence are orthonormal and for some  $\tau \geq 2$ , we have that

$$\tau \gamma_\tau^{gap} / 2 > \gamma_1^{gap}. \quad (7)$$

604 Set  $h(\mathbf{x}) = \mathbf{v}^\top \mathbf{x} - \lambda \|\mathbf{x}\|^2$  where  $\tau \gamma_\tau^{gap} / 2 > \lambda > \gamma_1^{gap}$ ,  $\ell(x) = -x$ , and  $Y = 1$ . Let  $\Delta_T$  denote the  
 605  $T$ -dimensional simplex. Define the unconstrained softmax optimization associated to the objective  $h$   
 606 where we make  $\mathbf{s} := \mathbb{S}(\mathbf{X}\mathbf{W}\mathbf{z})$  a free variable, namely,

$$\min_{\mathbf{s} \in \Delta_T} \ell(h(\mathbf{X}\mathbf{s})) = \min_{\mathbf{s} \in \Delta_T} \lambda \|\mathbf{X}^\top \mathbf{s}\|^2 - \mathbf{v}^\top \mathbf{X}^\top \mathbf{s}. \quad (8)$$

607 Then, the optimal solution  $\mathbf{s}^*$  contains at least 2 and at most  $\tau$  nonzero entries.

608 Figure 5 presents experimental findings concerning Lemma 3 across random problem instances. For  
 609 this experiment, we set  $n = 1$ ,  $T = 10$ , and  $d = 10$ . The results are averaged over 100 random  
 610 trials, with each trial involving the generation of randomly orthonormal vectors  $\mathbf{x}_{1t}$  and the random  
 611 sampling of vector  $\mathbf{v}$  from the unit sphere. Similar to the processing step in Figure 4, and following  
 612 Figure 4 (lower) which illustrates that smaller softmax outputs over masked sets correspond to higher  
 613 correlation coefficients, we define the selected and masked token sets. Specifically, tokens with  
 614 softmax outputs  $> 10^{-3}$  are considered selected, while tokens with softmax outputs  $< 10^{-8}$  are  
 615 masked. Instances with softmax outputs between  $10^{-8}$  and  $10^{-3}$  are filtered out.

616 Figure 5(a) shows that the number of selected tokens grows alongside  $\lambda$ , a prediction consistent with  
 617 Lemma 3. When  $\lambda = 0$ , the head  $h(\mathbf{x}) = \mathbf{v}^\top \mathbf{x}$  is linear, resulting in the selection of only one token  
 618 per input. Conversely, as  $\lambda$  exceeds a certain threshold (e.g.,  $\lambda > 2.0$  based on our criteria), the  
 619 optimization consistently selects all tokens. Figure 5(b) and 5(c) delve into the predictivity of attention  
 620 SVM solutions for varying  $\lambda$  and different numbers of selected tokens. The dotted curves in both  
 621 figures represent  $1 - \text{corr\_coef}(\mathbf{W}, \mathbf{W}^{mm})$ , while solid curves indicate  $1 - \text{corr\_coef}(\mathbf{W}, \mathbf{W}^{SVM^{eq}})$ ,



622 where  $\mathbf{W}$  denotes the GD solution. Overall, the SVM-equivalence demonstrates a strong correlation  
 623 with the GD solution (consistently above 0.95). However, selecting more tokens (aligned with larger  
 624  $\lambda$  values) leads to reduced predictivity.

625 To sum up, we have showcased the predictive capacity of the generalized SVM equivalence regarding  
 626 the inductive bias of 1-layer transformers with nonlinear heads. Nevertheless, it's important to  
 627 acknowledge that this section represents an initial approach to a complex problem, with certain  
 628 caveats requiring further investigation (e.g., the use of filtering in Figures 4 and 5, and the presence of  
 629 imperfect correlations). We aspire to conduct a more comprehensive investigation, both theoretically  
 630 and empirically, in forthcoming work.

## 631 B.2 Proof of Lemma 3

632 Suppose  $\tau$  described by (7) exists and set  $\lambda$  accordingly. Let  $\mathcal{S} \subset [T]$  denote the top  $\tau$  indices of  
 633  $\boldsymbol{\gamma}$  with largest scores. Denote  $\mathbf{X}^1 \in \mathbb{R}^{\tau \times d}$  to be the sequence corresponding to  $\mathcal{S}$  and  $\mathbf{X}^2 \in \mathbb{R}^{(T-\tau) \times d}$   
 634 to be the sequence corresponding to  $[T] - \mathcal{S}$ . Similarly, denote the subvectors  $\boldsymbol{\gamma}_1, \mathbf{s}^{(1)} \in \mathbb{R}^\tau$  and  
 635  $\boldsymbol{\gamma}_2, \mathbf{s}^{(2)} \in \mathbb{R}^{T-\tau}$  and define the probability over  $\mathcal{S}$  as  $S_1 = \sum_{i \in \mathcal{S}} s_i$ . The orthogonality and unit norm  
 636 assumption on the tokens imply

$$1 \geq \|\mathbf{X}^\top \mathbf{s}\|^2 = \sum_{i=1}^T s_i^2 \geq S_1^2/\tau + (1 - S_1)^2/(T - \tau).$$

637 Also note that  $\mathbf{v}^\top \mathbf{X} \mathbf{s} = \boldsymbol{\gamma}_1^\top \mathbf{s}^{(1)} + \boldsymbol{\gamma}_2^\top \mathbf{s}^{(2)}$ . With these, we can write the objective  $\mathcal{L}(\mathbf{s}) := \ell(h(\mathbf{X} \mathbf{s}))$  as  
 638 follows

$$\mathcal{L}(\mathbf{s}) = \lambda \sum_{i=1}^T s_i^2 - \boldsymbol{\gamma}_1^\top \mathbf{s}^{(1)} - \boldsymbol{\gamma}_2^\top \mathbf{s}^{(2)}.$$

639 Note that, for fixed  $\boldsymbol{\gamma}$  and over all permutations of entries of  $\mathbf{s}$ ,  $\boldsymbol{\gamma}^\top \mathbf{s}$  is maximized when  $\mathbf{s}$  and  $\boldsymbol{\gamma}$  are  
 640 aligned namely, when the entries of  $\mathbf{s}$  are sorted according to the entries of  $\boldsymbol{\gamma}$ . Otherwise, we could  
 641 swap two unsorted entries of  $\mathbf{s}$  (i.e. with unaligned  $\boldsymbol{\gamma}$  entries) to a sorted position to obtain a strictly  
 642 better optimal (where we also used the fact that  $\mathbf{s}$  has nonnegative entries). Thus, we can assume the  
 643 entries of  $\mathbf{s}^*$  are sorted according to  $\boldsymbol{\gamma}$ . Specifically, the largest  $\tau$  entries of  $\mathbf{s}^*$  lie on the set  $\mathcal{S}$ .

644 • **We first show that  $\mathbf{s} := \mathbf{s}^*$  cannot have more than  $\tau$  entries.** To prove this, we compare  $\mathbf{s}$  against  
 645 the baseline  $\bar{\mathbf{s}}$  where  $\bar{\mathbf{s}}^1 = \mathbf{s}^{(1)}/S_1$  and  $\bar{\mathbf{s}}^2 = 0$  so that  $\bar{\mathbf{s}}$  is  $\tau$ -sparse. In this scenario,  $\bar{\mathbf{s}}$  yields the  
 646 objective

$$\mathcal{L}(\bar{\mathbf{s}}) = \frac{\lambda}{S_1^2} \sum_{i \in \mathcal{S}} s_i^2 - \frac{1}{S_1} \boldsymbol{\gamma}_1^\top \mathbf{s}^{(1)}.$$

647 We claim that  $\mathcal{L}(\bar{\mathbf{s}}) < \mathcal{L}(\mathbf{s})$ . To see this, we first observe that  $\boldsymbol{\gamma}_1^\top \mathbf{s}^{(1)}/S_1 \geq \boldsymbol{\gamma}_2^\top \mathbf{s}^{(2)}/(1 - S_1) + \boldsymbol{\gamma}_\tau^{\text{gap}}$ .  
 648 This implies

$$(1/S_1 - 1)\boldsymbol{\gamma}_1^\top \mathbf{s}^{(1)} - \boldsymbol{\gamma}_2^\top \mathbf{s}^{(2)} \geq (1 - S_1)\boldsymbol{\gamma}_\tau^{\text{gap}}.$$

649 Recalling  $\sum_{i \in \mathcal{S}} s_i^2 \leq S_1^2/\tau$ , we can now utilize the following chain of implications

$$\begin{aligned} & \mathcal{L}(\bar{\mathbf{s}}) < \mathcal{L}(\mathbf{s}) \\ \iff & \frac{\lambda}{S_1^2} \sum_{i \in \mathcal{S}} s_i^2 - \frac{1}{S_1} \boldsymbol{\gamma}_1^\top \mathbf{s}^{(1)} < \lambda \sum_{i=1}^T s_i^2 - \boldsymbol{\gamma}_1^\top \mathbf{s}^{(1)} - \boldsymbol{\gamma}_2^\top \mathbf{s}^{(2)} \\ \iff & \lambda(1/S_1^2 - 1) \sum_{i \in \mathcal{S}} s_i^2 < (1/S_1 - 1)\boldsymbol{\gamma}_1^\top \mathbf{s}^{(1)} - \boldsymbol{\gamma}_2^\top \mathbf{s}^{(2)} \\ \iff & \lambda(1/S_1^2 - 1) \sum_{i \in \mathcal{S}} s_i^2 < (1 - S_1)\boldsymbol{\gamma}_\tau^{\text{gap}} \\ \iff & \lambda(1 - S_1^2)/\tau < (1 - S_1)\boldsymbol{\gamma}_\tau^{\text{gap}} \\ \iff & \lambda(1 + S_1)/\tau < \boldsymbol{\gamma}_\tau^{\text{gap}} \\ \iff & 2\lambda/\tau < \boldsymbol{\gamma}_\tau^{\text{gap}} \\ \iff & \lambda < \tau \boldsymbol{\gamma}_\tau^{\text{gap}}/2. \end{aligned}$$

650 • **We next prove that there are at least two nonzeros in the optimal solution.** Denote the largest  
651 and second largest entry of  $\boldsymbol{\gamma}$  by  $\bar{\gamma}_1$  and  $\bar{\gamma}_2$  respectively. For  $\mathbf{s}^{\text{one}} \in \Delta_T$  containing a single nonzero  
652 (i.e. one-hot vector), the best achievable risk is given by

$$\mathcal{L}(\mathbf{s}^{\text{one}}) = \lambda - \bar{\gamma}_1.$$

653 On the other hand consider the 2-sparse reference solution  $\mathbf{s}^{\text{ref}}$  which assigns equal likelihood over  
654 the top two entries. This achieves

$$\mathcal{L}(\mathbf{s}^{\text{ref}}) = \frac{\lambda}{2} - \boldsymbol{\gamma}^\top \mathbf{s}^{\text{ref}} \leq \frac{\lambda}{2} - \frac{\bar{\gamma}_1 + \bar{\gamma}_2}{2}.$$

655 The latter is superior as soon as

$$\frac{\lambda}{2} - \frac{\bar{\gamma}_1 + \bar{\gamma}_2}{2} < \lambda - \bar{\gamma}_1 \iff \lambda > \boldsymbol{\gamma}_1^{\text{gap}}.$$

656 Thus, we conclude with the statement by selecting  $\tau \boldsymbol{\gamma}_\tau^{\text{gap}} / 2 > \lambda > \boldsymbol{\gamma}_1^{\text{gap}}$ . ■

## 657 C Auxiliary Lemmas

### 658 C.1 Proof of Lemma 1

659 Suppose the claim is wrong and row space of  $\mathbf{W}_\diamond^{mm}$  does not lie within  $\mathcal{S} = \text{span}(\{\mathbf{z}_i\}_{i=1}^n)$ . Let  
660  $\mathbf{W} = \Pi_{\mathcal{S}}(\mathbf{W}_\diamond^{mm})$  denote the matrix obtained by projecting the rows of  $\mathbf{W}_\diamond^{mm}$  on  $\mathcal{S}$ . Observe that  $\mathbf{W}$   
661 satisfies all SVM constraints since  $\mathbf{W}\mathbf{z}_i = \mathbf{W}_\diamond^{mm}\mathbf{z}_i$  for all  $i \in [n]$ . For Frobenius norm, using  $\mathbf{W}_\diamond^{mm} \neq \mathbf{W}$ ,  
662 we obtain a contradiction via  $\|\mathbf{W}_\diamond^{mm}\|_F^2 = \|\mathbf{W}\|_F^2 + \|\mathbf{W}_\diamond^{mm} - \mathbf{W}\|_F^2 > \|\mathbf{W}\|_F^2$ . For nuclear norm, we can  
663 write  $\mathbf{W} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^\top$  with  $\boldsymbol{\Sigma} \in \mathbb{R}^{r \times r}$  where  $r$  is dimension of  $\mathcal{S}$  and  $\text{column\_span}(\mathbf{V}) = \mathcal{S}$ .

664 To proceed, we split the problem into two scenarios.

665 **Scenario 1:** Let  $\mathbf{U}_\perp, \mathbf{V}_\perp$  be orthogonal complements of  $\mathbf{U}, \mathbf{V}$  – viewing matrices with orthonormal  
666 columns as subspaces. Suppose  $\mathbf{U}_\perp^\top \mathbf{W}_\diamond^{mm} \mathbf{V}_\perp \neq 0$ . Then, singular value inequalities (which were  
667 also used in earlier works on nuclear norm analysis [RXH11, OH10, OMFH11]) guarantee that  
668  $\|\mathbf{W}_\diamond^{mm}\|_\star \geq \|\mathbf{U}_\perp^\top \mathbf{W}_\diamond^{mm} \mathbf{V}\|_\star + \|\mathbf{U}_\perp^\top \mathbf{W}_\diamond^{mm} \mathbf{V}_\perp\|_\star > \|\mathbf{W}\|_\star$ .

669 **Scenario 2:** Now suppose  $\mathbf{U}_\perp^\top \mathbf{W}_\diamond^{mm} \mathbf{V}_\perp = 0$ . Since  $\mathbf{W}_\diamond^{mm} \mathbf{V}_\perp \neq 0$ , this implies  $\mathbf{U}^\top \mathbf{W}_\diamond^{mm} \mathbf{V}_\perp \neq 0$ . Let  
670  $\mathbf{W}' = \mathbf{U}\mathbf{U}^\top \mathbf{W}_\diamond^{mm}$  which is a rank- $r$  matrix. Since  $\mathbf{W}'$  is a subspace projection, we have  $\|\mathbf{W}'\|_\star \leq$   
671  $\|\mathbf{W}_\diamond^{mm}\|_\star$ . Next, observe that  $\|\mathbf{W}\|_\star = \text{trace}(\mathbf{U}^\top \mathbf{W} \mathbf{V}) = \text{trace}(\mathbf{U}^\top \mathbf{W}' \mathbf{V})$ . On the other hand,  
672  $\text{trace}(\mathbf{U}^\top \mathbf{W}' \mathbf{V}) < \|\mathbf{W}'\|_\star$  because the equality in *von Neumann's trace inequality* happens if and  
673 only if the two matrices we are inner-producting, namely  $(\mathbf{W}', \mathbf{U}\mathbf{V}^\top)$ , share a joint set of singular  
674 vectors [Car21]. However, this is not true as the row space of  $\mathbf{W}_\diamond^{mm}$  does not lie within  $\mathcal{S}$ . Thus, we  
675 obtain  $\|\mathbf{W}\|_\star < \|\mathbf{W}'\|_\star \leq \|\mathbf{W}_\diamond^{mm}\|_\star$  concluding the proof via contradiction. ■

### 676 C.2 Proof of Lemma 2

677 **Lemma 4 (Lemma 2 restated)** Under Assumption A,  $\nabla \mathcal{L}(\mathbf{W})$ ,  $\nabla_{\mathbf{K}} \mathcal{L}(\mathbf{K}, \mathbf{Q})$ , and  $\nabla_{\mathbf{Q}} \mathcal{L}(\mathbf{K}, \mathbf{Q})$  are  $L_{\mathbf{W}}$ ,  
678  $L_{\mathbf{K}}$ ,  $L_{\mathbf{Q}}$ -Lipschitz continuous, respectively, where  $a_i = \|\mathbf{v}\| \|\mathbf{z}_i\|^2 \|\mathbf{X}_i\|^3$ ,  $b_i = M_0 \|\mathbf{v}\| \|\mathbf{X}_i\| + 3M_1$  for all  
679  $i \in [n]$ ,

$$L_{\mathbf{W}} := \frac{1}{n} \sum_{i=1}^n a_i b_i, \quad L_{\mathbf{K}} := \|\mathbf{Q}\| L_{\mathbf{W}}, \quad \text{and} \quad L_{\mathbf{Q}} := \|\mathbf{K}\| L_{\mathbf{W}}. \quad (9)$$

680 **Proof.** Let

$$\boldsymbol{\gamma}_i = Y_i \cdot \mathbf{X}_i \mathbf{v}, \quad \mathbf{h}_i = \mathbf{X}_i \mathbf{W} \mathbf{z}_i. \quad (10)$$

681 From Assumption A,  $\ell : \mathbb{R} \rightarrow \mathbb{R}$  is differentiable. Hence, the gradient evaluated at  $\mathbf{W}$  is given by

$$\nabla \mathcal{L}(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \ell'(\boldsymbol{\gamma}_i^\top \mathcal{S}(\mathbf{h}_i)) \cdot \mathbf{X}_i^\top \mathcal{S}'(\mathbf{h}_i) \boldsymbol{\gamma}_i \mathbf{z}_i^\top, \quad (11)$$

682 where

$$\mathcal{S}'(\mathbf{h}) = \text{diag}(\mathcal{S}(\mathbf{h})) - \mathcal{S}(\mathbf{h})\mathcal{S}(\mathbf{h})^\top \in \mathbb{R}^{T \times T}. \quad (12)$$

683 Note that

$$\|\mathbb{S}'(\mathbf{h})\| \leq \|\mathbb{S}'(\mathbf{h})\|_F \leq 1. \quad (13)$$

684 Hence, for any  $\mathbf{W}, \dot{\mathbf{W}} \in \mathbb{R}^{d \times d}$ ,  $i \in [n]$ , we have

$$\|\mathbb{S}(\mathbf{h}_i) - \mathbb{S}(\dot{\mathbf{h}}_i)\| \leq \|\mathbf{h}_i - \dot{\mathbf{h}}_i\| \leq \|\mathbf{X}_i\| \|\mathbf{z}_i\| \|\mathbf{W} - \dot{\mathbf{W}}\|_F, \quad (14a)$$

685 where  $\dot{\mathbf{h}}_i = \mathbf{X}_i \dot{\mathbf{W}} \mathbf{z}_i$ .

686 Similarly,

$$\begin{aligned} \|\mathbb{S}'(\mathbf{h}_i) - \mathbb{S}'(\dot{\mathbf{h}}_i)\|_F &\leq \|\mathbb{S}(\mathbf{h}_i) - \mathbb{S}(\dot{\mathbf{h}}_i)\| + \|\mathbb{S}(\mathbf{h}_i)\mathbb{S}(\mathbf{h}_i)^\top - \mathbb{S}(\dot{\mathbf{h}}_i)\mathbb{S}(\dot{\mathbf{h}}_i)^\top\|_F \\ &\leq 3\|\mathbf{X}_i\| \|\mathbf{z}_i\| \|\mathbf{W} - \dot{\mathbf{W}}\|_F. \end{aligned} \quad (14b)$$

687 Next, for any  $\mathbf{W}, \dot{\mathbf{W}} \in \mathbb{R}^{d \times d}$ , we get

$$\begin{aligned} \|\nabla \mathcal{L}(\mathbf{W}) - \nabla \mathcal{L}(\dot{\mathbf{W}})\|_F &\leq \frac{1}{n} \sum_{i=1}^n \left\| \ell'(\boldsymbol{\gamma}_i^\top \mathbb{S}(\mathbf{h}_i)) \cdot \mathbf{z}_i \boldsymbol{\gamma}_i^\top \mathbb{S}'(\mathbf{h}_i) \mathbf{X}_i - \ell'(\boldsymbol{\gamma}_i^\top \mathbb{S}(\dot{\mathbf{h}}_i)) \cdot \mathbf{z}_i \boldsymbol{\gamma}_i^\top \mathbb{S}'(\dot{\mathbf{h}}_i) \mathbf{X}_i \right\|_F \\ &\leq \frac{1}{n} \sum_{i=1}^n \|\mathbf{z}_i \boldsymbol{\gamma}_i^\top \mathbb{S}'(\dot{\mathbf{h}}_i) \mathbf{X}_i\|_F \left| \ell'(\boldsymbol{\gamma}_i^\top \mathbb{S}(\mathbf{h}_i)) - \ell'(\boldsymbol{\gamma}_i^\top \mathbb{S}(\dot{\mathbf{h}}_i)) \right| \\ &\quad + \frac{1}{n} \sum_{i=1}^n \left| \ell'(\boldsymbol{\gamma}_i^\top \mathbb{S}(\mathbf{h}_i)) \right| \|\mathbf{z}_i \boldsymbol{\gamma}_i^\top \mathbb{S}'(\mathbf{h}_i) \mathbf{X}_i - \mathbf{z}_i \boldsymbol{\gamma}_i^\top \mathbb{S}'(\dot{\mathbf{h}}_i) \mathbf{X}_i\|_F \\ &\leq \frac{1}{n} \sum_{i=1}^n M_0 \|\boldsymbol{\gamma}_i\|^2 \|\mathbf{z}_i\| \|\mathbf{X}_i\| \|\mathbb{S}(\mathbf{h}_i) - \mathbb{S}(\dot{\mathbf{h}}_i)\| \\ &\quad + \frac{1}{n} \sum_{i=1}^n M_1 \|\boldsymbol{\gamma}_i\| \|\mathbf{z}_i\| \|\mathbf{X}_i\| \|\mathbb{S}'(\mathbf{h}_i) - \mathbb{S}'(\dot{\mathbf{h}}_i)\|_F, \end{aligned} \quad (15)$$

688 where the second inequality follows from the fact that  $|ab - cd| \leq |d||a - c| + |a||b - d|$  and the third  
689 inequality uses Assumption A and (13).

690 Substituting (14a) and (14b) into (15), we get

$$\begin{aligned} \|\nabla \mathcal{L}(\mathbf{W}) - \nabla \mathcal{L}(\dot{\mathbf{W}})\|_F &\leq \frac{1}{n} \sum_{i=1}^n \left( M_0 \|\boldsymbol{\gamma}_i\|^2 \|\mathbf{z}_i\|^2 \|\mathbf{X}_i\|^2 + 3M_1 \|\boldsymbol{\gamma}_i\| \|\mathbf{z}_i\| \|\mathbf{X}_i\|^2 \right) \|\mathbf{W} - \dot{\mathbf{W}}\|_F \\ &\leq \frac{1}{n} \sum_{i=1}^n \left( M_0 \|\boldsymbol{\nu}\|^2 \|\mathbf{z}_i\|^2 \|\mathbf{X}_i\|^4 + 3M_1 \|\boldsymbol{\nu}\| \|\mathbf{z}_i\|^2 \|\mathbf{X}_i\|^3 \right) \|\mathbf{W} - \dot{\mathbf{W}}\|_F \\ &\leq L_W \|\mathbf{W} - \dot{\mathbf{W}}\|_F, \end{aligned}$$

691 where  $L_W$  is defined in (9).

692 Let  $\mathbf{g}_i = \mathbf{X}_i \mathbf{K} \boldsymbol{\mathcal{Q}}^\top \mathbf{z}_i$ . We have

$$\nabla_{\mathbf{K}} \mathcal{L}(\mathbf{K}, \boldsymbol{\mathcal{Q}}) = \frac{1}{n} \sum_{i=1}^n \ell'(\boldsymbol{\gamma}_i^\top \mathbb{S}(\mathbf{g}_i)) \cdot \mathbf{z}_i \boldsymbol{\gamma}_i^\top \mathbb{S}'(\mathbf{g}_i) \mathbf{X}_i \boldsymbol{\mathcal{Q}}, \quad (16a)$$

$$\nabla_{\boldsymbol{\mathcal{Q}}} \mathcal{L}(\mathbf{K}, \boldsymbol{\mathcal{Q}}) = \frac{1}{n} \sum_{i=1}^n \ell'(\boldsymbol{\gamma}_i^\top \mathbb{S}(\mathbf{g}_i)) \cdot \mathbf{X}_i^\top \mathbb{S}'(\mathbf{g}_i) \boldsymbol{\gamma}_i \mathbf{z}_i^\top \mathbf{K}. \quad (16b)$$

693 By the similar argument as in (15), for any  $\boldsymbol{\mathcal{Q}}$  and  $\dot{\boldsymbol{\mathcal{Q}}} \in \mathbb{R}^{d \times m}$ , we have

$$\begin{aligned} \|\nabla_{\boldsymbol{\mathcal{Q}}} \mathcal{L}(\mathbf{K}, \boldsymbol{\mathcal{Q}}) - \nabla_{\boldsymbol{\mathcal{Q}}} \mathcal{L}(\mathbf{K}, \dot{\boldsymbol{\mathcal{Q}}})\|_F &\leq \frac{\|\mathbf{K}\|}{n} \sum_{i=1}^n \left\| \ell'(\boldsymbol{\gamma}_i^\top \mathbb{S}(\mathbf{h}_i)) \cdot \mathbf{z}_i \boldsymbol{\gamma}_i^\top \mathbb{S}'(\mathbf{h}_i) \mathbf{X}_i - \ell'(\boldsymbol{\gamma}_i^\top \mathbb{S}(\dot{\mathbf{h}}_i)) \cdot \mathbf{z}_i \boldsymbol{\gamma}_i^\top \mathbb{S}'(\dot{\mathbf{h}}_i) \mathbf{X}_i \right\|_F \\ &\leq L_W \|\mathbf{K}\| \|\boldsymbol{\mathcal{Q}} - \dot{\boldsymbol{\mathcal{Q}}}\|_F. \end{aligned} \quad (17)$$

Similarly, for any  $\mathbf{K}, \dot{\mathbf{K}} \in \mathbb{R}^{d \times m}$ , we get

$$\|\nabla_{\mathbf{K}} \mathcal{L}(\mathbf{K}, \boldsymbol{\mathcal{Q}}) - \nabla_{\mathbf{K}} \mathcal{L}(\dot{\mathbf{K}}, \boldsymbol{\mathcal{Q}})\|_F \leq L_W \|\boldsymbol{\mathcal{Q}}\| \|\mathbf{K} - \dot{\mathbf{K}}\|_F.$$

694

■

695 **C.3 Useful Lemmas**

696 **Lemma 5 (Optimal Tokens Minimize Training Loss)** *Suppose Assumption A (i)-(ii) hold, and not*  
 697 *all tokens are optimal per Definition 1. Then, training risk obeys  $\mathcal{L}(\mathbf{W}) > \mathcal{L}_\star := \frac{1}{n} \sum_{i=1}^n \ell(\boldsymbol{\gamma}_{i\text{opt}_i})$ .*  
 698 *Additionally, suppose there are optimal indices  $(\text{opt}_i)_{i=1}^n$  for which (Att-SVM) is feasible, i.e. there*  
 699 *exists a  $\mathbf{W}$  separating optimal tokens. This  $\mathbf{W}$  choice obeys  $\lim_{R \rightarrow \infty} \mathcal{L}(R \cdot \mathbf{W}) = \mathcal{L}_\star$ .*

700 The result presented in Lemma 5 originates from the observation that the output tokens of the attention  
 701 layer constitute a convex combination of the input tokens. Consequently, when subjected to a strictly  
 702 decreasing loss function, attention optimization inherently leans towards the selection of a singular  
 703 token, specifically, the optimal token  $(\text{opt}_i)_{i=1}^n$ .

704 **Proof.** The token at the output of the attention layer is given by  $\mathbf{a}_i = \mathbf{X}_i^\top \mathbb{S}(\mathbf{X}_i \mathbf{W} \mathbf{z}_i)$ . Here,  $\mathbf{a}_i$  can be  
 705 written as  $\mathbf{a}_i = \sum_{t \in [T]} c_{it} \mathbf{x}_{it}$  where  $c_{it} \geq 0$  and  $\sum_{t \in [T]} c_{it} = 1$ . Note that, for any finite  $\mathbf{W}$ ,  $c_{it}$  as softmax  
 706 probabilities are strictly positive. To proceed, using the linearity of  $h(\mathbf{x}) = \mathbf{v}^\top \mathbf{x}$  and strictly-decreasing  
 707 nature of the loss  $\ell$ , we find that

$$\mathcal{L}(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i \cdot h(\mathbf{a}_i)) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i \cdot \sum_{t \in [T]} c_{it} h(\mathbf{x}_{it})) \geq \frac{1}{n} \sum_{i=1}^n \ell(Y_i \cdot h(\mathbf{x}_{i\text{opt}_i})) = \mathcal{L}_\star,$$

708 which implies that  $\mathcal{L}(\mathbf{W}) \geq \mathcal{L}_\star$  for any  $\mathbf{W}$ .

709 On the other hand, since not all tokens are optimal, there exists a token index  $(i, t)$  for which  
 710  $Y_i \cdot h(\mathbf{x}_{it}) < Y_i \cdot h(\mathbf{x}_{i\text{opt}_i})$ . Since all softmax entries obey  $c_{it} > 0$  for finite  $\mathbf{W}$ , this implies the strict  
 711 inequality  $\ell(Y_i \cdot h(\mathbf{a}_i)) > \ell(Y_i \cdot h(\mathbf{x}_{i\text{opt}_i}))$ . This leads to the desired conclusion  $\mathcal{L}(\mathbf{W}) > \mathcal{L}_\star$ .

712 Secondly, suppose (Att-SVM) is feasible i.e. there exists a  $\mathbf{W}$  separating some optimal indices  
 713  $(\text{opt}_i)_{i=1}^n$  from the other tokens. Note that, this does not exclude the existence of other optimal  
 714 indices. This implies that, letting  $\lim_{R \rightarrow \infty} \mathbb{S}(\mathbf{X}_i(R \cdot \mathbf{W}) \mathbf{z}_i)$  saturates the softmax and will be equal to the  
 715 indicator function at  $\text{opt}_i$  for all inputs  $i \in [n]$ . Thus,  $c_{it} \rightarrow 0$  for  $t \neq \text{opt}_i$  and  $c_{it} \rightarrow 1$  for  $t = \text{opt}_i$ .  
 716 Using  $M_1$ -Lipschitzness of  $\ell$ , we can write

$$\left| \ell(Y_i \cdot h(\mathbf{x}_{i\text{opt}_i})) - \ell(Y_i \cdot h(\mathbf{a}_i)) \right| \leq M_1 \left| h(\mathbf{a}_i) - h(\mathbf{x}_{i\text{opt}_i}) \right|.$$

717 Since  $h$  is linear, it is  $\|\mathbf{v}\|$ -Lipschitz implying

$$\left| \ell(Y_i \cdot h(\mathbf{x}_{i\text{opt}_i})) - \ell(Y_i \cdot h(\mathbf{a}_i)) \right| \leq M_1 \|\mathbf{v}\| \cdot \|\mathbf{a}_i - \mathbf{x}_{i\text{opt}_i}\|.$$

718 Since  $\mathbf{a}_i \rightarrow \mathbf{x}_{i\text{opt}_i}$  as  $R \rightarrow \infty$ , we conclude with the advertised result. ■ ■

719 **Lemma 6** *For any  $\mathbf{X} \in \mathbb{R}^{T \times d}$ ,  $\mathbf{W}, \mathbf{V} \in \mathbb{R}^{d \times d}$  and  $\mathbf{z}, \mathbf{v} \in \mathbb{R}^d$ , let  $\mathbf{a} = \mathbf{XVz}$ ,  $\mathbf{s} = \mathbb{S}(\mathbf{XWz})$ , and  $\boldsymbol{\gamma} = \mathbf{Xv}$ .*  
 720 *Set*

$$\Gamma = \sup_{t, \tau \in [T]} |\gamma_t - \gamma_\tau| \quad \text{and} \quad A = \sup_{t \in [T]} \|\mathbf{a}_t\|.$$

721 *We have that*

$$\left| \mathbf{a}^\top \text{diag}(\mathbf{s}) \boldsymbol{\gamma} - \mathbf{a}^\top \mathbf{s} \mathbf{s}^\top \boldsymbol{\gamma} - \sum_{t \geq 2}^T (\mathbf{a}_1 - \mathbf{a}_t) s_t (\gamma_1 - \gamma_t) \right| \leq 2\Gamma A (1 - s_1)^2.$$

722 **Proof.** The proof is similar to [TLZ023, Lemma 4], but for the sake of completeness, we provide it  
 723 here. Set  $\bar{\boldsymbol{\gamma}} = \sum_{t=1}^T \boldsymbol{\gamma}_t s_t$ . We have

$$\boldsymbol{\gamma}_1 - \bar{\boldsymbol{\gamma}} = \sum_{t \geq 2}^T (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) s_t, \quad \text{and} \quad |\boldsymbol{\gamma}_1 - \bar{\boldsymbol{\gamma}}| \leq \Gamma (1 - s_1).$$

724 Then,

$$\begin{aligned} \mathbf{a}^\top \text{diag}(\mathbf{s}) \boldsymbol{\gamma} - \mathbf{a}^\top \mathbf{s} \mathbf{s}^\top \boldsymbol{\gamma} &= \sum_{t=1}^T \mathbf{a}_t \boldsymbol{\gamma}_t s_t - \sum_{t=1}^T \mathbf{a}_t s_t \sum_{t=1}^T \boldsymbol{\gamma}_t s_t \\ &= \mathbf{a}_1 s_1 (\boldsymbol{\gamma}_1 - \bar{\boldsymbol{\gamma}}) - \sum_{t \geq 2}^T \mathbf{a}_t s_t (\bar{\boldsymbol{\gamma}} - \boldsymbol{\gamma}_t). \end{aligned} \tag{18}$$

Since

$$\left| \sum_{t \geq 2}^T \mathbf{a}_t s_t (\bar{\gamma} - \gamma_t) - \sum_{t \geq 2}^T \mathbf{a}_t s_t (\gamma_1 - \gamma_t) \right| \leq A\Gamma(1 - s_1)^2,$$

725 we obtain<sup>1</sup>

$$\begin{aligned} \mathbf{a}^\top \text{diag}(s)\boldsymbol{\gamma} - \mathbf{a}^\top s s^\top \boldsymbol{\gamma} &= \mathbf{a}_1 s_1 (\gamma_1 - \bar{\gamma}) - \sum_{t \geq 2}^T \mathbf{a}_t s_t (\gamma_1 - \gamma_t) \pm A\Gamma(1 - s_1)^2 \\ &= \mathbf{a}_1 s_1 \sum_{t \geq 2}^T (\gamma_1 - \gamma_t) s_t - \sum_{t \geq 2}^T \mathbf{a}_t s_t (\gamma_1 - \gamma_t) \pm A\Gamma(1 - s_1)^2 \\ &= \sum_{t \geq 2}^T (\mathbf{a}_1 s_1 - \mathbf{a}_t) s_t (\gamma_1 - \gamma_t) \pm A\Gamma(1 - s_1)^2 \\ &= \sum_{t \geq 2}^T (\mathbf{a}_1 - \mathbf{a}_t) s_t (\gamma_1 - \gamma_t) \pm 2A\Gamma(1 - s_1)^2. \end{aligned}$$

726 Here,  $\pm$  on the right handside uses the fact that

$$\left| \sum_{t \geq 2}^T (\mathbf{a}_1 s_1 - \mathbf{a}_t) s_t (\gamma_1 - \gamma_t) \right| \leq (1 - s_1) A\Gamma \sum_{t \geq 2}^T s_t = (1 - s_1)^2 A\Gamma.$$

727

■

## 728 D Global Regularization Path

### 729 D.1 Proof of Theorem 1

730 Throughout  $\diamond$  denotes either Frobenius norm or nuclear norm. We will prove that  $\bar{\mathbf{W}}(R)$  asymptotically  
731 aligns with the set of globally-optimal directions and also  $\|\bar{\mathbf{W}}(R)\|_\diamond \rightarrow \infty$ .  $\mathcal{R}_m \subseteq \mathbb{R}^{d \times d}$  denote the  
732 manifold of rank  $\leq m$  matrices.

733 **Step 1:** Let us first prove that  $\bar{\mathbf{W}}(R)$  achieves the optimal risk as  $R \rightarrow \infty$  – rather than problem  
734 having finite optima. Define  $\Xi_\diamond = 1/\|\mathbf{W}^{mm}\|_\diamond$  and norm-normalized  $\bar{\mathbf{W}}^{mm} = \Xi_\diamond \mathbf{W}^{mm}$ . Note that  $\mathbf{W}^{mm}$   
735 separates tokens opt from rest of the tokens for each  $i \in [n]$ . Thus, we have that

$$\lim_{R \rightarrow \infty} \mathcal{L}(\bar{\mathbf{W}}(R)) \leq \lim_{R \rightarrow \infty} \mathcal{L}(R \cdot \bar{\mathbf{W}}^{mm}) := \mathcal{L}_\star = \frac{1}{n} \sum_{i=1}^n \ell(\gamma_i^{\text{opt}}). \quad (19)$$

736 On the other hand, for any  $\mathbf{W} \in \mathcal{R}_m$ , define the softmax probabilities  $\mathbf{s}^{(i)} = \mathbb{S}(X_i \mathbf{W} z_i)$  and attention  
737 features  $\mathbf{x}_i^{\mathbf{W}} = \sum_{t=1}^T s_t^{(i)} \mathbf{x}_t$ . Decompose  $\mathbf{x}_i^{\mathbf{W}}$  as  $\mathbf{x}_i^{\mathbf{W}} = \mathbf{s}_{\text{opt}_i}^{(i)} \mathbf{x}_{i\text{opt}_i} + \sum_{t \neq \text{opt}_i} s_t^{(i)} \mathbf{x}_{it}$ . Set  $\gamma_{it}^{\text{gap}} = \gamma_i^{\text{opt}} - \gamma_{it} =$   
738  $Y_i \cdot \mathbf{v}^\top (\mathbf{x}_{i\text{opt}_i} - \mathbf{x}_{it}) > 0$ , and define

$$B := \max_{i \in [n]} \max_{t, \tau \in [T]} \|\mathbf{v}\| \cdot \|\mathbf{x}_{it} - \mathbf{x}_{i\tau}\| \geq \gamma_{it}^{\text{gap}}. \quad (20)$$

739 Define  $c_{\text{opt}} = \min_{i \in [n], t \neq \text{opt}_i} \gamma_{it}^{\text{gap}} > 0$  and  $\boldsymbol{\gamma}_i^{\mathbf{W}} = Y_i \cdot \mathbf{v}^\top \mathbf{x}_i^{\mathbf{W}}$ . We obtain the following score inequalities

$$\begin{aligned} \boldsymbol{\gamma}_i^{\mathbf{W}} &\leq \boldsymbol{\gamma}_i^{\text{opt}} - c_{\text{opt}}(1 - \mathbf{s}_{\text{opt}_i}^{(i)}) < \boldsymbol{\gamma}_i^{\text{opt}}, \\ |\boldsymbol{\gamma}_i^{\mathbf{W}} - \boldsymbol{\gamma}_i^{\text{opt}}| &\leq \|\mathbf{v}\| \cdot \|\mathbf{x}_i^{\mathbf{W}} - \mathbf{x}_i^{\alpha}\| \leq \|\mathbf{v}\| \sum_{t \neq \text{opt}_i} s_t^{(i)} \|\mathbf{x}_{it} - \mathbf{x}_i^{\alpha}\| \leq B(1 - \mathbf{s}_{\text{opt}_i}^{(i)}). \end{aligned} \quad (21)$$

740 We will use the  $\boldsymbol{\gamma}_i^{\mathbf{W}} - \boldsymbol{\gamma}_i^{\text{opt}}$  term in (21) to evaluate  $\mathbf{W}$  against the reference loss  $\mathcal{L}_\star$  of (19). Using the  
741 strictly-decreasing nature of  $\ell$ , we conclude with the fact that for all (finite)  $\mathbf{W} \in \mathcal{R}_m$ ,

$$\mathcal{L}(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \ell(\boldsymbol{\gamma}_i^{\mathbf{W}}) > \mathcal{L}_\star = \frac{1}{n} \sum_{i=1}^n \ell(\boldsymbol{\gamma}_i^{\text{opt}}),$$

<sup>1</sup>For simplicity, we use  $\pm$  on the right hand side to denote the upper and lower bounds.

742 which implies  $\|\bar{\mathbf{W}}(R)\|_\diamond \rightarrow \infty$  together with (19).

743 **Step 2:** To proceed, we show that  $\bar{\mathbf{W}}(R)$  converges in direction to  $\mathcal{W}^{mm}$ , which denotes the set of  
 744 SVM minima. Suppose this is not the case and convergence fails. We will obtain a contradiction by  
 745 showing that  $\bar{\mathbf{W}}_R^{mm} = R \cdot \bar{\mathbf{W}}^{mm}$  achieves a strictly superior loss compared to  $\bar{\mathbf{W}}(R)$ . Let us introduce  
 746 the normalized parameters  $\bar{\mathbf{W}}_0(R) = \frac{\bar{\mathbf{W}}(R)}{R\Xi_\circ}$  and  $\mathbf{W}' = \frac{\bar{\mathbf{W}}(R)}{\|\bar{\mathbf{W}}(R)\|_\diamond \Xi_\circ}$ . Note that  $\bar{\mathbf{W}}_0(R)$  is obtained by scaling  
 747 down  $\mathbf{W}'$  since  $\|\bar{\mathbf{W}}(R)\|_\diamond \leq R$  and  $\mathbf{W}'$  obeys  $\|\mathbf{W}'\|_\diamond = \|\mathbf{W}^{mm}\|_\diamond$ . Since  $\bar{\mathbf{W}}_0(R)$  fails to converge to  $\mathcal{W}^{mm}$ ,  
 748 for some  $\delta > 0$ , there exists arbitrarily large  $R > 0$  such that  $\text{dist}(\bar{\mathbf{W}}_0(R), \mathcal{W}^{mm}) \geq \delta$ . This translates  
 749 to the suboptimality in terms of the margin constraints as follows: First, since nuclear norm dominates  
 750 Frobenius, distance with respect to the  $\diamond$ -norm obeys  $\text{dist}_\diamond(\bar{\mathbf{W}}_0(R), \mathcal{W}^{mm}) \geq \delta$ . Secondly, using  
 751 triangle inequality,

$$\text{this implies that either } \|\bar{\mathbf{W}}_0(R)\|_\diamond \leq \|\mathbf{W}^{mm}\|_\diamond - \delta/2 \text{ or } \text{dist}_\diamond(\mathbf{W}', \mathcal{W}^{mm}) \geq \delta/2.$$

752 In either scenario,  $\bar{\mathbf{W}}_0(R)$  strictly violates one of the margin constraints of (Att-SVM) ( $\diamond = F$ ) or  
 753 (Att-SVM $\star$ ) ( $\diamond = \star$ ): If  $\|\bar{\mathbf{W}}_0(R)\|_\diamond \leq \|\mathbf{W}^{mm}\|_\diamond - \delta/2$ , then, since the optimal SVM objective is  
 754  $\|\mathbf{W}^{mm}\|_\diamond$ , there exists a constraint  $i, t \neq \text{opt}_i$  for which  $\langle (\mathbf{x}_i^{\text{opt}} - \mathbf{x}_{it})\mathbf{z}_i^\top, \bar{\mathbf{W}}_0(R) \rangle \leq 1 - \frac{\delta}{2\|\mathbf{W}^{mm}\|_\diamond}$ . If  
 755  $\text{dist}_\diamond(\mathbf{W}', \mathcal{W}^{mm}) \geq \delta/2$ , then,  $\mathbf{W}'$  has the same SVM objective but it is strictly bounded away from  
 756 the solution set. Thus, for some  $\epsilon := \epsilon(\delta) > 0$ ,  $\mathbf{W}'$  and its scaled down version  $\bar{\mathbf{W}}_0(R)$  strictly violate  
 757 an SVM constraint achieving margin  $\leq 1 - \epsilon$ . Without losing generality, suppose  $\bar{\mathbf{W}}_0(R)$  violates the  
 758 first constraint  $i = 1$ . Thus, for a properly updated  $\delta > 0$  (that is function of the initial  $\delta > 0$ ) and for  
 759  $i = 1$  and some support index  $\tau \in \mathcal{T}_1$ ,

$$\langle (\mathbf{x}_1^{\text{opt}} - \mathbf{x}_{1\tau})\mathbf{z}_1^\top, \bar{\mathbf{W}}_0(R) \rangle \leq 1 - \delta. \quad (22)$$

760 Now, we will argue that this leads to a contradiction by proving  $\mathcal{L}(\bar{\mathbf{W}}_R^{mm}) < \mathcal{L}(\bar{\mathbf{W}}(R))$  for sufficiently  
 761 large  $R$ .

762 To obtain the result, we establish a refined softmax probability control as in Step 1 by studying  
 763 distance to  $\mathcal{L}_\star$ . Following (21), denote the score function at  $\bar{\mathbf{W}}(R)$  via  $\gamma_i^R := \gamma_i^{\bar{\mathbf{W}}(R)}$ . Similarly,  
 764 let  $s_i^R = \mathbb{S}(\mathbf{a}_i^R)$  with  $\mathbf{a}_i^R = \mathbf{X}_i \bar{\mathbf{W}}(R) \mathbf{z}_i$ . Set the corresponding notation for the reference parameter  
 765  $\bar{\mathbf{W}}_R^{mm}$  as  $\gamma_i^\star, s_i^\star, \mathbf{a}_i^\star$ . Recall that  $R \geq \|\bar{\mathbf{W}}(R)\|_\diamond$  and  $\Xi_\circ := 1/\|\mathbf{W}^{mm}\|_\diamond$ . We note the following softmax  
 766 inequalities

$$\begin{aligned} s_{i\text{opt}_i}^\star &\geq \frac{1}{1 + T e^{-R\Xi_\circ}} \geq 1 - T e^{-R\Xi_\circ} \quad \text{for all } i \in [n], \\ s_{i\text{opt}_i}^R &\leq \frac{1}{1 + e^{-(1-\delta)\|\bar{\mathbf{W}}(R)\|_\diamond \Xi_\circ}} \leq \frac{1}{1 + e^{-(1-\delta)R\Xi_\circ}} \quad \text{for } i = 1. \end{aligned} \quad (23)$$

767 The former inequality is thanks to  $\mathbf{W}^{mm}$  achieving  $\geq 1$  margins on all tokens  $[T] - \text{opt}_i$  and the latter  
 768 arises from the  $\delta$ -margin violation of  $\bar{\mathbf{W}}(R)$  at  $i = 1$  i.e. Eq. (22). Since  $\ell$  is strictly decreasing with  
 769 Lipschitz derivative and the scores are upper/lower bounded by an absolute constant (as tokens are  
 770 bounded and fixed), we have that  $c_{\text{up}} \geq -\ell'(\gamma_i^W) \geq c_{\text{dn}}$  for some constants  $c_{\text{up}} > c_{\text{dn}} > 0$ . Thus,  
 771 following Eq. (20), the score decomposition (21), and (23) we can write

$$\begin{aligned} \mathcal{L}(\bar{\mathbf{W}}(R)) - \mathcal{L}_\star &\geq \frac{1}{n} [\ell(\gamma_1^{\bar{\mathbf{W}}(R)}) - \ell(\gamma_1^{\text{opt}})] \geq \frac{c_{\text{dn}}}{n} (\gamma_1^{\text{opt}} - \gamma_1^{\bar{\mathbf{W}}(R)}) \\ &\geq \frac{c_{\text{dn}}}{n} c_{\text{opt}} (1 - s_{1\text{opt}_1}^R). \\ &\geq \frac{c_{\text{dn}} c_{\text{opt}}}{n} \frac{1}{1 + e^{(1-\delta)R\Xi_\circ}}. \end{aligned} \quad (24)$$

772 Conversely, we upper bound the difference between  $\mathcal{L}(\bar{\mathbf{W}}_R^{mm})$  and  $\mathcal{L}_\star$  as follows. Define the worst-  
 773 case loss difference for  $\bar{\mathbf{W}}(R)$  as  $j = \arg \max_{i \in [n]} [\ell(\gamma_i^\star) - \ell(\gamma_i^{\text{opt}})]$ . Using (21)&(23), we write

$$\begin{aligned} \mathcal{L}(\bar{\mathbf{W}}_R^{mm}) - \mathcal{L}_\star &\leq \max_{i \in [n]} [\ell(\gamma_i^\star) - \ell(\gamma_i^{\text{opt}})] \leq c_{\text{up}} \cdot (\gamma_j^{\text{opt}} - \gamma_j^\star) \\ &\leq c_{\text{up}} \cdot (1 - s_{j\text{opt}_j}^\star) B \\ &\leq c_{\text{up}} \cdot T e^{-R\Xi_\circ} B. \end{aligned}$$

774 Combining the last inequality and (24), we conclude that  $\mathcal{L}(\bar{\mathbf{W}}_R^{mm}) < \mathcal{L}(\bar{\mathbf{W}}(R))$  whenever

$$c_{\text{up}}T \cdot e^{-R\Xi_0} B < \frac{c_{\text{dn}} \cdot c_{\text{opt}}}{n} \frac{1}{1 + e^{(1-\delta)R\Xi_0}} \iff \frac{e^{R\Xi_0}}{1 + e^{(1-\delta)R\Xi_0}} > \frac{c_{\text{up}}TnB}{c_{\text{dn}}c_{\text{opt}}}.$$

775 The left hand-side inequality holds for all sufficiently large  $R$ : Specifically, as soon as  $R$  obeys  
776  $R > \frac{1}{\delta\Xi_0} \log\left(\frac{2c_{\text{up}}TnB}{c_{\text{dn}}c_{\text{opt}}}\right)$ . This completes the proof of the theorem by contradiction since we obtained  
777  $\mathcal{L}(\bar{\mathbf{W}}(R)) > \mathcal{L}(\bar{\mathbf{W}}_R^{mm})$ .

## 778 E Convergence of Gradient Descent

779 **Optimization problem definition.** Recap the problem, where we use a linear head  $h(\mathbf{x}) = \mathbf{v}^\top \mathbf{x}$  for  
780 most of our theoretical exposition. Given dataset  $(Y_i, \mathbf{X}_i, \mathbf{z}_i)_{i=1}^n$ , we minimize the empirical risk of an  
781 1-layer transformer using combined weights  $\mathbf{W} \in \mathbb{R}^{d \times d}$  or individual weights  $\mathbf{K}, \mathbf{Q} \in \mathbb{R}^{d \times m}$  for a fixed  
782 head and decreasing loss function:

$$\mathcal{L}(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i \cdot \mathbf{v}^\top \mathbf{X}_i^\top \mathbb{S}(\mathbf{X}_i \mathbf{W} \mathbf{z}_i)), \quad (\text{W-ERM})$$

$$\mathcal{L}(\mathbf{K}, \mathbf{Q}) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i \cdot \mathbf{v}^\top \mathbf{X}_i^\top \mathbb{S}(\mathbf{X}_i \mathbf{K} \mathbf{Q}^\top \mathbf{z}_i)). \quad (\text{KQ-ERM})$$

783 We can recover the self-attention model by setting  $\mathbf{z}_i$  to be the first token of  $\mathbf{X}_i$ , i.e.,  $\mathbf{z}_i \leftarrow \mathbf{x}_{i1}$ .

### 784 E.1 Divergence of norm of the iterates $\mathbf{W}(k)$

785 The next lemma establishes the descent property of gradient descent for  $\mathcal{L}(\mathbf{W})$  under Assumption A.

786 **Lemma 7 (Descent Lemma)** *Under Assumption A, if  $\eta \leq 1/L_{\mathbf{W}}$ , then for any initialization  $\mathbf{W}(0)$ ,*  
787 *Algorithm W-GD satisfies:*

$$\mathcal{L}(\mathbf{W}(k+1)) - \mathcal{L}(\mathbf{W}(k)) \leq -\frac{\eta}{2} \|\nabla \mathcal{L}(\mathbf{W}(k))\|_F^2, \quad (25)$$

788 *for all  $k \geq 0$ . Additionally, it holds that  $\sum_{k=0}^{\infty} \|\nabla \mathcal{L}(\mathbf{W}(k))\|_F^2 < \infty$ , and  $\lim_{k \rightarrow \infty} \|\nabla \mathcal{L}(\mathbf{W}(k))\|_F^2 = 0$ .*

789 **Proof.** The proof is similar to [TLZO23, Lemma 5]. ■

790 The lemma below reveals that the correlation between the training loss's gradient at any arbitrary matrix  $\mathbf{W}$  and the attention SVM solution  $\mathbf{W}^{mm}$  is negative. Consequently, for any finite  $\mathbf{W}$ ,  
791  $\langle \nabla \mathcal{L}(\mathbf{W}), \mathbf{W}^{mm} \rangle$  cannot be equal to zero.  
792

793 **Lemma 8** *Let  $\mathbf{W}^{mm}$  be the SVM solution of (Att-SVM). Suppose Assumptions A and B hold. Then,*  
794 *for all  $\mathbf{W} \in \mathbb{R}^{d \times d}$ , the training loss (W-ERM) obeys  $\langle \nabla \mathcal{L}(\mathbf{W}), \mathbf{W}^{mm} \rangle \leq -c < 0$ , for some constant*  
795  *$c > 0$  (see (34)) depending on the data, the head  $\mathbf{v}$ , and a loss derivative bound.*

796 **Proof.** Let

$$\bar{\mathbf{h}}_i = \mathbf{X}_i \mathbf{W}^{mm} \mathbf{z}_i, \quad \boldsymbol{\gamma}_i = Y_i \cdot \mathbf{X}_i \mathbf{v}, \quad \text{and} \quad \mathbf{h}_i = \mathbf{X}_i \mathbf{W} \mathbf{z}_i. \quad (26)$$

797 Let us recall the gradient evaluated at  $\mathbf{W}$  which is given by

$$\nabla \mathcal{L}(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \ell'(\boldsymbol{\gamma}_i^\top \mathbb{S}(\mathbf{h}_i)) \cdot \mathbf{X}_i^\top \mathbb{S}'(\mathbf{h}_i) \boldsymbol{\gamma}_i \mathbf{z}_i^\top, \quad (27)$$



798 which implies that

$$\begin{aligned}
\langle \nabla \mathcal{L}(\mathbf{W}), \mathbf{W}^{mm} \rangle &= \frac{1}{n} \sum_{i=1}^n \ell'_i \left( \boldsymbol{\gamma}_i^\top \mathbb{S}(\mathbf{h}_i) \right) \cdot \left\langle \mathbf{X}_i^\top \mathbb{S}'(\mathbf{h}_i) \boldsymbol{\gamma}_i \mathbf{z}_i^\top, \mathbf{W}^{mm} \right\rangle \\
&= \frac{1}{n} \sum_{i=1}^n \ell'_i \cdot \text{trace} \left( (\mathbf{W}^{mm})^\top \mathbf{X}_i^\top \mathbb{S}'(\mathbf{h}_i) \boldsymbol{\gamma}_i \mathbf{z}_i^\top \right) \\
&= \frac{1}{n} \sum_{i=1}^n \ell'_i \cdot \bar{\mathbf{h}}_i^\top \mathbb{S}'(\mathbf{h}_i) \boldsymbol{\gamma}_i \\
&= \frac{1}{n} \sum_{i=1}^n \ell'_i \cdot \left( \bar{\mathbf{h}}_i^\top \text{diag}(\mathbf{s}_i) \boldsymbol{\gamma}_i - \bar{\mathbf{h}}_i^\top \mathbf{s}_i \mathbf{s}_i^\top \boldsymbol{\gamma}_i \right).
\end{aligned} \tag{28}$$

799 Here, let  $\ell'_i := \ell'(\boldsymbol{\gamma}_i^\top \mathbb{S}(\mathbf{h}_i))$ ,  $\mathbf{s}_i = \mathbb{S}(\mathbf{h}_i)$  and the third equality uses  $\text{trace}(\mathbf{b}\mathbf{a}^\top) = \mathbf{a}^\top \mathbf{b}$ .

800 In order to move forward, we will establish the following result, with a focus on the equal score  
801 condition (Assumption (i)): Let  $\gamma = \gamma_{i \geq 2}$  be a constant, and let  $\boldsymbol{\gamma}_1$  and  $\bar{\mathbf{h}}_1$  represent the largest indices  
802 of vectors  $\boldsymbol{\gamma}$  and  $\bar{\mathbf{h}}$  respectively. For any vector  $\mathbf{s}$  that satisfies  $\sum_{t \in [T]} s_t = 1$  and  $s_t > 0$ , we aim to  
803 prove that  $\bar{\mathbf{h}}^\top \text{diag}(\mathbf{s}) \boldsymbol{\gamma} - \bar{\mathbf{h}}^\top \mathbf{s} \mathbf{s}^\top \boldsymbol{\gamma} > 0$ . To demonstrate this, we proceed by writing the following:

$$\begin{aligned}
\bar{\mathbf{h}}^\top \text{diag}(\mathbf{s}) \boldsymbol{\gamma} - \bar{\mathbf{h}}^\top \mathbf{s} \mathbf{s}^\top \boldsymbol{\gamma} &= \sum_{t=1}^T \bar{\mathbf{h}}_t \boldsymbol{\gamma}_t s_t - \sum_{t=1}^T \bar{\mathbf{h}}_t s_t \sum_{i=1}^T \boldsymbol{\gamma}_i s_i \\
&= \left( \bar{\mathbf{h}}_1 \boldsymbol{\gamma}_1 s_1 + \gamma \sum_{t \geq 2}^T \bar{\mathbf{h}}_t s_t \right) - (\boldsymbol{\gamma}_1 s_1 + \gamma(1 - s_1)) \left( \bar{\mathbf{h}}_1 s_1 + \sum_{t \geq 2}^T \bar{\mathbf{h}}_t s_t \right) \\
&= \bar{\mathbf{h}}_1 (\boldsymbol{\gamma}_1 - \gamma) s_1 (1 - s_1) - (\boldsymbol{\gamma}_1 - \gamma) s_1 \sum_{t \geq 2}^T \bar{\mathbf{h}}_t s_t \\
&= (\boldsymbol{\gamma}_1 - \gamma) (1 - s_1) s_1 \left[ \bar{\mathbf{h}}_1 - \frac{\sum_{t \geq 2}^T \bar{\mathbf{h}}_t s_t}{\sum_{t \geq 2}^T s_t} \right] \\
&\geq (\boldsymbol{\gamma}_1 - \gamma) (1 - s_1) s_1 (\bar{\mathbf{h}}_1 - \max_{t \geq 2} \bar{\mathbf{h}}_t).
\end{aligned} \tag{29}$$

804 To proceed, define

$$\gamma_{\text{gap}}^i = \boldsymbol{\gamma}_{i \text{opt}_i} - \max_{t \neq \text{opt}_i} \boldsymbol{\gamma}_{it} \quad \text{and} \quad \bar{h}_{\text{gap}}^i = \bar{\mathbf{h}}_{i \text{opt}_i} - \max_{t \neq \text{opt}_i} \bar{\mathbf{h}}_{it}.$$

805 With these, we obtain

$$\bar{\mathbf{h}}_i^\top \text{diag}(\mathbf{s}_i) \boldsymbol{\gamma}_i - \bar{\mathbf{h}}_i^\top \mathbf{s}_i \mathbf{s}_i^\top \boldsymbol{\gamma}_i \geq \gamma_{\text{gap}}^i \bar{h}_{\text{gap}}^i (1 - s_{i \text{opt}_i}) s_{i \text{opt}_i}. \tag{30}$$

806 Note that

$$\begin{aligned}
\bar{h}_{\text{gap}}^i &= \min_{t \neq \text{opt}_i} (\mathbf{x}_{i \text{opt}_i} - \mathbf{x}_{it})^\top \mathbf{W}^{mm} \mathbf{z}_i \geq 1, \\
\gamma_{\text{gap}}^i &= \min_{t \neq \text{opt}_i} \boldsymbol{\gamma}_{i \text{opt}_i} - \boldsymbol{\gamma}_{it} > 0, \\
s_{i \text{opt}_i} (1 - s_{i \text{opt}_i}) &> 0.
\end{aligned}$$

807 Hence,

$$c_0 := \min_{i \in [n]} \left\{ \left( \min_{t \neq \text{opt}_i} (\mathbf{x}_{i \text{opt}_i} - \mathbf{x}_{it})^\top \mathbf{W}^{mm} \mathbf{z}_i \right) \cdot \left( \min_{t \neq \text{opt}_i} \boldsymbol{\gamma}_{i \text{opt}_i} - \boldsymbol{\gamma}_{it} \right) \cdot s_{i \text{opt}_i} (1 - s_{i \text{opt}_i}) \right\} > 0. \tag{31}$$

808 It follows from (30) and (31) that

$$\min_{i \in [n]} \left\{ \bar{\mathbf{h}}_i^\top \text{diag}(\mathbf{s}_i) \boldsymbol{\gamma}_i - \bar{\mathbf{h}}_i^\top \mathbf{s}_i \mathbf{s}_i^\top \boldsymbol{\gamma}_i \right\} \geq c_0 > 0. \tag{32}$$

809 Further, by our assumption  $\ell'_i < 0$ . Since by Assumption A,  $\ell'$  is continuous and the domain is  
810 bounded, the maximum is attained and negative, and thus

$$-c_1 = \max_x \ell'(x), \quad \text{for some } c_1 > 0. \tag{33}$$

811 Hence, using (32) and (33) in (28), we obtain

$$\langle \nabla \mathcal{L}(\mathbf{W}), \mathbf{W}^{mm} \rangle \leq -c < 0, \quad \text{where } c = c_1 \cdot c_0. \quad (34)$$

812 In the scenario that Assumption B(ii) holds (all tokens are support),  $\bar{\mathbf{h}}_t = \mathbf{x}_{it}^\top \mathbf{W}^{mm} \mathbf{z}_i$  is constant for all  
813  $t \geq 2$ . Hence, following similar steps as in (29) completes the proof. ■

814 **Theorem 4** Suppose Assumption A on the loss function  $\ell$  and Assumption B on the tokens hold.  
815 Then,

- 816 • There is no  $\mathbf{W} \in \mathbb{R}^{d \times d}$  satisfying  $\nabla \mathcal{L}(\mathbf{W}) = 0$ .
- 817 • Algorithm W-GD with the step size  $\eta \leq 1/L_W$  and any starting point  $\mathbf{W}(0)$  satisfies  
818  $\lim_{k \rightarrow \infty} \|\mathbf{W}(k)\|_F = \infty$ .

819 **Proof.** It follows from Lemma 7 that under Assumption A,  $\eta \leq 1/L_W$ , and for any initialization  $\mathbf{W}(0)$ ,  
820 the gradient descent sequence  $\mathbf{W}(k+1) = \mathbf{W}(k) - \eta \nabla \mathcal{L}(\mathbf{W}(k))$  satisfies  $\lim_{k \rightarrow \infty} \|\nabla \mathcal{L}(\mathbf{W}(k))\|_F^2 = 0$ .

821 Further, it follows from Lemma 8 that  $\langle \nabla \mathcal{L}(\mathbf{W}), \mathbf{W}^{mm} \rangle < 0$  for all  $\mathbf{W} \in \mathbb{R}^{d \times d}$ . Hence, for any finite  $\mathbf{W}$ ,  
822  $\langle \nabla \mathcal{L}(\mathbf{W}), \mathbf{W}^{mm} \rangle$  cannot be equal to zero. Therefore, there are no finite critical points  $\mathbf{W}$ , for which  
823  $\nabla \mathcal{L}(\mathbf{W}) = 0$  which contradicts Lemma 7. This implies that  $\|\mathbf{W}(k)\| \rightarrow \infty$ . ■

## 824 E.2 Global Convergence of Gradient Descent

825 The following lemma illustrates that when non-optimal tokens within an input share the same scores,  
826 the negative gradient of the loss function at  $\mathbf{W}$  becomes more correlated with the max-margin solution  
827 ( $\mathbf{W}^{mm}$ ) than with  $\mathbf{W}$  itself.

828 **Lemma 9** Let  $\mathbf{W}^{mm}$  be the SVM solution of (Att-SVM). Suppose Assumption (i) on the tokens'  
829 score hold and  $\ell(\cdot)$  is strictly decreasing and differentiable. For any choice of  $\pi > 0$ , there exists  
830  $R := R_\pi$  such that, for any  $\mathbf{W}$  with  $\|\mathbf{W}\|_F \geq R$ , we have

$$\left\langle \nabla \mathcal{L}(\mathbf{W}), \frac{\mathbf{W}}{\|\mathbf{W}\|_F} \right\rangle \geq (1 + \pi) \left\langle \nabla \mathcal{L}(\mathbf{W}), \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F} \right\rangle.$$

831 **Proof.** Let  $\bar{\mathbf{W}} = \|\mathbf{W}^{mm}\|_F \mathbf{W} / \|\mathbf{W}\|_F$ ,  $M = \sup_{i,t} \|\mathbf{x}_{it} \mathbf{z}_i^\top\|$ ,  $\Theta = 1 / \|\mathbf{W}^{mm}\|_F$ ,  $\mathbf{s}_i = \mathbb{S}(\mathbf{X}_i \mathbf{W} \mathbf{z}_i)$ ,  $\mathbf{h}_i = \mathbf{X}_i \bar{\mathbf{W}} \mathbf{z}_i$ ,  
832  $\bar{\mathbf{h}}_i = \mathbf{X}_i \mathbf{W}^{mm} \mathbf{z}_i$ , and  $\gamma_i = \gamma_{i,t \geq 2}$ . Without losing generality assume  $\alpha_i = \text{opt}_i = 1$  for all  $i \in [n]$ .  
833 Repeating the proof of Lemma 8 yields

$$\begin{aligned} \langle \nabla \mathcal{L}(\mathbf{W}), \mathbf{W}^{mm} \rangle &= \frac{1}{n} \sum_{i=1}^n \ell'_i \cdot (\gamma_{i1} - \gamma_i)(1 - s_{i1}) s_{i1} \left[ \bar{\mathbf{h}}_{i1} - \frac{\sum_{t \geq 2}^T \bar{\mathbf{h}}_{it} s_{it}}{\sum_{t \geq 2}^T s_{it}} \right], \\ \langle \nabla \mathcal{L}(\mathbf{W}), \bar{\mathbf{W}} \rangle &= \frac{1}{n} \sum_{i=1}^n \ell'_i \cdot (\gamma_{i1} - \gamma_i)(1 - s_{i1}) s_{i1} \left[ \mathbf{h}_{i1} - \frac{\sum_{t \geq 2}^T \mathbf{h}_{it} s_{it}}{\sum_{t \geq 2}^T s_{it}} \right]. \end{aligned}$$

834 Focusing on a single example  $i \in [n]$  with  $\mathbf{s}, \mathbf{h}, \bar{\mathbf{h}}$  vectors (dropping subscript  $i$ ), given  $\pi$ , for  
835 sufficiently large  $R$ , we wish to show that

$$\left[ \mathbf{h}_1 - \frac{\sum_{t \geq 2}^T \mathbf{h}_t s_t}{\sum_{t \geq 2}^T s_t} \right] \leq (1 + \pi) \cdot \left[ \bar{\mathbf{h}}_1 - \frac{\sum_{t \geq 2}^T \bar{\mathbf{h}}_t s_t}{\sum_{t \geq 2}^T s_t} \right]. \quad (35)$$

836 We consider two scenarios.

837 **Scenario 1:**  $\|\bar{\mathbf{W}} - \mathbf{W}^{mm}\|_F \leq \epsilon := \pi / (2M)$ . In this scenario, for any token, we find that

$$|\mathbf{h}_t - \bar{\mathbf{h}}_t| = |\mathbf{x}_t^\top (\bar{\mathbf{W}} - \mathbf{W}^{mm}) \mathbf{z}_t| \leq M \|\bar{\mathbf{W}} - \mathbf{W}^{mm}\|_F \leq M \epsilon.$$

838 Consequently, we obtain

$$\bar{\mathbf{h}}_1 - \frac{\sum_{t \geq 2}^T \bar{\mathbf{h}}_t s_t}{\sum_{t \geq 2}^T s_t} \geq \mathbf{h}_1 - \frac{\sum_{t \geq 2}^T \mathbf{h}_t s_t}{\sum_{t \geq 2}^T s_t} - 2M \epsilon = \mathbf{h}_1 - \frac{\sum_{t \geq 2}^T \mathbf{h}_t s_t}{\sum_{t \geq 2}^T s_t} - \pi.$$

839 Also noticing  $\bar{\mathbf{h}}_1 - \frac{\sum_{t \geq 2}^T \bar{\mathbf{h}}_t s_t}{\sum_{t \geq 2}^T s_t} \geq 1$  (thanks to  $\mathbf{W}^{mm}$  satisfying  $\geq 1$  margin), this implies (35).

840 **Scenario 2:**  $\|\bar{\mathbf{W}} - \mathbf{W}^{mm}\|_F \geq \epsilon := \pi/(2M)$ . In this scenario, for some  $\delta = \delta(\epsilon)$  and  $\tau \geq 2$ , we have  
841 that

$$\mathbf{h}_1 - \mathbf{h}_\tau \leq 1 - 2\delta.$$

842 Recall that  $s = \mathbb{S}(\bar{\mathbf{R}}\mathbf{h})$  where  $\bar{\mathbf{R}} = \|\mathbf{W}\|_F / \|\mathbf{W}^{mm}\|_F$ . To proceed, split the tokens into two groups: Let  
843  $\mathcal{N}$  be the group of tokens obeying  $(\mathbf{x}_1 - \mathbf{x}_t)^\top \bar{\mathbf{W}}\mathbf{z} \geq 1 - \delta$  for  $t \in \mathcal{N}$  and  $[T] - \mathcal{N}$  be the rest. Observe  
844 that

$$\frac{\sum_{t \in \mathcal{N}} s_t}{\sum_{t \geq 2}^T s_t} \leq \frac{\sum_{t \in \mathcal{N}} s_t}{s_\tau} \leq T \frac{e^{\delta \bar{\mathbf{R}}}}{e^{2\delta \bar{\mathbf{R}}}} = T e^{-\bar{\mathbf{R}}\delta}.$$

845 Set  $\bar{M} = M/\Theta$  and note that  $\|\mathbf{h}_t\| \leq \|\mathbf{W}^{mm}\|_F \cdot \|\mathbf{x}_t \mathbf{z}^\top\| \leq \bar{M}$ . Using  $(\mathbf{x}_1 - \mathbf{x}_t)^\top \bar{\mathbf{W}}\mathbf{z} < 1 - \delta$  over  
846  $t \in [T] - \mathcal{N}$  and plugging in the above bound, we obtain

$$\begin{aligned} \frac{\sum_{t \geq 2}^T (\mathbf{h}_1 - \mathbf{h}_t) s_t}{\sum_{t \geq 2}^T s_t} &= \frac{\sum_{t \in [T] - \mathcal{N}} (\mathbf{h}_1 - \mathbf{h}_t) s_t}{\sum_{t \geq 2}^T s_t} + \frac{\sum_{t \in \mathcal{N}} (\mathbf{h}_1 - \mathbf{h}_t) s_t}{\sum_{t \geq 2}^T s_t} \\ &\leq (1 - \delta) + 2\bar{M}T e^{-\bar{\mathbf{R}}\delta}. \end{aligned}$$

847 Using the fact that  $\bar{\mathbf{h}}_1 - \frac{\sum_{t \geq 2}^T \bar{\mathbf{h}}_t s_t}{\sum_{t \geq 2}^T s_t} \geq 1$ , the above implies (35) with  $\pi' = 2\bar{M}T e^{-\bar{\mathbf{R}}\delta} - \delta$ . To proceed,  
848 choose

$$R_\pi = \delta^{-1} \Theta^{-1} \log\left(\frac{2\bar{M}T}{\pi}\right) \quad \text{to ensure} \quad \pi' \leq \pi. \quad (36)$$

849 ■

## 850 E.2.1 Proof of Theorem 2.

851 The proof is similar to [TLZO23, Theorem 2]. Given any  $\epsilon \in (0, 1)$ , let  $\pi = \epsilon/(1 - \epsilon)$ . It follows from  
852 Theorem 4 that  $\lim_{k \rightarrow \infty} \|\mathbf{W}(k)\|_F = \infty$ . Hence, we can choose  $k_\epsilon$  such that for any  $k \geq k_\epsilon$ , it holds that  
853  $\|\mathbf{W}(k)\|_F > R_\epsilon \vee 1/2$  for some parameter  $R_\epsilon$ . Now for any  $k \geq k_\epsilon$ , it follows from Lemma 9 that

$$\left\langle -\nabla \mathcal{L}(\mathbf{W}(k)), \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F} \right\rangle \geq (1 - \epsilon) \left\langle -\nabla \mathcal{L}(\mathbf{W}(k)), \frac{\mathbf{W}(k)}{\|\mathbf{W}(k)\|_F} \right\rangle.$$

854 Multiplying both sides by the stepsize  $\eta$  and using the gradient descent update, we get

$$\begin{aligned} \left\langle \mathbf{W}(k+1) - \mathbf{W}(k), \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F} \right\rangle &\geq (1 - \epsilon) \left\langle \mathbf{W}(k+1) - \mathbf{W}(k), \frac{\mathbf{W}(k)}{\|\mathbf{W}(k)\|_F} \right\rangle \\ &= \frac{(1 - \epsilon)}{2\|\mathbf{W}(k)\|_F} \left( \|\mathbf{W}(k+1)\|_F^2 - \|\mathbf{W}(k)\|_F^2 - \|\mathbf{W}(k+1) - \mathbf{W}(k)\|_F^2 \right) \\ &\geq (1 - \epsilon) \left( \frac{1}{2\|\mathbf{W}(k)\|_F} \left( \|\mathbf{W}(k+1)\|_F^2 - \|\mathbf{W}(k)\|_F^2 \right) - \|\mathbf{W}(k+1) - \mathbf{W}(k)\|_F^2 \right) \\ &\geq (1 - \epsilon) \left( \|\mathbf{W}(k+1)\|_F - \|\mathbf{W}(k)\|_F - \|\mathbf{W}(k+1) - \mathbf{W}(k)\|_F^2 \right) \\ &\geq (1 - \epsilon) \left( \|\mathbf{W}(k+1)\|_F - \|\mathbf{W}(k)\|_F - 2\eta (\mathcal{L}(\mathbf{W}(k)) - \mathcal{L}(\mathbf{W}(k+1))) \right). \end{aligned} \quad (37)$$

855 Here, the second inequality is obtained from  $\|\mathbf{W}(k)\|_F \geq 1/2$ ; the third inequality follows since for  
856 any  $a, b > 0$ , we have  $(a^2 - b^2)/(2b) - (a - b) \geq 0$ ; and the last inequality uses Lemma 7.

857 Summing the above inequality over  $k \geq k_\epsilon$  gives

$$\left\langle \frac{\mathbf{W}(k)}{\|\mathbf{W}(k)\|_F}, \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F} \right\rangle \geq 1 - \epsilon + \frac{C(\epsilon, \eta)}{\|\mathbf{W}(k)\|_F},$$

858 for some finite constant  $C(\epsilon, \eta)$  defined as

$$C(\epsilon, \eta) := \left\langle \mathbf{W}(k_\epsilon), \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F} \right\rangle - (1 - \epsilon) \|\mathbf{W}(k_\epsilon)\|_F - 2\eta(1 - \epsilon)(\mathcal{L}(\mathbf{W}(k_\epsilon)) - \mathcal{L}_\star), \quad (38)$$

859 where  $\mathcal{L}_\star \leq \mathcal{L}(\mathbf{W}(k))$  for all  $k \geq 0$ .

860 Since  $\|\mathbf{W}(k)\| \rightarrow \infty$ , we get

$$\liminf_{k \rightarrow \infty} \left\langle \frac{\mathbf{W}(k)}{\|\mathbf{W}(k)\|_F}, \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F} \right\rangle \geq 1 - \epsilon. \quad (39)$$

861 Given that  $\epsilon$  is arbitrary, we can consider the limit as  $\epsilon$  approaches zero. Thus,  $\mathbf{W}(k)/\|\mathbf{W}(k)\|_F \rightarrow$   
862  $\mathbf{W}^{mm}/\|\mathbf{W}^{mm}\|_F$ . ■

### 863 E.3 Local Convergence of Gradient Descent

864 To provide a basis for discussing local convergence of GD, we establish a cone centered around  $\mathbf{W}_\alpha^{mm}$   
865 using the following construction. For parameters  $\mu \in (0, 1)$  and  $R > 0$ , we define  $C_{\mu,R}(\mathbf{W}_\alpha^{mm})$  as the  
866 set of matrices  $\mathbf{W} \in \mathbb{R}^{d \times d}$  such that  $\|\mathbf{W}\|_F \geq R$  and the correlation coefficient between  $\mathbf{W}$  and  $\mathbf{W}_\alpha^{mm}$  is  
867 at least  $1 - \mu$ :

$$\mathcal{S}_\mu(\mathbf{W}_\alpha^{mm}) := \left\{ \mathbf{W} \in \mathbb{R}^{d \times d} : \left\langle \frac{\mathbf{W}}{\|\mathbf{W}\|_F}, \frac{\mathbf{W}_\alpha^{mm}}{\|\mathbf{W}_\alpha^{mm}\|_F} \right\rangle \geq 1 - \mu \right\}, \quad (40a)$$

$$C_{\mu,R}(\mathbf{W}_\alpha^{mm}) := \mathcal{S}_\mu(\mathbf{W}_\alpha^{mm}) \cap \left\{ \mathbf{W} \in \mathbb{R}^{d \times d} : \|\mathbf{W}\|_F \geq R \right\}. \quad (40b)$$

868 **Lemma 10** Suppose Assumption A on the loss function  $\ell$  holds, and let  $\alpha = (\alpha_i)_{i=1}^n$  be locally optimal  
869 tokens according to Definition 2. Let  $\mathbf{W}^{mm} = \mathbf{W}_\alpha^{mm}$  denote the SVM solution obtained via (Att-SVM)  
870 by applying the Frobenius norm and replacing  $(\text{opt}_i)_{i=1}^n$  with  $\alpha = (\alpha_i)_{i=1}^n$ . To provide a basis for  
871 discussing the local convergence of gradient descent, we establish a cone centered around  $\mathbf{W}^{mm}$  using  
872 the following construction. There exists a scalar  $\mu = \mu(\alpha) > 0$  such that for sufficiently large  $\bar{R}_\mu$ :

873 **L1.** There is no stationary point within  $C_{\mu,\bar{R}_\mu}(\mathbf{W}^{mm})$ .

874 **L2.** For all  $\mathbf{V} \in \mathcal{S}_\mu(\mathbf{W}^{mm})$  with  $\|\mathbf{V}\|_F = \|\mathbf{W}^{mm}\|_F$  and  $\mathbf{W} \in C_{\mu,\bar{R}_\mu}(\mathbf{W}^{mm})$ , there exist dataset dependent  
875 constants  $C, c > 0$  such that

$$C \cdot \frac{1}{n} \sum_{i=1}^n (1 - s_{i\alpha_i}) \geq -\langle \nabla \mathcal{L}(\mathbf{W}), \mathbf{V} \rangle \geq c \cdot \frac{1}{n} \sum_{i=1}^n (1 - s_{i\alpha_i}) > 0, \quad (41a)$$

$$\|\nabla \mathcal{L}(\mathbf{W})\|_F \leq \bar{A} C \cdot \frac{1}{n} \sum_{i=1}^n (1 - s_{i\alpha_i}), \quad (41b)$$

$$-\left\langle \frac{\mathbf{V}}{\|\mathbf{V}\|_F}, \frac{\nabla \mathcal{L}(\mathbf{W})}{\|\nabla \mathcal{L}(\mathbf{W})\|_F} \right\rangle \geq \frac{c}{C} \cdot \frac{\Theta}{\bar{A}} > 0. \quad (41c)$$

876 Here,  $s_{i\alpha_i} = (\mathbb{S}(\mathbf{X}_i \mathbf{W} \mathbf{z}_i))_{\alpha_i}$ ,  $\bar{A} = \max_{i \in [n], t, \tau \in [T]} \|(\mathbf{x}_{it} - \mathbf{x}_{i\tau})\| \|\mathbf{z}_i\|$ , and  $\Theta = 1/\|\mathbf{W}^{mm}\|_F$ .

877 **Proof.** Let  $R = \bar{R}_\mu$ ,  $(\mathcal{T}_i)_{i=1}^n$  be the set of all support indices per Definition 2. Let  $\bar{\mathcal{T}}_i = [T] - \mathcal{T}_i - \{\alpha_i\}$   
878 be the non-support indices. Let

$$\begin{aligned} \Theta &= 1/\|\mathbf{W}^{mm}\|_F, \\ \delta &= \frac{1}{2} \min_{i \in [n]} \min_{t \in \mathcal{T}_i, \tau \in \bar{\mathcal{T}}_i} (\mathbf{x}_{it} - \mathbf{x}_{i\tau})^\top \mathbf{W}^{mm} \mathbf{z}_i, \\ A &= \max_{i \in [n], t \in [T]} \frac{\|\mathbf{x}_{it} \mathbf{z}_i^\top\|_F}{\Theta}, \\ \mu &\leq \mu(\delta) = \frac{1}{8} \left( \frac{\min(0.5, \delta)}{A} \right)^2. \end{aligned} \quad (42)$$

879 Since  $\mathbf{W}^{mm}$  is the max-margin model ensuring  $(\mathbf{x}_{i\alpha_i} - \mathbf{x}_{it})^\top \mathbf{W}^{mm} \mathbf{z}_i \geq 1$ , the following inequalities hold  
880 for all  $\mathbf{W} \in \mathcal{S}_\mu(\mathbf{W}^{mm})$ ,  $\|\mathbf{W}\|_F = \|\mathbf{W}^{mm}\|_F$  and all  $i \in [n]$ ,  $t \in \mathcal{T}_i$ ,  $\tau \in \bar{\mathcal{T}}_i$ :

$$\begin{aligned} (\mathbf{x}_{it} - \mathbf{x}_{i\tau})^\top \mathbf{W} \mathbf{z}_i &\geq \delta > 0, \\ (\mathbf{x}_{i\alpha_i} - \mathbf{x}_{i\tau})^\top \mathbf{W} \mathbf{z}_i &\geq 1 + \delta, \\ \frac{3}{2} &\geq (\mathbf{x}_{i\alpha_i} - \mathbf{x}_{it})^\top \mathbf{W} \mathbf{z}_i \geq \frac{1}{2}. \end{aligned} \quad (43)$$

881 Here, we used  $\|\mathbf{W} - \mathbf{W}^{mm}\|_F^2 / \|\mathbf{W}^{mm}\|_F^2 \leq 2\mu$  which implies  $\|\mathbf{W} - \mathbf{W}^{mm}\|_F \leq \sqrt{2\mu}/\Theta$ .

882 To proceed, we write the gradient correlation following (11) and (29)

$$\langle \nabla \mathcal{L}(\mathbf{W}), \mathbf{V} \rangle = \frac{1}{n} \sum_{i=1}^n \ell'_i \cdot \mathbf{h}_i^\top \mathbb{S}'(\tilde{\mathbf{h}}_i) \boldsymbol{\gamma}_i, \quad (44)$$

883 where we denoted  $\ell'_i = \ell'(Y_i \cdot \mathbf{v}^\top \mathbf{X}_i^\top \mathbb{S}(\tilde{\mathbf{h}}_i))$ ,  $\mathbf{h}_i = \mathbf{X}_i \mathbf{V} \mathbf{z}_i$ ,  $\tilde{\mathbf{h}}_i = \mathbf{X}_i \mathbf{W} \mathbf{z}_i$ ,  $s_i = \mathbb{S}(\tilde{\mathbf{h}}_i)$ .

884 Using (43), for all  $t \in \mathcal{T}_i$ ,  $\tau \in \bar{\mathcal{T}}_i$ , for all  $\mathbf{W} \in C_{\mu,R}(\mathbf{W}^{mm})$ , we have that

$$\begin{aligned} \tilde{\mathbf{h}}_{it} - \tilde{\mathbf{h}}_{i\tau} &\geq R\Theta\delta, \\ \tilde{\mathbf{h}}_{i\alpha_i} - \tilde{\mathbf{h}}_{i\tau} &\geq R\Theta(1 + \delta), \\ \tilde{\mathbf{h}}_{i\alpha_i} - \tilde{\mathbf{h}}_{it} &\geq R\Theta/2. \end{aligned}$$

885 Consequently, we can bound the softmax probabilities  $s_i = \mathbb{S}(\tilde{\mathbf{h}}_i)$  over non-support indices as follows:

886 For all  $i \in [n]$  and any  $t_i \in \mathcal{T}_i$

$$S_i := \sum_{\tau \in \bar{\mathcal{T}}_i} s_{i\tau} \leq T e^{-R\Theta/2} s_{i\alpha_i} \leq T e^{-R\Theta/2}, \quad (45a)$$

$$Q_i := \sum_{\tau \in \bar{\mathcal{T}}_i} s_{i\tau} \leq T e^{-R\Theta\delta} s_{it_i} \leq T e^{-R\Theta\delta} S_i. \quad (45b)$$

887 Recall scores  $\boldsymbol{\gamma}_{it} = Y_i \cdot \mathbf{v}^\top \mathbf{x}_{it}$ . Define the score gaps over support indices:

$$\boldsymbol{\gamma}_i^{gap} = \boldsymbol{\gamma}_{i\alpha_i} - \max_{t \in \mathcal{T}_i} \boldsymbol{\gamma}_{it} \quad \text{and} \quad \tilde{\boldsymbol{\gamma}}_i^{gap} = \boldsymbol{\gamma}_{i\alpha_i} - \min_{t \in \mathcal{T}_i} \boldsymbol{\gamma}_{it}.$$

888 It follows from (42) that

$$A = \max_{i \in [n], t \in [T]} \frac{\|\mathbf{x}_{it} \mathbf{z}_i^\top\|_F}{\Theta} \geq \max_{i \in [n], t \in [T]} \|\mathbf{h}_{it}\|.$$

889 Define the  $\alpha$ -dependent global scalar  $\Gamma = \sup_{i \in [n], t, \tau \in [T]} |\boldsymbol{\gamma}_{it} - \boldsymbol{\gamma}_{i\tau}|$ .

890 Let us focus on a fixed datapoint  $i \in [n]$ , assume (without losing generality)  $\alpha_i = 1$ , and drop  
891 subscripts  $i$ . Directly applying Lemma 6, we obtain

$$\left| \mathbf{h}^\top \text{diag}(s) \boldsymbol{\gamma} - \mathbf{h}^\top s s^\top \boldsymbol{\gamma} - \sum_{t \geq 2} (\mathbf{h}_1 - \mathbf{h}_t) s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) \right| \leq 2\Gamma A (1 - s_1)^2.$$

892 To proceed, let us decouple the non-support indices within  $\sum_{t \geq 2} (\mathbf{h}_1 - \mathbf{h}_t) s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t)$  via

$$\left| \sum_{t \in \bar{\mathcal{T}}} (\mathbf{h}_1 - \mathbf{h}_t) s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) \right| \leq 2Q\Gamma A.$$

893 Aggregating these, we found

$$\left| \mathbf{h}^\top \text{diag}(s) \boldsymbol{\gamma} - \mathbf{h}^\top s s^\top \boldsymbol{\gamma} - \sum_{t \in \mathcal{T}} (\mathbf{h}_1 - \mathbf{h}_t) s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) \right| \leq 2\Gamma A ((1 - s_1)^2 + Q). \quad (46)$$

894 To proceed, let us upper/lower bound the gradient correlation. We use two bounds depending on  
895  $\mathbf{V} \in \mathcal{S}_\mu(\mathbf{W}^{mm})$  (**Case 1**) or general  $\mathbf{V} \in \mathbb{R}^{d \times d}$  (**Case 2**).

896 • **Case 1:**  $\mathbf{V} \in \mathcal{S}_\mu(\mathbf{W}^{mm})$ . Since  $1.5 \geq \mathbf{h}_1 - \mathbf{h}_t \geq 0.5$  following (43), we find

$$1.5 \cdot S \cdot \tilde{\boldsymbol{\gamma}}^{gap} \geq \sum_{t \in \mathcal{T}} (\mathbf{h}_1 - \mathbf{h}_t) s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) \geq 0.5 \cdot S \cdot \boldsymbol{\gamma}^{gap},$$

897 where recall the definition of  $S$  (having dropped subscripts) in (45a).

• **Case 2:**  $\mathbf{V} \in \mathbb{R}^{d \times d}$  and  $\|\mathbf{V}\|_F = \|\mathbf{W}^{mm}\|_F$ . Define  $\bar{A} = \max_{i \in [n], t, \tau \in [T]} \|\mathbf{x}_{it} - \mathbf{x}_{i\tau}\| \|\mathbf{z}_i\|$ . For any  $\|\mathbf{V}\|_F = \|\mathbf{W}^{mm}\|$ , we use the fact that

$$\|\mathbf{h}_1 - \mathbf{h}_t\| \leq \|(\mathbf{x}_{it} - \mathbf{x}_{i\tau}) \mathbf{z}_i^\top\|_F \cdot \|\mathbf{V}\|_F \leq \frac{\bar{A}}{\Theta}.$$

898 Note that by definition  $\frac{\bar{A}}{\Theta} \geq 1$ . To proceed, we can upper bound

$$\frac{\bar{A}}{\Theta} \cdot S \cdot \bar{\gamma}^{gap} \geq \sum_{i \in \mathcal{T}} (\mathbf{h}_1 - \mathbf{h}_i) s_i (\gamma_1 - \gamma_i). \quad (47)$$

899 Next we claim that for both cases,  $S$  dominates  $((1 - s_1)^2 + Q)$  for large  $R$ . Specifically, we wish for

$$\frac{S \cdot \gamma^{gap}}{4} \geq 4\Gamma A \max((1 - s_1)^2, Q) \iff S \geq 16 \frac{\Gamma A}{\gamma^{gap}} \max((1 - s_1)^2, Q). \quad (48)$$

900 Now choose  $R \geq \delta^{-1} \log(T)/\Theta$  to ensure  $Q \leq S$  since  $Q \leq T e^{-R\Theta\delta} S$  from (45a). Consequently

$$(1 - s_1)^2 = (Q + S)^2 \leq 4S^2 \leq 4ST e^{-R\Theta/2}.$$

901 Combining these, what we wish is ensured by guaranteeing

$$S \geq 16 \frac{\Gamma A}{\gamma^{gap}} \max(4ST e^{-R\Theta/2}, T e^{-R\Theta\delta} S). \quad (49)$$

902 This in turn is ensured for all inputs  $i \in [n]$  by choosing

$$R \geq \frac{\max(2, \delta^{-1})}{\Theta} \log \left( \frac{64T\Gamma A}{\gamma_{\min}^{gap}} \right), \quad (50)$$

903 where  $\gamma_{\min}^{gap} = \min_{i \in [n]} \gamma_i^{gap}$  is the global scalar which is the worst case score gap over all inputs.

904 • **Case 1:**  $V \in \mathcal{S}_\mu(\mathbf{W}^{mm})$ . With the above choice of  $R$ , we guaranteed

$$2(1 - s_1) \cdot \bar{\gamma}^{gap} \geq 2 \cdot S \cdot \bar{\gamma}^{gap} \geq \mathbf{h}^\top \text{diag}(s) \boldsymbol{\gamma} - \mathbf{h}^\top s s^\top \boldsymbol{\gamma} \geq \frac{S \cdot \gamma^{gap}}{4} \geq \frac{(1 - s_1) \gamma^{gap}}{8}.$$

905 via (48) and (46).

906 Since this holds over all inputs, going back to the gradient correlation (44) and averaging above  
907 over all inputs  $i \in [n]$  and plugging back the indices  $i$ , we obtain the advertised bound by setting  
908  $q_i = 1 - s_{i\alpha_i}$  (where we set  $\alpha_i = 1$  above without losing generality)

$$\frac{2}{n} \sum_{i \in [n]} -\ell'_i \cdot q_i \cdot \bar{\gamma}_i^{gap} \geq -\langle \nabla \mathcal{L}(\mathbf{W}), \mathbf{V} \rangle \geq \frac{1}{8n} \sum_{i \in [n]} -\ell'_i \cdot q_i \cdot \gamma_i^{gap}. \quad (51)$$

909 Let  $-\ell'_{\min/\max}$  be the min/max values negative loss derivative admits over the ball  $[-A, A]$  and note  
910 that  $\max_{i \in [n]} \bar{\gamma}_i^{gap} > 0$  and  $\min_{i \in [n]} \gamma_i^{gap} > 0$  are dataset dependent constants. Then, we declare the  
911 constants  $C = -2\ell'_{\max} \cdot \max_{i \in [n]} \bar{\gamma}_i^{gap} > 0$ ,  $c = -(1/8)\ell'_{\min} \cdot \min_{i \in [n]} \gamma_i^{gap} > 0$  to obtain the bound (41a).

• **Case 2:**  $V \in \mathbb{R}^{d \times d}$  and  $\|\mathbf{V}\|_F = \|\mathbf{W}^{mm}\|_F$ . Next, we show (41b) and (41c). For any  $V \in \mathbb{R}^{d \times d}$  satisfying  $\|\mathbf{V}\|_F = \|\mathbf{W}^{mm}\|_F$ , using (47) and the choice of  $R$  in (50) similarly guarantees

$$\frac{2\bar{A}}{\Theta} (1 - s_1) \bar{\gamma}^{gap} \geq \mathbf{h}^\top \text{diag}(s) \boldsymbol{\gamma} - \mathbf{h}^\top s s^\top \boldsymbol{\gamma},$$

912 for fixed input. Going back to the gradient correlation (44) and averaging above over all inputs  $i \in [n]$ ,  
913 with the same definition of  $C > 0$ , we obtain

$$\frac{\bar{A}C}{\Theta n} \sum_{i \in [n]} q_i \geq -\langle \nabla \mathcal{L}(\mathbf{W}), \mathbf{V} \rangle. \quad (52)$$

914 To proceed, since (52) holds for any  $V \in \mathbb{R}^{d \times d}$ , we observe that when setting  $V = \frac{\|\mathbf{W}^{mm}\|_F}{\|\nabla \mathcal{L}(\mathbf{W})\|_F} \cdot \nabla \mathcal{L}(\mathbf{W})$ ,  
915 this implies that

$$\langle \nabla \mathcal{L}(\mathbf{W}), \mathbf{V} \rangle = \|\nabla \mathcal{L}(\mathbf{W})\|_F \cdot \|\mathbf{W}^{mm}\|_F \leq \frac{\bar{A}C}{\Theta n} \sum_{i \in [n]} q_i.$$

916 Simplifying  $\Theta = 1/\|\mathbf{W}^{mm}\|_F$  on both sides gives (41b).

917 Combining the above inequality with (51), we obtain that for all  $V, \mathbf{W} \in \mathcal{S}_\mu(\mathbf{W}^{mm})$

$$-\left\langle \frac{\mathbf{V}}{\|\mathbf{V}\|_F}, \frac{\nabla \mathcal{L}(\mathbf{W})}{\|\nabla \mathcal{L}(\mathbf{W})\|_F} \right\rangle \geq \frac{c\Theta}{C\bar{A}},$$

918 which gives (41c).

919

■

920 **Lemma 11** Suppose Assumption A on the loss function  $\ell$  holds, and let  $\alpha = (\alpha_i)_{i=1}^n$  be locally optimal  
 921 tokens according to Definition 2. Let  $\mathbf{W}^{mm} = \mathbf{W}_\alpha^{mm}$  denote the SVM solution obtained via (Att-SVM)  
 922 by replacing  $(\text{opt}_i)_{i=1}^n$  with  $\alpha = (\alpha_i)_{i=1}^n$ . Let  $\mu = \mu(\alpha) > 0$  and  $\bar{R}_\mu$  be defined as in Lemma 10. For any  
 923 choice of  $\pi > 0$ , there exists  $R_\pi \geq \bar{R}_\mu$  such that, for any  $\mathbf{W} \in C_{\mu, R_\pi}(\mathbf{W}^{mm})$ , we have

$$\left\langle \nabla \mathcal{L}(\mathbf{W}), \frac{\mathbf{W}}{\|\mathbf{W}\|_F} \right\rangle \geq (1 + \pi) \left\langle \nabla \mathcal{L}(\mathbf{W}), \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F} \right\rangle.$$

924 **Proof.** Let  $R = R_\pi$ ,  $\bar{\mathbf{W}} = \|\mathbf{W}^{mm}\|_F \mathbf{W} / \|\mathbf{W}\|_F$ ,  $\mathbf{h}_i = X_i \bar{\mathbf{W}} \mathbf{z}_i$ , and  $\bar{\mathbf{h}}_i = X_i \mathbf{W}^{mm} \mathbf{z}_i$ . To establish the result,  
 925 we will prove that, for sufficiently large  $R$ , for any  $\mathbf{W} \in C_{\mu, R}(\mathbf{W}^{mm})$  and for any  $i \in [n]$ ,

$$\langle \mathbf{h}_i, \mathbb{S}'(X_i \mathbf{W} \mathbf{z}_i) \boldsymbol{\gamma}_i \rangle \leq (1 + \pi) \langle \bar{\mathbf{h}}_i, \mathbb{S}'(X_i \mathbf{W} \mathbf{z}_i) \boldsymbol{\gamma}_i \rangle. \quad (53)$$

926 Once (53) holds for all  $i$ , the same conclusion will hold for the gradient correlations via (44). Moving  
 927 forward, we shall again focus on a single point  $i \in [n]$  and drop all subscripts  $i$ . Also, assume  
 928  $\alpha = \alpha_i = 1$  without losing generality (same as above).

929 Following (46), for all  $\mathbf{W} \in S_\mu(\mathbf{W}^{mm})$  with  $\|\mathbf{W}\|_F = \|\mathbf{W}^{mm}\|_F$  and  $\tilde{\mathbf{h}} = X \mathbf{W} \mathbf{z}$ , and  $s = \mathbb{S}(\tilde{\mathbf{h}})$ , we have  
 930 found

$$\left| \tilde{\mathbf{h}}^\top \text{diag}(s) \boldsymbol{\gamma} - \tilde{\mathbf{h}}^\top s s^\top \boldsymbol{\gamma} - \sum_{t \in \mathcal{T}} (\tilde{\mathbf{h}}_1 - \tilde{\mathbf{h}}_t) s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) \right| \leq 2\Gamma A((1 - s_1)^2 + Q), \quad (54)$$

931 where  $\mathcal{T}$  is the set of support indices. Plugging in  $\mathbf{h}, \bar{\mathbf{h}}$  in the bound above and assuming  $\pi \leq 1$   
 932 (w.l.o.g.), (53) is implied by the following stronger inequality

$$\begin{aligned} 6\Gamma A((1 - s_1)^2 + Q) + \sum_{t \in \mathcal{T}} (\mathbf{h}_1 - \mathbf{h}_t) s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) &\leq (1 + \pi) \sum_{t \in \mathcal{T}} (\bar{\mathbf{h}}_1 - \bar{\mathbf{h}}_t) s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) \\ &= (1 + \pi) \sum_{t \in \mathcal{T}} s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t). \end{aligned}$$

933 First, we claim that  $0.5\pi \sum_{t \in \mathcal{T}} s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) \geq 6\Gamma A((1 - s_1)^2 + Q)$ . The proof of this claim directly  
 934 follows the earlier argument, namely, following (48), (50) and (49) which leads to the choice

$$R \geq \frac{\max(2, \delta^{-1})}{\Theta} \log \left( \frac{C_0 \cdot T\Gamma A}{\pi \gamma_{\min}^{\text{gap}}} \right), \quad (55)$$

935 for some constant  $C_0 > 0$ . Using (50), we choose  $C_0 \geq 64\pi$  to guarantee  $R = R_\pi \geq \bar{R}_\mu$ .

936 Following this control over the perturbation term  $6\Gamma A((1 - s_1)^2 + Q)$ , to conclude with the result,  
 937 what remains is proving the comparison

$$\sum_{t \in \mathcal{T}} (\mathbf{h}_1 - \mathbf{h}_t) s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) \leq (1 + 0.5\pi) \sum_{t \in \mathcal{T}} s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t). \quad (56)$$

938 To proceed, we split the problem into two scenarios.

939 **Scenario 1:**  $\|\bar{\mathbf{W}} - \mathbf{W}^{mm}\|_F \leq \epsilon = \frac{\pi}{4A\Theta}$  for some  $\epsilon > 0$ . In this scenario, for any token, we find that

$$|\mathbf{h}_t - \bar{\mathbf{h}}_t| \leq A\Theta\epsilon = \pi/4.$$

940 Consequently, we obtain

$$\mathbf{h}_1 - \mathbf{h}_t \leq \bar{\mathbf{h}}_1 - \bar{\mathbf{h}}_t + 2A\Theta\epsilon = 1 + 0.5\pi.$$

941 Similarly,  $\mathbf{h}_1 - \mathbf{h}_t \geq 1 - 0.5\pi \geq 0.5$ . Since all terms  $\mathbf{h}_1 - \mathbf{h}_t, s_t, \boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t$  in (56) are nonnegative and  
 942  $(\mathbf{h}_1 - \mathbf{h}_t) s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) \leq (1 + 0.5\pi) s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t)$ , the above implies the desired result (56).

943 **Scenario 2:**  $\|\bar{\mathbf{W}} - \mathbf{W}^{mm}\|_F \geq \epsilon = \frac{\pi}{4A\Theta}$ . Since  $\bar{\mathbf{W}}$  is not (locally) max-margin, in this scenario, for  
 944 some  $\nu = \nu(\epsilon) > 0$  and  $\tau \in \mathcal{T}$ , we have that  $\mathbf{h}_1 - \mathbf{h}_\tau \leq 1 - 2\nu$ . Here  $\tau = \arg \max_{\tau \in \mathcal{T}} \mathbf{x}_\tau \bar{\mathbf{W}} \mathbf{z}$  denotes the  
 945 nearest point to  $\mathbf{h}_1$  (along the  $\bar{\mathbf{W}}$  direction). Note that a non-support index  $\tau \in \bar{\mathcal{T}}$  cannot be closest  
 946 because  $\mathbf{W} \in C_\mu$  and (43) holds. Recall that  $s = \mathbb{S}(\bar{R}\mathbf{h})$  where  $\bar{R} = \|\mathbf{W}\|_F \Theta \geq R\Theta$ . To proceed, split  
 947 the tokens into two groups: Let  $\mathcal{N}$  be the group of tokens obeying  $(\mathbf{x}_1 - \mathbf{x}_\tau) \mathbf{W} \mathbf{z} \leq 1 - \nu$  and  $\mathcal{T} - \mathcal{N}$   
 948 be the rest of the support indices. Observe that

$$\frac{\sum_{t \in \mathcal{T} - \mathcal{N}} s_t}{\sum_{t \in \mathcal{T}} s_t} \leq \frac{\sum_{t \in \mathcal{T} - \mathcal{N}} s_t}{\sum_{t = \tau} s_t} \leq T \frac{e^{\nu \bar{R}}}{e^{2\nu \bar{R}}} = T e^{-\bar{R}\nu}.$$



949 Thus, using  $|\mathbf{h}_1 - \mathbf{h}_t| \leq 2A$  and recalling the definition of  $\gamma^{\text{gap}}$ , observe that

$$\sum_{t \in \mathcal{T} - \mathcal{N}} (\mathbf{h}_1 - \mathbf{h}_t) s_t(\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) \leq \frac{2\Gamma A T e^{-\bar{R}\nu}}{\gamma^{\text{gap}}} \sum_{t \in \mathcal{T}} s_t(\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t).$$

950 Plugging this into (56), we obtain

$$\begin{aligned} \sum_{t \in \mathcal{T}} (\mathbf{h}_1 - \mathbf{h}_t) s_t(\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) &= \sum_{t \in \mathcal{N}} (\mathbf{h}_1 - \mathbf{h}_t) s_t(\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) + \sum_{t \in \mathcal{T} - \mathcal{N}} (\mathbf{h}_1 - \mathbf{h}_t) s_t(\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) \\ &\leq \sum_{t \in \mathcal{N}} (1 - \nu) s_t(\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) + \sum_{t \in \mathcal{T} - \mathcal{N}} 2\Gamma T e^{-\bar{R}\nu} \\ &\leq \left(1 - \nu + \frac{2\Gamma A T e^{-\bar{R}\nu}}{\gamma^{\text{gap}}}\right) \sum_{t \in \mathcal{T}} s_t(\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) \\ &\leq \left(1 + \frac{2\Gamma A T e^{-\bar{R}\nu}}{\gamma^{\text{gap}}}\right) \sum_{t \in \mathcal{T}} s_t(\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t). \end{aligned}$$

951 Consequently, the proof boils down to ensuring the perturbation term  $\frac{2\Gamma A T e^{-\bar{R}\nu}}{\gamma^{\text{gap}}} \leq 0.5\pi$ . Recalling  
952  $\bar{R} \geq R\Theta$ , this is guaranteed for all inputs  $i \in [n]$  by recalling  $\gamma_{\min}^{\text{gap}} = \min_{i \in [n]} \gamma_i^{\text{gap}}$  and choosing

$$R \geq \frac{1}{\nu\Theta} \log \left( \frac{4\Gamma A T}{\gamma_{\min}^{\text{gap}} \pi} \right),$$

953 where  $\nu = \nu(\frac{\pi}{4A\Theta})$  depends only on  $\pi$  and global problem variables.

954 Combining this with the prior  $R$  lower bound of (55) (by taking maximum), we conclude with the  
955 statement.  $\blacksquare$

### 956 E.3.1 Proof of Theorem 3

957 **Theorem 5 (Theorem 3 restated)** Suppose Assumption A on the loss  $\ell$  holds, and let  $\boldsymbol{\alpha} = (\alpha_i)_{i=1}^n$  be  
958 locally optimal tokens according to Definition 2. Let  $\mathbf{W}_{\boldsymbol{\alpha}}^{\text{mm}}$  denote the SVM solution obtained via  
959 (Att-SVM) by replacing  $(\text{opt}_i)_{i=1}^n$  with  $\boldsymbol{\alpha} = (\alpha_i)_{i=1}^n$ . Then,

- 960 • There exist parameters  $\mu = \mu(\boldsymbol{\alpha}) \in (0, 1)$  and  $R > 0$  such that  $C_{\mu, R}(\mathbf{W}_{\boldsymbol{\alpha}}^{\text{mm}})$  does not contain  
961 any stationary points.
- 962 • Algorithm **W-GD** with  $\eta \leq 1/L_{\mathbf{W}}$  and any  $\mathbf{W}(0) \in C_{\mu, R}(\mathbf{W}_{\boldsymbol{\alpha}}^{\text{mm}})$  satisfies  $\lim_{k \rightarrow \infty} \|\mathbf{W}(k)\|_F = \infty$   
963 and  $\lim_{k \rightarrow \infty} \frac{\mathbf{W}(k)}{\|\mathbf{W}(k)\|_F} = \frac{\mathbf{W}_{\boldsymbol{\alpha}}^{\text{mm}}}{\|\mathbf{W}_{\boldsymbol{\alpha}}^{\text{mm}}\|_F}$ .

964 The proof of this theorem follows the proof of [TLZO23, Theorem 3]. Let us denote the initialization  
965 lower bound as  $R_{\mu}^0 := R$ , where  $R$  is given in the Theorem 3's statement. Consider an arbitrary value  
966 of  $\epsilon \in (0, \mu/2)$  and let  $1/(1 + \pi) = 1 - \epsilon$ . We additionally denote  $R_{\epsilon} \leftarrow R_{\pi} \vee 1/2$  where  $R_{\pi}$  was defined  
967 in Lemma 11. At initialization  $\mathbf{W}(0)$ , we set  $\epsilon = \mu/2$  to obtain  $R_{\mu}^0 = R_{\mu/2}$ , and provide the proof in  
968 four steps:

969 **Step 1: There are no stationary points within  $C_{\mu, R_{\mu}^0}(\mathbf{W}^{\text{mm}})$ .** We begin by proving that there are  
970 no stationary points within  $C_{\mu, R_{\mu}^0}(\mathbf{W}^{\text{mm}})$ . Let  $(\mathcal{T}_i)_{i=1}^n$  denote the sets of support indices as defined in  
971 Definition 2. We define  $\bar{\mathcal{T}}_i = [T] - \mathcal{T}_i - \{\alpha_i\}$  as the tokens that are non-support indices. Additionally,  
972 let  $\mu$  be defined as in (42). Then, since  $R_{\mu}^0 \geq \bar{R}_{\mu}$  per Lemma 11, we can apply Lemma 10 to find  
973 that: For all  $\mathbf{V}, \mathbf{W} \in \mathcal{S}_{\mu}(\mathbf{W}^{\text{mm}})$  with  $\|\mathbf{W}\|_F \neq 0$  and  $\|\mathbf{W}\|_F \geq R_{\mu}^0$ , we have that  $-\langle \mathbf{V}, \nabla \mathcal{L}(\mathbf{W}) \rangle$  is strictly  
974 positive.

975 **Step 2:** It follows from Lemma 11 that, there exists  $R_{\epsilon} \geq \bar{R}_{\mu} \vee 1/2$  such that all  $\mathbf{W} \in C_{\mu, R_{\epsilon}}(\mathbf{W}^{\text{mm}})$   
976 satisfy

$$\left\langle -\nabla \mathcal{L}(\mathbf{W}), \frac{\mathbf{W}^{\text{mm}}}{\|\mathbf{W}^{\text{mm}}\|_F} \right\rangle \geq (1 - \epsilon) \left\langle -\nabla \mathcal{L}(\mathbf{W}), \frac{\mathbf{W}}{\|\mathbf{W}\|_F} \right\rangle. \quad (57)$$

977 The argument below applies to a general  $\epsilon \in (0, \mu/2)$ . However, at initialization  $\mathbf{W}(0)$ , we set  $\epsilon = \mu/2$   
 978 and, recalling above, initialization lower bound was defined as  $R_\mu^0 := R_{\mu/2}$ . To proceed, for any  
 979  $\epsilon \in (0, \mu/2)$ , we will show that after gradient descent enters the conic set  $C_{\mu, R_\epsilon}(\mathbf{W}^{mm})$  for the first  
 980 time, it will never leave the set. Let  $t_\epsilon$  be the first time gradient descent enters  $C_{\mu, R_\epsilon}(\mathbf{W}^{mm})$ . In **Step 4**,  
 981 we will prove that such  $t_\epsilon$  is guaranteed to exist. Additionally, for  $\epsilon \leftarrow \mu/2$ , note that  $t_\epsilon = 0$  i.e. the  
 982 point of initialization.

983 **Step 3: Updates remain inside the cone  $C_{\mu, R_\epsilon}(\mathbf{W}^{mm})$ .** By leveraging the results from **Step 1** and  
 984 **Step 2**, we demonstrate that the gradient iterates, with an appropriate constant step size, starting from  
 985  $\mathbf{W}(k_\epsilon) \in C_{\mu, R_\epsilon}(\mathbf{W}^{mm})$ , remain within this cone.

986 We proceed by induction. Suppose that the claim holds up to iteration  $k \geq k_\epsilon$ . This implies that  
 987  $\mathbf{W}(k) \in C_{\mu, R_\epsilon}(\mathbf{W}^{mm})$ . Hence, recalling cone definition, there exists scalar  $\mu = \mu(\alpha) \in (0, 1)$  and  $R$  such  
 988 that  $\|\mathbf{W}(k)\|_F \geq R$ , and

$$\left\langle \frac{\mathbf{W}(k)}{\|\mathbf{W}(k)\|_F}, \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F} \right\rangle \geq 1 - \mu.$$

989 For all  $k \geq 1$ , let

$$\rho(k) := -\frac{1}{1 - \epsilon} \left\langle \nabla \mathcal{L}(\mathbf{W}(k)), \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F} \right\rangle. \quad (58)$$

990 Note that  $\rho(k) > 0$  due to **Step 1**. This together with the gradient descent update rule gives

$$\begin{aligned} \left\langle \frac{\mathbf{W}(k+1)}{\|\mathbf{W}(k+1)\|_F}, \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F} \right\rangle &= \left\langle \frac{\mathbf{W}(k)}{\|\mathbf{W}(k)\|_F} - \frac{\eta}{\|\mathbf{W}(k)\|_F} \nabla \mathcal{L}(\mathbf{W}(k)), \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F} \right\rangle \\ &\geq 1 - \mu - \frac{\eta}{\|\mathbf{W}(k)\|_F} \left\langle \nabla \mathcal{L}(\mathbf{W}(k)), \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F} \right\rangle \\ &\geq 1 - \mu + \frac{\eta \rho(k)(1 - \epsilon)}{\|\mathbf{W}(k)\|_F}. \end{aligned} \quad (59a)$$

991 Note that from Lemma 10, we have  $\langle \nabla \mathcal{L}(\mathbf{W}(k)), \mathbf{W}(k) \rangle < 0$  which implies that  $\|\mathbf{W}(k+1)\|_F \geq$   
 992  $\|\mathbf{W}(k)\|_F$ . This together with  $R_\epsilon$  definition and  $\|\mathbf{W}(k)\|_F \geq 1/2$  implies that

$$\begin{aligned} \|\mathbf{W}(k+1)\|_F &\leq \frac{1}{2\|\mathbf{W}(k)\|_F} \left( \|\mathbf{W}(k+1)\|_F^2 + \|\mathbf{W}(k)\|_F^2 \right) \\ &= \frac{1}{2\|\mathbf{W}(k)\|_F} \left( 2\|\mathbf{W}(k)\|_F^2 - 2\eta \langle \nabla \mathcal{L}(\mathbf{W}(k)), \mathbf{W}(k) \rangle + \eta^2 \|\nabla \mathcal{L}(\mathbf{W}(k))\|_F^2 \right) \\ &\leq \|\mathbf{W}(k)\|_F - \frac{\eta}{\|\mathbf{W}(k)\|_F} \langle \nabla \mathcal{L}(\mathbf{W}(k)), \mathbf{W}(k) \rangle + \eta^2 \|\nabla \mathcal{L}(\mathbf{W}(k))\|_F^2, \end{aligned}$$

993 which gives

$$\begin{aligned} \frac{\|\mathbf{W}(k+1)\|_F}{\|\mathbf{W}(k)\|_F} &\leq 1 - \frac{\eta}{\|\mathbf{W}(k)\|_F} \left\langle \nabla \mathcal{L}(\mathbf{W}(k)), \frac{\mathbf{W}(k)}{\|\mathbf{W}(k)\|_F} \right\rangle + \eta^2 \frac{\|\nabla \mathcal{L}(\mathbf{W}(k))\|_F^2}{\|\mathbf{W}(k)\|_F} \\ &\leq 1 - \frac{\eta}{(1 - \epsilon)\|\mathbf{W}(k)\|_F} \left\langle \nabla \mathcal{L}(\mathbf{W}(k)), \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F} \right\rangle + \eta^2 \frac{\|\nabla \mathcal{L}(\mathbf{W}(k))\|_F^2}{\|\mathbf{W}(k)\|_F} \\ &\leq 1 + \frac{\eta \rho(k)}{\|\mathbf{W}(k)\|_F} + \frac{\eta^2 \|\nabla \mathcal{L}(\mathbf{W}(k))\|_F^2}{\|\mathbf{W}(k)\|_F} =: C_1(\rho(k), \eta). \end{aligned} \quad (59b)$$

994 Here, the second inequality follows from (57) and (58).

995 Now, it follows from (59a) and (59b) that

$$\begin{aligned}
\left\langle \frac{\mathbf{W}(k+1)}{\|\mathbf{W}(k+1)\|}, \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|} \right\rangle &\geq \frac{1}{C_1(\rho(k), \eta)} \left( 1 - \mu + \frac{\eta\rho(k)(1-\epsilon)}{\|\mathbf{W}(k)\|_F} \right) \\
&= 1 - \mu + \frac{1}{C_1(\rho(k), \eta)} \left( (1-\mu)(1 - C_1(\rho(k), \eta)) + \frac{\eta\rho(k)(1-\epsilon)}{\|\mathbf{W}(k)\|_F} \right) \\
&= 1 - \mu + \frac{\eta}{C_1(\rho(k), \eta)} \left( (\mu-1) \left( \frac{\rho(k)}{\|\mathbf{W}(k)\|_F} + \frac{\eta\|\nabla\mathcal{L}(\mathbf{W}(k))\|^2}{\|\mathbf{W}(k)\|_F} \right) + \frac{\rho(k)(1-\epsilon)}{\|\mathbf{W}(k)\|_F} \right) \\
&= 1 - \mu + \frac{\eta}{C_1(\rho(k), \eta)} \left( \frac{\rho(k)(\mu-\epsilon)}{\|\mathbf{W}(k)\|_F} - \eta(1-\mu) \frac{\|\nabla\mathcal{L}(\mathbf{W}(k))\|^2}{\|\mathbf{W}(k)\|_F} \right) \\
&\geq 1 - \mu,
\end{aligned} \tag{60}$$

996 where the last inequality uses our choice of stepsize  $\eta \leq 1/L_W$  in Theorem 3's statement. Specifically,  
997 we need  $\eta$  to be small to ensure the last inequality. We will guarantee this by choosing a proper  $R_\epsilon$  in  
998 Lemma 11. Specifically, Lemma 11 leaves the choice of  $C_0$  in  $R_\epsilon$  lower bound of (55) open (it can  
999 always be chosen larger). Here, by choosing  $C_0 \geq 1/L_W$  will ensure  $\eta \leq 1/L_W$  works well.

$$\begin{aligned}
\eta &\leq \frac{\mu}{2(1-\mu)(1-\frac{\mu}{2})} \frac{c}{C} \frac{\Theta}{\bar{A}} \frac{1}{\bar{A}CT} e^{R_\mu^0 \Theta/2} \\
&\leq \frac{\mu-\epsilon}{1-\mu} \cdot \frac{1}{1-\epsilon} \cdot \frac{c}{C} \cdot \frac{\Theta}{\bar{A}} \cdot \frac{1}{\bar{A}CT} e^{R_\mu^0 \Theta/2} \leq \frac{(\mu-\epsilon)}{1-\mu} \frac{\rho(k)}{\|\nabla\mathcal{L}(\mathbf{W}(k))\|_F^2}.
\end{aligned} \tag{61}$$

1000 Here, the first inequality uses our choice of  $\epsilon \in (0, \mu/2)$  (see Step 2), and the last inequality is  
1001 obtained from Lemma 10 since

$$\begin{aligned}
\frac{\rho(k)}{\|\nabla\mathcal{L}(\mathbf{W}(k))\|_F} &= -\frac{1}{1-\epsilon} \left\langle \frac{\nabla\mathcal{L}(\mathbf{W}(k))}{\|\nabla\mathcal{L}(\mathbf{W}(k))\|_F}, \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F} \right\rangle \geq \frac{1}{1-\epsilon} \cdot \frac{c}{C} \cdot \frac{\Theta}{\bar{A}}, \\
\frac{1}{\|\nabla\mathcal{L}(\mathbf{W}(k))\|_F} &\geq \frac{1}{\bar{A}C \cdot \frac{1}{n} \sum_{i=1}^n (1-s_{i\alpha_i})} \geq \frac{1}{\bar{A}CT e^{-R_\mu^0 \Theta/2}}
\end{aligned}$$

1002 for some data dependent constants  $c$  and  $C$ ,  $\bar{A} = \max_{i \in [n], t, \tau \in [T]} \|\mathbf{x}_{it} - \mathbf{x}_{i\tau}\| \|\mathbf{z}_i\|$ , and  $\Theta = 1/\|\mathbf{W}^{mm}\|_F$ .

1003 Next, we will demonstrate that the choice of  $\eta$  in (61) does indeed meet our step size condition as  
1004 stated in the theorem, i.e.,  $\eta \leq 1/L_W$ . Recall that  $1/(1+\pi) = 1-\epsilon$ , which implies that  $\pi = \epsilon/(1-\epsilon)$ .

1005 Combining this with (55), we obtain:

$$R_\pi \geq \frac{\max(2, \delta^{-1})}{\Theta} \log \left( \frac{C_0 T \Gamma A}{\pi \gamma_{\min}^{gap}} \right), \quad \text{where } C_0 \geq 64\pi. \tag{62}$$

$$\Rightarrow R_\epsilon \geq \frac{\max(2, \delta^{-1})}{\Theta} \log \left( \frac{(1-\epsilon)C_0 T \Gamma A}{\epsilon \gamma_{\min}^{gap}} \right), \quad \text{where } C_0 \geq 64 \frac{\epsilon}{1-\epsilon}. \tag{63}$$

1006 On the other hand, at the initialization, we have  $\epsilon = \mu/2$  which implies that

$$R_\mu^0 \geq \frac{\max(2, \delta^{-1})}{\Theta} \log \left( \frac{(2-\mu)C_0 T \Gamma A}{\mu \gamma_{\min}^{gap}} \right), \quad \text{where } C_0 \geq 64 \frac{\mu}{2(1-\frac{\mu}{2})}. \tag{64}$$

1007 In the following, we will determine a lower bound on  $C_0$  such that our step size condition in  
1008 Theorem 3's statement, i.e.,  $\eta \leq 1/L_W$ , is satisfied. Note that for the choice of  $\eta$  in (61) to meet the  
1009 condition  $\eta \leq 1/L_W$ , the following condition must hold:

$$\frac{1}{L_W} \leq \frac{\mu}{(2-\mu)} \frac{1}{C_2 T} e^{R_\mu^0 \Theta/2} \Rightarrow R_\mu^0 \geq \frac{2}{\Theta} \log \left( \frac{1}{L_W} \frac{2-\mu}{\mu} C_2 T \right). \tag{65}$$

1010 where  $C_2 = (1-\mu) \frac{\bar{A}^2 C^2}{\Theta c}$ .

1011 This together with (64) implies that

$$\frac{C_0 \Gamma A}{\gamma_{\min}^{gap}} \geq (1-\mu) \frac{C_2}{L_W} \Rightarrow C_0 \geq \max \left( \frac{(1-\mu)C_2 \gamma_{\min}^{gap}}{L_W \Gamma A}, \frac{64\mu}{2-\mu} \right). \tag{66}$$

1012 Therefore, with this lower bound on  $C_0$ , the step size bound in (61) is sufficiently large to ensure that  
 1013  $\eta \leq 1/L_W$  guarantees (60).

1014 Hence, it follows from (60) that  $\mathbf{W}(k+1) \in C_{\mu, R_\epsilon}(\mathbf{W}^{mm})$ .

1015 **Step 4: The correlation of  $\mathbf{W}(k)$  and  $\mathbf{W}^{mm}$  increases over  $k$ .** The remainder is similar to the proof  
 1016 of Theorem 2. From Step 3, we have that all iterates remain within the initial conic set i.e.  $\mathbf{W}(k) \in$   
 1017  $C_{\mu, R_\mu^0}(\mathbf{W}^{mm})$  for all  $k \geq 0$ . Note that it follows from Lemma 10 that  $\langle \nabla \mathcal{L}(\mathbf{W}), \mathbf{W}^{mm} / \|\mathbf{W}^{mm}\|_F \rangle < 0$ ,  
 1018 for any finite  $\mathbf{W} \in C_{\mu, R_\mu^0}(\mathbf{W}^{mm})$ . Hence, there are no finite critical points  $\mathbf{W} \in C_{\mu, R_\mu^0}(\mathbf{W}^{mm})$ , for which  
 1019  $\nabla \mathcal{L}(\mathbf{W}) = 0$ . Now, based on Lemma 7, which guarantees that  $\nabla \mathcal{L}(\mathbf{W}(k)) \rightarrow 0$ , this implies that  
 1020  $\|\mathbf{W}(t)\|_F \rightarrow \infty$ . Consequently, for any choice of  $\epsilon \in (0, \mu/2)$  there is an iteration  $k_\epsilon$  such that, for all  
 1021  $k \geq k_\epsilon$ ,  $\mathbf{W}(k) \in C_{\mu, R_\epsilon}(\mathbf{W}^{mm})$ . Once within  $C_{\mu, R_\epsilon}(\mathbf{W}^{mm})$ , following similar steps in (37) and (38), for  
 1022 any  $k \geq k_\epsilon$ ,

$$\left\langle \frac{\mathbf{W}(k)}{\|\mathbf{W}(k)\|_F}, \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F} \right\rangle \geq 1 - \epsilon + \frac{C(\epsilon, \eta)}{\|\mathbf{W}(k)\|_F}, \quad \mathbf{W}(k) \in C_{\mu, R_\epsilon}(\mathbf{W}^{mm}),$$

1023 for some finite constant  $C(\epsilon, \eta)$  (that depends only on  $\eta, \epsilon, \|\mathbf{W}(k_\epsilon)\|_F$ ).

1024 Consequently, as  $k \rightarrow \infty$

$$\liminf_{k \rightarrow \infty} \left\langle \frac{\mathbf{W}(k)}{\|\mathbf{W}(k)\|_F}, \frac{\mathbf{W}^{mm}}{\|\mathbf{W}^{mm}\|_F} \right\rangle \geq 1 - \epsilon, \quad \mathbf{W}(k) \in C_{\mu, R_\epsilon}(\mathbf{W}^{mm}).$$

1025 Since  $\epsilon \in (0, \mu/2)$  is arbitrary, we get  $\mathbf{W}(k)/\|\mathbf{W}(k)\|_F \rightarrow \mathbf{W}^{mm}/\|\mathbf{W}^{mm}\|_F$ . ■

## 1026 F Supporting Experiments

1027 In this section, we introduce implementation details and additional experiments. We create a 1-layer  
 1028 self-attention using PyTorch, training it with the SGD optimizer and a learning rate of  $\eta = 0.1$ . We  
 1029 apply normalized gradient descent to ensure divergence of attention weights. The attention weight  $\mathbf{W}$   
 1030 is then updated through

$$\mathbf{W}(k+1) = \mathbf{W}(k) - \eta \frac{\nabla \mathcal{L}(\mathbf{W}(k))}{\|\nabla \mathcal{L}(\mathbf{W}(k))\|_F}.$$

1031 In the setting of  $(\mathbf{K}, \mathbf{Q})$ -parameterization, we noted that with extended training iterations, the norm  
 1032 of the combined parameter  $\mathbf{K}\mathbf{Q}^\top$  consistently rises, despite the gradient being treated as zero due to  
 1033 computational limitations. To tackle this issue, we introduce a minor regularization penalty to the  
 1034 loss function, ensuring that the norms of  $\mathbf{K}$  and  $\mathbf{Q}$  remain within reasonable bounds. This adjustment  
 1035 involves

$$\tilde{\mathcal{L}}(\mathbf{K}, \mathbf{Q}) = \mathcal{L}(\mathbf{K}, \mathbf{Q}) + \lambda(\|\mathbf{K}\|_F^2 + \|\mathbf{Q}\|_F^2).$$

1036 Here, we set  $\lambda$  to be the the smallest representable number, e.g. computed as  $1 + \lambda = 1$  in Python,  
 1037 which is around  $2.22 \times 10^{-16}$ . Therefore,  $\mathbf{K}, \mathbf{Q}$  parameters are updated as follows.

$$\mathbf{K}(k+1) = \mathbf{K}(k) - \eta \frac{\nabla \tilde{\mathcal{L}}_{\mathbf{K}}(\mathbf{K}(k), \mathbf{Q}(k))}{\|\nabla \tilde{\mathcal{L}}_{\mathbf{K}}(\mathbf{K}(k), \mathbf{Q}(k))\|_F}, \quad \mathbf{Q}(k+1) = \mathbf{Q}(k) - \eta \frac{\nabla \tilde{\mathcal{L}}_{\mathbf{Q}}(\mathbf{K}(k), \mathbf{Q}(k))}{\|\nabla \tilde{\mathcal{L}}_{\mathbf{Q}}(\mathbf{K}(k), \mathbf{Q}(k))\|_F}.$$

1038 • As observed in previous work [TLZO23], and due to the exponential expression of softmax  
 1039 nonlinearity and computation limitation, PyTorch has no guarantee to select optimal tokens when  
 1040 the score gap is too small. Therefore in Figures 2, 9 and 10, we generate random tokens making sure  
 1041 that  $\min_{i \in [n], i \neq \text{opt}_i} \gamma_{\text{opt}_i} - \gamma_{it} \geq \underline{\gamma}$  and we choose  $\underline{\gamma} = 0.1$  in our experiments.

1042 **Rank sensitivity of  $(\mathbf{K}, \mathbf{Q})$ -parameterization (Figures 6&7).** In Lemma 1, we have theoretically  
 1043 established that the rank of the SVM solution, denoted as  $\mathbf{W}^{mm}$  in (Att-SVM) or  $\mathbf{W}_\star^{mm}$  in (Att-SVM $\star$ ),  
 1044 is at most rank  $\max(n, d)$ . To further verify it, Figure 6 illustrates rank range of  $\mathbf{W}^{mm}$  and  $\mathbf{W}_\star^{mm}$ ,  
 1045 solved using optimal tokens  $(\text{opt}_i)_{i=1}^n$  and setting  $m = d$  (the rank constraint is eliminated). Each  
 1046 result is averaged over 100 trials, and for each trial,  $\mathbf{x}_{it}, \mathbf{z}_i$ , and linear head  $\mathbf{v}$  are randomly sampled  
 1047 from the unit sphere. In Fig. 6(a), we fix  $T = 5$  and vary  $n$  across  $\{5, 10, 15\}$ . Conversely, in Fig. 6(b),  
 1048 we keep  $n = 5$  constant and alter  $T$  across  $\{5, 10, 15\}$ . Both figures confirm rank of  $\mathbf{W}^{mm}$  and  $\mathbf{W}_\star^{mm}$   
 1049 are bounded by  $\max(n, d)$ , validating Lemma 1.

1050 Now, moving to Figure 7, we delve into GD performance across various dimensions of  $\mathbf{K}, \mathbf{Q} \in \mathbb{R}^{d \times m}$   
 1051 while keeping  $d = 20$  fixed and varying  $m$  from 1 to 10. In the upper subfigure, we maintain a constant

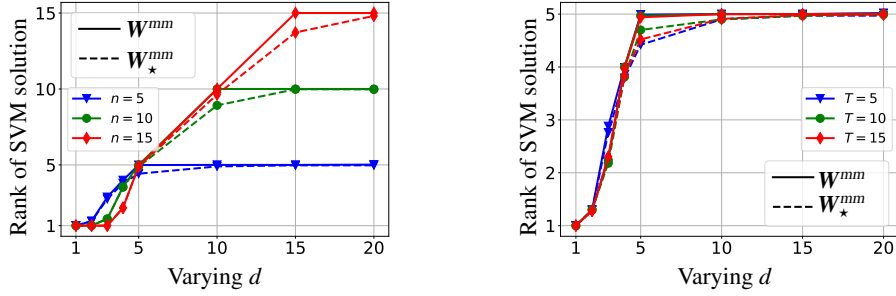
(a) Rank of SVM solutions with fixed  $T = 5$ (b) Rank of SVM solutions with fixed  $n = 5$ 

Figure 6: Rank range of solutions for (Att-SVM) and (Att-SVM<sub>\*</sub>), denoted as  $W^{mm}$  and  $W_{*}^{mm}$ , solved using optimal tokens  $(\text{opt}_i)_{i=1}^n$  and setting  $m = d$  (the rank constraint is eliminated). Both figures confirm ranks of  $W^{mm}$  and  $W_{*}^{mm}$  are bounded by  $\max(n, d)$ , validating Lemma 1.

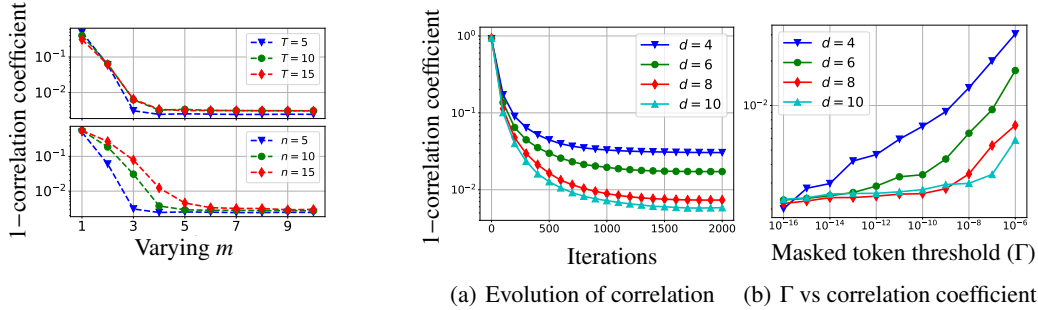
(a) Evolution of correlation (b)  $\Gamma$  vs correlation coefficient

Figure 7: Convergence behavior of GD when training  $(K, Q) \in \mathbb{R}^{d \times m}$  with varying  $m$ . The misalignment,  $1 - \text{corr\_coef}(W, KQ^T)$ , is studied, where  $W_{*,\alpha}^{mm}$  is from (Att-SVM<sub>\*</sub>) with  $\text{opt}$  replaced by  $\alpha$  and  $m = d$ . Subfigures with fixed  $n = 5$  (upper) and  $T = 5$  (lower) show that as  $m$  approaches or exceeds  $n$ ,  $KQ^T$  aligns more with  $W_{*,\alpha}^{mm}$ .

Figure 8: Behavior of GD with nonlinear nonconvex prediction head and multi-token compositions. (a): Blue, green, red and teal curves represent the evolution of  $1 - \text{corr\_coef}(W, W^{\text{SVMeq}})$  for  $d = 4, 6, 8$  and  $10$  respectively, which have been displayed in Figure 4(upper). (b): Over the 500 random instances as discussed in Figure 4, we filter different instances by constructing masked set with tokens whose softmax output  $< \Gamma$  and vary  $\Gamma$  from  $10^{-16}$  to  $10^{-6}$ . The corresponding results of  $1 - \text{corr\_coef}(W, W^{\text{SVMeq}})$  are displayed in blue, green, red and teal curves.

1052  $n = 5$  and vary  $T$  within  $\{5, 10, 15\}$ , while in the lower subfigure,  $T$  is fixed at 5 and  $n$  changes  
 1053 within  $\{5, 10, 15\}$ . Results are depicted using blue, green, and red dashed curves, with both y-axes  
 1054 representing  $1 - \text{corr\_coef}(W, W_{*,\alpha}^{mm})$ , where  $W$  represents the GD solution and  $W_{*,\alpha}^{mm}$  is obtained  
 1055 from (Att-SVM<sub>\*</sub>) by employing token indices  $\alpha$  selected via GD and setting the rank limit to  $m = d$ .  
 1056 Observing both subfigures, we note that a larger  $n$  necessitates a larger  $m$  for attention weights  $KQ^T$   
 1057 to accurately converge to the SVM solution (Figure 7(lower)). Meanwhile, performances remain  
 1058 consistent across varying  $T$  values (Figure 7(upper)). This observation further validates Lemma 1.  
 1059 Furthermore, the results demonstrate that  $W$  converges directionally towards  $W_{*,\alpha}^{mm}$  as long as  $m \geq n$ .

1060 **Global Convergence via overparameterization (Figures 9&10).** The trend depicted in Figure  
 1061 9, where the percentage of global convergence (red bars) approaches 100% and Assumption B(ii)  
 1062 holds with higher probability (green bars) as  $d$  grows, reinforces this insight. Specifically, Fig. 9(a)  
 1063 is same as Figure 2, and Fig. 9(b) displays the same evaluation over  $(K, Q)$ -parameterization setting. In  
 1064 both experiments, and for each chosen  $d$  value, a total of 500 random instances are conducted under  
 1065 the conditions of  $n = T = 5$ . The outcomes are reported in terms of the percentages of Not Local,  
 1066 Local, and Global convergence, represented by the teal, blue, and red bars, respectively. We validate  
 1067 Assumption B(ii) as follows: Given a problem instance, we compute the average margin over all  
 1068 non-optimal tokens of all inputs and declare that problem satisfies Assumption B(ii), if the average

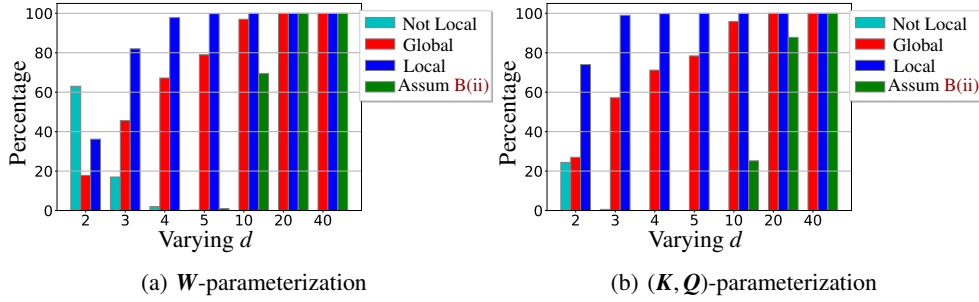


Figure 9: Percentage of different convergence types of GD when training cross-attention weights (a):  $W$  or (b):  $(K, Q)$  with varying  $d$ . In both figures, red, blue, and teal bars represent the percentages of Global, Local (including Global), and Not Local convergence, respectively. The green bar corresponds to Assumption B(ii) where all tokens act as support vectors. Larger overparameterization ( $d$ ) relates to a higher percentage of globally-optimal SVM convergence.

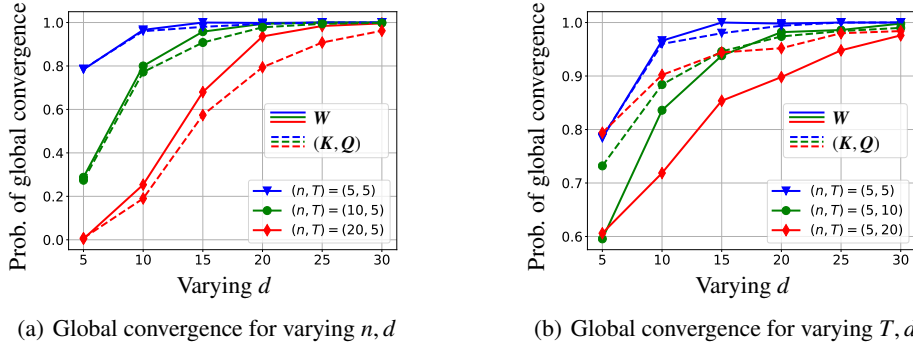


Figure 10: Global convergence behavior of GD when training cross-attention weights  $W$  (solid) or  $(K, Q)$  (dashed) with random data. The blue, green, and red curves represent the probabilities of global convergence for (a): fixing  $T = 5$  and varying  $n \in \{5, 10, 20\}$  and (b): fixing  $n = 5$  and varying  $T \in \{5, 10, 20\}$ . Results demonstrate that for both attention models, as  $d$  increases (due to over-parameterization), attention weights tend to select optimal tokens  $(\text{opt}_i)_{i=1}^n$ .

1069 margin is below 1.1 (where 1 is the minimum). Here, recall that margin of a non-optimal token is  
 1070 defined as  $(\mathbf{x}_{i_{\text{opt}_t}} - \mathbf{x}_{it})^\top \mathbf{W}^{mm} \mathbf{z}_i$  or  $(\mathbf{x}_{i_{\text{opt}_t}} - \mathbf{x}_{it})^\top \mathbf{W}_\star^{mm} \mathbf{z}_i$  for  $t \neq \text{opt}_i$ .

1071 Furthermore, the observations in Figure 10 regarding the percentages of achieving global convergence  
 1072 reaching 100 with larger  $d$  reaffirm that overparameterization leads the attention weights to converge  
 1073 directionally towards the optimal max-margin direction outlined by (Att-SVM) and (Att-SVM $_\star$ ).

1074 **Behavior of GD with nonlinear nonconvex prediction head and multi-token compositions**  
 1075 **(Figure 8).** To better investigate how correlation changes with data dimension  $d$ , we collect the  
 1076 solid curves in Figure 4(upper) and construct as Figure 8(a). Moreover, Figure 8(b) displays the  
 1077 average correlation of instances (refer to scatters in Figure 4 (lower)), considering masked tokens  
 1078 with softmax probability  $< \Gamma$ . Both findings highlight that higher  $d$  enhances alignment. For  $d \geq 8$  or  
 1079  $\Gamma \leq 10^{-9}$ , the GD solution  $W$  achieves a correlation of  $> 0.99$  with the SVM-equivalence  $W^{\text{SVMeq}}$ ,  
 1080 defined in Section B.

1081 **Investigation of Lemma 3 over different  $\tau$  selections (Figure 11).** Consider the setting of Sec-  
 1082 tion B.1 and Lemma 3. Figure 5 explores the influence of  $\lambda$  on the count of tokens selected by  
 1083 GD-derived attention weights. As  $\lambda$  increases, the likelihood of selecting more tokens also increases.  
 1084 Shifting focus to Figure 11, we examine the effect of  $\tau$ . For each outcome, we generate random  
 1085  $\lambda$  values, retaining pairs  $(\lambda, X)$  satisfying  $\tau$  constraints, with averages derived from 100 successful  
 1086 trials. The results indicate a positive correlation among  $\tau$ ,  $\lambda$ , and the number of selected tokens.

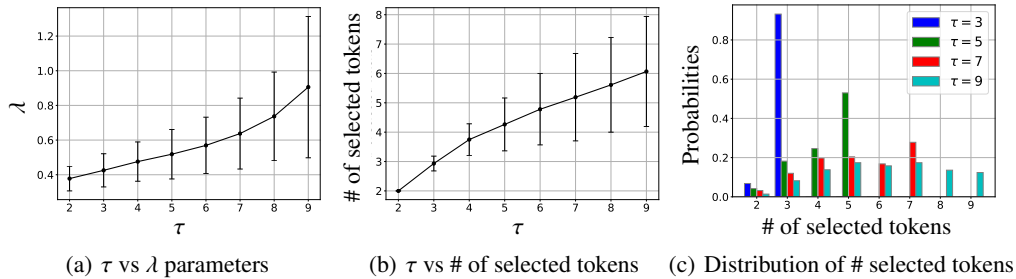


Figure 11: Behavior of GD when selecting multiple tokens.

1087 Moreover, Figure 11(c) provides a precise distribution of selected token counts across various  $\tau$   
 1088 values (specifically  $\tau \in \{3, 5, 7, 9\}$ ). The findings confirm that the number of selected tokens remains  
 1089 within the limit of  $\tau$ , thus validating the assertion made in Lemma 3.

## 1090 G Discussion, Future Directions, and Open Problems

1091 Our optimization-theoretic characterization of the self-attention model provides a comprehensive  
 1092 understanding of its underlying principles. The developed framework, along with the research  
 1093 presented in [TLZO23], introduces new avenues for studying transformers and language models. The  
 1094 key findings include:

- 1095 ✓ The optimization geometry of self-attention exhibits a fascinating connection to hard-margin SVM  
 1096 problems. By leveraging linear constraints formed through outer products of token pairs, optimal  
 1097 input tokens can be effectively separated from non-optimal ones.
- 1098 ✓ When gradient descent is employed without early-stopping, implicit regularization and conver-  
 1099 gence of self-attention naturally occur. This convergence leads to the maximum margin solution  
 1100 when minimizing specific requirements using logistic loss, exp-loss, or other smooth decreasing loss  
 1101 functions. Moreover, this implicit bias is unaffected by the step size, as long as it is sufficiently small  
 1102 for convergence, and remains independent of the initialization process.

1103 The fact that gradient descent leads to a maximum margin solution may not be surprising to those  
 1104 who are familiar with the relationship between regularization path and gradient descent in linear and  
 1105 nonlinear neural networks [SHN<sup>+</sup>18, GLSS18, NLG<sup>+</sup>19, JT21, MWG<sup>+</sup>20, JT20]. However, there is  
 1106 a lack of prior research or discussion regarding this connection to the attention mechanism. Moreover,  
 1107 there has been no rigorous analysis or investigation into the exactness and independence of this bias  
 1108 with respect to the initialization and step size. Thus, we believe our findings and insights deepen  
 1109 our understanding of transformers and language models, paving the way for further research in this  
 1110 domain. Below, we discuss some notable directions and highlight open problems that are not resolved  
 1111 by the existing theory.

- 1112 • **Convergence Rates:** The current paper establishes asymptotic convergence of gradient  
 1113 descent; nonetheless, there is room for further exploration to characterize non-asymptotic  
 1114 convergence rates. Indeed, such an exploration can also provide valuable insights into the  
 1115 choice of learning rate, initialization, and the optimization method.
- 1116 • **Gradient descent on  $(K, Q)$  parameterization:** We find it remarkable that regularization  
 1117 path analysis was able to predict the implicit bias of gradient descent. Complete analysis  
 1118 of gradient descent is inherently connected to the fundamental question of low-rank factor-  
 1119 ization [GWB<sup>+</sup>17, LMZ18]. We believe formalizing the implicit bias of gradient descent  
 1120 under margin constraints presents an exciting open research direction for further research.
- 1121 • **Generalization analysis:** An important direction is the generalization guarantees for  
 1122 gradient-based algorithms. The established connection to hard-margin SVM can facilitate  
 1123 this because the SVM problem is amenable to statistical analysis. This would be akin to  
 1124 how kernel/NTK analysis for deep nets enabled a rich literature on generalization analysis  
 1125 for traditional deep learning.



- 1126
- 1127
- 1128
- 1129
- 1130
- **Realistic architectures:** Naturally, we wish to explore whether max-margin equivalence can be extended to more realistic settings: Can the theory be expanded to handle multi-head attention, multi-layer architectures, and MLP nonlinearities? We believe the results in Section B take an important step towards this direction by including analytical formulae for the implicit bias of the attention layer under nonlinear prediction heads.
- 1131
- **Jointly optimizing attention and prediction head:** It would be interesting to study the joint optimization dynamics of attention weights and prediction head  $h(\cdot)$ . This problem can be viewed as a novel low-rank factorization type problem where  $h(\cdot)$  and  $\mathbf{W}$  are factors of the optimization problem, only, here,  $\mathbf{W}$  passes through the softmax nonlinearity. To this aim, [TLZO23] provides a preliminary geometric characterization of the implicit bias for a simpler attention model using regularization path analysis. Such findings can potentially be generalized to the analysis of gradient methods and full transformer block.
- 1132
- 1133
- 1134
- 1135
- 1136
- 1137