
Margin Maximization in Attention Mechanism

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Attention mechanism is a central component of the transformer architecture which
2 led to the phenomenal success of large language models. However, the theoretical
3 principles underlying the attention mechanism are poorly understood, especially its
4 nonconvex optimization dynamics. In this work, we explore the seminal softmax-
5 attention model $f(X) = \mathbf{v}^\top X^\top \text{softmax}(XW^\top \mathbf{p})$, where, X is the tokenized input, \mathbf{v}
6 is the value weights, W is the key-query weights, and \mathbf{p} is a tunable token/prompt.
7 We prove that running gradient descent on \mathbf{p} , or equivalently W , converges to a max-
8 margin solution that separates locally-optimal tokens from non-optimal ones. We
9 also develop regularization path analysis that generalizes these findings to nonlinear
10 classifier head – rather than linear \mathbf{v} . When optimizing \mathbf{v} and \mathbf{p} simultaneously with
11 logistic loss, we identify conditions under which the regularization paths converge
12 to their respective max-margin solutions where \mathbf{v} separates the input features based
13 on their labels. Finally, we verify our results through numerical insights.

14 1 Introduction

15 Since its introduction in the seminal work [1], attention mechanism has played an influential role in
16 the advances in natural language processing, and more recently, large language models [2, 3, 4, 5].
17 Attention is initially introduced for encoder-decoder RNN architectures in order to allow the decoder
18 to use the most relevant parts of the input sequence, rather than relying on a fixed-length hidden
19 state. Attention mechanism has taken the center stage in the transformers [6] where the self-attention
20 layer – which calculates softmax similarities between input tokens – forms the backbone of the
21 architecture. Since their inception, transformers have revolutionized natural language processing
22 (from BERT to ChatGPT [7, 8]) and they have also become the architecture of choice for foundation
23 models [9] to address diverse challenges in generative modeling [3, 10], computer vision [11, 12],
24 and reinforcement learning [13, 14, 15].

25 The prominence of the attention mechanism motivate a fundamental theoretical understanding of its
26 role in optimization and learning. While it is well-known that attention enables the model to focus on
27 the relevant parts of the input sequence, the precise mechanism by which this is achieved is far from
28 clear. To this end, we ask

29 **Q:** What are the optimization dynamics and inductive biases of the attention mechanism?

30 We study this question using the fundamental attention model $f(X) = \mathbf{v}^\top X^\top \mathbb{S}(XW^\top \mathbf{p})$. Here, X is the
31 sequence of input tokens, \mathbf{v} is the classifier head, W is the trainable key-query weights, and \mathbb{S} denotes
32 the softmax nonlinearity. For transformers, \mathbf{p} corresponds to the [CLS] token or tunable prompt [16]
33 whereas for RNN architectures [1], \mathbf{p} corresponds to the hidden state.

34 Given training data $(Y_i, X_i)_{i=1}^n$ with labels $Y_i \in \{-1, 1\}$ and inputs $X_i \in \mathbb{R}^{T \times d}$, we consider the empirical
35 risk minimization with a decreasing loss function $\ell(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$,

$$\mathcal{L}(\mathbf{v}, \mathbf{p}, W) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i \cdot f(X_i)).$$

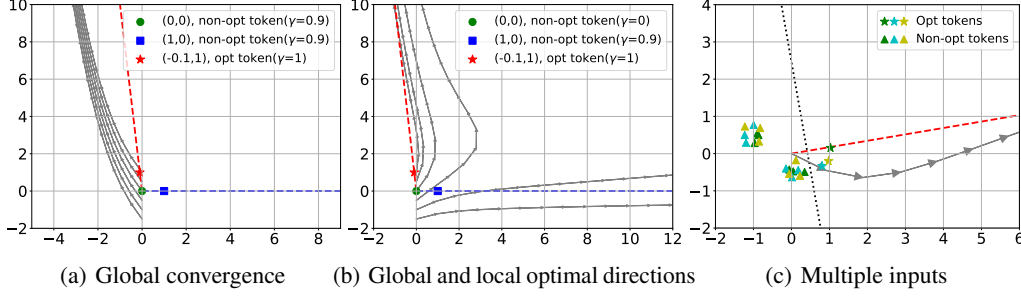


Figure 1: The convergence behavior of the gradient descent on the attention weights \mathbf{p} using the logistic loss in (ERM). Here, (---) and (---) denote the global and local max-margin solutions. γ denotes the *score* of a token per Definition 1. Discussion is provided under Theorem 1.

We operate under the assumption that the most relevant tokens within each input are separable from the rest through softmax nonlinearity. Our main contributions are as follows:

• **Optimize \mathbf{p} or \mathbf{W} for fixed \mathbf{v} (Sec. 2):** We first prove that gradient iterations of \mathbf{p} and \mathbf{W} admit a one-to-one mapping, thus we focus on optimizing \mathbf{p} without losing generality. We prove that gradient iterates of \mathbf{p} converges to a max-margin solution (namely (ATT-SVM)) that separates locally-optimal tokens from non-optimal ones. The notion of *relevant tokens* is clearly quantified in terms of scores $\gamma_t = \mathbf{Y} \cdot \mathbf{v}^\top \mathbf{x}_t$ where \mathbf{x}_t is the t 'th token of input \mathbf{X} . The *locally-optimal* tokens are those with higher scores than their nearest neighbors determined by the SVM solution. These are illustrated in Figure 1.

• **Optimize (\mathbf{v}, \mathbf{p}) jointly (Sec. 3):** We study the joint problem under logistic loss function. We use regularization path analysis where (ERM) is solved under ridge constraints and we study the solution trajectory as the constraints are relaxed. Since the problem is linear in \mathbf{v} , if the attention features $\mathbf{x}_i^{\text{att}} = \mathbb{S}(\mathbf{X}_i \mathbf{W}^\top \mathbf{p})$ are separable based on their labels Y_i , \mathbf{v} would implement a max-margin classifier. Building on this, we prove that \mathbf{p} and \mathbf{v} converges to their respective max-margin solutions under certain margin conditions. Relaxing these conditions, we obtain a more general solution where margin constraints on \mathbf{p} are relaxed on the inputs whose attention features are not support vectors of \mathbf{v} . Figure 2 illustrates these outcomes.

In Sec. 4, we extend the ideas in Sec. 2 to the more general model $f(\mathbf{X}) = \psi(\mathbf{X}^\top \mathbb{S}(\mathbf{X} \mathbf{W}^\top \mathbf{p}))$ with nonlinear head ψ . Overall, our results clearly formalize the role of the attention mechanism as a token-selection/context-discovery mechanism and lay the groundwork for future research by connecting it to the implicit bias literature and max-margin SVM formulation.

Next section introduces preliminaries, Section 5 discusses related literature, and Section 6 provides a discussion of limitations and future work.

1.1 Preliminaries

Notations. For any integer $N \geq 1$, let $[N] = \{1, \dots, N\}$. We use lower-case and upper-case bold letters (e.g. \mathbf{a} and \mathbf{A}) to represent vectors and matrices, respectively. The entries of \mathbf{a} are denoted as a_i . We use $\sigma_{\max}(\mathbf{A})$ and $\sigma_{\min}(\mathbf{A})$ to denote the maximum and minimum singular values of \mathbf{A} , respectively. We denote the minimum of two numbers a, b as $a \wedge b$, and the maximum $a \vee b$. We use the standard big-Oh notation $O(\cdot)$ to hide universal constants.

Optimization. Given a function $\mathcal{L} : \mathbb{R}^d \rightarrow \mathbb{R}$ and an ℓ_2 -norm bound R , the regularized solution is defined as

$$\bar{\mathbf{p}}(R) := \arg \min_{\|\mathbf{p}\| \leq R} \mathcal{L}(\mathbf{p}). \quad (1)$$

Note that $\bar{\mathbf{p}}(R)$ is not unique in general. For gradient descent, we assume the objective $\mathcal{L}(\mathbf{p})$ is smooth and describe the gradient descent process as

$$\mathbf{p}(t+1) = \mathbf{p}(t) - \eta(t) \nabla \mathcal{L}(\mathbf{p}(t)), \quad (2)$$

where $\eta(t)$ is the stepsize at time t and $\nabla \mathcal{L}(\mathbf{p}(t))$ is the gradient of \mathcal{L} at $\mathbf{p}(t)$.

Attention in Transformers. We now describe how our model relates to the attention mechanism in transformers. Our exposition follows the recent work [17] which focuses on the theoretical properties of prompt-tuning.

72 • **Self-attention** is the core building block of transformers [6]. Given an input consisting of T
73 tokens $X = [\mathbf{x}_1, \dots, \mathbf{x}_T]^\top \in \mathbb{R}^{T \times d}$, self-attention with key-query matrix $W \in \mathbb{R}^{d \times d}$, and value matrix
74 $V \in \mathbb{R}^{d \times v}$, the self-attention model is defined as follows:

$$f_{sa}(X) = \mathbb{S}(XWX^\top)XV. \quad (3)$$

75 Here, $\mathbb{S}(\cdot)$ is the softmax nonlinearity that applies row-wise on the similarity matrix XWX^\top .

76 • **Tunable tokens: [CLS] and prompt-tuning.** In practice, we append additional tokens to the raw
77 input features X : For instance, a [CLS] token is used for classification purposes [7] and prompt
78 vectors can be appended for adapting a pretrained model to new tasks [16, 18]. Let $\mathbf{p} \in \mathbb{R}^d$ be the
79 tunable token ([CLS] or prompt vector) and concatenate it to X to obtain $X_p := [\mathbf{p} X^\top]^\top \in \mathbb{R}^{(T+1) \times d}$.
80 Consider the cross-attention features obtained from X_p and X given by

$$\begin{bmatrix} f_{cls}^\top(X) \\ f_{sa}^\top(X) \end{bmatrix} = \mathbb{S}(X_p W X^\top) X V = \begin{bmatrix} \mathbb{S}(\mathbf{p}^\top W X^\top) \\ \mathbb{S}(X W X^\top) \end{bmatrix} X V.$$

81 The beauty of cross-attention is that it isolates the contribution of \mathbf{p} under the upper term $f_{cls}(X) =$
82 $V^\top X^\top \mathbb{S}(X W^\top \mathbf{p}) \in \mathbb{R}^v$. In this work, we use the value weights for classification, thus we set $v = 1$,
83 and denote $\mathbf{v} = V \in \mathbb{R}^d$. This brings us to our attention model of interest:

$$f(X, \Theta) = \mathbf{v}^\top X^\top \mathbb{S}(K \mathbf{p}) \quad \text{where} \quad K = X W^\top. \quad (4)$$

84 Above $\Theta = (\mathbf{v}, W, \mathbf{p})$ are the tunable model parameters and K is the key embeddings. Note that
85 W and \mathbf{p} are playing the same role within softmax, thus, it is intuitive that they exhibit similar
86 optimization dynamics. Confirming this, the next lemma shows that gradient iterations of \mathbf{p} (after
87 setting $W \leftarrow \text{Identity}$) and W admit a one-to-one mapping.

88 **Lemma 1** Fix $\mathbf{a} \in \mathbb{R}^d$. Let $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$ and $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ be differentiable functions. On the same
89 training data, define $\mathcal{L}(\mathbf{p}) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, \psi(X_i^\top \mathbb{S}(X_i \mathbf{p})))$ and $\mathcal{L}(W) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, \psi(X_i^\top \mathbb{S}(X_i W^\top \mathbf{a})))$.
90 Starting from $\mathbf{p}(0)$ and $W(0) = \mathbf{a} \mathbf{p}(0)^\top / \|\mathbf{a}\|^2$ consider the gradient descent iterations with stepsize η :

$$\begin{aligned} \mathbf{p}(t+1) &= \mathbf{p}(t) - \eta \nabla \mathcal{L}(\mathbf{p}_t), \\ W(t+1) &= W(t) - \eta \|\mathbf{a}\|^{-2} \nabla \mathcal{L}(W(t)). \end{aligned}$$

91 We have that $W_t = \mathbf{a} \mathbf{p}_t^\top / \|\mathbf{a}\|^2$ for all $t \geq 0$.

92 Thanks to this lemma, W 's optimization dynamics is directly characterized by \mathbf{p} 's dynamics, since we
93 can always reconstruct W from \mathbf{p} using the relationship between their gradient iterations. Hence, in
94 what follows, we fix W , and focus on optimizing \mathbf{p} in Sec 2 and joint optimization of (\mathbf{v}, \mathbf{p}) in Sec 3.
95

Problem definition: Throughout, $(Y_i, X_i)_{i=1}^n$ denotes our training dataset where $Y_i \in \{-1, 1\}$ and
 $X_i \in \mathbb{R}^{T \times d}$. We denote the key embeddings of X_i via $K_i = X_i W^\top$ and explore the training risk

$$\mathcal{L}(\mathbf{v}, \mathbf{p}) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i \cdot \mathbf{v}^\top X_i^\top \mathbb{S}(K_i \mathbf{p})). \quad (\text{ERM})$$

Importantly, our results apply to general tuples (Y_i, X_i, K_i) and do not assume that (X_i, K_i) are
tied via W . Finally, the t^{th} tokens of X_i, K_i are denoted by $\mathbf{x}_{it}, \mathbf{k}_{it} \in \mathbb{R}^d$ respectively for $t \in [T]$.

97 2 Global and Local Margin Maximization with Attention

98 In this section, we establish the main results of this paper (Theorems 1 and 3) by characterizing the
99 implicit bias of gradient descent on learning $\mathbf{p} \in \mathbb{R}^d$ for a fixed $\mathbf{v} \in \mathbb{R}^d$ choice. A notable feature of our
100 results is that they apply to general decreasing loss functions. The underlying reason is that margin
101 maximization arises from the exponentially-tailed nature of the softmax within attention rather than ℓ .
102 Throughout, we make the following assumption on the loss function:

103 **Assumption A (Well-Behaved Losses)** Over any bounded set in \mathbb{R} , $\ell : \mathbb{R} \rightarrow \mathbb{R}$ obeys

104 **A1.** ℓ is strictly decreasing and bounded from below.

105 **A2.** ℓ' is M_0 -Lipschitz continuous and $|\ell'(u)| \leq M_1$.

Assumption A includes many common loss functions, including the logistic loss $\ell(u) = \log(1 + e^{-u})$, exponential loss $\ell(u) = e^{-u}$, and correlation loss $\ell(u) = -u$. Assumption A implies that $\mathcal{L}(\mathbf{p})$ is a $O((M_0 + M_1)\bar{\sigma}_{\max}^4)$ -smooth function (see supplementary), where $\bar{\sigma}_{\max} := 1/n \sum_{i=1}^n \sigma_{\max}(\mathbf{X}_i)$. The central feature of this assumption is that *we do not require convexity for the loss function*.

We now introduce a convex hard-margin SVM problem that separates one token of the input sequence from the rest, jointly solved over all inputs. We will show that this problem captures the optimization properties of softmax-attention. Fix indices $\alpha = (\alpha_i)_{i=1}^n$ and consider

$$\mathbf{p}^{mm}(\alpha) = \arg \min_{\mathbf{p}} \|\mathbf{p}\| \quad \text{subject to} \quad \min_{i \neq \alpha_i} \mathbf{p}^\top (\mathbf{k}_{i\alpha_i} - \mathbf{k}_{it}) \geq 1 \quad \text{for all} \quad 1 \leq i \leq n. \quad (\text{ATT-SVM})$$

We are ready to introduce our main results by characterizing global and local convergence of the attention weights \mathbf{p} in the direction of (ATT-SVM) solutions.

2.1 Global Convergence of the Attention Weights \mathbf{p}

We first identify the conditions that guarantees the global convergence of gradient descent for \mathbf{p} . The intuition is that, in order for attention to exhibit implicit bias, the softmax nonlinearity should be forced to select the *optimal token within each input sequence*. Fortunately, the optimal tokens that achieve the smallest training objective under decreasing loss function $\ell(\cdot)$ have a clear definition.

Definition 1 (Token Score and Optimality) The score of token \mathbf{x}_{it} of input \mathbf{X}_i is defined as $\gamma_{it} := \mathbf{Y}_i \cdot \mathbf{v}^\top \mathbf{x}_{it}$. The optimal tokens for input \mathbf{X}_i are those tokens with highest scores given by

$$\text{opt}_i \in \arg \max_{i \in [T]} \gamma_{it}.$$

We denote the solution of (ATT-SVM) with optimal indices $(\text{opt}_i)_{i=1}^n$ by \mathbf{p}^{mm*} . Note that multiple tokens within an \mathbf{X}_i might attain same score, thus opt_i and \mathbf{p}^{mm*} may not be unique.

To proceed with our global convergence analysis, we need to make the assumption that all non-optimal tokens have equal scores. In other words, if a potential solution includes tokens that do not appear in the final optimal solution, all of these tokens are assumed to have the same score value.

Assumption B For all $i \in [n]$ and $t, \tau \neq \text{opt}_i$, the scores per Def. 1 obey $\gamma_{it} = \gamma_{i\tau} < \gamma_{i\text{opt}_i}$.

Theorem 1 (Global Convergence of Gradient Descent) Suppose Assumption A on the loss function ℓ and Assumption B on the tokens' score hold. Then the gradient descent iterates $\mathbf{p}(t+1) = \mathbf{p}(t) - \eta \nabla \mathcal{L}_{\mathbf{p}}(\mathbf{p}(t))$ on (ERM), with the step size $\eta \leq O(\bar{\sigma}_{\max}^{-4}/(M_0 + M_1))$ and any starting point $\mathbf{p}(0)$ satisfies $\lim_{t \rightarrow \infty} \mathbf{p}(t)/\|\mathbf{p}(t)\| = \mathbf{p}^{mm*}/\|\mathbf{p}^{mm*}\|$.

Theorem 1 shows that gradient descent dynamics of the normalized predictor $\mathbf{p}(t)/\|\mathbf{p}(t)\|$ converges towards $\mathbf{p}^{mm*}/\|\mathbf{p}^{mm*}\|$, effectively separating globally optimal tokens from non-optimal ones. To illustrate this theorem, we have conducted synthetic experiments. Let us first explain the setup used in Figure 1 and 2(a). We set $d = 3$ with each token having three entries $\mathbf{x} = [x_1, x_2, x_3]$. We reserve the first two coordinates as key embeddings $\mathbf{k} = [x_1, x_2, 0]$ by setting $\mathbf{W} = \text{diag}([1, 1, 0])$. This is what we display in our figures as token positions. Finally, in order to assign scores to the tokens we use the last coordinate by setting $\mathbf{v} = [0, 0, 1]$. This way score becomes $\mathbf{Y} \cdot \mathbf{v}^\top \mathbf{x} = \mathbf{Y} \cdot x_3$, allowing us to assign any score (regardless of key embedding).

In Figure 1(a), the gray paths represent gradient descent trajectories initiated from different points, while the points (0, 0) and (1, 0) correspond to non-optimal tokens, and (-0.1, 1) represents the optimal token. Notably, gradient descent iterates with various starting points converge towards the direction of the max-margin solution \mathbf{p}^{mm*} (depicted by - - -). Moreover, as the iteration count t increases, the inner product $\langle \mathbf{p}(t)/\|\mathbf{p}(t)\|, \mathbf{p}^{mm*}/\|\mathbf{p}^{mm*}\| \rangle$ consistently increases. Figure 1(c) also illustrates the result of Theorem 1 on multiple inputs (gray dot line is the separating hyperplane). These observations emphasize the gradual alignment between the evolving predictor and the max-margin solution throughout optimization.

Transient optimization dynamics and the role of loss function. While asymptotic direction of gradient descent is determined by \mathbf{p}^{mm} , intuitively transient dynamics can exhibit bias towards tokens with extreme scores. We aim to capture this intuition in Figure 2(a) which depicts the gradient

trajectories for different scores and loss functions. We have two optimal tokens (\star) with scores $\gamma_1 = 1, \gamma_2 = C$ for varying C and we consider correlation loss $\ell(x) = -x$ and exponential loss $\ell(x) = e^{-x}$. In a nutshell, as C grows, it can be seen that $\ell(x) = -x$ is biased towards token with high-score whereas $\ell(x) = e^{-x}$ is biased towards the low-score token. The underlying reason can be seen from the gradient of individual inputs: $\nabla \mathcal{L}_i(\mathbf{p}) = \ell'_i \cdot \mathbf{K}_i^\top \mathbb{S}'(\mathbf{X}\mathbf{p})\mathbf{X}\mathbf{v}$ where $\mathbb{S}'(\cdot)$ is the softmax derivative and $\ell'_i := \ell'(Y_i \cdot \mathbf{v}^\top \mathbf{X}_i^\top \mathbb{S}(\mathbf{X}_i\mathbf{p}))$. Assuming \mathbf{p} (approximately) selects the optimal tokens, this would simplify to $\ell'_i \approx \ell'(\gamma_i)$ and $\|\nabla \mathcal{L}_i(\mathbf{p})\| \propto |\ell'(\gamma_i)| \cdot \gamma_i$. Now, with correlation loss $|\ell'| = 1$, thus, $\|\nabla \mathcal{L}_i(\mathbf{p})\| \propto \gamma_i$ and larger score induces larger gradient. Whereas with exponential loss $|\ell'| = e^{-u}$, thus, $\|\nabla \mathcal{L}_i(\mathbf{p})\| \propto \gamma_i e^{-\gamma_i}$ and smaller score induces larger gradient explaining the empirical behavior.

We next provide the regularization path analysis that requires relaxed assumptions on both loss function and tokens' score.

Theorem 2 (Regularization Path) *Suppose Assumption A on the loss function holds, and for all $i \in [n]$ and $t \neq \text{opt}_i$, scores obey $\gamma_{it} < \gamma_{i\text{opt}_i}$. Then the regularization path $\bar{\mathbf{p}}(R) = \arg \min_{\|\mathbf{p}\| \leq R} \mathcal{L}(\mathbf{p})$ satisfies $\lim_{R \rightarrow \infty} \bar{\mathbf{p}}(R)/R = \mathbf{p}^{\text{mm}\star} / \|\mathbf{p}^{\text{mm}\star}\|$.*

Theorem 2 reveals that as we loosen the regularization strength R to achieve ridgeless optimization with $\min_{\mathbf{p}} \mathcal{L}(\mathbf{p})$, the optimal direction $\bar{\mathbf{p}}(R)$ gradually aligns with the max-margin solution $\mathbf{p}^{\text{mm}\star}$. A central feature of this theorem is its ability to handle non-optimal tokens that possess different arbitrary scores. Thus, it demonstrates that max-margin convergence is a global feature of attention mechanism. As we shall see in the next section, due to nonconvex landscape and nonlinearity of softmax, convergence of regularization path without Assumption B does not imply that Theorem 1 can avoid this condition.

2.2 Local convergence of the attention weights \mathbf{p}

Theorem 1 on the global convergence of gradient descent serves as a prelude to the general behavior of the optimization. Once we relax Assumption B by allowing for arbitrary token scores, we will show that \mathbf{p} can converge (in direction) to a locally-optimal solution. However, this locally-optimal solution is still characterized in terms of (ATT-SVM) which separates *locally-optimal* tokens from the rest. Our theory builds on two new concepts: locally-optimal tokens and neighbors of these tokens.

Definition 2 (SVM-Neighbor and Locally-Optimal Tokens) *Fix token indices $\alpha = (\alpha_i)_{i=1}^n$. Solve (ATT-SVM) to obtain $\mathbf{p}^{\text{mm}} = \mathbf{p}^{\text{mm}}(\alpha)$. Consider tokens $\mathcal{T}_i \subset [T]$ such that $(\mathbf{k}_{i\alpha_i} - \mathbf{k}_{it})^\top \mathbf{p}^{\text{mm}} = 1$ for all $t \in \mathcal{T}_i$. We refer to \mathcal{T}_i as SVM-neighbors of $\mathbf{k}_{i\alpha_i}$. Additionally, tokens $\alpha = (\alpha_i)_{i=1}^n$ are called locally-optimal if for all $i \in [n]$, $t \in \mathcal{T}_i$ scores per Def. 1 obey $\gamma_{i\alpha_i} > \gamma_{it}$.*

To provide a basis for discussing local convergence, we establish a cone centered around \mathbf{p}^{mm} using the following construction. Let μ be a positive scalar, and define the cone as:

$$\text{cone}_\mu(\mathbf{p}^{\text{mm}}) := \left\{ \mathbf{p} \in \mathbb{R}^d \mid \left\langle \frac{\mathbf{p}}{\|\mathbf{p}\|}, \frac{\mathbf{p}^{\text{mm}}}{\|\mathbf{p}^{\text{mm}}\|} \right\rangle \geq 1 - \mu \right\}. \quad (5)$$

In the subsequent theorem, we demonstrate the existence of a scalar $\mu = \mu(\alpha) > 0$ and a radius R such that when R is sufficiently large, there are no stationary points within the intersection of $\text{cone}_\mu(\mathbf{p}^{\text{mm}})$ and the set $\{\mathbf{p} \mid \|\mathbf{p}\| \geq R\}$. Further, the gradient descent initialized within this intersection converges in direction to $\mathbf{p}^{\text{mm}}/\|\mathbf{p}^{\text{mm}}\|$.

Theorem 3 (Local Convergence of Gradient Descent) *Suppose Assumption A on the loss function ℓ holds and assume $\alpha = (\alpha_i)_{i=1}^n$ are locally-optimal tokens per Definition 2. Then, there exists a scalar $\mu = \mu(\alpha) \in (0, 1)$ and a radius $R > 0$ such that $\text{cone}_\mu(\mathbf{p}^{\text{mm}}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| \geq R\}$ does not contain any stationary points. Further, the gradient descent iterates $\mathbf{p}(t+1) = \mathbf{p}(t) - \eta \nabla \mathcal{L}(\mathbf{p}(t))$ on (ERM) with*

$$\eta \leq O\left(\min\left(\frac{1}{(M_0 + M_1)\bar{\sigma}_{\max}^4}, \frac{\mu - \epsilon}{(1 - \mu)}\right)\right), \quad (6)$$

for any $\epsilon \in (0, \min(\mu, 1))$, and any starting point $\mathbf{p}(0) \in \text{cone}_\mu(\mathbf{p}^{\text{mm}}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| \geq R\}$ satisfies $\lim_{t \rightarrow \infty} \mathbf{p}(t)/\|\mathbf{p}(t)\| = \mathbf{p}^{\text{mm}}/\|\mathbf{p}^{\text{mm}}\|$.

Proof sketch. We provide the proof in four steps:

Step 1. We begin by proving that there are no stationary points within $\text{cone}_\mu(\mathbf{p}^{\text{mm}}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| \geq R_\mu\}$

for a specific radius R_μ . Let $(\mathcal{T}_i)_{i=1}^n$ denote the set of SVM-neighbors as defined in Definition 2. We define $\tilde{\mathcal{T}}_i = [T] - \mathcal{T}_i - \alpha_i$ as the tokens that are non-SVM neighbors. Additionally, let

$$\delta := \frac{1}{2} \min_{i \in [n]} \min_{t \in \mathcal{T}_i, \tau \in \tilde{\mathcal{T}}_i} (\mathbf{k}_{it} - \mathbf{k}_{i\tau})^\top \mathbf{p}^{mm}, \quad A := \max_{i \in [n], t \in [T]} \|\mathbf{k}_{it}\| \cdot \|\mathbf{p}^{mm}\|, \quad \mu := \frac{1}{8} \left(\frac{\min(0.5, \delta)}{A} \right)^2.$$

For all $\mathbf{q}, \mathbf{p} \in \text{cone}_\mu(\mathbf{p}^{mm})$ with $\|\mathbf{q}\| = \|\mathbf{p}^{mm}\|$, we establish the existence of R_μ such that $-\mathbf{q}^\top \nabla \mathcal{L}(\mathbf{p})$ is strictly positive for $\|\mathbf{p}\| \geq R_\mu$. Specifically, we show the existence of positive constants C and c satisfying:

$$C \cdot \max_{i \in [n]} q_i \geq -\langle \nabla \mathcal{L}(\mathbf{p}), \mathbf{q} \rangle \geq c \cdot \min_{i \in [n]} q_i > 0.$$

Here, $q_i = 1 - \mathbb{S}(\mathbf{K}_i \mathbf{p})_{\alpha_i}$ and $\alpha = (\alpha_i)_{i=1}^n$ are locally-optimal tokens per Definition 2.

Step 2. We demonstrate that for any $\epsilon \in (0, \min(\mu, 1))$, there exists R_ϵ such that all $\mathbf{p} \in \text{cone}_\mu(\mathbf{p}^{mm}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| \geq R_\epsilon\}$ satisfy

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

Step 3. By leveraging the results from Step 1 and Step 2, we can demonstrate that the gradient iterates, with an appropriate step size, starting from $\mathbf{p}(0) \in \text{cone}_\mu(\mathbf{p}^{mm}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| \geq R\}$, remain within this cone. Specifically, if $\mathbf{p}(t) \in \text{cone}_\mu(\mathbf{p}^{mm}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| \geq R\}$, then $\|\mathbf{p}(t+1)\| \geq \|\mathbf{p}(t)\|$, and

$$\left\langle \frac{\mathbf{p}(t+1)}{\|\mathbf{p}(t+1)\|}, \frac{\mathbf{p}^{mm}}{\|\mathbf{p}^{mm}\|} \right\rangle \geq 1 - \mu + O(\eta(\mu - \epsilon) - \eta^2(1 - \mu)),$$

which implies that, with the step size η satisfying (6), $\mathbf{p}(t+1) \in \text{cone}_\mu(\mathbf{p}^{mm}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| \geq R\}$.

Step 4. The remaining part of the proof follows the same reasoning as the proof of Theorem 1 and is provided in the supplementary material. ■

To further illustrate Theorem 3, we can consider Figure 1(b) where $n = 1$ and $T = 3$. In this figure, the point $(0, 0)$ represents the non-optimal tokens, while $(1, 0)$ represents the locally optimal token. Additionally, the gray paths represent the trajectories of gradient descent initiated from different points. By observing the figure, we can see that gradient descent, when properly initialized, converges towards the direction of \mathbf{p}^{mm} (depicted by - - -). This direction of convergence effectively separates the locally optimal tokens $(1, 0)$ from the non-optimal token $(0, 0)$.

2.3 Tightness of the locally-optimal token definition

An important question is whether our definition of locally-optimal tokens (Def. 2) covers all token configurations $\alpha = (\alpha_i)_{i=1}^n$ that can be selected by the attention mechanism asymptotically (as $\|\mathbf{p}\| \rightarrow \infty$). The following theorem essentially establishes the tightness of our definition: It shows that, given $\alpha = (\alpha_i)_{i=1}^n$, if any of the α_i 's have an SVM-neighbor with a higher score, then regularization path will not prefer the $\mathbf{p}^{mm}(\alpha)$ direction.

Theorem 4 Fix indices $\alpha = (\alpha_i)_{i=1}^n$ with SVM-neighbors $(\mathcal{T}_i)_{i=1}^n$. Set $\mathbf{p}^{mm} := \mathbf{p}^{mm}(\alpha)$. Suppose that:

- For some $j \in [n]$, there exists $\beta \in \mathcal{T}_j$ with higher score than α_j , i.e., $Y_j \cdot \mathbf{v}^\top \mathbf{x}_{j\beta} > Y_j \cdot \mathbf{v}^\top \mathbf{x}_{j\alpha_j}$.
- For all $i \in [n]$ and $t \in \mathcal{T}_i$, the vectors $\mathbf{k}_{i\alpha_i} - \mathbf{k}_{it}$ are linearly independent (We note that this holds for almost all datasets).

For any $\epsilon > 0$, there exists $R_\epsilon > 0$ as follows: Consider the neighborhood of \mathbf{p}^{mm} : $C_\epsilon = \text{cone}_\epsilon(\mathbf{p}^{mm}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| \geq R_\epsilon\}$. Define the local path $\tilde{\mathbf{p}}(R) = \min_{\mathbf{p} \in C_\epsilon, \|\mathbf{p}\| \leq R} \mathcal{L}(\mathbf{p})$. Then $\lim_{R \rightarrow \infty} \frac{\tilde{\mathbf{p}}(R)}{\|\tilde{\mathbf{p}}(R)\|} \neq \frac{\mathbf{p}^{mm}}{\|\mathbf{p}^{mm}\|}$.

3 Joint Convergence of Head \mathbf{v} and Attention Weights \mathbf{p}

In this section, we extend the preceding results to the general case of joint optimization of head \mathbf{v} and attention weights \mathbf{p} using a logistic loss function. To this aim, we focus on regularization path analysis, which involves solving (ERM) under ridge constraints and examining the solution trajectory as the constraints are relaxed.

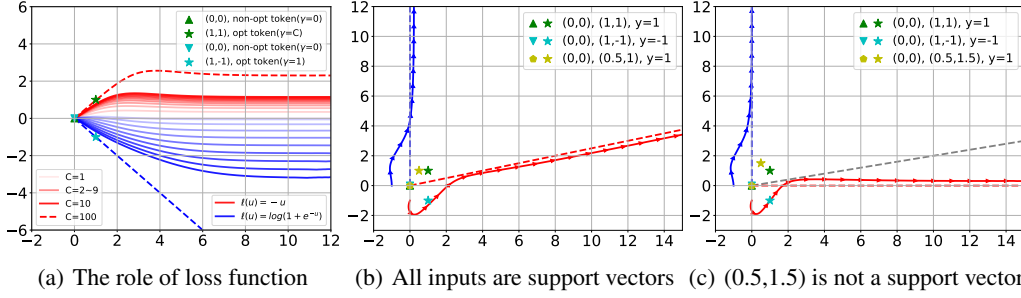


Figure 2: (a) Global convergence of p with different loss functions and scores. (b)&(c) Joint convergence of attention weights p and classifier head v to max-margin directions.

High-level intuition. Since the prediction is linear as a function of v , logistic regression in v can exhibit its own implicit bias to a max-margin solution. Concretely, define the attention features $x_i^p = X_i^T \mathbb{S}(K_i p)$ and define the dataset $\mathcal{S}^p = (Y_i, x_i^p)$. If the dataset \mathcal{S}^p is separable by v , then optimizing only v will converge in the direction of the max-margin classifier by setting $r_i \leftarrow x_i^p$:

$$v^{mm} = \arg \min_{v \in \mathbb{R}^d} \|v\| \quad \text{subject to} \quad Y_i \cdot v^T r_i \geq 1. \quad (\text{CLS-SVM})$$

This motivates a clear question: *Under what conditions, optimizing v, p jointly will converge to their respective max-margin solutions?* We study this question in two steps. Loosely speaking, we will first assume that when solving (CLS-SVM), all inputs $i \in [n]$ are also the support vectors. We will then relax this condition to uncover a more general implicit bias for p . Throughout we assume that the joint problem is separable and there exists (v, p) asymptotically achieving zero training loss.

3.1 When all attention features are support vectors

In (CLS-SVM), define *label margin* to be $1/\|v^{mm}\|$. Our first insight in quantifying joint implicit bias is that, optimal tokens admit a natural definition: Those that maximize the downstream label margin when selected. This is formalized below where we assume that: (1) selecting the token indices $\alpha = (\alpha_i)_{i=1}^n$ from each input data achieves the largest label margin. (2) The optimality of the α choice is strict in the sense that mixing other tokens will shrink the label margin in (CLS-SVM).

Assumption C Let $\Gamma > 0$ be the label margin when solving (CLS-SVM) with $r_i \leftarrow x_{i\alpha_i}$. There exists $\nu > 0$ such that for all p , solving (CLS-SVM) with $r_i \leftarrow x_i^p$ results in a label margin of at most $\Gamma - \nu \cdot \max_{i \in [n]} (1 - s_{i\alpha_i})$ where $s_i = \mathbb{S}(K_i p)$.

Example: To gain intuition, let us fix $v_\star \in \mathbb{R}^d$ and consider the dataset obeying $x_{i1} = Y_i \cdot v_\star$ and $\|x_{it}\| < \|v_\star\|$ for all $t \geq 2$ and all $i \in [n]$. For this dataset, we can choose $\alpha_i = 1$, $\Gamma = \|v_\star\|$ and $\nu = \|v_\star\| - \sup_{i \in [n], t \geq 2} \|x_{it}\|$.

Theorem 5 Consider the ridge-constrained solutions (v_r, p_R) of (ERM) defined as

$$v_r, p_R = \arg \min_{\|v\| \leq r, \|p\| \leq R} \mathcal{L}(v, p).$$

Suppose Assumption C holds for some $\Gamma, \nu > 0$. As $r, R \rightarrow \infty$, the joint regularization path (v_r, p_R) converges as follows: $\frac{p_R}{R} \rightarrow \frac{p^{mm}}{\|p^{mm}\|}$ where p^{mm} is the solution of (ATT-SVM). $\frac{v_r}{r} \rightarrow \frac{v^{mm}}{\|v^{mm}\|}$ where v^{mm} is the solution of (CLS-SVM) with $r_i = x_{i1}$.

As further discussion, consider Figure 2(b) where we set $n = 3, T = d = 2$ and $W = \text{Identity}$. All three inputs share the point (0, 0) which corresponds to their non-optimal tokens. The optimal tokens (denoted by \star) are all support vectors of the (CLS-SVM) since $v^{mm} = [0, 1]$ is the optimal classifier direction (blue color). Because of this, p^{mm} will separate optimal tokens from (0, 0) coordinate via (ATT-SVM) which results in the red direction (yellow and teal \star are the support tokens).

3.2 General solution when selecting one token per input

Can we relax Assumption C, and if so, what is the resulting behavior? Consider the scenario where the optimal p diverges to ∞ and ends up selecting one token per input. Suppose this p selects some

coordinates $\alpha = (\alpha_i)_{i=1}^n$. Let $\mathcal{N} \subset [n]$ be the set of indices where the associated token $\mathbf{x}_{i\alpha_i}$ is **not** a support vector when solving (CLS-SVM). Our intuition is as follows: Even if we slightly perturb this \mathbf{p} choice and mix other tokens $t \neq \alpha_i$ over the input set $\mathcal{N} \subset [n]$, since \mathcal{N} is not support vector for (CLS-SVM), we can preserve the label margin (by only preserving the support vectors $[n] - \mathcal{N}$). This means that \mathbf{p} may not have to enforce *max-margin* constraint over inputs $i \in \mathcal{N}$, instead, it suffices to just select these tokens (asymptotically). This results in the following **relaxed SVM** problem:

$$\mathbf{p}^{\text{relax}} = \min_{\mathbf{p}} \|\mathbf{p}\| \quad \text{such that} \quad \mathbf{p}^\top (\mathbf{k}_{i\alpha_i} - \mathbf{k}_{it}) \geq \begin{cases} 1 & \text{for all } t \neq \alpha_i, i \in [n] - \mathcal{N} \\ 0 & \text{for all } t \neq \alpha_i, i \in \mathcal{N} \end{cases}. \quad (7)$$

Here, $\mathbf{p}^\top \mathbf{x}_{i\alpha_i} \geq 0$ corresponds to the *selection* idea. Building on this intuition, the following theorem captures the generalized behavior of the joint regularization path.

Theorem 6 Consider the path of $(\mathbf{v}_r, \mathbf{p}_R)$ as $r, R \rightarrow \infty$ as in Theorem 5. Suppose $\mathbb{S}(\mathbf{K}_i \mathbf{p}_R)_{\alpha_i} \rightarrow 1$, i.e., the tokens $(\alpha_i)_{i=1}^n$ are asymptotically selected. Let \mathbf{v}^{mm} be the solution of (CLS-SVM) with $\mathbf{r}_i = \mathbf{x}_{i\alpha_i}$ and \mathcal{N} be its set of non-support indices. Suppose Assumption C holds over the support vectors $[n] - \mathcal{N}$. Then, $\frac{\mathbf{v}_r}{r} \rightarrow \frac{\mathbf{v}^{\text{mm}}}{\|\mathbf{v}^{\text{mm}}\|}$ and $\frac{\mathbf{p}_R}{R} \rightarrow \frac{\mathbf{p}^{\text{relax}}}{\|\mathbf{p}^{\text{relax}}\|}$ where $\mathbf{p}^{\text{relax}}$ is the solution of (7) with α_i choices.

To illustrate this numerically, consider Figure 2(c) which modifies Figure 2(b) by pushing the yellow \star to the northern position (0.5, 1.5). We still have $\mathbf{v}^{\text{mm}} = [0, 1]$ however the yellow \star is no longer a support vector of (CLS-SVM). Thus, \mathbf{p} solves the relaxed problem which separates green and teal \star 's by enforcing the max-margin constraint on \mathbf{p} (which is the red direction). Instead, yellow \star only needs to achieve positive correlation with \mathbf{p} (unlike Figure 2(c) where it dictates the direction).

4 Regularization Path of Attention with Nonlinear Head

So far our discussion has focused on the attention model with linear head. However, the conceptual ideas on optimal token selection via margin maximization also extends to a general nonlinear model under mild assumptions. The aim of this section is showcasing this generalization. Specifically, we consider the prediction model $f(\mathbf{X}) = \psi(\mathbf{X}^\top \mathbb{S}(\mathbf{K}\mathbf{p}))$ where $\psi(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ generalizes the linear head \mathbf{v} of our attention model. For instance, following exposition in Section 1.1, $\psi(\cdot)$ can represent a multilayer transformer with \mathbf{p} being a tunable prompt at the input layer. Recall that $\mathcal{S} = (\mathbf{X}_i, \mathbf{K}_i, \mathbf{Y}_i)_{i=1}^n$ is the dataset of the input-key-label tuples. We consider the training risk

$$\mathcal{L}(\mathbf{p}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{Y}_i, \psi(\mathbf{X}_i^\top \mathbf{s}_i^{\mathbf{p}})) \quad \text{where} \quad \mathbf{s}_i^{\mathbf{p}} = \mathbb{S}(\mathbf{K}_i \mathbf{p}) \in \mathbb{R}^T. \quad (8)$$

The challenge with nonlinear $\psi(\cdot)$ is that, we lack a clear score function (Def. 1) unlike the previous sections. The assumption below introduces a generic condition that splits the tokens of each \mathbf{X}_i into an *optimal* set \mathcal{O}_i and *non-optimal* set $\bar{\mathcal{O}}_i = [T] - \mathcal{O}_i$. In words, non-optimal tokens are those that strictly increase the training risk $\mathcal{L}(\mathbf{p})$ if they are not fully suppressed by attention probabilities $\mathbf{s}_i^{\mathbf{p}}$.

Assumption D (Mixing non-optimal tokens hurt) There exists sets $(\mathcal{O}_i)_{i=1}^n \subset [T]$ as follows. Let $q_i^{\mathbf{p}} = \sum_{t \in \bar{\mathcal{O}}_i} s_{it}^{\mathbf{p}}$ be the sum of softmax similarities over the non-optimal set for \mathbf{p} . Set $q_{\max}^{\mathbf{p}} = \max_{i \in [n]} q_i^{\mathbf{p}}$. For any $\Delta > 0$, there exists $\rho < 0$ such that:

$$\text{For all } \mathbf{p}, \mathbf{p}' \in \mathbb{R}^d, \text{ if } \log(q_{\max}^{\mathbf{p}}) \leq (1 + \Delta) \log(q_{\max}^{\mathbf{p}'}) \wedge \rho, \text{ then } \mathcal{L}(\mathbf{p}) < \mathcal{L}(\mathbf{p}').$$

This assumption is titled *mixing hurts* because the attention output $\mathbf{X}_i^\top \mathbf{s}_i^{\mathbf{p}}$ is mixing the tokens of \mathbf{X}_i and our condition is that, to achieve optimal risk, this mixture should not contain any non-optimal tokens. In particular, we require that, a model \mathbf{p} that contains *exponentially less non-optimality* (quantified via $\log(q_{\max})$) compared to \mathbf{p}' is strictly preferable. As we discuss in the supplementary material, Theorem 2 is in fact a concrete instance (with linear head \mathbf{v}) satisfying this condition.

Before stating our generic theorem, we need to introduce the max-margin separator towards which regularization path of attention will converge. This is a slightly general version of Section 2's (ATT-SVM) problem where we allow for a set of optimal tokens \mathcal{O}_i for each input.

$$\mathbf{p}^{\text{mm}} = \arg \min_{\mathbf{p}} \|\mathbf{p}\| \quad \text{subject to} \quad \max_{\alpha \in \mathcal{O}_i} \min_{\beta \in \bar{\mathcal{O}}_i} \mathbf{p}^\top (\mathbf{k}_{i\alpha} - \mathbf{k}_{i\beta}) \geq 1 \quad \text{for all } i \in [n]. \quad (\text{ATT-SVM}')$$

Unlike (ATT-SVM), this problem is not necessarily convex when the optimal set \mathcal{O}_i is not a singleton. To see this, imagine $n = d = 1$ and $T = 3$: Set the two optimal tokens as $\mathbf{k}_1 = 1$ and $\mathbf{k}_2 = -1$ and the

non-optimal token as $k_3 = 0$. The solution set of (ATT-SVM') is $\mathbf{p}^{mm} \in \{-1, 1\}$ whereas their convex combination $\mathbf{p} = 0$ violates the constraints. To proceed, our final result establishes the convergence of regularization path to the solution set of (ATT-SVM') under Assumption D.

Theorem 7 Let \mathcal{P}^{mm} be the set of global minima of (ATT-SVM'). Suppose its objective $\Gamma := \|\mathbf{p}^{mm}\|$ is finite and Assumption D holds. Let $\text{dist}(\cdot, \cdot)$ denote the ℓ_2 -distance between a vector and a set. Following (8), define $\bar{\mathbf{p}}(R) = \arg \min_{\|\mathbf{p}\| \leq R} \mathcal{L}(\mathbf{p})$. We have that $\lim_{R \rightarrow \infty} \text{dist}\left(\Gamma \frac{\bar{\mathbf{p}}(R)}{R}, \mathcal{P}^{mm}\right) = 0$.

We note that Theorem 2 is a corollary of this result where the set \mathcal{P}^{mm} is a singleton.

5 Related Work

Implicit Regularization. The implicit bias of gradient descent in classification tasks involving separable data has been extensively examined by [19, 20, 21, 22, 23, 24]. These works typically use logistic loss or, more generally, exponentially-tailed losses to make connections to margin maximization. These results are also extended to non-separable data by [25, 26, 27]. Furthermore, there have been notable investigations into the implicit bias in regression problems/losses utilizing techniques such as mirror descent [28, 20, 29, 30, 31, 32]. In addition, several papers have explored the implicit bias of stochastic gradient descent [33, 34, 35, 36, 37, 38], as well as adaptive and momentum-based methods [39, 40, 41, 42]. Although there are similarities between our optimization approach for \mathbf{v} and existing works, the optimization of \mathbf{p} stands out as significantly different. Firstly, our optimization problem is nonconvex, introducing new challenges and complexities. Secondly, it necessitates the introduction of novel concepts such as locally-optimal tokens and requires a fresh analysis specifically tailored to the cones surrounding them.

Attention Mechanism. Transformers, introduced by [6], revolutionized the field of NLP and machine translation, with earlier works on self-attention by [43, 44, 45, 46]. Self-attention differs from traditional models like MLPs and CNNs by leveraging global interactions for feature representations, showing exceptional empirical performance. However, the underlying mechanisms and learning processes of the attention layer remain unknown. Recent studies such as [47, 48, 49, 50, 51] have focused on specific aspects like representing sparse functions, convex-relaxations, and expressive power. [52, 53] have developed initial results to characterize the optimization and generalization dynamics of attention. [17] is another closely related work where the authors analyze the same attention model (ERM) as us. However, all of these works make stringent assumptions on the data, namely, tokens are tightly clusterable or can be clearly split into clear relevant and irrelevant sets. Additionally [53] requires assumptions on initialization and [52] considers a simplified attention structure where the attention matrix is not directly parameterized with respect to the input. Our work offers a comprehensive optimization-theoretic analysis of the attention model by establishing a formal connection to max-margin problems. Notably, our work presents the first theoretical understanding of the implicit bias exhibited by gradient descent methods in the context of the attention model.

6 Discussion

We have provided a thorough optimization-theoretic characterization of the fundamental attention model $f(\mathbf{X}) = \mathbf{v}^\top \mathbf{X}^\top \mathbb{S}(\mathbf{X}\mathbf{W}\mathbf{p})$ by formally connecting it to max-margin problems. We first established the convergence of gradient descent on \mathbf{p} (or equivalently \mathbf{W}) in isolation. We also explored joint convergence of (\mathbf{v}, \mathbf{p}) via regularization path which revealed surprising implicit biases such as (7). These findings motivate several exciting avenues for future research. An immediate open problem is characterizing the (local) convergence of gradient descent for joint optimization of (\mathbf{v}, \mathbf{p}) . Another major direction is to extend similar analysis to study self-attention layer (3) or to allow for multiple tunable tokens (where \mathbf{p} becomes a matrix). Either setting will enrich the problem by allowing the attention to discover multiple hyperplanes to separate tokens. While we assumed the tokens to be separable, it would be interesting to relax this assumption by leveraging results developed for logistic regression analysis [26, 19]. Ideas from these results can also be useful for characterizing the non-asymptotic behavior of how gradient descent aligns with the max-margin direction.

References

- [1] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *The International Conference on Learning Representations*,

2015.

[2] 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.

[3] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.

[4] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, and et al. Language models are unsupervised multitask learners. *arXiv preprint arXiv:1911.05786*, 2019.

[5] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.

[6] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, volume 30, 2017.

[7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.

[8] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[9] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.

[10] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021.

[11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.

[12] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.

[13] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.

[14] 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.

[15] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al. A generalist agent. *arXiv preprint arXiv:2205.06175*, 2022.

[16] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, 2021.

- 408 [17] Samet Oymak, Ankit Singh Rawat, Mahdi Soltanolkotabi, and Christos Thrampoulidis. On the
409 role of attention in prompt-tuning. In *ICLR 2023 Workshop on Mathematical and Empirical*
410 *Understanding of Foundation Models*, 2023.
- 411 [18] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation.
412 *arXiv preprint arXiv:2101.00190*, 2021.
- 413 [19] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The
414 implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*,
415 19(1):2822–2878, 2018.
- 416 [20] Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias
417 in terms of optimization geometry. In *International Conference on Machine Learning*, pages
418 1832–1841. PMLR, 2018.
- 419 [21] Mor Shpigel Nacson, Jason Lee, Suriya Gunasekar, Pedro Henrique Pamplona Savarese, Nathan
420 Srebro, and Daniel Soudry. Convergence of gradient descent on separable data. In *The 22nd*
421 *International Conference on Artificial Intelligence and Statistics*, pages 3420–3428. PMLR,
422 2019.
- 423 [22] Ziwei Ji and Matus Telgarsky. Characterizing the implicit bias via a primal-dual analysis. In
424 *Algorithmic Learning Theory*, pages 772–804. PMLR, 2021.
- 425 [23] Edward Moroshko, Blake E Woodworth, Suriya Gunasekar, Jason D Lee, Nati Srebro, and
426 Daniel Soudry. Implicit bias in deep linear classification: Initialization scale vs training accuracy.
427 *Advances in neural information processing systems*, 33:22182–22193, 2020.
- 428 [24] Ziwei Ji and Matus Telgarsky. Directional convergence and alignment in deep learning. In
429 H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural*
430 *Information Processing Systems*, volume 33, pages 17176–17186. Curran Associates, Inc.,
431 2020.
- 432 [25] Ziwei Ji and Matus Telgarsky. Risk and parameter convergence of logistic regression. *arXiv*
433 *preprint arXiv:1803.07300*, 2018.
- 434 [26] Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In
435 *Conference on Learning Theory*, pages 1772–1798. PMLR, 2019.
- 436 [27] Ziwei Ji, Miroslav Dudík, Robert E Schapire, and Matus Telgarsky. Gradient descent follows
437 the regularization path for general losses. In *Conference on Learning Theory*, pages 2109–2136.
438 PMLR, 2020.
- 439 [28] Blake Woodworth, Suriya Gunasekar, Jason D Lee, Edward Moroshko, Pedro Savarese, Itay
440 Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models.
441 In *Conference on Learning Theory*, pages 3635–3673. PMLR, 2020.
- 442 [29] Chulhee Yun, Shankar Krishnan, and Hossein Mobahi. A unifying view on implicit bias in
443 training linear neural networks. *arXiv preprint arXiv:2010.02501*, 2020.
- 444 [30] Tomas Vaskevicius, Varun Kanade, and Patrick Rebeschini. Implicit regularization for optimal
445 sparse recovery. *Advances in Neural Information Processing Systems*, 32:2972–2983, 2019.
- 446 [31] Ehsan Amid and Manfred K Warmuth. Winnowing with gradient descent. In *Conference on*
447 *Learning Theory*, pages 163–182. PMLR, 2020.
- 448 [32] Ehsan Amid and Manfred KK Warmuth. Reparameterizing mirror descent as gradient descent.
449 *Advances in Neural Information Processing Systems*, 33:8430–8439, 2020.
- 450 [33] Yuanzhi Li, Colin Wei, and Tengyu Ma. Towards explaining the regularization effect of initial
451 large learning rate in training neural networks. *arXiv preprint arXiv:1907.04595*, 2019.
- 452 [34] Guy Blanc, Neha Gupta, Gregory Valiant, and Paul Valiant. Implicit regularization for deep
453 neural networks driven by an ornstein-uhlenbeck like process. In *Conference on learning theory*,
454 pages 483–513. PMLR, 2020.

- [35] Jeff Z HaoChen, Colin Wei, Jason D Lee, and Tengyu Ma. Shape matters: Understanding the implicit bias of the noise covariance. *arXiv preprint arXiv:2006.08680*, 2020.
- [36] Zhiyuan Li, Tianhao Wang, and Sanjeev Arora. What happens after SGD reaches zero loss? –a mathematical framework. In *International Conference on Learning Representations*, 2022.
- [37] Alex Damian, Tengyu Ma, and Jason Lee. Label noise sgd provably prefers flat global minimizers. *arXiv preprint arXiv:2106.06530*, 2021.
- [38] Difan Zou, Jingfeng Wu, Vladimir Braverman, Quanquan Gu, Dean P Foster, and Sham Kakade. The benefits of implicit regularization from sgd in least squares problems. *Advances in Neural Information Processing Systems*, 34:5456–5468, 2021.
- [39] Qian Qian and Xiaoyuan Qian. The implicit bias of adagrad on separable data. *Advances in Neural Information Processing Systems*, 32, 2019.
- [40] Bohan Wang, Qi Meng, Huishuai Zhang, Ruoyu Sun, Wei Chen, and Zhi-Ming Ma. Momentum doesn’t change the implicit bias. *arXiv preprint arXiv:2110.03891*, 2021.
- [41] Bohan Wang, Qi Meng, Wei Chen, and Tie-Yan Liu. The implicit bias for adaptive optimization algorithms on homogeneous neural networks. In *International Conference on Machine Learning*, pages 10849–10858. PMLR, 2021.
- [42] Ziwei Ji, Nathan Srebro, and Matus Telgarsky. Fast margin maximization via dual acceleration. In *International Conference on Machine Learning*, pages 4860–4869. PMLR, 2021.
- [43] 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.
- [44] Ankur Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. A decomposable attention model for natural language inference. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2249–2255, Austin, Texas, November 2016. Association for Computational Linguistics.
- [45] Romain Paulus, Caiming Xiong, and Richard Socher. A deep reinforced model for abstractive summarization. In *International Conference on Learning Representations*, 2018.
- [46] Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. A structured self-attentive sentence embedding. In *International Conference on Learning Representations*, 2017.
- [47] Benjamin L Edelman, Surbhi Goel, Sham Kakade, and Cyril Zhang. Inductive biases and variable creation in self-attention mechanisms. In *International Conference on Machine Learning*, pages 5793–5831. PMLR, 2022.
- [48] Arda Sahiner, Tolga Ergen, Batu Ozturkler, John Pauly, Morteza Mardani, and Mert Pilanci. Unraveling attention via convex duality: Analysis and interpretations of vision transformers. In *International Conference on Machine Learning*, pages 19050–19088. PMLR, 2022.
- [49] Tolga Ergen, Behnam Neyshabur, and Harsh Mehta. Convexifying transformers: Improving optimization and understanding of transformer networks. *arXiv preprint arXiv:2211.11052*, 2022.
- [50] Pierre Baldi and Roman Vershynin. The quarks of attention. *arXiv preprint arXiv:2202.08371*, 2022.
- [51] Yihe Dong, Jean-Baptiste Cordonnier, and Andreas Loukas. Attention is not all you need: Pure attention loses rank doubly exponentially with depth. In *International Conference on Machine Learning*, pages 2793–2803. PMLR, 2021.
- [52] Samy Jelassi, Michael Eli Sander, and Yuanzhi Li. Vision transformers provably learn spatial structure. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.

- [53] Hongkang Li, Meng Wang, Sijia Liu, and Pin-Yu Chen. A theoretical understanding of shallow vision transformers: Learning, generalization, and sample complexity. *arXiv preprint arXiv:2302.06015*, 2023.
- [54] Olvi L Mangasarian. *Nonlinear programming*. SIAM, 1994.
- [55] Vladimir Vapnik. *Estimation of dependences based on empirical data*. Springer Science & Business Media, 2006.
- [56] Peter Bartlett. For valid generalization the size of the weights is more important than the size of the network. *Advances in neural information processing systems*, 9, 1996.
- [57] Albert B Novikoff. On convergence proofs for perceptrons. Technical report, STANFORD RESEARCH INST MENLO PARK CA, 1963.
- [58] Peter Bartlett, Yoav Freund, Wee Sun Lee, and Robert E Schapire. Boosting the margin: A new explanation for the effectiveness of voting methods. *The annals of statistics*, 26(5):1651–1686, 1998.
- [59] Tong Zhang and Bin Yu. Boosting with early stopping: Convergence and consistency. *Annals of Statistics*, page 1538, 2005.
- [60] Matus Telgarsky. Margins, shrinkage, and boosting. In *International Conference on Machine Learning*, pages 307–315. PMLR, 2013.
- [61] Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. *arXiv preprint arXiv:1810.02032*, 2018.
- [62] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [63] Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks. *arXiv preprint arXiv:1906.05890*, 2019.
- [64] Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss. In *Conference on Learning Theory*, pages 1305–1338. PMLR, 2020.
- [65] Spencer Frei, Gal Vardi, Peter L Bartlett, and Nathan Srebro. Benign overfitting in linear classifiers and leaky relu networks from kkt conditions for margin maximization. *arXiv e-prints*, pages arXiv–2303, 2023.
- [66] Gal Vardi, Ohad Shamir, and Nati Srebro. On margin maximization in linear and relu networks. *Advances in Neural Information Processing Systems*, 35:37024–37036, 2022.
- [67] Gal Vardi and Ohad Shamir. Implicit regularization in relu networks with the square loss. In *Conference on Learning Theory*, pages 4224–4258. PMLR, 2021.
- [68] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [69] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. Deep interest network for click-through rate prediction. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1059–1068, 2018.
- [70] Qiwei Chen, Huan Zhao, Wei Li, Pipei Huang, and Wenwu Ou. Behavior sequence transformer for e-commerce recommendation in alibaba. In *Proceedings of the 1st International Workshop on Deep Learning Practice for High-Dimensional Sparse Data*, pages 1–4, 2019.
- [71] Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin, Wenwu Ou, and Peng Jiang. Bert4rec: Sequential recommendation with bidirectional encoder representations from transformer. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, pages 1441–1450, 2019.

- [72] Mia Xu Chen, Orhan Firat, Ankur Bapna, Melvin Johnson, Wolfgang Macherey, George Foster, Llion Jones, Mike Schuster, Noam Shazeer, Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Zhifeng Chen, Yonghui Wu, and Macduff Hughes. The best of both worlds: Combining recent advances in neural machine translation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 76–86, Melbourne, Australia, July 2018. Association for Computational Linguistics.
- [73] Michael Janner, Qiyang Li, and Sergey Levine. Reinforcement learning as one big sequence modeling problem. In *ICML 2021 Workshop on Unsupervised Reinforcement Learning*, 2021.
- [74] Qinqing Zheng, Amy Zhang, and Aditya Grover. Online decision transformer. In *Proceedings of the 39th International Conference on Machine Learning*, 2022.
- [75] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.
- [76] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021.
- [77] Zi-Hang Jiang, Qibin Hou, Li Yuan, Daquan Zhou, Yujun Shi, Xiaojie Jin, Anran Wang, and Jiashi Feng. All tokens matter: Token labeling for training better vision transformers. In *Advances in Neural Information Processing Systems*, volume 34, pages 18590–18602, 2021.
- [78] Bingbin Liu, Jordan T Ash, Surbhi Goel, Akshay Krishnamurthy, and Cyril Zhang. Transformers learn shortcuts to automata. *arXiv preprint arXiv:2210.10749*, 2022.
- [79] Zizheng Pan, Bohan Zhuang, Jing Liu, Haoyu He, and Jianfei Cai. Scalable vision transformers with hierarchical pooling. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 377–386, 2021.
- [80] Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dynamicvit: Efficient vision transformers with dynamic token sparsification. In *Advances in Neural Information Processing Systems*, volume 34, pages 13937–13949, 2021.
- [81] Youwei Liang, GE Chongjian, Zhan Tong, Yibing Song, Jue Wang, and Pengtao Xie. Not all patches are what you need: Expediting vision transformers via token reorganizations. In *International Conference on Learning Representations*, 2022.
- [82] Yehui Tang, Kai Han, Yunhe Wang, Chang Xu, Jianyuan Guo, Chao Xu, and Dacheng Tao. Patch slimming for efficient vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12165–12174, 2022.
- [83] Hongxu Yin, Arash Vahdat, Jose M Alvarez, Arun Mallya, Jan Kautz, and Pavlo Molchanov. A-vit: Adaptive tokens for efficient vision transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10809–10818, 2022.
- [84] Hyunjik Kim, George Papamakarios, and Andriy Mnih. The lipschitz constant of self-attention. In *International Conference on Machine Learning*, pages 5562–5571. PMLR, 2021.
- [85] Jiri Hron, Yasaman Bahri, Jascha Sohl-Dickstein, and Roman Novak. Infinite attention: Nngp and ntk for deep attention networks. In *International Conference on Machine Learning*, pages 4376–4386. PMLR, 2020.
- [86] Greg Yang. Tensor programs ii: Neural tangent kernel for any architecture. *arXiv preprint arXiv:2006.14548*, 2020.
- [87] Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Lukasz Kaiser. Universal transformers. In *International Conference on Learning Representations*, 2018.
- [88] Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar. Are transformers universal approximators of sequence-to-sequence functions? In *International Conference on Learning Representations*, 2019.

- 598 [89] Satwik Bhattamishra, Kabir Ahuja, and Navin Goyal. On the ability and limitations of trans-
599 formers to recognize formal languages. *arXiv preprint arXiv:2009.11264*, 2020.
- 600 [90] Valerii Likhoshesterov, Krzysztof Choromanski, and Adrian Weller. On the expressive power of
601 self-attention matrices. *arXiv preprint arXiv:2106.03764*, 2021.
- 602 [91] Jean-Baptiste Cordonnier, Andreas Loukas, and Martin Jaggi. On the relationship between self-
603 attention and convolutional layers. In *International Conference on Learning Representations*,
604 2019.
- 605 [92] Yoav Levine, Noam Wies, Or Sharir, Hofit Bata, and Amnon Shashua. Limits to depth
606 efficiencies of self-attention. In *Advances in Neural Information Processing Systems*, volume 33,
607 pages 22640–22651, 2020.
- 608 [93] Charlie Snell, Ruiqi Zhong, Dan Klein, and Jacob Steinhardt. Approximating how single head
609 attention learns. *arXiv preprint arXiv:2103.07601*, 2021.
- 610 [94] Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan
611 Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv*
612 *preprint arXiv:2109.01652*, 2021.
- 613 [95] Tan Minh Nguyen, Tam Minh Nguyen, Nhat Ho, Andrea L Bertozzi, Richard Baraniuk, and
614 Stanley Osher. A primal-dual framework for transformers and neural networks. In *The Eleventh*
615 *International Conference on Learning Representations*, 2023.

Roadmap. The appendix is organized as follows: Section A provides basic facts about the training risk. Section B presents the proof of local and global gradient descent and regularized path for learning $\mathbf{p} \in \mathbb{R}^d$ with a fixed $\mathbf{v} \in \mathbb{R}^d$ choice. Section C provides the proof of regularized path applied to the general case of joint optimization of head \mathbf{v} and attention weights \mathbf{p} using a logistic loss function. Section D presents the proof for the regularized path applied to a more general model $f(\mathbf{X}) = \psi(\mathbf{X}^\top \mathbb{S}(\mathbf{X}\mathbf{W}^\top \mathbf{p}))$ with a nonlinear head ψ . Section E provides implementation details. Finally, Section F discusses additional related work on implicit bias and self-attention.

Corrections and Refinements. We have made the following changes to the main submission.

- In the first bullet point of Theorem 4, we corrected indices i into j . This was a typo.
- In the statement of Theorem 4, we now include the norm lower bound R_ε over the conic neighborhood. Note that, this is consistent with the setting of success guarantee Theorem 3 and the main message on the tightness of local optimality remains intact.
- In Theorem 6, we corrected the statement from $\mathbf{X}_i^\top \mathbb{S}(\mathbf{K}_i \mathbf{p}_R) \rightarrow \mathbf{x}_{i\alpha_i}$ to $\mathbb{S}(\mathbf{K}_i \mathbf{p}_R)_{\alpha_i} \rightarrow 1$. Note that, the former statement does not actually imply token index α_i is selected because combination of other tokens can still add up to $\mathbf{x}_{i\alpha_i}$. Instead, the new statement says softmax probability fully concentrates over α_i .

633

634

Table of Contents

636	A Addendum to Section 1	17
637	A.1 Preliminaries on the Training Risk	17
638	A.2 Proof of Lemma 1	17
639	B Addendum to Section 2	18
640	B.1 Local Gradient Condition	18
641	B.2 Descent and Gradient Correlation Conditions	21
642	B.3 Proof of Theorem 1	25
643	B.4 Proof of Theorem 2	26
644	B.5 Proof of Theorem 3	26
645	B.6 Proof of Theorem 4: Regularization Path Fails for Non-Locally-Optimal Tokens .	27
646	C Addendum to Section 3	30
647	C.1 Proof of Theorem 5	30
648	C.2 Proof of Theorem 6	31
649	D Addendum to Section 4	33
650	D.1 Proof of Theorem 7	33
651	D.2 Application to Linearly-mixed Labels	34
652	E Experimental Details	35
653	F Addendum to Section 5	35
654	F.1 Related Work on Implicit Regularization	35
655	F.2 Related Work on Attention Mechanism	36

656

657

658

659 A Addendum to Section 1

660 A.1 Preliminaries on the Training Risk

661 Recall the objective

$$\mathcal{L}(\mathbf{v}, \mathbf{p}, \mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i \cdot f(\mathbf{X}_i)). \quad (9)$$

662 with the generic prediction model $f(\mathbf{X}) = \psi(\mathbf{X}^\top \mathbb{S}(\mathbf{K}\mathbf{p}))$ and $\mathbf{K} = \mathbf{X}\mathbf{W}^\top$. Here, we write down the
 663 gradients of \mathbf{W} and \mathbf{p} in (9) to highlight the connection. Set $\mathbf{q} := \mathbf{W}^\top \mathbf{p}$, $\mathbf{z}\{\mathbf{X}\} := \mathbf{X}^\top \mathbb{S}(\mathbf{K}\mathbf{p})$ and
 664 $\mathbf{a}\{\mathbf{X}\} := \mathbf{K}\mathbf{p}$. Given \mathbf{X} and using $\mathbf{K} = \mathbf{X}\mathbf{W}^\top$, we have that

$$\nabla_{\mathbf{q}} f_\psi(\Theta) = \mathbf{X}^\top \mathbb{S}'(\mathbf{a}\{\mathbf{X}\}) \mathbf{X} \cdot \nabla \psi(\mathbf{z}\{\mathbf{X}\}) \quad (10a)$$

$$\nabla_{\mathbf{p}} f_\psi(\Theta) = \mathbf{W} \nabla_{\mathbf{q}} f_\psi(\Theta), \quad (10b)$$

$$\nabla_{\mathbf{W}} f_\psi(\Theta) = \mathbf{p} \nabla_{\mathbf{q}}^\top f_\psi(\Theta). \quad (10c)$$

665 Setting $\psi(\mathbf{z}) = \mathbf{v}^\top \mathbf{z}$ and recalling the score definition $\boldsymbol{\gamma} = \mathbf{X}\mathbf{v}$, for linear head, we obtain

$$\nabla_{\mathbf{q}} f_\psi(\Theta) = \mathbf{X}^\top \mathbb{S}'(\mathbf{a}\{\mathbf{X}\}) \boldsymbol{\gamma} \quad (11a)$$

$$\nabla_{\mathbf{p}} f_\psi(\Theta) = \mathbf{W} \nabla_{\mathbf{q}} f_\psi(\Theta) = \mathbf{K}^\top \mathbb{S}'(\mathbf{a}\{\mathbf{X}\}) \boldsymbol{\gamma}, \quad (11b)$$

$$\nabla_{\mathbf{W}} f_\psi(\Theta) = \mathbf{p} \nabla_{\mathbf{q}}^\top f_\psi(\Theta) = \mathbf{p} \boldsymbol{\gamma}^\top \mathbb{S}'(\mathbf{a}\{\mathbf{X}\}) \mathbf{X}. \quad (11c)$$

666 Note that the gradient of \mathbf{W} is rank-1 with fixed left singular direction. The proof of Lemma 1 below
 667 shows that solutions induced by matrix \mathbf{W} and vectors \mathbf{q}, \mathbf{p} can be mapped to each other exactly.

668 A.2 Proof of Lemma 1

669 **Proof.** Let us prove the result for a general step size sequence $(\eta_t)_{t \geq 0}$. By our assumption $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$
 670 and $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ are differentiable functions. Recall $\mathcal{L}(\mathbf{p}) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, \psi(\mathbf{X}_i^\top \mathbb{S}(\mathbf{X}_i \mathbf{p})))$ and
 671 $\mathcal{L}(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, \psi(\mathbf{X}_i^\top \mathbb{S}(\mathbf{X}_i \mathbf{W}^\top \mathbf{a})))$ for fixed \mathbf{a} . Suppose claim is true till iteration t . For iteration
 672 $t + 1$, using $\mathbf{W}_t^\top \mathbf{a} = \mathbf{p}_t$, define and observe that

$$\mathbf{S}_i' = \mathbb{S}'(\mathbf{X}_i \mathbf{W}_t^\top \mathbf{a}) = \mathbb{S}'(\mathbf{X}_i \mathbf{p}_t) \quad (12)$$

$$\mathbf{s}_i' = \mathbb{S}(\mathbf{X}_i \mathbf{W}_t^\top \mathbf{a}) = \mathbb{S}(\mathbf{X}_i \mathbf{p}_t) \quad (13)$$

$$\mathbf{z}\{\mathbf{X}_i\} := \mathbf{X}_i^\top \mathbb{S}(\mathbf{X}_i \mathbf{p}_t) = \mathbf{X}_i^\top \mathbb{S}(\mathbf{X}_i \mathbf{W}_t^\top \mathbf{a}) \quad (14)$$

673 for all $i \in [n]$. Thus, recalling (10a) and (10c), and defining $\ell_i' = \ell'(Y_i, \psi(\mathbf{z}\{\mathbf{X}_i\}))$ we have that

$$\nabla_{\mathbf{p}} \ell(Y_i, \psi(\mathbf{X}_i^\top \mathbb{S}(\mathbf{X}_i \mathbf{p}_t))) = \ell_i' \cdot \mathbf{X}_i^\top \mathbf{S}_i' \mathbf{X}_i \cdot \nabla \psi(\mathbf{z}\{\mathbf{X}_i\}), \quad (15)$$

$$\nabla_{\mathbf{W}} \ell(Y_i, \psi(\mathbf{X}_i^\top \mathbb{S}(\mathbf{X}_i \mathbf{W}_t^\top \mathbf{a}))) = \mathbf{a} \left(\ell_i' \cdot \mathbf{X}_i^\top \mathbf{S}_i' \mathbf{X}_i \cdot \nabla \psi(\mathbf{z}\{\mathbf{X}_i\}) \right)^\top. \quad (16)$$

674 Consequently, we found that gradient is rank-1 with left singular space equal given by \mathbf{a}

$$\nabla_{\mathbf{W}} \mathcal{L}_{\mathbf{W}}(\mathbf{W}_t) = \mathbf{a} \nabla_{\mathbf{p}}^\top \mathcal{L}_{\mathbf{q}}(\mathbf{p}_t).$$

675 Since \mathbf{W}_t 's left singular space is guaranteed to be in \mathbf{a} (including \mathbf{W}_0 by initialization), we only need
 676 to study the right singular vector. Using the induction till t , this yields

$$\mathbf{W}_{t+1}^\top \mathbf{a} = \mathbf{W}_t^\top \mathbf{a} - \eta_t \|\mathbf{a}\|^{-2} \nabla_{\mathbf{W}}^\top \mathcal{L}_{\mathbf{W}}(\mathbf{W}_t) \mathbf{a} \quad (17)$$

$$= \mathbf{p}_t - \eta_t \|\mathbf{a}\|^{-2} \mathbf{a}^\top \mathbf{a} \nabla_{\mathbf{p}} \mathcal{L}_{\mathbf{q}}(\mathbf{p}_t) \quad (18)$$

$$= \mathbf{q}_{t+1}. \quad (19)$$

677 This concludes the induction. ■

B Addendum to Section 2

B.1 Local Gradient Condition

Lemma 2 (Key lemma) Let $\mathbf{p}, \mathbf{q} \in \mathbb{R}^d$, $\mathbf{a} = \mathbf{K}\mathbf{q}$, $\mathbf{s} = \mathbb{S}(\mathbf{K}\mathbf{p})$, $\boldsymbol{\gamma} = \mathbf{X}\mathbf{v}$. Set $\Gamma = \sup_{t, \tau \in [T]} |\gamma_t - \gamma_\tau|$ and $A = \sup_{t \in [T]} \|\mathbf{k}_t\| \cdot \|\mathbf{q}\|$. We have that

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

Proof. Set $\bar{\gamma} = \sum_{t=1}^T \gamma_t s_t$. $\gamma_1 - \bar{\gamma} = \sum_{t \geq 2}^T (\gamma_1 - \gamma_t) s_t$. Also note that $|\bar{\gamma} - \gamma_1| \leq \Gamma(1 - s_1)$.

Proceeding,

$$\begin{aligned} \mathbf{a}^\top \text{diag}(\mathbf{s}) \boldsymbol{\gamma} - \mathbf{a}^\top \mathbf{s} \mathbf{s}^\top \boldsymbol{\gamma} &= \sum_{t=1}^T \mathbf{a}_t \gamma_t s_t - \sum_{t=1}^T \mathbf{a}_t s_t \sum_{t=1}^T \gamma_t s_t \\ &= \mathbf{a}_1 s_1 (\gamma_1 - \bar{\gamma}) - \sum_{t \geq 2}^T \mathbf{a}_t s_t (\bar{\gamma} - \gamma_t). \end{aligned} \quad (20)$$

Now using $|\sum_{t \geq 2}^T \mathbf{a}_t s_t (\bar{\gamma} - \gamma_t) - \sum_{t \geq 2}^T \mathbf{a}_t s_t (\gamma_1 - \gamma_t)| \leq A\Gamma(1 - s_1)^2$, we obtain¹

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

Above, in the last inequality (i.e., \pm on the right handside), we used the fact that

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

686

This lemma will play a key role in the following lemma.

Lemma 3 (Local Gradient Condition) Let $\alpha = (\alpha_i)_{i=1}^n$ be locally-optimal tokens per Definition 2. Define $\text{cone}_\mu(\mathbf{p}^{\text{mm}})$ to be the set of vectors obeying $\text{corr}(\mathbf{p}, \mathbf{p}^{\text{mm}}) \geq 1 - \mu$. There exists a scalar $\mu = \mu(\alpha) > 0$ such that for sufficiently large $R = R_\mu$:

- There is no stationary point within $\text{cone}_\mu(\mathbf{p}^{\text{mm}}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| \geq R\}$.
- Let $q_i = 1 - \mathbb{S}(\mathbf{K}_i \mathbf{p})_{\alpha_i}$ and $\ell'_i = \ell'(Y_i \cdot \mathbf{v}^\top \mathbf{X}_i^\top \mathbb{S}(\mathbf{K}_i \mathbf{p})) < 0$, $\gamma_i^{\text{gap}} = \min_{t \in \mathcal{T}_i} Y_i \cdot (\mathbf{x}_{i\alpha_i} - \mathbf{x}_{it})^\top \mathbf{v}$, $\bar{\gamma}_i^{\text{gap}} = \max_{t \in \mathcal{T}_i} Y_i \cdot (\mathbf{x}_{i\alpha_i} - \mathbf{x}_{it})^\top \mathbf{v}$. For all $\mathbf{q}, \mathbf{p} \in \text{cone}_\mu(\mathbf{p}^{\text{mm}})$ with $\|\mathbf{q}\| = \|\mathbf{p}^{\text{mm}}\|$, we have
$$\frac{2}{n} \sum_{i \in [n]} \ell'_i \cdot q_i \cdot \bar{\gamma}_i^{\text{gap}} \geq \langle \nabla \mathcal{L}(\mathbf{p}), \mathbf{q} \rangle \geq \frac{1}{8n} \sum_{i \in [n]} \ell'_i \cdot q_i \cdot \gamma_i^{\text{gap}}. \quad (21)$$

Note that above $-\ell'_i$ and $\gamma_i^{\text{gap}}, \bar{\gamma}_i^{\text{gap}}$ are upper/lower bounded by positive dataset-dependent constants. The only term that can vanish (as $\|\mathbf{p}\| \rightarrow \infty$) is q_i . Consequently, there exists constants $C, c > 0$ such that,

$$C \cdot \max_{i \in [n]} q_i \geq -\langle \nabla \mathcal{L}(\mathbf{p}), \mathbf{q} \rangle \geq c \cdot \min_{i \in [n]} q_i > 0. \quad (22)$$

Note that, the identical bound holds by setting $\mathbf{q} = \mathbf{p}^{\text{mm}}$ or $\mathbf{q} = \|\mathbf{p}^{\text{mm}}\| \mathbf{p} / \|\mathbf{p}\|$.

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

698 • Denote $\bar{\mathbf{p}} = \|\mathbf{p}^{mm}\| \mathbf{p} / \|\mathbf{p}\|$. For any $\pi > 0$, there exists $R := R_\pi$ such that all $\mathbf{p} \in \text{cone}_\mu(\mathbf{p}^{mm})$
699 with $\|\mathbf{p}\| \geq R$ obeys

$$\langle \nabla \mathcal{L}(\mathbf{p}), \bar{\mathbf{p}} \rangle \geq (1 + \pi) \langle \nabla \mathcal{L}(\mathbf{p}), \mathbf{p}^{mm} \rangle,$$

700 **Proof.** Let $\mathbf{p}^{mm} = \mathbf{p}^{mm}(\alpha)$ be the solution of (ATT-SVM). Define $\text{cone}_{\mu,R}(\mathbf{p}^{mm}) = \{\mathbf{p} \in$
701 $\mathbb{R}^d \mid \text{corr}(\mathbf{p}, \mathbf{p}^{mm}) \geq 1 - \mu, \|\mathbf{p}\| \geq R\}$. Let $(\mathcal{T}_i)_{i=1}^n$ be the set of all SVM-neighbors per Defini-
702 tion 2. Let $\bar{\mathcal{T}}_i = [T] - \mathcal{T}_i - \{\alpha_i\}$ be the non-SVM-neighbor tokens. Introduce the notation

$$\Theta = 1 / \|\mathbf{p}^{mm}\|, \quad (23)$$

$$\delta = 0.5 \min_{i \in [n]} \min_{t \in \mathcal{T}_i, \tau \in \bar{\mathcal{T}}_i} (\mathbf{k}_{it} - \mathbf{k}_{i\tau})^\top \mathbf{p}^{mm}, \quad (24)$$

$$A = \max_{i \in [n], t \in [T]} \|\mathbf{k}_{it}\| / \Theta, \quad (25)$$

$$\mu = \mu(\delta) = \frac{1}{8} \left(\frac{\min(0.5, \delta)}{A} \right)^2. \quad (26)$$

703 Since \mathbf{p}^{mm} is the max-margin model ensuring $(\mathbf{k}_{i\alpha_i} - \mathbf{k}_{i\tau})^\top \mathbf{p}^{mm} \geq 1$, the following inequalities hold
704 for all $\mathbf{p} \in \text{cone}_\mu(\mathbf{p}^{mm})$, $\|\mathbf{p}\| = \|\mathbf{p}^{mm}\|$ and all $i \in [n]$, $t \in \mathcal{T}_i$, $\tau \in \bar{\mathcal{T}}_i$:

$$(\mathbf{k}_{it} - \mathbf{k}_{i\tau})^\top \mathbf{p} \geq \delta > 0, \quad (27)$$

$$(\mathbf{k}_{i\alpha_i} - \mathbf{k}_{i\tau})^\top \mathbf{p} \geq 1 + \delta, \quad (28)$$

$$3/2 \geq (\mathbf{k}_{i\alpha_i} - \mathbf{k}_{i\tau})^\top \mathbf{p} \geq 1/2. \quad (29)$$

705 Above we used $\|\mathbf{p} - \mathbf{p}^{mm}\|^2 / \|\mathbf{p}^{mm}\|^2 \leq 2\mu$ which implies $\|\mathbf{p} - \mathbf{p}^{mm}\| \leq \sqrt{2\mu} / \Theta$.

706 **Proving Steps 1 and 2: No stationary point and $-\mathbf{q}^\top \nabla \mathcal{L}(\mathbf{p}) > 0$ within cone.** Now that the choice
707 of local cone is determined, we need to prove the main claims. We will lower bound $-\mathbf{q}^\top \nabla \mathcal{L}(\mathbf{p})$
708 and establish its strict positivity for $\|\mathbf{p}\| \geq R$. This will show that there is no stationary point as a by
709 product. Given any $\mathbf{p} \in \text{cone}_{\mu,R}(\mathbf{p}^{mm})$, denote $\bar{\mathbf{p}} = (\|\mathbf{p}^{mm}\| / \|\mathbf{p}\|) \mathbf{p}$ and recall $\|\mathbf{q}\| = \|\mathbf{p}^{mm}\|$. To proceed,
710 we write the gradient correlation following (44) and (46)

$$\langle \nabla \mathcal{L}(\mathbf{p}), \mathbf{q} \rangle = \frac{1}{n} \sum_{i=1}^n \ell'_i \cdot \langle \mathbf{a}_i, \mathbb{S}'(\mathbf{a}'_i) \gamma_i \rangle. \quad (30)$$

711 where we denoted $\ell'_i = \ell'(Y_i \cdot \mathbf{v}^\top X_i^\top \mathbb{S}(\mathbf{K}_i \mathbf{p}))$, $\mathbf{a}_i = \mathbf{K}_i \mathbf{q}$, $\mathbf{a}'_i = \mathbf{K}_i \mathbf{p}$, $s_i = \mathbb{S}(\mathbf{K}_i \mathbf{p})$. Using (27), for all
712 $t \in \mathcal{T}_i$, $\tau \in \bar{\mathcal{T}}_i$, for all $\mathbf{p} \in \text{cone}_{\mu,R}(\mathbf{p}^{mm})$, we have that

$$\mathbf{a}'_{i\alpha_i} - \mathbf{a}'_{i\tau} \geq R\Theta(1 + \delta), \quad \mathbf{a}'_{it} - \mathbf{a}'_{i\tau} \geq R\Theta\delta$$

713 Consequently, we can bound the softmax probabilities $s_i = \mathbb{S}(\mathbf{K}_i \mathbf{p})$ over non-neighbors as follows:
714 For all $i \in [n]$ and any $t_i \in \mathcal{T}_i$

$$S_i := \sum_{\tau \in \mathcal{T}_i} s_{i\tau} \leq \sum_{\tau \neq \alpha_i} s_{i\tau} \leq T e^{-R\Theta/2} s_{i\alpha_i} \leq T e^{-R\Theta/2}, \quad (31)$$

$$Q_i := \sum_{\tau \in \bar{\mathcal{T}}_i} s_{i\tau} \leq T e^{-R\Theta\delta} s_{i\alpha_i} \leq T e^{-R\Theta\delta} S_i. \quad (32)$$

715 Recall scores $\gamma_{it} = Y_i \cdot \mathbf{v}^\top \mathbf{x}_{it}$. Define the score gaps over neighbors: $\gamma_i^{\text{gap}} = \gamma_{i\alpha_i} - \max_{t \in \mathcal{T}_i} \gamma_{it}$,
716 $\bar{\gamma}_i^{\text{gap}} = \gamma_{i\alpha_i} - \min_{t \in \mathcal{T}_i} \gamma_{it}$. Recall that $A := \max_{i \in [n], t \in [T]} \|\mathbf{k}_{it}\| / \Theta \geq \max_{i, t \in [T]} \|\mathbf{a}_{it}\| = \|\mathbf{k}_{it} \mathbf{q}\|$. Define the
717 α -dependent global scalar $\Gamma = \sup_{i \in [n], t, \tau \in [T]} |\gamma_{it} - \gamma_{i\tau}|$.

718 Let us focus on a fixed datapoint $i \in [n]$, assume (without losing generality) $\alpha := \alpha_i = 1$, and drop
719 subscripts i , that is, $\alpha := \alpha_i$, $X := X_i$, $Y := Y_i$, $\mathbf{K} := \mathbf{K}_i$, $\mathbf{a}' = \mathbf{K} \mathbf{p}$, $\mathbf{a} = \mathbf{K} \mathbf{q}$, $s = \mathbb{S}(\mathbf{K} \mathbf{p})$, $\gamma = Y \cdot X \mathbf{v}$,
720 $\gamma^{\text{gap}} := \gamma_i^{\text{gap}}$. Directly applying Lemma 2, we obtain

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

721 To proceed, let us decouple the non-neighbors within $\sum_{t \geq 2}^T (\mathbf{a}_1 - \mathbf{a}_t) s_t (\gamma_1 - \gamma_t)$ via

$$|\sum_{t \in \bar{\mathcal{T}}} (\mathbf{a}_1 - \mathbf{a}_t) s_t (\gamma_1 - \gamma_t)| \leq 2Q\Gamma A.$$

722 Aggregating these, we found

$$\left| \mathbf{a}^\top \text{diag}(\mathbf{s}) \boldsymbol{\gamma} - \mathbf{a}^\top \mathbf{s} \mathbf{s}^\top \boldsymbol{\gamma} - \sum_{t \in \mathcal{T}} (\mathbf{a}_1 - \mathbf{a}_t) s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) \right| \leq 2\Gamma A((1 - s_1)^2 + Q). \quad (33)$$

723 To proceed, let us upper/lower bound the gradient correlation. Since $1.5 \geq \mathbf{a}_1 - \mathbf{a}_t \geq 0.5$, we find

$$1.5 \cdot S \cdot \bar{\gamma}^{\text{gap}} \sum_{t \in \mathcal{T}} (\mathbf{a}_1 - \mathbf{a}_t) s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) \geq 0.5 \cdot S \cdot \gamma^{\text{gap}}.$$

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

$$S \cdot \gamma^{\text{gap}}/4 \geq 4\Gamma A \max((1 - s_1)^2, Q) \iff S \geq 16 \frac{\Gamma A}{\gamma^{\text{gap}}} \max((1 - s_1)^2, Q). \quad (34)$$

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

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

726 Combining these, what we wish is ensured by guaranteeing

$$S \geq 16 \frac{\Gamma A}{\gamma^{\text{gap}}} \max(4ST e^{-R\Theta/2}, T e^{-R\Theta\delta} S). \quad (35)$$

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

$$R = \frac{\max(2, \delta^{-1})}{\Theta} \log\left(\frac{64T\Gamma A}{\gamma_{\min}^{\text{gap}}}\right), \quad (36)$$

728 where $\gamma_{\min}^{\text{gap}} = \sup_{i \in [n]} \gamma_i^{\text{gap}}$ is the global scalar which is the worst case score gap over all inputs. With
729 the above choice of R , we guaranteed

$$2(1 - s_1) \cdot \bar{\gamma}^{\text{gap}} \geq 2 \cdot S \cdot \bar{\gamma}^{\text{gap}} \geq \sum_{t \in \mathcal{T}} (\mathbf{a}_1 - \mathbf{a}_t) s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) \geq \frac{S \cdot \gamma^{\text{gap}}}{4} \geq \frac{(1 - s_1) \gamma^{\text{gap}}}{8}.$$

730 via (34) and (39). Since this holds over all inputs, going back to the gradient correlation (30) and
731 averaging above over all inputs $i \in [n]$ and plugging back the indices i , we obtain the advertised
732 bound by setting $q_i = 1 - s_{i\alpha_i}$ (where we set $\alpha_i = 1$ above without losing generality)

$$\frac{2}{n} \sum_{i \in [n]} \ell'_i \cdot q_i \cdot \bar{\gamma}_i^{\text{gap}} \geq \langle \nabla \mathcal{L}(\mathbf{p}), \mathbf{q} \rangle \geq \frac{1}{8n} \sum_{i \in [n]} \ell'_i \cdot q_i \cdot \gamma_i^{\text{gap}}. \quad (37)$$

733 **Proving Step 3: Establishing gradient correlation.** Our final goal is establishing gradient compari-
734 son between $\mathbf{p}, \mathbf{p}^{\text{mm}}$ for the same choice of $\mu > 0$ provided in (23). Define $\bar{\mathbf{p}} = \|\mathbf{p}^{\text{mm}}\| \mathbf{p} / \|\mathbf{p}\|$ to be the
735 normalized vector. Set notations $\mathbf{a}_i = \mathbf{K}_i \bar{\mathbf{p}}$, $\bar{\mathbf{a}}_i = \mathbf{K}_i \mathbf{p}^{\text{mm}}$, and $s_i = \mathbb{S}(\mathbf{K}_i \mathbf{p})$. To establish the result, we
736 will prove that, for sufficiently large $R = R_\pi$, for any $\mathbf{p} \in \text{cone}_{\mu, R}(\mathbf{p}^{\text{mm}})$ and for any $i \in [n]$,

$$\langle \mathbf{a}_i, \mathbb{S}'(\mathbf{a}_i) \boldsymbol{\gamma}_i \rangle \leq (1 + \pi) \langle \bar{\mathbf{a}}_i, \mathbb{S}'(\mathbf{a}_i) \boldsymbol{\gamma}_i \rangle. \quad (38)$$

737 Once (38) holds for all i , the same conclusion will hold for the gradient correlations via (30). Moving
738 forward, we shall again focus on a single point $i \in [n]$ and drop all subscripts i . Also assume
739 $\alpha = \alpha_i = 1$ without losing generality (same as above).

740 Following (39), for all $\mathbf{q} \in \text{cone}_\mu$ with $\|\mathbf{q}\| = \|\mathbf{p}^{\text{mm}}\|$ and $\mathbf{a}' = \mathbf{K} \mathbf{q}$, we have found

$$\left| \mathbf{a}'^\top \text{diag}(\mathbf{s}) \boldsymbol{\gamma} - \mathbf{a}'^\top \mathbf{s} \mathbf{s}^\top \boldsymbol{\gamma} - \sum_{t \in \mathcal{T}} (\mathbf{a}'_1 - \mathbf{a}'_t) s_t (\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_t) \right| \leq 2\Gamma A((1 - s_1)^2 + Q). \quad (39)$$

741 Plugging in $\mathbf{a}, \bar{\mathbf{a}}$ in the bound above and assuming $\pi \leq 1$ (w.l.o.g.), (38) is implied by the following
742 stronger inequality

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

743 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
744 follows the earlier argument, namely, following (34), (36) and (35) which leads to the choice

745 $R_\pi = \frac{\max(2, \delta^{-1})}{\Theta} \log\left(\frac{C \cdot T \Gamma A}{\pi \gamma_{\min}^{\text{gap}}}\right)$ for some constant $C > 0$.

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

$$\sum_{t \in \mathcal{T}} (a_1 - a_t) s_t (\gamma_1 - \gamma_t) \leq (1 + 0.5\pi) \sum_{t \in \mathcal{T}} s_t (\gamma_1 - \gamma_t) \quad (40)$$

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

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

$$|a_t - \bar{a}_t| = |\mathbf{k}_t^\top (\bar{\mathbf{p}} - \mathbf{p}^{mm})| \leq A\Theta\varepsilon = \pi/4.$$

750 Consequently, we obtain

$$a_1 - a_t \leq \bar{a}_1 - \bar{a}_t + 2A\Theta\varepsilon \leq 1 + 0.5\pi.$$

751 Similarly, $a_1 - a_t \geq 1 - 0.5\pi \geq 0.5$. Since all terms $a_1 - a_t, s_t, \gamma_1 - \gamma_t$ in (40) are nonnegative and
 752 $(a_1 - a_t)s_t(\gamma_1 - \gamma_t) \leq (1 + 0.5\pi)s_t(\gamma_1 - \gamma_t)$, above implies the desired result (40).

753 **Scenario 2:** $\|\bar{\mathbf{p}} - \mathbf{p}^{mm}\| \geq \varepsilon = \frac{\pi}{4A\Theta}$. Since $\bar{\mathbf{p}}$ is not (locally) max-margin, in this scenario, for some
 754 $\nu = \nu(\varepsilon) > 0$ and $\tau \in \mathcal{T}$, we have that $\bar{\mathbf{p}}^\top(\mathbf{k}_1 - \mathbf{k}_\tau) = a_1 - a_\tau \leq 1 - 2\nu$. Here $\tau = \arg \max_{\tau \in \mathcal{T}} \bar{\mathbf{p}}^\top \mathbf{k}_\tau$
 755 denotes the nearest point to \mathbf{k}_1 (along the $\bar{\mathbf{p}}$ direction). Note that a non-neighbor $\tau \in \bar{\mathcal{T}}$ cannot be
 756 nearest because $\mathbf{p} \in C_\mu$ and (27) holds. Recall that $s = \mathbb{S}(\bar{R}\mathbf{a})$ where $\bar{R} = R\Theta$. To proceed, split the
 757 tokens into two groups: Let \mathcal{N} be the group of tokens obeying $\mathbf{p}^\top(\mathbf{k}_1 - \mathbf{k}_\tau) \leq 1 - \nu$ and $\mathcal{T} - \mathcal{N}$ be the
 758 rest of the neighbors. Observe that

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

759 Thus, using $|a_1 - a_t| \leq 2A$ and recalling the definition of γ^{gap} , observe that

$$\sum_{t \in \mathcal{T} - \mathcal{N}} (a_1 - a_t) s_t (\gamma_1 - \gamma_t) \leq \frac{2\Gamma A T e^{-\bar{R}\nu}}{\gamma^{gap}} \sum_{t \in \mathcal{N}} s_t (\gamma_1 - \gamma_t).$$

760 Plugging this into (40), we obtain

$$\begin{aligned} \sum_{t \in \mathcal{T}} (a_1 - a_t) s_t (\gamma_1 - \gamma_t) &= \sum_{t \in \mathcal{N}} (a_1 - a_t) s_t (\gamma_1 - \gamma_t) + \sum_{t \in \mathcal{T} - \mathcal{N}} (a_1 - a_t) s_t (\gamma_1 - \gamma_t) \\ &\leq \sum_{t \in \mathcal{N}} (1 - \nu) s_t (\gamma_1 - \gamma_t) + \sum_{t \in \mathcal{T} - \mathcal{N}} 2\Gamma A T e^{-\bar{R}\nu} \\ &\leq (1 - \nu + \frac{2\Gamma A T e^{-\bar{R}\nu}}{\gamma^{gap}}) \sum_{t \in \mathcal{T}} s_t (\gamma_1 - \gamma_t) \end{aligned} \quad (41)$$

$$\leq (1 + \frac{2\Gamma A T e^{-\bar{R}\nu}}{\gamma^{gap}}) \sum_{t \in \mathcal{T}} s_t (\gamma_1 - \gamma_t). \quad (42)$$

$$(43)$$

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

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

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

764 Combining this with the prior R_π choice (by taking maximum), we conclude with the statement. ■

765 B.2 Descent and Gradient Correlation Conditions

766 The lemma below identifies conditions under which $\mathbf{p}^{mm\star}$ is a global descent direction for $\mathcal{L}(\mathbf{p})$.

767 **Lemma 4 (Global descent conditions)** Suppose $\ell(\cdot)$ is a strictly decreasing loss function and either
 768 of the following two conditions holds

- **Scores of non-optimal tokens are same:** For all $i \in [n]$ and $t_1, t_2 \neq \text{opt}_i$, $\mathbf{v}^\top \mathbf{x}_{it_1} = \mathbf{v}^\top \mathbf{x}_{it_2}$.
- **All tokens are support vectors:** Consider (ATT-SVM) with optimal indices $(\text{opt}_i)_{i=1}^n$. $(\mathbf{k}_i^{\text{opt}} - \mathbf{k}_{it})^\top \mathbf{p}^{\text{mm}\star} = 1$ for all $t \neq \text{opt}_i, i \in [n]$.

Define

- $\mathbf{a}_{\text{gap}}^i := 1 = \inf_{t \neq \text{opt}_i} (\mathbf{k}_i^{\text{opt}} - \mathbf{k}_{it})^\top \mathbf{p}^{\text{mm}\star}$,
- $\gamma_{\text{gap}}^i = \inf_{t \neq \text{opt}_i} Y_i \cdot (\mathbf{x}_i^{\text{opt}} - \mathbf{x}_{it})^\top \mathbf{v}$,
- $\text{lg}t'_i = p_{\text{opt}}(1 - p_{\text{opt}})$ where $p_{\text{opt}} = \mathbb{S}(\mathbf{K}_i \mathbf{p})_{\text{opt}_i}$,
- $\ell'_i = \ell'(Y_i \cdot \mathbf{v}^\top \mathbf{X}_i^\top \mathbb{S}(\mathbf{K}_i \mathbf{p})) < 0$.

Then, for all $\mathbf{p} \in \mathbb{R}^d$, the training loss (ERM) obeys

$$-\langle \nabla \mathcal{L}(\mathbf{p}), \mathbf{p}^{\text{mm}\star} \rangle \geq \min_{i \in [n]} \{ -\ell'_i \cdot \text{lg}t'_i \cdot \mathbf{a}_{\text{gap}}^i \cdot \gamma_{\text{gap}}^i \} > 0.$$

Proof. Set $\bar{\mathbf{a}}_i = \mathbf{K}_i \mathbf{p}^{\text{mm}\star}$ to obtain In order to show this result, let us recall the gradient evaluated at \mathbf{p} which is given by

$$\nabla \mathcal{L}(\mathbf{p}) = \frac{1}{n} \sum_{i=1}^n \ell'_i \cdot \mathbf{K}_i^\top \mathbb{S}'(\mathbf{a}_i) \gamma_i. \quad (44)$$

Here $\gamma_i = Y_i \cdot \mathbf{X}_i \mathbf{v}$, $\mathbf{a}_i = \mathbf{K}_i \mathbf{p}$, and $\ell'_i = \ell'(Y_i \cdot \mathbf{v}^\top \mathbf{X}_i^\top \mathbb{S}(\mathbf{K}_i \mathbf{p}))$. This implies that

$$\langle \nabla \mathcal{L}(\mathbf{p}), \mathbf{p}^{\text{mm}\star} \rangle = \frac{1}{n} \sum_{i=1}^n \ell'_i \cdot \langle \bar{\mathbf{a}}_i, \mathbb{S}'(\mathbf{a}_i) \gamma_i \rangle.$$

To proceed, we will prove that individual summands are all strictly negative. To show that, without losing generality, let us focus on the first input and drop the subscript i for cleaner notation. This yields

$$\langle \bar{\mathbf{a}}, \mathbb{S}'(\mathbf{a}) \gamma \rangle = \bar{\mathbf{a}}^\top \text{diag}(\mathbb{S}(\mathbf{a})) \gamma - \bar{\mathbf{a}}^\top \mathbb{S}(\mathbf{a}) \mathbb{S}(\mathbf{a})^\top \gamma. \quad (45)$$

Without losing generality, assume optimal token is the first one. The lemma has two scenarios. In the first scenario (same non-optimal scores), γ_t is a constant for all $t \geq 2$. In the second scenario (all tokens are support), $\bar{\mathbf{a}}_t = \mathbf{k}_t \mathbf{p}^{\text{mm}\star}$ is constant for all $t \geq 2$. Since $\bar{\mathbf{a}}, \gamma$ vectors are represented symmetrically in the gradient correlation, verifying these two conditions are equivalent.

To proceed, we will prove the following (focusing on the first condition): Suppose $\gamma = \gamma_{t \geq 2}$ is constant, $\gamma_1, \bar{\mathbf{a}}_1$ are the largest indices of $\gamma, \bar{\mathbf{a}}$. Then, for any s obeying $\sum_{t \in [T]} s_t = 1, s_t \geq 0$, we have that $\bar{\mathbf{a}}^\top \text{diag}(s) \gamma - \bar{\mathbf{a}}^\top s s^\top \gamma > 0$. To see this, we write

$$\bar{\mathbf{a}}^\top \text{diag}(s) \gamma - \bar{\mathbf{a}}^\top s s^\top \gamma = \sum_{t=1}^T \bar{\mathbf{a}}_t \gamma_t s_t - \sum_{t=1}^T \bar{\mathbf{a}}_t s_t \sum_{t=1}^T \gamma_t s_t \quad (46)$$

$$= (\bar{\mathbf{a}}_1 \gamma_1 s_1 + \gamma \sum_{t \geq 2} \bar{\mathbf{a}}_t s_t) - (\gamma_1 s_1 + \gamma(1 - s_1))(\bar{\mathbf{a}}_1 s_1 + \sum_{t \geq 2} \bar{\mathbf{a}}_t s_t) \quad (47)$$

$$= \bar{\mathbf{a}}_1(\gamma_1 - \gamma)s_1(1 - s_1) + (\gamma - (\gamma_1 s_1 + \gamma(1 - s_1))) \sum_{t \geq 2} \bar{\mathbf{a}}_t s_t \quad (48)$$

$$= \bar{\mathbf{a}}_1(\gamma_1 - \gamma)s_1(1 - s_1) - (\gamma_1 - \gamma)s_1 \sum_{t \geq 2} \bar{\mathbf{a}}_t s_t \quad (49)$$

$$= (\gamma_1 - \gamma)(1 - s_1)s_1 \left[\bar{\mathbf{a}}_1 - \frac{\sum_{t \geq 2} \bar{\mathbf{a}}_t s_t}{\sum_{t \geq 2} s_t} \right]. \quad (50)$$

To proceed, recall the definitions $\gamma_{\text{gap}} = \gamma_1 - \max_{t \geq 2} \gamma_t$ and $\mathbf{a}_{\text{gap}} = \bar{\mathbf{a}}_1 - \max_{t \geq 2} \bar{\mathbf{a}}_t$. With these, we obtain

$$\bar{\mathbf{a}}^\top \text{diag}(s) \gamma - \bar{\mathbf{a}}^\top s s^\top \gamma \geq \mathbf{a}_{\text{gap}} \gamma_{\text{gap}} s_1 (1 - s_1),$$

which is the advertised result after noticing $s_1(1 - s_1)$ is the logistic derivative and infimum'ing over all inputs and multiplying by ℓ'_i . ■

Lemma 5 (Gradient correlation conditions) Fix indices $\alpha = (\alpha_{i=1}^n)$ and let $\mathbf{p}^{mm} = \mathbf{p}^{mm}(\alpha)$ be the SVM solution separating α_i from remaining tokens of input \mathbf{X}_i for $i \in [n]$. Suppose for all $i \in [n]$ and $t_1, t_2 \neq \alpha_i$, $\mathbf{v}^\top \mathbf{x}_{it_1} = \mathbf{v}^\top \mathbf{x}_{it_2} < \mathbf{v}^\top \mathbf{x}_{i\alpha_i}$ and $\ell(\cdot)$ is strictly decreasing. Let $\bar{\mathbf{p}} = \|\mathbf{p}^{mm}\| \mathbf{p} / \|\mathbf{p}\|$. $M = \sup_{i,t} \|\mathbf{k}_i\|$ and $\Xi = 1/\|\mathbf{p}^{mm}\|$. For any choice of $\pi > 0$, there exists $R := R_\pi$ such that, for any \mathbf{p} with $\|\mathbf{p}\| \geq R$, we have

$$\langle \nabla \mathcal{L}(\mathbf{p}), \bar{\mathbf{p}} \rangle \geq (1 + \pi) \langle \nabla \mathcal{L}(\mathbf{p}), \mathbf{p}^{mm} \rangle.$$

Above, observe that as $R \rightarrow \infty$, we eventually get to set $\pi = 0$.

Proof. The proof is similar to Lemma 4 at a high-level. However, we also need to account for the impact of \mathbf{p} besides \mathbf{p}^{mm} in the gradient correlation. The main goal is showing that \mathbf{p}^{mm} is the near-optimal descent direction, thus, \mathbf{p} cannot significantly outperform it.

To proceed, set $s_i = \mathbb{S}(\mathbf{K}_i \mathbf{p})$, $\mathbf{a}_i = \mathbf{K}_i \bar{\mathbf{p}}$, $\bar{\mathbf{a}}_i = \mathbf{K}_i \mathbf{p}^{mm}$. Without losing generality assume $\alpha_i = 1$ for all $i \in [n]$. Set $\text{lg}\mathbf{t}'_i = s_{i1}/(1 - s_{i1})$. Repeating the proof of Lemma 4 yields

$$\langle \nabla \mathcal{L}(\mathbf{p}), \mathbf{p}^{mm} \rangle = \frac{1}{n} \sum_{i=1}^n \ell'_i \cdot \text{lg}\mathbf{t}'_i \cdot (\gamma_{i1} - \gamma_i) \left[\bar{\mathbf{a}}_{i1} - \frac{\sum_{t \geq 2}^T \bar{\mathbf{a}}_{it} s_{it}}{\sum_{t \geq 2} s_{it}} \right] \quad (51)$$

$$\langle \nabla \mathcal{L}(\mathbf{p}), \bar{\mathbf{p}} \rangle = \frac{1}{n} \sum_{i=1}^n \ell'_i \cdot \text{lg}\mathbf{t}'_i \cdot (\gamma_{i1} - \gamma_i) \left[\mathbf{a}_{i1} - \frac{\sum_{t \geq 2}^T \mathbf{a}_{it} s_{it}}{\sum_{t \geq 2} s_{it}} \right] \quad (52)$$

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

$$\left[\mathbf{a}_1 - \frac{\sum_{t \geq 2}^T \mathbf{a}_t s_t}{\sum_{t \geq 2} s_t} \right] \leq (1 + \pi) \cdot \left[\bar{\mathbf{a}}_1 - \frac{\sum_{t \geq 2}^T \bar{\mathbf{a}}_t s_t}{\sum_{t \geq 2} s_t} \right]. \quad (53)$$

We consider two scenarios. Let $M = \max_{i \in [n], t \in [T]} \|\mathbf{k}_{it}\|$.

Scenario 1: $\|\bar{\mathbf{p}} - \mathbf{p}^{mm}\| \leq \varepsilon := \pi/2M$. In this scenario, for any token, we find that

$$|\mathbf{a}_t - \bar{\mathbf{a}}_t| = |\mathbf{k}_t^\top (\bar{\mathbf{p}} - \mathbf{p}^{mm})| \leq M \|\bar{\mathbf{p}} - \mathbf{p}^{mm}\| \leq M\varepsilon.$$

Consequently, we obtain

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

Also noticing $\bar{\mathbf{a}}_1 - \frac{\sum_{t \geq 2}^T \bar{\mathbf{a}}_t s_t}{\sum_{t \geq 2} s_t} \geq 1$ (thanks to \mathbf{p}^{mm} satisfying ≥ 1 margin), this implies (53).

Scenario 2: $\|\bar{\mathbf{p}} - \mathbf{p}^{mm}\| \geq \varepsilon := \pi/2M$. In this scenario, for some $\delta = \delta(\varepsilon)$ and $\tau \geq 2$, we have that $\mathbf{p}^\top (\mathbf{k}_1 - \mathbf{k}_\tau) = \mathbf{a}_1 - \mathbf{a}_\tau \leq 1 - 2\delta$. Here $\tau = \arg \max_{t \geq 2} \mathbf{p}^\top \mathbf{k}_t$ denotes the nearest point to \mathbf{k}_1 . Recall that $s = \mathbb{S}(\bar{R}\mathbf{a})$ where $\bar{R} = R\Xi = R/\|\mathbf{p}^{mm}\|$. To proceed, split the tokens into two groups: Let \mathcal{N} be the group of tokens obeying $\mathbf{p}^\top (\mathbf{k}_1 - \mathbf{k}_\tau) \geq 1 - \delta$ and $[T] - \mathcal{N}$ be the rest. Observe that

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

Set $\tilde{M} = M/\Xi$ and note that $\|\mathbf{a}_t\| \leq \|\mathbf{p}^{mm}\| \cdot \|\mathbf{k}_t\| \leq \tilde{M}$. Using $\mathbf{p}^\top (\mathbf{k}_1 - \mathbf{k}_\tau) < 1 - \delta$ over $\tau \in [T] - \mathcal{N}$ and plugging in the above bound, we obtain

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

Using the fact that $\bar{\mathbf{a}}_1 - \frac{\sum_{t \geq 2}^T \bar{\mathbf{a}}_t s_t}{\sum_{t \geq 2} s_t} \geq 1$, the above implies (53) with $\pi' = (1 - \delta) + 2\tilde{M} T e^{-\bar{R}\delta}$. To proceed, choose $R_\pi = \delta^{-1} \Xi^{-1} \log(2\tilde{M} T / \pi)$ to ensure $\pi' \leq \pi$. ■

The following lemma states the descent property of gradient descent for $\mathcal{L}(\mathbf{p})$ under Assumption A. It is important to note that although the infimum of the optimization problem is \mathcal{L}^* , it is not achieved at any finite \mathbf{p} . Additionally, there are no finite critical points \mathbf{p} .

823 **Lemma 6** Under Assumption A, the objective $\mathcal{L}(\mathbf{p})$ is L_p -smooth, where

$$L_p := \frac{1}{n} \sum_{i=1}^n \left(M_0 \|\mathbf{v}\|^2 \|\mathbf{W}\|^2 + M_1 \|\mathbf{v}\| \|\mathbf{W}\|^3 \right) \|\mathbf{X}_i\|^4. \quad (54)$$

824 Further, if $\eta \leq 2/L_p$, then, for any initialization $\mathbf{p}(0)$, with the GD sequence $\mathbf{p}(t+1) = \mathbf{p}(t) - \eta \nabla \mathcal{L}(\mathbf{p}(t))$,
825 we have

$$\mathcal{L}(\mathbf{p}(t+1)) - \mathcal{L}(\mathbf{p}(t)) \leq -\frac{\eta}{2} \|\nabla \mathcal{L}(\mathbf{p}(t))\|^2, \quad (55)$$

826 for all $t \geq 0$, $\sum_{t=0}^{\infty} \|\nabla \mathcal{L}(\mathbf{p}(t))\|^2 < \infty$ and $\lim_{t \rightarrow \infty} \|\nabla \mathcal{L}(\mathbf{p}(t))\|^2 = 0$.

827 **Proof.** Recall that we defined $\gamma_i = Y_i \cdot \mathbf{X}_i \mathbf{v}$, $\mathbf{a}_i = \mathbf{K}_i \mathbf{p}$, and $\ell'_i = \ell'(Y_i \cdot \mathbf{v}^\top \mathbf{X}_i^\top \mathbb{S}(\mathbf{K}_i \mathbf{p}))$. The gradient
828 evaluated at \mathbf{p} is given by

$$\nabla \mathcal{L}(\mathbf{p}) = \frac{1}{n} \sum_{i=1}^n \ell'_i \cdot \mathbf{K}_i^\top \mathbb{S}'(\mathbf{a}_i) \gamma_i.$$

829 Now, for any $\mathbf{p}, \dot{\mathbf{p}} \in \mathbb{R}^d$, we have

$$\begin{aligned} \|\nabla \mathcal{L}(\mathbf{p}) - \nabla \mathcal{L}(\dot{\mathbf{p}})\| &\leq \frac{1}{n} \sum_{i=1}^n \left\| \ell'(\gamma_i^\top \mathbb{S}(\mathbf{K}_i \mathbf{p})) \cdot \mathbf{K}_i^\top \mathbb{S}'(\mathbf{K}_i \mathbf{p}) \gamma_i - \ell'(\gamma_i^\top \mathbb{S}(\mathbf{K}_i \dot{\mathbf{p}})) \cdot \mathbf{K}_i^\top \mathbb{S}'(\mathbf{K}_i \dot{\mathbf{p}}) \gamma_i \right\| \\ &\leq \frac{1}{n} \sum_{i=1}^n \left\| \mathbf{K}_i^\top \mathbb{S}'(\mathbf{K}_i \dot{\mathbf{p}}) \gamma_i \right\| \left\| \ell'(\gamma_i^\top \mathbb{S}(\mathbf{K}_i \mathbf{p})) - \ell'(\gamma_i^\top \mathbb{S}(\mathbf{K}_i \dot{\mathbf{p}})) \right\| \\ &\quad + \frac{1}{n} \sum_{i=1}^n \left\| \ell'(\gamma_i^\top \mathbb{S}(\mathbf{K}_i \mathbf{p})) \right\| \left\| \mathbf{K}_i^\top \mathbb{S}'(\mathbf{K}_i \mathbf{p}) \gamma_i - \mathbf{K}_i^\top \mathbb{S}'(\mathbf{K}_i \dot{\mathbf{p}}) \gamma_i \right\| \\ &\leq \frac{1}{n} \sum_{i=1}^n M_0 \|\gamma_i\|^2 \|\mathbf{K}_i\| \|\mathbb{S}(\mathbf{K}_i \mathbf{p}) - \mathbb{S}(\mathbf{K}_i \dot{\mathbf{p}})\| + M_1 \|\mathbf{K}_i\| \|\gamma_i\| \|\mathbb{S}'(\mathbf{K}_i \mathbf{p}) - \mathbb{S}'(\mathbf{K}_i \dot{\mathbf{p}})\|, \end{aligned}$$

830 where the second inequality follows from the fact that $|ab - cd| \leq |d||a - c| + |a||b - d|$ and the third
831 inequality uses Assumption A.

832 Note also that for any $\mathbf{p}, \dot{\mathbf{p}} \in \mathbb{R}^d$,

$$\|\mathbb{S}(\mathbf{K}_i \mathbf{p}) - \mathbb{S}(\mathbf{K}_i \dot{\mathbf{p}})\| \leq \|\mathbf{K}_i\| \|\mathbf{p} - \dot{\mathbf{p}}\| \quad \text{and} \quad \|\mathbb{S}'(\mathbf{K}_i \mathbf{p}) - \mathbb{S}'(\mathbf{K}_i \dot{\mathbf{p}})\| \leq \|\mathbf{K}_i\|^2 \|\mathbf{p} - \dot{\mathbf{p}}\|.$$

833 Hence,

$$\begin{aligned} \|\nabla \mathcal{L}(\mathbf{p}) - \nabla \mathcal{L}(\dot{\mathbf{p}})\| &\leq \frac{1}{n} \sum_{i=1}^n \left(M_0 \|\gamma_i\|^2 \|\mathbf{K}_i\|^2 + M_1 \|\mathbf{K}_i\|^3 \|\gamma_i\| \right) \|\mathbf{p} - \dot{\mathbf{p}}\| \\ &\leq \frac{1}{n} \sum_{i=1}^n \left(M_0 \|\mathbf{v}\|^2 \|\mathbf{W}\|^2 \|\mathbf{X}_i\|^4 + M_1 \|\mathbf{v}\| \|\mathbf{W}\|^3 \|\mathbf{X}_i\|^4 \right) \|\mathbf{p} - \dot{\mathbf{p}}\| \\ &\leq L_p \|\mathbf{p} - \dot{\mathbf{p}}\|, \end{aligned}$$

834 where L_p is defined in (54).

835 The reminder of proof is similar to the proof of [19, Lemma 10]. Since $\mathcal{L}(\mathbf{p})$ is L_p -smooth, we get

$$\begin{aligned} \mathcal{L}(\mathbf{p}(t+1)) &\leq \mathcal{L}(\mathbf{p}(t)) + \nabla \mathcal{L}(\mathbf{p}(t))^\top (\mathbf{p}(t+1) - \mathbf{p}(t)) + \frac{L_p}{2} \|\mathbf{p}(t+1) - \mathbf{p}(t)\|^2 \\ &= \mathcal{L}(\mathbf{p}(t)) - \eta \|\nabla \mathcal{L}(\mathbf{p}(t))\|^2 + \frac{L_p \eta^2}{2} \|\nabla \mathcal{L}(\mathbf{p}(t))\|^2 \\ &= \mathcal{L}(\mathbf{p}(t)) - \eta \left(1 - \frac{L_p \eta}{2} \right) \|\nabla \mathcal{L}(\mathbf{p}(t))\|^2 \\ &= \mathcal{L}(\mathbf{p}(t)) - \frac{\eta}{2} \|\nabla \mathcal{L}(\mathbf{p}(t))\|^2, \end{aligned}$$

836 where the last inequality follows from our assumption on the stepsize.

837 The above inequality implies that

$$\sum_{t=0}^{\infty} \|\nabla \mathcal{L}(\mathbf{p}(t))\|^2 \leq \frac{2}{\eta} (\mathcal{L}(\mathbf{p}(0)) - \mathcal{L}^*).$$

838 Here, the right hand side is upper bounded by a finite constant, since by Assumption A, $\mathcal{L}(\mathbf{p}(0)) < \infty$
 839 and $\mathcal{L}^* \leq \mathcal{L}(\mathbf{p}(t))$. This implies $\sum_{t=0}^{\infty} \|\nabla \mathcal{L}(\mathbf{p}(t))\|^2 < \infty$ and therefore $\|\nabla \mathcal{L}(\mathbf{p}(t))\|^2 \rightarrow 0$. ■

840 B.3 Proof of Theorem 1

841 **Proof.** We first show that $\lim_{t \rightarrow \infty} \|\mathbf{p}(t)\| = \infty$. From Lemma 4, we have

$$\langle \nabla \mathcal{L}(\mathbf{p}), \mathbf{p}^{mm\star} \rangle = \frac{1}{n} \sum_{i=1}^n \ell'_i \cdot \langle \mathbf{K}_i \mathbf{p}^{mm\star}, \mathbb{S}'(\mathbf{a}_i) \gamma_i \rangle,$$

842 where $\gamma_i = Y_i \cdot \mathbf{X}_i \mathbf{v}$, $\mathbf{a}_i = \mathbf{K}_i \mathbf{p}$, and $\ell'_i = \ell'(Y_i \cdot \mathbf{v}^\top \mathbf{X}_i^\top \mathbb{S}(\mathbf{K}_i \mathbf{p}))$.

843 It follows from Lemma 4 that $\langle \nabla \mathcal{L}(\mathbf{p}), \mathbf{p}^{mm\star} \rangle < 0$ for all $\mathbf{p} \in \mathbb{R}^d$. Hence, for any finite \mathbf{p} ,
 844 $\langle \nabla \mathcal{L}(\mathbf{p}), \mathbf{p}^{mm\star} \rangle$ cannot be equal to zero, as a sum of negative terms. Therefore, there are no finite
 845 critical points \mathbf{p} , for which $\nabla \mathcal{L}(\mathbf{p}) = 0$ which contradicts Lemma 6. This implies that $\|\mathbf{p}(t)\| \rightarrow \infty$.

846 Now, given any $\epsilon \in (0, 1)$, let $\pi = \epsilon/(1 - \epsilon)$. Since $\lim_{t \rightarrow \infty} \|\mathbf{p}(t)\| = \infty$, we can choose t_0 such that for
 847 any $t \geq t_0$, it holds that $\|\mathbf{p}(t)\| > R_\epsilon \vee 1/2$ for some radius R_ϵ . Now for any $t \geq t_0$, it follows from
 848 Lemma 5 that

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

849 Multiplying both sides by the stepsize η and using the gradient descent update, we get

$$\begin{aligned} \left\langle \mathbf{p}(t+1) - \mathbf{p}(t), \frac{\mathbf{p}^{mm\star}}{\|\mathbf{p}^{mm\star}\|} \right\rangle &\geq (1 - \epsilon) \left\langle \mathbf{p}(t+1) - \mathbf{p}(t), \frac{\mathbf{p}(t)}{\|\mathbf{p}(t)\|} \right\rangle \\ &= \frac{(1 - \epsilon)}{2\|\mathbf{p}(t)\|} \left(\|\mathbf{p}(t+1)\|^2 - \|\mathbf{p}(t)\|^2 - \|\mathbf{p}(t+1) - \mathbf{p}(t)\|^2 \right) \\ &\geq (1 - \epsilon) \left(\|\mathbf{p}(t+1)\|^2 - \|\mathbf{p}(t)\|^2 - \|\mathbf{p}(t+1) - \mathbf{p}(t)\|^2 \right) \\ &\geq (1 - \epsilon) \left(\|\mathbf{p}(t+1)\| - \|\mathbf{p}(t)\| - \|\mathbf{p}(t+1) - \mathbf{p}(t)\|^2 \right) \\ &\geq (1 - \epsilon) \left(\|\mathbf{p}(t+1)\| - \|\mathbf{p}(t)\| - 2\eta(\mathcal{L}(\mathbf{p}(t)) - \mathcal{L}(\mathbf{p}(t+1))) \right). \end{aligned} \quad (56)$$

850 Here, the last inequality uses Lemma 6.

851 Summing the above inequality over $t \geq t_0$ gives

$$\left\langle \frac{\mathbf{p}(t)}{\|\mathbf{p}(t)\|}, \frac{\mathbf{p}^{mm\star}}{\|\mathbf{p}^{mm\star}\|} \right\rangle \geq 1 - \epsilon + \frac{C(\epsilon, \eta)}{\|\mathbf{p}(t)\|},$$

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

$$C(\epsilon, \eta) := \left\langle \mathbf{p}(t_0), \frac{\mathbf{p}^{mm\star}}{\|\mathbf{p}^{mm\star}\|} \right\rangle - (1 - \epsilon)\|\mathbf{p}(t_0)\| - 2\eta(1 - \epsilon)(\mathcal{L}(\mathbf{p}(t_0)) - \mathcal{L}^*), \quad (57)$$

853 where $\mathcal{L}^* \leq \mathcal{L}(\mathbf{p}(t))$ for all $t \geq 0$.

854 Since $\|\mathbf{p}(t)\| \rightarrow \infty$, we get

$$\liminf_{t \rightarrow \infty} \left\langle \frac{\mathbf{p}(t)}{\|\mathbf{p}(t)\|}, \frac{\mathbf{p}^{mm\star}}{\|\mathbf{p}^{mm\star}\|} \right\rangle \geq 1 - \epsilon.$$

855 Given that ϵ is arbitrary, we can consider the limit as ϵ approaches zero. Thus, we have: $\mathbf{p}(t)/\|\mathbf{p}(t)\| \rightarrow$
 856 $\mathbf{p}^{mm\star}/\|\mathbf{p}^{mm\star}\|$. ■

857 B.4 Proof of Theorem 2

858 This proof is a direct corollary of Lemma 9 which itself is a special case of the nonlinear head
 859 Theorem 7. Let us verify that $f(X) = \mathbf{v}^\top X^\top \mathbb{S}(X\mathbf{p})$ satisfies the assumptions of Lemma 9 where
 860 we replace the nonlinear head with linear \mathbf{v} . To see this, set the optimal sets to be the singletons
 861 $O_i = \{\text{opt}_i\}$, given (X_i, Y_i) and defining $s_i = \mathbb{S}(K_i \mathbf{p})$ and $q_i := q_i^p = \sum_{t \neq \text{opt}_i} \gamma_{it} s_{it}$. Recalling score
 862 definition $\gamma_i = Y_i \cdot X_i \mathbf{v}$ and setting $v_i := \gamma_{i \text{opt}_i}$ and $Z_i := \sum_{t \neq \text{opt}_i} \gamma_{it} s_{it}$, a particular prediction can be
 863 written as

$$Y_i \cdot \mathbf{v}^\top X_i^\top \mathbb{S}(X_i \mathbf{p}) = \gamma_i^\top s_i = \gamma_{i \text{opt}_i} (1 - q_i) + \sum_{t \neq \text{opt}_i} \gamma_{it} s_{it} \quad (58)$$

$$= v_i (1 - q_i) + Z_i. \quad (59)$$

864 To proceed, we demonstrate the choices for $C, \varepsilon > 0$. Let $C := -\min_{i \in [n], t \in [T]} \gamma_{it} \wedge 0$ and $q_{\max} =$
 865 $\max_{i \in [n]} q_i$. Note that $Z_i \geq \sum_{t \neq \text{opt}_i} \gamma_{it} s_{it} \geq q_i \gamma_{\min} \geq -C q_{\max}$. Now, using strict score optimality of
 866 opt_i 's for all $i \in [n]$, we set

$$\varepsilon := 1 - \sup_{i \in [n]} \frac{\sum_{t \neq \text{opt}_i} \gamma_{it} s_{it}}{v_i q_i} \geq 1 - \sup_{i \in [n]} \frac{\sup_{t \neq \text{opt}_i} \gamma_{it}}{\gamma_{i \text{opt}_i}} > 0.$$

867 We conclude by observing $Z_i \leq v_i q_i \frac{\sum_{t \neq \text{opt}_i} \gamma_{it} s_{it}}{v_i q_i} \leq v_i q_i \varepsilon$ as desired.

868 B.5 Proof of Theorem 3

869 **Proof.** We provide the proof in four steps:

870 **Step 1: There are no stationary points within the cone.** We begin by proving that there are no
 871 stationary points within $\text{cone}_\mu(\mathbf{p}^{\text{mm}}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| \geq R_\mu\}$ for a specific radius R_μ . Let $(\mathcal{T}_i)_{i=1}^n$ denote the
 872 set of SVM-neighbors as defined in Definition 2. We define $\bar{\mathcal{T}}_i = [T] - \mathcal{T}_i - \alpha_i$ as the tokens that
 873 are non-SVM neighbors. Additionally, let μ be defined as in (23). For all $\mathbf{q}, \mathbf{p} \in \text{cone}_\mu(\mathbf{p}^{\text{mm}})$ with
 874 $\|\mathbf{q}\| = \|\mathbf{p}^{\text{mm}}\|$, it follows from Lemma 3 that there exists R_μ such that $-\mathbf{q}^\top \nabla \mathcal{L}(\mathbf{p})$ is strictly positive
 875 for $\|\mathbf{p}\| \geq R_\mu$.

876 **Step 2:** Let $\epsilon \in (0, \min(\mu, 1))$, $1/(1 + \pi) = 1 - \epsilon$. It follows from Lemma 5 that, there exists R_ϵ such
 877 that all $\mathbf{p} \in \text{cone}_\mu(\mathbf{p}^{\text{mm}}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| \geq R_\epsilon\}$ satisfy

$$\left\langle -\nabla \mathcal{L}(\mathbf{p}), \frac{\mathbf{p}^{\text{mm}}(\alpha)}{\|\mathbf{p}^{\text{mm}}(\alpha)\|} \right\rangle \geq (1 - \epsilon) \left\langle -\nabla \mathcal{L}(\mathbf{p}), \frac{\mathbf{p}}{\|\mathbf{p}\|} \right\rangle. \quad (60)$$

878
 879 **Step 3: Updates remain inside the cone.** By leveraging the results from Step 1 and Step
 880 2, we show that that the gradient iterates, with an appropriate step size, starting from $\mathbf{p}(0) \in$
 881 $\text{cone}_\mu(\mathbf{p}^{\text{mm}}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| \geq R\}$, remain within this cone.

882 We proceed by induction. Suppose that the claim holds up to iteration $t \geq 0$. This implies that
 883 $\mathbf{p}(t) \in \text{cone}_\mu(\mathbf{p}^{\text{mm}}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| \geq R\}$. Hence, there exists scalar $\mu = \mu(\alpha) \in (0, 1]$ and R_μ such that
 884 $\text{corr}(\mathbf{p}(t), \mathbf{p}^{\text{mm}}(\alpha)) \geq 1 - \mu$ and $\|\mathbf{p}(t)\| \geq R_\mu$. Let $\rho := -(1/(1 - \epsilon)) \langle \nabla \mathcal{L}(\mathbf{p}(t)), \frac{\mathbf{p}^{\text{mm}}(\alpha)}{\|\mathbf{p}^{\text{mm}}(\alpha)\|} \rangle > 0$. We have

$$\begin{aligned} \left\langle \frac{\mathbf{p}(t+1)}{\|\mathbf{p}(t+1)\|}, \frac{\mathbf{p}^{\text{mm}}(\alpha)}{\|\mathbf{p}^{\text{mm}}(\alpha)\|} \right\rangle &= \left\langle \frac{\mathbf{p}(t)}{\|\mathbf{p}(t)\|} - \frac{\eta}{\|\mathbf{p}(t)\|} \nabla \mathcal{L}(\mathbf{p}(t)), \frac{\mathbf{p}^{\text{mm}}(\alpha)}{\|\mathbf{p}^{\text{mm}}(\alpha)\|} \right\rangle \\ &\geq 1 - \mu - \frac{\eta}{\|\mathbf{p}(t)\|} \left\langle \nabla \mathcal{L}(\mathbf{p}(t)), \frac{\mathbf{p}^{\text{mm}}(\alpha)}{\|\mathbf{p}^{\text{mm}}(\alpha)\|} \right\rangle \\ &\geq 1 - \mu + \frac{\eta \rho (1 - \epsilon)}{\|\mathbf{p}(t)\|}. \end{aligned} \quad (61a)$$

885 Note that from Lemma 3, we have $\langle \nabla f(\mathbf{p}(t)), \mathbf{p}(t) \rangle < 0$ which implies that $\|\mathbf{p}(t+1)\| \leq \|\mathbf{p}(t)\| -$
 886 $\frac{\eta}{\|\mathbf{p}(t)\|} \langle \nabla f(\mathbf{p}(t)), \mathbf{p}(t) \rangle + \eta^2 \|\nabla f(\mathbf{p}(t))\|^2$. Hence, $\|\mathbf{p}(t+1)\| \geq \|\mathbf{p}(t)\|$, and

$$\begin{aligned} \frac{\|\mathbf{p}(t+1)\|}{\|\mathbf{p}(t)\|} &\leq 1 - \eta \left\langle \nabla f(\mathbf{p}(t)), \frac{\mathbf{p}(t)}{\|\mathbf{p}(t)\|} \right\rangle + \eta^2 \frac{\|\nabla \mathcal{L}(\mathbf{p}(t))\|^2}{\|\mathbf{p}(t)\|} \\ &\leq 1 - \frac{\eta}{1 - \epsilon} \left\langle \nabla \mathcal{L}(\mathbf{p}(t)), \frac{\mathbf{p}^{\text{mm}}(\alpha)}{\|\mathbf{p}^{\text{mm}}(\alpha)\|} \right\rangle + \eta^2 \frac{\|\nabla \mathcal{L}(\mathbf{p}(t))\|^2}{\|\mathbf{p}(t)\|} \\ &\leq 1 + \frac{\eta \rho}{\|\mathbf{p}(t)\|} + \frac{\eta^2 \|\nabla \mathcal{L}(\mathbf{p}(t))\|^2}{\|\mathbf{p}(t)\|} =: C(\eta, \rho). \end{aligned} \quad (61b)$$

887 Here, the second inequality follows from (60).

888 Now, it follows from (61a) and (61b) that

$$\begin{aligned}
\left\langle \frac{\mathbf{p}(t+1)}{\|\mathbf{p}(t+1)\|}, \frac{\mathbf{p}^{mm}(\alpha)}{\|\mathbf{p}^{mm}(\alpha)\|} \right\rangle &\geq \frac{1}{C(\eta, \rho)} \left(1 - \mu + \frac{\eta\rho(1-\epsilon)}{\|\mathbf{p}(t)\|} \right) \\
&\geq \frac{1}{C(\eta, \rho)} \left(1 - \mu + \frac{\eta\rho(1-\epsilon)}{\|\mathbf{p}(t)\|} \right) \\
&\geq 1 - \mu + \frac{\eta}{C(\eta, \rho)} \left(\frac{\rho(\mu - \epsilon)}{\|\mathbf{p}(t)\|} - \eta(1 - \mu) \frac{\|\nabla \mathcal{L}(\mathbf{p}(t))\|^2}{\|\mathbf{p}(t)\|} \right) \\
&\geq 1 - \mu,
\end{aligned} \tag{62}$$

889 where the last inequality uses $\eta \leq \frac{(\mu-\epsilon)\rho}{1-\mu} \frac{1}{\|\nabla f(\mathbf{p}(t))\|^2}$.

890 Hence, $\mathbf{p}(t+1) \in \text{cone}_\mu(\mathbf{p}^{mm}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| \geq R_\mu\}$.

891 **Step 4: The correlation of $\mathbf{p}(t)$ and $\mathbf{p}^{mm}(\alpha)$ increases over t .** The reminder is similar to the proof of
892 Theorem 1. Note that it follows from Lemma 4 that $\langle \nabla \mathcal{L}(\mathbf{p}), \mathbf{p}^{mm}(\alpha)/\|\mathbf{p}^{mm}(\alpha)\| \rangle < 0$, for any finite \mathbf{p} .
893 Hence, there are no finite critical points \mathbf{p} , for which $\nabla \mathcal{L}(\mathbf{p}) = 0$ which contradicts Lemma 6. This
894 implies that $\|\mathbf{p}(t)\| \rightarrow \infty$. Hence, we can choose t_0 such that for any $t \geq t_0$, it holds that $\|\mathbf{p}(t)\| > R$
895 for some $R \geq R_\mu \vee R_\epsilon \vee 1/2$. Now, following similar steps in (56) and (57), we obtain

$$\left\langle \frac{\mathbf{p}(t)}{\|\mathbf{p}(t)\|}, \frac{\mathbf{p}^{mm}(\alpha)}{\|\mathbf{p}^{mm}(\alpha)\|} \right\rangle \geq 1 - \epsilon + \frac{C(\epsilon, \eta)}{\|\mathbf{p}(t)\|},$$

896 for some finite constant $C(\epsilon, \eta)$.

897 Consequently,

$$\liminf_{t \rightarrow \infty} \left\langle \frac{\mathbf{p}(t)}{\|\mathbf{p}(t)\|}, \frac{\mathbf{p}^{mm}(\alpha)}{\|\mathbf{p}^{mm}(\alpha)\|} \right\rangle \geq 1 - \epsilon.$$

898 Since $\epsilon \in (0, \min(\mu, 1))$ is arbitrary, we get $\mathbf{p}(t)/\|\mathbf{p}(t)\| \rightarrow \mathbf{p}^{mm}(\alpha)/\|\mathbf{p}^{mm}(\alpha)\|$. ■

899 B.6 Proof of Theorem 4: Regularization Path Fails for Non-Locally-Optimal Tokens

900 The theorem below is essentially a restatement of Theorem 4 and shows that regularization path does
901 not converge to the max-margin solution if token indices α does not satisfy Definition 2. The only
902 difference is that, Theorem 4 replaces the second condition below with a cleaner statement which
903 assumes the linear-independence of the support vectors.

904 **Theorem 8 (Failure of Local Regularization Path)** Fix token indices $\alpha = (\alpha_i)_{i=1}^n$ with
905 SVM-neighbors $(\mathcal{T}_i)_{i=1}^n$. Suppose for some $j \in [n]$, there exists an SVM-neighbor $\beta \in \mathcal{T}_j$ sat-
906 isfying the following:

- 907 • $\mathbf{x}_{j\beta}$ has a higher score than $\mathbf{x}_{j\alpha_j}$: $Y_j \cdot \mathbf{v}^\top \mathbf{x}_{j\beta} > Y_j \cdot \mathbf{v}^\top \mathbf{x}_{j\alpha_j}$.
- 908 • Recall $\mathbf{p}^{mm} = \mathbf{p}^{mm}(\alpha)$ be the solution of (ATT-SVM) and let \mathbf{p}^β be the solution of
909 (ATT-SVM) where the constraint $(\mathbf{k}_{j\alpha_j} - \mathbf{k}_{j\beta})^\top \mathbf{p} \geq 1$ is not enforced. β is an active
910 SVM-neighbor in the sense that \mathbf{p}^β violates the constraint i.e. $(\mathbf{k}_{j\alpha_j} - \mathbf{k}_{j\beta})^\top \mathbf{p}^\beta < 1$.

911 For any $\epsilon > 0$, there exists $R_\epsilon > 0$ as follows: Consider the neighborhood of \mathbf{p}^{mm} : $C_\epsilon = \text{cone}_\epsilon(\mathbf{p}^{mm})$
912 $\cap \{\mathbf{p} \mid \|\mathbf{p}\| \geq R_\epsilon\}$. Define the local path $\tilde{\mathbf{p}}(R) = \min_{\mathbf{p} \in C_\epsilon, \|\mathbf{p}\| \leq R} \mathcal{L}(\mathbf{p})$. Then $\lim_{R \rightarrow \infty} \frac{\tilde{\mathbf{p}}(R)}{\|\tilde{\mathbf{p}}(R)\|} \neq \frac{\mathbf{p}^{mm}}{\|\mathbf{p}^{mm}\|}$.

913 **Proof of Theorem 4:** Using the above theorem we can now conclude the proof of Theorem 4 by
914 showing that, second bullet of Theorem 4 implies the second bullet of Theorem 8. We are given
915 solution \mathbf{p}^{mm} and \mathbf{p}^β . Suppose that \mathbf{p}^β in Theorem 4 does not violate the constraint $(\mathbf{k}_{j\alpha_j} - \mathbf{k}_{j\beta})^\top \mathbf{p}^\beta \geq 1$.
916 Then, it would imply that $\mathbf{p}^\beta = \mathbf{p}^{mm}$ because \mathbf{p}^β satisfies all margin constraints and $\|\mathbf{p}^\beta\| \leq \|\mathbf{p}^{mm}\|$
917 (because it solves the problem with less constraints), thus, if $\mathbf{p}^\beta \neq \mathbf{p}^{mm}$, it would contradict with the
918 optimality of \mathbf{p}^{mm} . Since the active constraints are linearly independent, their Lagrange multipliers
919 are unique. Since \mathbf{p}^β is missing a linearly independent constraint, the solution \mathbf{p}^β expressed in terms
920 of Lagrange-weighted constraints cannot equate to the solution \mathbf{p}^{mm} expressed in terms of its own
921 Lagrange-weighted constraints that also include the constraint induced by $\mathbf{k}_{j\alpha_j} - \mathbf{k}_{j\beta}$.

922 B.6.1 Proof of Theorem 8

923 **Proof strategy:** Without losing generality, let us prove the result for 2ε (to simplify the downstream
 924 notation). To accomplish the proof, we will follow the following strategy. Fix $\bar{\mathbf{p}}_\varepsilon^{mm} = \frac{\varepsilon \mathbf{p}^\beta + (1-\varepsilon) \mathbf{p}^{mm}}{\|\varepsilon \mathbf{p}^\beta + (1-\varepsilon) \mathbf{p}^{mm}\|}$
 925 and $\bar{\mathbf{p}}^{mm} = \mathbf{p}^{mm} / \|\mathbf{p}^{mm}\|$. Using $\|\mathbf{p}^\beta\| \leq \|\mathbf{p}^{mm}\|$, we observe that $\bar{\mathbf{p}}_\varepsilon^{mm}$ obeys the correlation inequality

$$(\bar{\mathbf{p}}_\varepsilon^{mm})^\top \bar{\mathbf{p}}^{mm} \geq \frac{(1-\varepsilon)\|\mathbf{p}^{mm}\|^2 - \varepsilon\|\mathbf{p}^{mm}\|\|\mathbf{p}^\beta\|}{\|\mathbf{p}^{mm}\|((1-\varepsilon)\|\mathbf{p}^{mm}\| + \varepsilon\|\mathbf{p}^\beta\|)} \geq \frac{(1-2\varepsilon)\|\mathbf{p}^{mm}\|^2}{\|\mathbf{p}^{mm}\|^2} \geq 1-2\varepsilon.$$

926 This establishes that $r \cdot \bar{\mathbf{p}}_\varepsilon^{mm} \in \text{cone}_{2\varepsilon}(\mathbf{p}^{mm})$. Thus, we will use $\bar{\mathbf{p}}_\varepsilon^{mm}$ to show that it is a superior
 927 direction to \mathbf{p}^{mm} . Concretely, for all $R \geq R_\varepsilon$, suppose that, there exists $\delta = \delta(\varepsilon)$ such that,

$$\mathcal{L}(R \cdot \bar{\mathbf{p}}_\varepsilon^{mm}) < \inf_{\|\mathbf{p}\|=R, \mathbf{p} \in \text{cone}_\delta(\mathbf{p}^{mm})} \mathcal{L}(\mathbf{p}). \quad (63)$$

928 In words, suppose that $R \cdot \bar{\mathbf{p}}_\varepsilon^{mm}$ achieves strictly better loss than all points of ℓ_2 -norm R within
 929 $\text{cone}_\delta(\mathbf{p}^{mm})$. Establishing this would imply the desired result $\lim_{R \rightarrow \infty} \frac{\bar{\mathcal{L}}(R)}{\|\bar{\mathbf{p}}(R)\|} \neq \frac{\mathbf{p}^{mm}}{\|\mathbf{p}^{mm}\|}$. Since for any choice
 930 of $R \geq R_\varepsilon$, (63) implies that the optimal direction $\frac{\bar{\mathbf{p}}(R)}{\|\bar{\mathbf{p}}(R)\|}$ is at least δ bounded away from $\frac{\mathbf{p}^{mm}}{\|\mathbf{p}^{mm}\|}$. In
 931 what follows, we will prove this by establishing (63).

932 First, let us establish the critical properties of \mathbf{p}^β . Set $K = n(T-1)$ and gather the set of margin
 933 equalities $\mathbf{p}^{mm} := \mathbf{p}^{mm}(\alpha)$ satisfies: These inequalities are given by vectors $(\mathbf{v}_{k=1}^K)$ where \mathbf{v}_k is the
 934 form $\mathbf{k}_{\alpha_i} - \mathbf{k}_t$ for $t \neq \alpha_i$. Also let $\mathbf{v}_1 = \mathbf{k}_{j_{\alpha_j}} - \mathbf{k}_{j_\beta}$ be the active constraint described in the theorem.

935 Note that $\|\mathbf{p}^\beta\| \leq \|\mathbf{p}^{mm}\|$ since \mathbf{p}^β is solving a max-margin problem with strictly less constraints (over
 936 $k \geq 2$). Secondly, we claim that \mathbf{p}^β achieves a strictly larger margin compared to \mathbf{p}^{mm} over $k \geq 2$,
 937 namely setting $\Gamma = \|\mathbf{p}^{mm}\|$ and $\Gamma_\beta = \|\mathbf{p}^\beta\|$

$$\min_{k \geq 2} \mathbf{v}_k^\top \mathbf{p}^\beta / \|\mathbf{p}^\beta\| = 1/\Gamma_\beta > \min_{k \geq 2} \mathbf{v}_k^\top \mathbf{p}^{mm} / \|\mathbf{p}^{mm}\| = 1/\Gamma.$$

938 If not, it would imply that $\|\mathbf{p}^\beta\| = \|\mathbf{p}^{mm}\|$ and that $\min_{k \geq 2} \mathbf{v}_k^\top \mathbf{p}^{mm} = \min_{k \geq 2} \mathbf{v}_k^\top \mathbf{p}^\beta$. Since theorem's
 939 statement guarantees $\mathbf{p}^{mm} \neq \mathbf{p}^\beta$, this contradicts with the unique optimality of \mathbf{p}^β when satisfying
 940 constraints $k \geq 2$ as \mathbf{p}^{mm} would achieve the same objective.

941 Finally, using same argument, we also note that, \mathbf{p}^β achieves strictly less margin over \mathbf{v}_1 , namely

$$\mathbf{v}_1^\top \mathbf{p}^\beta / \|\mathbf{p}^\beta\| < \mathbf{v}_1^\top \mathbf{p}^{mm} / \|\mathbf{p}^{mm}\| = 1/\Gamma. \quad (64)$$

942 If not, it would imply that \mathbf{p}^β achieves a better or equal margin at all constraints which would
 943 contradict with the optimality of \mathbf{p}^{mm} over constraints $k \geq 1$.

944 Now, let us define $\mathbf{p}_\varepsilon^{mm} = \varepsilon \mathbf{p}^\beta + (1-\varepsilon) \mathbf{p}^{mm}$ and observe that $\mathbf{p}_\varepsilon^{mm}$ also satisfies the discussion above.
 945 Namely, using $\|\mathbf{p}_\varepsilon^{mm}\| \leq \varepsilon \|\mathbf{p}^\beta\| + (1-\varepsilon) \|\mathbf{p}^{mm}\| < \Gamma$ we find

$$\frac{\min_{k \geq 2} \mathbf{v}_k^\top \mathbf{p}_\varepsilon^{mm}}{\|\mathbf{p}_\varepsilon^{mm}\|} \geq \frac{1}{\varepsilon \|\mathbf{p}^\beta\| + (1-\varepsilon) \|\mathbf{p}^{mm}\|} > \frac{1}{\Gamma}. \quad (65)$$

946 Similarly, on constraint \mathbf{v}_1 , we have that

$$\mathbf{v}_1^\top \mathbf{p}_\varepsilon^{mm} / \|\mathbf{p}_\varepsilon^{mm}\| := 1/\Gamma_\varepsilon < 1/\Gamma. \quad (66)$$

947 If not, it would imply that $\mathbf{p}_\varepsilon^{mm}$ achieves a better or equal margin on all constraints which would
 948 contradict with the unique optimality of \mathbf{p}^{mm} over constraints $k \geq 1$.

949 We will use $\Gamma_\varepsilon > \Gamma$, (66), and (65) to conclude that $\mathbf{p}_\varepsilon^{mm}$ is a strictly better direction compared to a
 950 $\delta = \delta(\varepsilon)$ conic neighborhood of $\bar{\mathbf{p}}^{mm} := \mathbf{p}^{mm} / \|\mathbf{p}^{mm}\|$. Pick the δ neighborhood of \mathbf{p}^{mm} such that, all
 951 \mathbf{p}_δ^{mm} it satisfies

$$\mathbf{v}_k^\top \mathbf{p}_\delta^{mm} / \|\mathbf{p}_\delta^{mm}\| \geq 1/\Gamma_\delta = 0.5(1/\Gamma_\varepsilon + 1/\Gamma) > 1/\Gamma_\varepsilon \quad \text{for all } k \in [K]. \quad (67)$$

952 In words, we choose a neighborhood with correlation profile dominated by $\mathbf{p}_\varepsilon^{mm}$ (on $k = 1$ and $k \geq 2$).
 953 We now lower bound the loss function over δ -neighborhood and upper bound over $\mathbf{p}_\varepsilon^{mm}$. Specifically,
 954 we will compare a \mathbf{p}_δ^{mm} within the δ neighborhood of \mathbf{p}^{mm} with $\|\mathbf{p}_\delta^{mm}\| = R$ and $\bar{\mathbf{p}}_\varepsilon^{mm} := R \cdot \bar{\mathbf{p}}_\varepsilon^{mm}$. To
 955 proceed, define:

$$q_i^* = 1 - \mathbb{S}(\mathbf{K}_i \bar{\mathbf{p}}_\varepsilon^{mm})_{\alpha_i}, \quad \hat{q}_i = 1 - \mathbb{S}(\mathbf{K}_i \mathbf{p}_\delta^{mm})_{\alpha_i}$$

956 Also define $q^\beta = \mathbb{S}(\mathbf{K}_j \tilde{\mathbf{p}}_\varepsilon^{mm})_\beta$ which is the j 'th softmax likelihood at token β . We will use the fact that
 957 margin at β is small for j 'th example to lower bound q^β carefully. We next bound these as follows
 958 based on (66), (65), (67) (e.g. following derivation of (109))²

$$\log(q_i^\star) \leq -(R/\Gamma) + \log T \quad \text{for all } i \neq j \quad (68)$$

$$\log(q^\beta) \geq -(R/\Gamma_\varepsilon) - \log T \quad \text{for all } i \neq j \quad (69)$$

$$\log(q_j^\star - q^\beta) \leq -(R/\Gamma) + \log T \quad \text{for all } i \neq j \quad (70)$$

$$\log(\hat{q}_i) \leq -(R/\Gamma_\delta) + \log T \quad \text{for all } i \in [n]. \quad (71)$$

959 • **Lower bounding $\mathcal{L}(\mathbf{p}_\delta^{mm})$:** Using the last inequality, on \mathbf{p}_δ^{mm} (within the δ neighborhood of \mathbf{p}^{mm}),
 960 we have the following lower bound: Set $\mathbf{x}_i^\delta = \mathbf{X}_i^\top \mathbb{S}(\mathbf{K}_i \mathbf{p}_\delta^{mm})$ and $M := \sup_{i \in [n], t, \tau \in [T]} \|\mathbf{x}_{it} - \mathbf{x}_{i\tau}\|$ and
 961 note that $\|\mathbf{x}_i^\delta - \mathbf{x}_{\alpha_i}\| \leq M \hat{q}_i$. Also let B and A be the lower and upper bound of $-\ell'$ over $[-M\|\mathbf{v}\|, M\|\mathbf{v}\|]$
 962 interval. Finally, define $\mathcal{L}_\star = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{v}^\top \mathbf{x}_{i\alpha_i})$. We find

$$|\mathcal{L}(\mathbf{p}_\delta^{mm}) - \mathcal{L}_\star| = \frac{1}{n} \sum_{i=1}^n |\ell(Y_i \cdot \mathbf{v}^\top \mathbf{x}_i^\delta) - \ell(Y_i \cdot \mathbf{v}^\top \mathbf{x}_{i\alpha_i})| \quad (72)$$

$$\leq B \hat{q}_{\max} M \|\mathbf{v}\| \quad (73)$$

$$\leq T B M \|\mathbf{v}\| e^{-R/\Gamma_\delta}. \quad (74)$$

963 This implies $\mathcal{L}(\mathbf{p}_\delta^{mm}) \geq \mathcal{L}_\star - T B M \|\mathbf{v}\| e^{-R/\Gamma_\delta}$. Note that this holds for all \mathbf{p}_δ^{mm} within the conic
 964 neighborhood $\text{cone}_\delta(\mathbf{p}^{mm}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| = R\}$ (defined above (63)).

965 • **Upper bounding $\mathcal{L}(\tilde{\mathbf{p}}_\varepsilon^{mm})$:** On $\mathbf{p}_\varepsilon^{mm}$, we upper bound the loss as follows. Define the loss $\mathcal{L}^{-j}(\mathbf{p}) =$
 966 $\frac{1}{n} \sum_{i \neq j} \ell(Y_i \cdot f(\mathbf{X}_i))$ i.e. loss over all training data except the j 'th one. Repeating the argument identical
 967 to (72), we find that

$$|\mathcal{L}^{-j}(\tilde{\mathbf{p}}_\varepsilon^{mm}) - \mathcal{L}_\star^{-j}| \leq T B M \|\mathbf{v}\| e^{-R/\Gamma}.$$

968 The critical term is the j 'th loss which we need to upper bound as follows. Set $\mathbf{x}_j^\varepsilon = \mathbf{X}_j^\top \mathbb{S}(\mathbf{K}_j \tilde{\mathbf{p}}_\varepsilon^{mm})$
 969 and define the score improvement by β to be $\gamma_{\text{gap}} = Y_j \cdot \mathbf{v}^\top (\mathbf{x}_{j\beta} - \mathbf{x}_{j\alpha_j}) > 0$. We note that

$$Y_j \cdot \mathbf{v}^\top \mathbf{x}_j^\varepsilon - Y_j \cdot \mathbf{v}^\top \mathbf{x}_{\alpha_j} = q^\beta Y_j \cdot \mathbf{v}^\top (\mathbf{x}_{j\beta} - \mathbf{x}_{\alpha_j}) + \sum_{t \notin \{\alpha_j, \beta\}} \mathbb{S}(\mathbf{K}_j \tilde{\mathbf{p}}_\varepsilon^{mm})_t Y_j \cdot \mathbf{v}^\top (\mathbf{x}_{it} - \mathbf{x}_{\alpha_j}) \quad (75)$$

$$\geq q^\beta \gamma_{\text{gap}} - (q_j^\star - q^\beta) M \|\mathbf{v}\| \quad (76)$$

$$\geq T^{-1} \gamma_{\text{gap}} e^{-R/\Gamma_\varepsilon} - T B M \|\mathbf{v}\| e^{-R/\Gamma}. \quad (77)$$

970 Combining these into $\mathcal{L}(\tilde{\mathbf{p}}_\varepsilon^{mm}) = \mathcal{L}^{-j}(\tilde{\mathbf{p}}_\varepsilon^{mm}) + n^{-1} \ell(Y_j \cdot \mathbf{v}^\top \mathbf{x}_j^\varepsilon)$ and using $A \leq -\ell' \leq B$, we obtain the
 971 lower bound

$$\mathcal{L}(\tilde{\mathbf{p}}_\varepsilon^{mm}) - \mathcal{L}_\star \leq 2 T B M \|\mathbf{v}\| e^{-R/\Gamma} - A n^{-1} T^{-1} \gamma_{\text{gap}} e^{-R/\Gamma_\varepsilon}. \quad (78)$$

972 In conclusion, we find that $\mathcal{L}(\tilde{\mathbf{p}}_\varepsilon^{mm}) > \mathcal{L}(\mathbf{p}_\delta^{mm})$ whenever

$$A n^{-1} T^{-1} \gamma_{\text{gap}} e^{-R/\Gamma_\varepsilon} > 2 T B M \|\mathbf{v}\| (e^{-R/\Gamma} + e^{-R/\Gamma_\delta}).$$

973 Using the relationship $\Gamma_\varepsilon \geq \Gamma_\delta \geq \Gamma$ and noticing $1/\Gamma_\delta - 1/\Gamma_\varepsilon = (1/\Gamma - 1/\Gamma_\varepsilon)/2$, this is implied by

$$e^{R(1/\Gamma_\delta - 1/\Gamma_\varepsilon)} > 4 T^2 n B M A^{-1} \|\mathbf{v}\| \iff R > \frac{2 \Gamma_\varepsilon \Gamma}{\Gamma_\varepsilon - \Gamma} \log(4 T^2 n B M A^{-1} \|\mathbf{v}\|).$$

974 Thus, as advertised, we found that, for any $\varepsilon > 0$, there exists R_ε such that, over the set
 975 $\text{cone}_\varepsilon(\mathbf{p}^{mm}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| \geq R_\varepsilon\}$, $\tilde{\mathbf{p}}_\varepsilon^{mm}$ with $\|\tilde{\mathbf{p}}_\varepsilon^{mm}\| = R > R_\varepsilon$ achieves smaller loss compared to $\mathcal{L}(\mathbf{p}_\delta^{mm})$
 976 for all $\mathbf{p}_\delta^{mm} \in \text{cone}_\delta(\mathbf{p}^{mm}) \cap \{\mathbf{p} \mid \|\mathbf{p}\| = R\}$. This in turn implies (63) for all $R > R_\varepsilon$ concluding the
 977 proof.

²We are essentially following identical arguments developed in the proofs of Theorem 7 or Theorem 5.

978 C Addendum to Section 3

979 C.1 Proof of Theorem 5

980 **Proof.** Suppose the claim is incorrect and either \mathbf{p}_R/R or \mathbf{v}_r/r fails to converge as R, r grows. Set
 981 $\Xi = 1/\|\mathbf{p}^{mm}\|$, $\tilde{\mathbf{p}}^{mm} = R\Xi\mathbf{p}^{mm}$ and $\tilde{\mathbf{v}}^{mm} = r\Gamma\mathbf{v}^{mm}$. The proof strategy is obtaining a contradiction by
 982 proving that $(\tilde{\mathbf{v}}^{mm}, \tilde{\mathbf{p}}^{mm})$ is a strictly better solution compared to $(\mathbf{v}_r, \mathbf{p}_R)$ for large R, r . Without losing
 983 generality, we will set $\alpha_i = 1$ for all $i \in [n]$ as the problem is invariant to tokens' permutation. Define
 984 $q_i^p = 1 - s_{i1}^p$ to be the amount of non-optimality (cumulative probability of non-first tokens) where
 985 $s_i^p = \mathbb{S}(\mathbf{K}_i\mathbf{p})$ is the softmax probabilities.

986 • **Case 1: \mathbf{p}_R/R does not converge.** Under this scenario there exists $\delta, \gamma = \gamma(\delta) > 0$ such that we can
 987 find arbitrarily large R with $\|\mathbf{p}_R/R - \tilde{\mathbf{p}}^{mm}/R\| \geq \delta$ and margin induced by \mathbf{p}_R/R is at most $\Xi(1 - \gamma)$
 988 (from strong convexity of (ATT-SVM)). Following q_i^p definition above, set $\hat{q}_{\max} = \sup_{i \in [n]} q_i^{p_R}$ to be
 989 worst non-optimality in \mathbf{p}_R and $q_{\max}^* = \sup_{i \in [n]} q_i^{\tilde{\mathbf{p}}^{mm}}$ to be the same for $\tilde{\mathbf{p}}^{mm}$. Repeating the identical
 990 argument in Theorem 7 (specifically (109)), we can bound the non-optimality amount q_i^* of \mathbf{p}_R^* as

$$q_i^* = \frac{\sum_{t \neq \alpha_i} \exp(\mathbf{k}_{it}^\top \mathbf{p}_R^*)}{\sum_{t \in [T]} \exp(\mathbf{k}_{it}^\top \mathbf{p}_R^*)} \leq \frac{\sum_{t \neq \alpha_i} \exp(\mathbf{k}_{it}^\top \mathbf{p}_R^*)}{\exp(\mathbf{k}_{i\alpha_i}^\top \mathbf{p}_R^*)} \leq T \exp(-R\Xi). \quad (79)$$

991 Thus, $q_{\max}^* = \max_{i \in [n]} q_i^* \leq T \exp(-R\Xi)$. Next without losing generality, assume first margin
 992 constraint is γ -violated by \mathbf{p}_R and $\min_{t \neq \alpha_1} (\mathbf{k}_{1\alpha_1} - \mathbf{k}_{1t})^\top \mathbf{p}_R \leq \Xi R(1 - \gamma)$. Denoting the amount of
 993 non-optimality of the first input as \hat{q}_1 , we find

$$\hat{q}_1 = \frac{\sum_{t \neq \alpha_1} \exp(\mathbf{k}_{1t}^\top \mathbf{p}_R)}{\sum_{t \in [T]} \exp(\mathbf{k}_{1t}^\top \mathbf{p}_R)} \geq \frac{1}{T} \frac{\sum_{t \neq \alpha_1} \exp(\mathbf{k}_{1t}^\top \mathbf{p}_R)}{\exp(\mathbf{k}_{1\alpha_1}^\top \mathbf{p}_R)} \geq T^{-1} \exp(-(1 - \gamma)R\Xi). \quad (80)$$

994 We similarly have $q_{\max}^* \geq T^{-1} \exp(-R\Xi)$ to find that

$$\log(\hat{q}_{\max}) \geq -(1 - \gamma)\Xi R - \log T, \quad (81)$$

$$-\Xi R - \log T \leq \log(q_{\max}^*) \leq -\Xi R + \log T. \quad (82)$$

995 In words, $\tilde{\mathbf{p}}^{mm}$ contains exponentially less non-optimality compared to \mathbf{p}_R as R grows. The remainder
 996 of the proof differs from Theorem 7 as we need to upper/lower bound the logistic loss of $(\tilde{\mathbf{v}}^{mm}, \tilde{\mathbf{p}}^{mm})$
 997 and $(\mathbf{v}_r, \mathbf{p}_R)$ respectively to conclude with the contradiction.

998 First, let us upper bound the logistic loss of $(\tilde{\mathbf{v}}^{mm}, \tilde{\mathbf{p}}^{mm})$. Set $\mathbf{r}_i = \mathbf{X}_i^\top \mathbb{S}(\mathbf{K}_i \tilde{\mathbf{p}}^{mm})$. Observe that
 999 if $\|\mathbf{r}_i - \mathbf{x}_{i1}\| \leq \varepsilon_i$, we have that \mathbf{v}^{mm} satisfies the SVM constraints on \mathbf{r}_i with $Y_i \cdot \mathbf{r}_i^\top \mathbf{v}^{mm} \geq 1 -$
 1000 ε_i/Γ . Consequently, setting $\varepsilon_{\max} = \sup_{i \in [n]} \varepsilon_i$, \mathbf{v}^{mm} achieves a label-margin of $\Gamma - \varepsilon_{\max}$ on the
 1001 dataset $(Y_i, \mathbf{r}_i)_{i \in [n]}$. With this, we upper bound the logistic loss of $(\tilde{\mathbf{v}}^{mm}, \tilde{\mathbf{p}}^{mm})$ as follows. Let
 1002 $M = \sup_{i \in [n], t \in [T]} \|\mathbf{x}_{it}\|$. In what follows, let us recall the fact (81) that worst-case perturbation is
 1003 $\varepsilon_{\max} \leq M \exp(-\Xi R + \log T) = MT \exp(-\Xi R)$.

$$\mathcal{L}(\tilde{\mathbf{v}}^{mm}, \tilde{\mathbf{p}}^{mm}) \leq \max_{i \in [n]} \log(1 + \exp(-Y_i \mathbf{r}_i^\top \tilde{\mathbf{v}}^{mm})). \quad (83)$$

$$\leq \max_{i \in [n]} \exp(-Y_i \mathbf{r}_i^\top \tilde{\mathbf{v}}^{mm}) \quad (84)$$

$$\leq \exp(-r\Gamma + r\varepsilon_{\max}) \quad (85)$$

$$\leq e^{rMT \exp(-\Xi R)} e^{-r\Gamma}. \quad (86)$$

1004 Conversely, we obtain a lower bound for $(\mathbf{v}_r, \mathbf{p}_R)$. Set $\mathbf{r}_i = \mathbf{X}_i^\top \mathbb{S}(\mathbf{K}_i \mathbf{p}_R)$. Using Assumption C, we find
 1005 that solving (CLS-SVM) on $(Y_i, \mathbf{r}_i)_{i \in [n]}$ achieves at most $\Gamma - \nu e^{-(1-\gamma)\Xi R}/T$ margin. Consequently, we
 1006 have

$$\mathcal{L}(\mathbf{v}_r, \mathbf{p}_R) \geq \frac{1}{n} \max_{i \in [n]} \log(1 + \exp(-Y_i \mathbf{r}_i^\top \mathbf{v}_r)) \quad (87)$$

$$\geq \frac{1}{2n} \max_{i \in [n]} \exp(-Y_i \mathbf{r}_i^\top \mathbf{v}_r) \wedge \log 2 \quad (88)$$

$$\geq \frac{1}{2n} \exp(-r(\Gamma - \nu e^{-(1-\gamma)\Xi R}/T)) \wedge \log 2 \quad (89)$$

$$\geq \frac{1}{2n} e^{r(\nu/T) \exp(-(1-\gamma)\Xi R)} e^{-r\Gamma} \wedge \log 2. \quad (90)$$

1007 Observe that, this lower bound dominates the previous upper bound when R is large, namely, when
 1008 (ignoring the multiplier $1/2n$ for brevity)

$$(\nu/T)e^{-(1-\gamma)\Xi R} \geq MT e^{-\Xi R} \iff R \geq R_0 := \frac{1}{\gamma\Xi} \log\left(\frac{MT^2}{\nu}\right).$$

1009 Thus, we indeed obtain the desired contradiction since such large R is guaranteed to exist when
 1010 $\mathbf{p}_R/R \rightarrow \mathbf{p}^{mm}$.

1011 • **Case 2: \mathbf{v}_r/r does not converge.** This is the simpler scenario: There exists $\delta > 0$ such that we
 1012 can find arbitrarily large r obeying $\|\mathbf{v}_r/r - \mathbf{v}^{mm}/\|\mathbf{v}^{mm}\|\| \geq \delta$. If $\|\mathbf{p}_R/R - \Xi\mathbf{p}^{mm}\| \rightarrow 0$, then ‘‘Case
 1013 1’’ applies. Otherwise, we have $\|\mathbf{p}_R/R - \Xi\mathbf{p}^{mm}\| \rightarrow 0$, thus we can assume $\|\mathbf{p}_R/R - \Xi\mathbf{p}^{mm}\| \leq \varepsilon$ for
 1014 arbitrary choice of $\varepsilon > 0$.

1015 On the other hand, due to the strong convexity of (CLS-SVM), for some $\gamma := \gamma(\delta) > 0$, \mathbf{v}_r achieves
 1016 a margin of at most $(1 - \gamma)\Gamma r$ on the dataset $(Y_i, \mathbf{x}_{i1})_{i \in [n]}$. Additionally, since $\|\mathbf{p}_R/R - \Xi\mathbf{p}^{mm}\| \leq \varepsilon$,
 1017 \mathbf{p}_R strictly separates all optimal tokens (for small enough $\varepsilon > 0$) and $\hat{q}_{\max} := f(\varepsilon) \rightarrow 0$ as $R \rightarrow \infty$.
 1018 Consequently, setting $\mathbf{r}_i = \mathbf{X}_i^\top \mathbb{S}(\mathbf{K}_i \mathbf{p}_R)$, for sufficiently large $R > 0$ setting $M = \sup_{i \in [n], t \in [T]} \|\mathbf{x}_{it}\|$, we
 1019 have that

$$\min_{i \in [n]} Y_i \mathbf{v}_r^\top \mathbf{r}_i \leq \min_{i \in [n]} Y_i \mathbf{v}_r^\top \mathbf{x}_{i1} + \sup_{i \in [n]} |\mathbf{r}_i - \mathbf{x}_{i1}| \mathbf{v}_r^\top \quad (91)$$

$$\leq (1 - \gamma)\Gamma r + Mf(\varepsilon)r \quad (92)$$

$$\leq (1 - \gamma/2)\Gamma r. \quad (93)$$

1020 This in turn implies that logistic loss is lower bounded by (following (90)),

$$\mathcal{L}(\mathbf{v}_r, \mathbf{p}_R) \geq \frac{1}{2n} e^{\gamma\Gamma r/2} e^{-\Gamma r} \wedge \log 2.$$

1021 Going back to (86), this exponentially dominates the upper bound of $(\tilde{\mathbf{p}}^{mm}, \tilde{\mathbf{v}}^{mm})$ whenever
 1022 $rMT \exp(-\Xi R) < r\gamma\Gamma/2$, (that is, whenever R, r are sufficiently large), again concluding the proof. ■

1023 C.2 Proof of Theorem 6

1024 We first restate Theorem 6 for ease of reference.

1025 **Theorem 9** Consider the path of $(\mathbf{v}_r, \mathbf{p}_R)$ as $r, R \rightarrow \infty$ as in Theorem 5. Suppose $\mathbb{S}(\mathbf{K}_i \mathbf{p}_R)_{\alpha_i} \rightarrow 1$, i.e.,
 1026 the tokens $(\alpha_i)_{i=1}^n$ are asymptotically selected. Then, $\mathbf{v}_r/r \rightarrow \mathbf{v}^{mm}/\|\mathbf{v}^{mm}\|$ where \mathbf{v}^{mm} is the solution
 1027 of (CLS-SVM) with $\mathbf{r}_i = \mathbf{x}_{i\alpha_i}$, \mathcal{N} is the set of non-support vectors for (CLS-SVM), and $\frac{\mathbf{p}_R}{R} \rightarrow \frac{\mathbf{p}^{relax}}{\|\mathbf{p}^{relax}\|}$
 1028 where \mathbf{p}^{relax} is the solution of (7) with α_i choices.

1029 We will prove this result in two steps. Our first claim restricts the optimization to the particular
 1030 quadrant induced by $\min_{t \neq \alpha_i} (\mathbf{k}_{i\alpha_i} - \mathbf{k}_{it}) \mathbf{p}_R$ under the theorem’s condition $\mathbb{S}(\mathbf{K}_i \mathbf{p}_R) \rightarrow \mathbf{e}_{\alpha_i}$.

1031 **Lemma 7** Suppose $\mathbb{S}(\mathbf{K}_i \mathbf{p}_R) \rightarrow \mathbf{e}_{\alpha_i}$. Then, there exists R_0 such that for all $R \geq R_0$, we have that,

$$\min_{t \neq \alpha_i} (\mathbf{k}_{i\alpha_i} - \mathbf{k}_{it}) \mathbf{p}_R \geq 0 \quad \text{for all } i \in [n]. \quad (94)$$

1032 **Proof.** Suppose the claim does not hold. Set $s_i^R = \mathbb{S}(\mathbf{K}_i \mathbf{p}_R)$. Fix R_0 such that $s_{i\alpha_i}^R \geq 0.9$ for all $R \geq R_0$.
 1033 On the other hand, there exists arbitrarily large R for which $(\mathbf{k}_{i\alpha_i} - \mathbf{k}_{it}) \mathbf{p}_R < 0$ for some $t \neq \alpha_i \in [T]$
 1034 and $i \in [n]$. At this (R, i, t) choices, we have that $s_{it}^R \geq s_{i\alpha_i}^R$. Since $s_{it}^R + s_{i\alpha_i}^R \leq 1$, we find $s_{i\alpha_i}^R < 0.5$
 1035 which contradicts with $s_{i\alpha_i}^R \geq 0.9$. ■

1036 Let \mathcal{Q} be the set of \mathbf{p} satisfying the quadrant constraint (94) – i.e. indices $(\alpha_i)_{i=1}^n$ are selected. Let
 1037 \mathbf{h}_R be the solution of regularization path of (\mathbf{v}, \mathbf{p}) subject to the constraint $\mathbf{p} \in \mathcal{Q}$. From Lemma
 1038 7, we know that, for some R_0 and all $R \geq R_0$, $\mathbf{h}_R = \mathbf{p}_R$. Thus, if the limit exists, we have that
 1039 $\lim_{R \rightarrow \infty} \mathbf{h}_R/R = \lim_{R \rightarrow \infty} \mathbf{p}_R/R$.

1040 To proceed, we will prove that $\lim_{R \rightarrow \infty} \mathbf{h}_R/R$ exists and is equal to $\mathbf{p}^{relax}/\|\mathbf{p}^{relax}\|$ and simultaneously
 1041 establish $\mathbf{v}_r/r \rightarrow \mathbf{v}^{mm}/\|\mathbf{v}^{mm}\|$.

1042 **Lemma 8** $\lim_{R,r \rightarrow \infty} \mathbf{h}_R/R = \mathbf{p}^{relax}/\|\mathbf{p}^{relax}\|$ and $\lim_{R,r \rightarrow \infty} \mathbf{v}_r/r = \mathbf{v}^{mm}/\|\mathbf{v}^{mm}\|$.

1043 **Proof.** The proof will be similar to that of Theorem 5. As usual, we aim to show that SVM-solutions
1044 constitute the most competitive direction. Set $\Xi = 1/\|\mathbf{p}^{relax}\|$.

1045 • **Case 1: \mathbf{h}_R/R does not converge.** Under this scenario there exists $\delta, \gamma = \gamma(\delta) > 0$ such that we can
1046 find arbitrarily large R with $\|\mathbf{h}_R/R - \Xi \mathbf{p}^{relax}\| \geq \delta$. This implies that margin induced by \mathbf{h}_R/R is at
1047 most $\Xi(1 - \gamma)$ over the support vectors $[n] - \mathcal{N}$ (from strong convexity of (7)). The reason is that, \mathbf{h}_R
1048 satisfies $\mathbf{h}_R^\top(\mathbf{k}_{i\alpha_i} - \mathbf{k}_{it}) \geq 0$ for all $t \neq \alpha_i$ by construction as $\mathbf{h}_R \in \mathcal{Q}$. Thus, a constraint over the
1049 support vectors have to be violated (when normalized to the same ℓ_2 norm as $\|\mathbf{p}^{relax}\| = 1/\Xi$).

1050 As usual, we will construct a solution strictly superior to \mathbf{h}_R and contradicts with its optimality.

1051 **Construction of competitor:** Rather than using \mathbf{p}^{relax} direction, we will choose a slightly deviating
1052 direction that ensures the selection of the correct tokens over non-supports \mathcal{N} . Specifically, consider
1053 the solution of (7) where we tighten the non-support constraints by arbitrarily small $\varepsilon > 0$.

$$\mathbf{p}^{\varepsilon-rlx} = \min_{\mathbf{p}} \|\mathbf{p}\| \quad \text{such that} \quad \mathbf{p}^\top(\mathbf{k}_{i\alpha_i} - \mathbf{k}_{it}) \geq \begin{cases} 1 & \text{for all } t \neq \alpha_i, i \in [n] - \mathcal{N} \\ \varepsilon & \text{for all } t \neq \alpha_i, i \in \mathcal{N} \end{cases} \quad (95)$$

1054 Let \mathbf{p}^{mm} be the solution of (ATT-SVM) with $\alpha = (\alpha_i)_{i=1}^n$ (which was assumed to be separable).
1055 Observe that $\mathbf{p}_\varepsilon^{mm} = \varepsilon \mathbf{p}^{mm} + (1 - \varepsilon) \mathbf{p}^{relax}$ satisfies the constraints of (95). Additionally, $\mathbf{p}_\varepsilon^{mm}$ would
1056 achieve a margin of $\frac{1}{(1-\varepsilon)\Xi + \varepsilon/\Delta} = \frac{\Delta\Xi}{\Delta + \varepsilon(\Xi - \Delta)}$ where $\Delta = 1/\|\mathbf{p}^{mm}\|$. Using optimality of $\mathbf{p}^{\varepsilon-rlx}$, this
1057 implies that the reduced margin $\Xi_\varepsilon = 1/\|\mathbf{p}^{\varepsilon-rlx}\|$ (by enforcing ε over non-support) over the support
1058 vectors is a Lipschitz function of ε . That is $\Xi_\varepsilon \geq \Xi - \varepsilon M$ for some $M \geq 0$. To proceed, choose an
1059 $\varepsilon > 0$ such that, it is strictly superior to margin induced by \mathbf{h}_R , that is,

$$\Xi_\varepsilon \geq \Xi(1 - \frac{\gamma}{2}).$$

1060 To proceed, set $\tilde{\mathbf{p}}^{\varepsilon-rlx} = R\Xi_\varepsilon \mathbf{p}^{\varepsilon-rlx}$. Let us recall the following notation from the proof of Theorem
1061 5: $s_i^p = \mathbb{S}(\mathbf{K}_i \mathbf{p})$ and $q_i^p = 1 - s_{i\alpha_i}$. Set $\hat{q}_{\max} = \max_i \hat{q}_{i \in [n] - \mathcal{N}}$ to be worst non-optimality of \mathbf{h}_R over
1062 **support set**. Similarly, define $q_{\max}^* = \max_{i \in [n] - \mathcal{N}} q_i^*$ to be the same for $\tilde{\mathbf{p}}^{\varepsilon-rlx}$. Repeating the identical
1063 arguments to (79), (80), (81), and using the fact that $\mathbf{p}^{\varepsilon-rlx}$ achieves a margin $\Xi(1 - \frac{\gamma}{2}) \leq \Xi_\varepsilon \leq \Xi$, we
1064 end up with the lines

$$\log(\hat{q}_{\max}) \geq -(1 - \gamma)\Xi R - \log T, \quad (96a)$$

$$-\Xi R - \log T \leq \log(q_{\max}^*) \leq -\Xi(1 - 0.5\gamma)R + \log T. \quad (96b)$$

1065 In what follows, we will prove that $\tilde{\mathbf{p}}^{\varepsilon-rlx}$ achieves a strictly smaller logistic loss contradicting with
1066 the optimality of \mathbf{p}_R (whenever $\|\mathbf{h}_R/R - \Xi \mathbf{p}^{relax}\| \geq \delta$).

1067 **Upper bounding logistic loss.** Let us now upper bound the logistic loss of $(\tilde{\mathbf{v}}^{mm}, \tilde{\mathbf{p}}^{\varepsilon-rlx})$ where
1068 $\tilde{\mathbf{v}}^{mm} = r\Gamma \mathbf{v}^{mm}$ with \mathbf{v}^{mm} being the solution of (CLS-SVM) with $\mathbf{r}_i \leftarrow \mathbf{x}_{i\alpha_i}$ and $\Gamma = 1/\|\mathbf{v}^{mm}\|$. Set
1069 $\mathbf{r}_i = \mathbf{X}_i^\top \mathbb{S}(\mathbf{K}_i \tilde{\mathbf{p}}^{\varepsilon-rlx})$. Set $\nu = \min_{i \in \mathcal{N}} Y_i \cdot \mathbf{x}_{i\alpha_i}^\top \mathbf{v}^{mm} - 1$ to be the additional margin buffer that non-support
1070 vectors have access to. Also set $M = \sup_{i \in [n], t, \tau \in [T]} \|\mathbf{x}_{it} - \mathbf{x}_{i\tau}\|$. Observe that we can write

$$\mathbf{x}_{i\alpha_i} - \mathbf{r}_i = \sum_{t \neq \alpha_i} s_{it}(\mathbf{x}_{i\alpha_i} - \mathbf{x}_{it}) \implies \|\mathbf{x}_{i\alpha_i} - \mathbf{r}_i\| \leq q_i M.$$

1071 Non-supports achieve strong label-margin: Using above and (95) for all $i \in \mathcal{N}$ and $t \neq \alpha_i$, we have
1072 that $s_{it} \leq e^{-\varepsilon \Xi_\varepsilon R} s_{i\alpha_i} \leq e^{-\varepsilon \Xi(1-\gamma/2)R} s_{i\alpha_i}$. Consequently, whenever $R \geq \bar{R}_0 := (\varepsilon \Xi(1 - \gamma/2))^{-1} \log(\frac{TM}{\Gamma\nu})$,

$$q_i^* \leq \frac{\sum_{t \neq \alpha_i} s_{it}}{s_{i\alpha_i}} \leq T e^{-\varepsilon \Xi(1-\gamma/2)R} \leq \frac{\Gamma\nu}{M}.$$

1073 This implies that, on $i \in \mathcal{N}$

$$Y_i \cdot \mathbf{r}_i^\top \mathbf{v}^{mm} \geq 1 + \nu + Y_i \cdot (\mathbf{r}_i - \mathbf{x}_{i\alpha_i})^\top \mathbf{v}^{mm} \geq 1 + \nu - q_i M \|\mathbf{v}^{mm}\| \geq 1. \quad (97)$$

1074 *In words:* Above a fixed \bar{R}_0 that only depends on $\gamma = \gamma(\delta)$, features \mathbf{r}_i induced by all non-support
1075 indices $i \in \mathcal{N}$ achieve margin at least 1. What remains is analyzing the margin shrinkage over the
1076 support vectors as in Theorem 5.

Controlling support margin and combining bounds: Over $[n] - \mathcal{N}$, suppose \mathbf{v}^{mm} satisfies the SVM constraints on \mathbf{r}_i with $Y_i \cdot \mathbf{r}_i^\top \mathbf{v}^{mm} \geq 1 - \varepsilon_i/\Gamma$. Consequently, setting $\varepsilon_{\max} = \sup_{i \in [n]} \varepsilon_i$, \mathbf{v}^{mm} achieves a label-margin of $\Gamma - \varepsilon_{\max}$ on the dataset $(Y_i, \mathbf{r}_i)_{i \in [n]}$. Next, we recall the fact (96b) that worst-case perturbation is $\varepsilon_{\max} \leq M \exp(-\Xi(1 - 0.5\gamma)R + \log T) = MT \exp(-\Xi(1 - 0.5\gamma)R)$. With this and (97), we upper bound the logistic loss of $(\tilde{\mathbf{v}}^{mm}, \tilde{\mathbf{p}}^{\varepsilon\text{-rlx}})$ as follows.

$$\mathcal{L}(\tilde{\mathbf{v}}^{mm}, \tilde{\mathbf{p}}^{mm}) \leq \max_{i \in [n]} \log(1 + \exp(-Y_i \mathbf{r}_i^\top \tilde{\mathbf{v}}^{mm})). \quad (98)$$

$$\leq \max_{i \in [n]} \exp(-Y_i \mathbf{r}_i^\top \tilde{\mathbf{v}}^{mm}) \quad (99)$$

$$\leq \exp(-r\Gamma + r\varepsilon_{\max}) \quad (100)$$

$$\leq e^{rMT \exp(-\Xi(1-0.5\gamma)R)} e^{-r\Gamma}. \quad (101)$$

Conversely, we obtain a lower bound for $(\mathbf{v}_r, \mathbf{h}_R)$. Set $\mathbf{r}_i = \mathbf{X}_i^\top \mathbb{S}(\mathbf{K}_i \mathbf{h}_R)$. Recall the lower bound (96a) over the support vector set $[n] - \mathcal{N}$. Combining this with our Assumption C over the support vectors of (CLS-SVM) implies that, solving (CLS-SVM) on $(Y_i, \mathbf{r}_i)_{i \in [n]}$ achieves at most $\Gamma - \nu e^{-(1-\gamma)\Xi R}/T$ margin. Consequently, we have

$$\mathcal{L}(\mathbf{v}_r, \mathbf{h}_R) \geq \frac{1}{n} \max_{i \in [n]} \log(1 + \exp(-Y_i \mathbf{r}_i^\top \mathbf{v}_r)) \quad (102)$$

$$\geq \frac{1}{2n} \max_{i \in [n]} \exp(-Y_i \mathbf{r}_i^\top \mathbf{v}_r) \wedge \log 2 \quad (103)$$

$$\geq \frac{1}{2n} \exp(-r(\Gamma - \nu e^{-(1-\gamma)\Xi R}/T)) \wedge \log 2 \quad (104)$$

$$\geq \frac{1}{2n} e^{r(v/T) \exp(-(1-\gamma)\Xi R)} e^{-r\Gamma} \wedge \log 2. \quad (105)$$

Observe that, this lower bound dominates the previous upper bound when R is large, namely, when (ignoring the multiplier $1/2n$ for brevity)

$$(v/T) e^{-(1-\gamma)\Xi R} \geq MT e^{-\Xi(1-0.5\gamma)R} \iff R \geq R_0 := \frac{2}{\gamma\Xi} \log\left(\frac{MT^2}{\nu}\right).$$

Thus, we obtain the desired contradiction since $\tilde{\mathbf{p}}^{\varepsilon\text{-rlx}}$ is a strictly better solution compared to $\mathbf{p}_R = \mathbf{h}_R$ (once R is sufficiently large).

• **Case 2: \mathbf{v}_r/r does not converge.** This is the simpler scenario: There exists $\delta > 0$ such that we can find arbitrarily large r obeying $\|\mathbf{v}_r/r - \mathbf{v}^{mm}/\|\mathbf{v}^{mm}\|\| \geq \delta$. First, note that, due to the strong convexity of (CLS-SVM), for some $\gamma := \gamma(\delta) > 0$, \mathbf{v}_r achieves a margin of at most $(\Gamma - \gamma)r$ on the dataset $(Y_i, \mathbf{x}_{i1})_{i \in [n]}$. By theorem's condition, we are provided that $\mathbb{S}(\mathbf{K}_i \mathbf{p}_R)_{\alpha_i} \rightarrow 1$. This immediately implies that, for any choice of $\varepsilon = \gamma/3 > 0$, above some sufficiently large (r_0, R_0) , we have that $\|\mathbf{x}_i^{\mathbf{p}_R} - \mathbf{r}_i\| \leq \varepsilon$. Following (101), this implies that, choosing $\tilde{\mathbf{v}}^{mm} = r\mathbf{v}^{mm}/\|\mathbf{v}^{mm}\|$ achieves a logistic loss of at most $e^{r\gamma/3} e^{-r\Gamma}$. Again using $\|\mathbf{x}_i^{\mathbf{p}_R} - \mathbf{r}_i\| \leq \varepsilon$, for sufficiently large (r, R) we have that

$$\min_{i \in [n]} Y_i \mathbf{v}_r^\top \mathbf{r}_i \leq \min_{i \in [n]} Y_i \mathbf{v}_r^\top \mathbf{x}_{i1} + \sup_{i \in [n]} |\mathbf{r}_i - \mathbf{x}_{i1}| \mathbf{v}_r^\top \quad (106)$$

$$\leq (\Gamma - \gamma)r + \varepsilon r \quad (107)$$

$$\leq (\Gamma - 2\gamma/3)r. \quad (108)$$

This in turn implies that logistic loss is lower bounded by (following (105)),

$$\mathcal{L}(\mathbf{v}_r, \mathbf{p}_R) \geq \frac{1}{2n} e^{2\gamma r/3} e^{-r\Gamma} \wedge \log 2.$$

This dominates the above upper bound $e^{r\gamma/3} e^{-r\Gamma}$ of $\tilde{\mathbf{v}}^{mm}$ whenever $\frac{1}{2n} e^{2\gamma r/3} > 1 \iff r > \frac{3}{\gamma} \log(2n)$, (that is, when r is sufficiently large), again concluding the proof. ■

1100 D Addendum to Section 4

1101 D.1 Proof of Theorem 7

1102 **Proof.** The key idea is showing that, thanks to the exponential tail of softmax-attention, (harmful)
1103 contribution of the non-optimal token with the minimum margin can dominate the contribution of

all other tokens as $R \rightarrow \infty$. This high-level approach is similar to earlier works on implicit bias of gradient descent with logistic loss.

Pick $\mathbf{p}^{mm} \in \mathcal{P}^{mm}$ and set $\mathbf{p}_R^* = R \frac{\mathbf{p}^{mm}}{\|\mathbf{p}^{mm}\|}$. This will be the baseline model that \mathbf{p}_R has to compete against. Also let $\bar{\mathbf{p}}_R = \Gamma \frac{\mathbf{p}_R}{R}$. Now suppose $\text{dist}(\bar{\mathbf{p}}_R, \mathcal{P}^{mm}) \rightarrow 0$ as $R \rightarrow \infty$. Then, there exists $\delta > 0$ such that, we can always find arbitrarily large R obeying $\text{dist}(\bar{\mathbf{p}}_R, \mathcal{P}^{mm}) \geq \delta$.

Since $\bar{\mathbf{p}}_R$ is $\delta > 0$ bounded away from \mathcal{P}^{mm} , $\bar{\mathbf{p}}_R$ and $\|\bar{\mathbf{p}}_R\| = \|\mathbf{p}^{mm}\|$, $\bar{\mathbf{p}}_R$ strictly violates at least one of the inequality constraints in (ATT-SVM'). Otherwise, we would have $\bar{\mathbf{p}}_R \in \mathcal{P}^{mm}$. Without losing generality, suppose $\bar{\mathbf{p}}_R$ violates the first margin constraint, that is, for some $\gamma := \gamma(\delta) > 0$, $\max_{\alpha \in \mathcal{O}_1} \min_{\beta \in \bar{\mathcal{O}}_1} \mathbf{p}^\top (\mathbf{k}_{1\alpha} - \mathbf{k}_{1\beta}) \leq 1 - \gamma$. Now, we will argue that this will lead to a contradiction as $R \rightarrow \infty$ since we will show that $\mathcal{L}(\mathbf{p}_R^*) < \mathcal{L}(\mathbf{p}_R)$ for sufficiently large R .

First, let us control $\mathcal{L}(\mathbf{p}_R^*)$. We study $\mathbf{s}_i^* = \mathbb{S}(\mathbf{K}_i \mathbf{p}_R^*)$ and let $\alpha_i \in \mathcal{O}_i$ be the index α in (ATT-SVM') for which $\text{margin}_i = \max_{\alpha \in \mathcal{O}_i} \min_{\beta \in \bar{\mathcal{O}}_i} (\mathbf{k}_{i\alpha} - \mathbf{k}_{i\beta})^\top \mathbf{p}^{mm} \geq 1$ is attained. Then, we bound the non-optimality amount q_i^* of \mathbf{p}_R^* as

$$q_i^* = \frac{\sum_{t \in \bar{\mathcal{O}}_i} \exp(\mathbf{k}_{it}^\top \mathbf{p}_R^*)}{\sum_{t \in [T]} \exp(\mathbf{k}_{it}^\top \mathbf{p}_R^*)} \leq \frac{\sum_{t \in \bar{\mathcal{O}}_i} \exp(\mathbf{k}_{it}^\top \mathbf{p}_R^*)}{\exp(\mathbf{k}_{i\alpha_i}^\top \mathbf{p}_R^*)} \leq T \exp(-R/\Gamma).$$

Thus, $q_{\max}^* = \max_{i \in [n]} q_i^* \leq T \exp(-R/\Gamma)$. Secondly, we wish to control $\mathcal{L}(\mathbf{p}_R)$ by lower bounding the non-optimality in \mathbf{p}_R . Focusing on the first margin constraint, let $\alpha \in \mathcal{O}_1$ be the index in (ATT-SVM') for which $\text{margin}_1 \leq 1 - \gamma$ is attained. Denoting the amount of non-optimality of the first input as \hat{q}_1 , we find³

$$\hat{q}_1 = \frac{\sum_{t \in \bar{\mathcal{O}}_1} \exp(\mathbf{k}_{1t}^\top \mathbf{p}_R)}{\sum_{t \in [T]} \exp(\mathbf{k}_{1t}^\top \mathbf{p}_R)} \geq \frac{1}{T} \frac{\sum_{t \in \bar{\mathcal{O}}_1} \exp(\mathbf{k}_{1t}^\top \mathbf{p}_R)}{\exp(\mathbf{k}_{1\alpha}^\top \mathbf{p}_R)} \geq T^{-1} \exp(-(R(1 - \gamma))/\Gamma).$$

We similarly have $q_{\max}^* \geq T^{-1} \exp(-R/\Gamma)$. In conclusion, for $\mathbf{p}_R, \mathbf{p}_R^*$, denoting maximum non-optimality by $\hat{q}_{\max} \geq \hat{q}_1$ and q_{\max}^* , we respectively obtained

$$\log(\hat{q}_{\max}) \geq -(1 - \gamma)(R/\Gamma) - \log T, \quad (109)$$

$$-(R/\Gamma) - \log T \leq \log(q_{\max}^*) \leq -(R/\Gamma) + \log T. \quad (110)$$

The above inequalities satisfy Assumption D as follows where $\mathbf{p} \leftarrow \mathbf{p}_R^*$ and $\mathbf{p}' \leftarrow \mathbf{p}_R$: Set $R_0 = 3\gamma^{-1}\Gamma \log T$ so that $\log T \leq \frac{\gamma R_0}{3\Gamma}$. Secondly, set $\rho_0 = -(R_0/\Gamma) - \log T$. This way, $\rho_0 \geq \log(q_{\max}^*)$ implies $R \geq R_0$ and $\log T \leq \frac{\gamma R}{3\Gamma}$. Using the latter inequality, we bound the $\log T$ terms to obtain

- $\log(\hat{q}_{\max}) \geq -(1 - 2\gamma/3)(R/\Gamma)$.
- $\log(q_{\max}^*) \leq -(1 - \gamma/3)(R/\Gamma)$.

To proceed, we pick $1 + \Delta = \frac{1 - \gamma/3}{1 - 2\gamma/3}$ implying $\Delta := \frac{\gamma}{3 - 2\gamma}$. Finally, for this Δ , there exists $\rho(\Delta)$ which we need to ensure $\log(\hat{q}_{\max}) \leq \rho(\Delta)$. This can be guaranteed by picking sufficiently large R that ensures $\log(q_{\max}^*) \leq -(1 - \gamma/3)(R/\Gamma) \leq \rho(\Delta)$ to satisfy all conditions of Assumption D. Since such large R exists by initial assumption $\text{dist}(\bar{\mathbf{p}}_R, \mathcal{P}^{mm}) \rightarrow 0$, Assumption D in turn implies that $\mathcal{L}(\mathbf{p}_R^*) < \mathcal{L}(\mathbf{p}_R)$ contradicting with the optimality of \mathbf{p}_R in (8). ■

D.2 Application to Linearly-mixed Labels

The following example shows that if non-optimal tokens result in reduced score (in terms of the alignment of prediction and label), Assumption D holds. The high-level idea behind this lemma is that, if the optimal risk is achieved by setting $q_{\max}^p = 0$, then, Assumption D will hold.

Lemma 9 (Linear label mixing) Recall $q_i^p = \sum_{t \in \bar{\mathcal{O}}_i} \mathbf{s}_{it}^p$ from Assumption D. Suppose $Y_i \in \{-1, 1\}$ and

$$Y_i \cdot \psi(\mathbf{X}_i^\top \mathbf{s}_i^p) = \nu_i(1 - q_i^p) + Z_i,$$

³Here, we assumed margin is non-negative i.e. $\mathbf{k}_{1\alpha}^\top \mathbf{p}_R \geq \sup_{t \in \bar{\mathcal{O}}_1} \mathbf{k}_{1t}^\top \mathbf{p}_R$. Otherwise, $\sup_{t \in [T]} \mathbf{k}_{1t}^\top \mathbf{p}_R$ is attained in $\bar{\mathcal{O}}_1$ which implies $\hat{q}_1 \geq T^{-1}$. Thus, we can still use the identical inequality (109) with the choice $\gamma = 1$.

for some $(v_i)_{i=1}^n > 0$. Here $Z_i = Z_i(\mathbf{p})$ is the contribution of non-optimal tokens to prediction. For some $C, \varepsilon > 0$ and for all $\mathbf{p} \in \mathbb{R}^d$, assume

$$-Cq_{\max}^{\mathbf{p}} \leq Z_i \leq (1 - \varepsilon)v_i q_i^{\mathbf{p}}. \quad (111)$$

Then, Assumption D holds for $\mathcal{L}(\mathbf{p}) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i \cdot \psi(\mathbf{X}_i^\top \mathbf{s}_i^{\mathbf{p}}))$ when $\ell(\cdot)$ is a strictly decreasing loss function with continuous derivative.

Proof. Recall the assumption $Y_i \cdot \psi(\mathbf{X}_i^\top \mathbf{s}_i^{\mathbf{p}}) = v_i(1 - q_i^{\mathbf{p}}) + Z_i$ with Z_i obeying (111). Let us also write the loss function

$$\mathcal{L}(\mathbf{p}) = \frac{1}{n} \sum_{i=1}^n \ell(v_i(1 - q_i^{\mathbf{p}}) + Z_i).$$

Define $q_{\max}^{\mathbf{p}} = \sup_{i \in [n]} q_i^{\mathbf{p}}$. Let M be the maximum absolute value of score over tokens. Let $B = \max_{|x| \leq M} -\ell'(x) \geq A = \min_{|x| \leq M} -\ell'(x) > 0$. Through Taylor's Theorem (integral remainder), we have that

$$B(q_i^{\mathbf{p}} v_i - Z_i) \geq \ell(v_i(1 - q_i^{\mathbf{p}}) + Z_i) - \ell(v_i) \geq A(q_i^{\mathbf{p}} v_i - Z_i) \geq \varepsilon A v_i q_i^{\mathbf{p}}.$$

Set $\mathcal{L}_\star = \frac{1}{n} \sum_{i=1}^n \ell(v_i)$. Set $C_+ = B(C + \max_{i \in [n]} v_i)$ and $C_- = n^{-1} A \varepsilon \min_{i \in [n]} v_i$. This also implies

$$C_+ q_{\max}^{\mathbf{p}} \geq \frac{1}{n} \sum_{i \in [n]} B(q_i^{\mathbf{p}} v_i - Z_i) \geq \mathcal{L}(\mathbf{p}) - \mathcal{L}_\star \geq \frac{1}{n} \sum_{i \in [n]} A(q_i^{\mathbf{p}} v_i - Z_i) \geq \frac{1}{n} \sum_{i \in [n]} \varepsilon A v_i q_i^{\mathbf{p}} \geq C_- q_{\max}^{\mathbf{p}}.$$

Thus, to prove $\mathcal{L}(\mathbf{p}') > \mathcal{L}(\mathbf{p})$, we simply need to establish the stronger statement $C_- q_{\max}^{\mathbf{p}'} > C_+ q_{\max}^{\mathbf{p}}$.

Going back to the condition of Assumption D, any $\log(q_{\max}^{\mathbf{p}}) \leq (1 + \Delta) \log(q_{\max}^{\mathbf{p}'})$ obeys $q_{\max}^{\mathbf{p}} \leq (q_{\max}^{\mathbf{p}'})^{1+\Delta}$ i.e. $q_{\max}^{\mathbf{p}'} \geq (q_{\max}^{\mathbf{p}})^{(1+\Delta)^{-1}}$. Following above, we wish to ensure $q_{\max}^{\mathbf{p}'} > \Theta q_{\max}^{\mathbf{p}}$ for such $(\mathbf{p}, \mathbf{p}')$ pairs where $\Theta = C_+/C_- > 1$. This is guaranteed by

$$(q_{\max}^{\mathbf{p}})^{(1+\Delta)^{-1}-1} > \Theta \iff \frac{\Delta}{1+\Delta} \log(q_{\max}^{\mathbf{p}}) < -\log(\Theta).$$

The above is satisfied by choosing a $\rho(\Delta) := -2(1 + \Delta^{-1}) \log(\Theta)$ in Assumption D. Thus, all \mathbf{p}, \mathbf{p}' with $\log(q_{\max}^{\mathbf{p}}) \leq \rho = \rho(\Delta)$ satisfies the condition of Assumption D finishing the proof. ■

E Experimental Details

In this section, we provide additional implementation details for the experiments.

1. We build one attention layer using PyTorch, and set input and embedding dimensions to be 3. During training, we use SGD optimizer with learning rate 1 in Figure 1 and 0.1 in Figure 2 and train the model for 1000 iterations. To better visualize the generalization path, we normalize the gradient of \mathbf{p} (and \mathbf{v}) at each iteration.
2. Next, given the solution $\hat{\mathbf{p}}$, we determine locally-optimal indices to be those with the highest softmax scores. Using these optimal indices, we utilize python package `cvxopt` to build and solve (ATT-SVM), and then get solution \mathbf{p}^{mm} . After obtaining \mathbf{p}^{mm} , we also verify that these indices satisfy our local-optimal definition. The examples we use in the paper are all trivial to verify (by construction).
3. In Fig. 2(b) and Fig. 2(c), \mathbf{v}^{mm} is solved using python package `sklearn.svm` based on the given label information, and \mathbf{p}^{mm} is the solution of (7) instead.

F Addendum to Section 5

We provide an overview of the current literature on implicit regularization and attention mechanism.

F.1 Related Work on Implicit Regularization

The introduction of Support Vector Machines (SVM), which utilize explicit regularization to choose maximum margin classifiers, represents one of the earliest relevant literature in this field [55]. The concept of maximizing the margin was later connected to generalization performance [56]. From a practical perspective, exponential losses with decaying regularization exhibit asymptotic behavior similar to SVMs, as demonstrated in [19]. While the analysis of the perceptron [57] originally introduced the concept of margins, the method itself does not possess an inherent bias as it terminates with zero classification error. However, establishing a meaningful lower bound for the attained margin is not possible. Initial empirical investigations highlighting the implicit bias of descent methods focused on ℓ_1 -regularization, revealing that coordinate descent, when combined with the exponential loss, exhibits an inherent inclination towards ℓ_1 -regularized solutions [58].

This work draws extensively from the literature on implicit bias and regularization, which has provided valuable techniques and inspiration. A common observation in these studies is the convergence to a specific optimal solution over the training set. This phenomenon has been observed in various approaches, including coordinate descent [59, 60], gradient descent [25, 19], deep linear networks [61, 62], ReLU networks [63, 64, 24, 65, 66, 67], mirror descent [20], and many others. The implicit bias of gradient descent in classification tasks involving separable data has been extensively examined by [19, 20, 21, 22, 23, 24]. These works typically utilize logistic loss or exponentially-tailed losses to establish connections to margin maximization. The results have also been extended to non-separable data by [25, 26, 27]. Furthermore, there have been notable investigations into the implicit bias in regression problems and losses, utilizing techniques such as mirror descent [28, 20, 29, 30, 31, 32]. Additionally, several papers have explored the implicit bias of stochastic gradient descent [33, 34, 35, 36, 37, 38], as well as adaptive and momentum-based methods [39, 40, 41, 42].

While there are some similarities between our optimization approach for \mathbf{v} and existing works, the optimization of \mathbf{p} presents notable differences. Firstly, our optimization problem is nonconvex and involves a composition of loss and softmax, which introduces new challenges and complexities. The presence of softmax adds a nonlinearity to the problem, requiring specialized techniques for analysis and optimization. Secondly, our analysis introduces the concept of locally-optimal tokens, which refers to tokens that achieve locally optimal solutions in their respective attention cones. This concept is crucial for understanding the behavior of the attention mechanism and its convergence properties. By focusing on the cones surrounding locally-optimal tokens, we provide a tailored analysis that captures the unique characteristics of the attention model. Overall, our work offers novel insights into the optimization of attention-based models and sheds light on the behavior of the attention mechanism during training.

F.2 Related Work on Attention Mechanism

As the backbone of Transformers [6], the self-attention mechanism [68] plays a crucial role in computing feature representations by globally modeling long-range interactions within the input. Transformers have achieved remarkable empirical success in various domains, including natural language processing [4, 2], recommendation systems [69, 70, 71], and reinforcement learning [72, 73, 74]. With the introduction of Vision Transformer (ViT) [75], Transformer-based models [76, 77, 78] have gradually replaced convolutional neural network (CNN) architectures and become prevalent in vision tasks. To train ViT efficiently, several techniques have been developed, among which token sparsification [79, 80, 81, 82, 83] remove redundant tokens (image patches) from the data, improving computational complexity while maintaining comparable learning performance.

However, the theoretical foundation of Transformers and self-attention mechanisms has remained largely unexplored. Some studies have established important results, including the Lipschitz constant of self-attention [84], properties of the neural tangent kernel [85, 86], and the expressive power and Turing-completeness of Transformers [87, 88, 89, 47, 51, 90, 91, 92] with statistical guarantees [93, 94].

Focusing on the self-attention component, Edelman et al. [47] theoretically proved that a single self-attention head can represent a sparse function of the input with a sample complexity for the generalization gap between the training loss and the test loss. However, they did not delve into the algorithmic aspects of training Transformers to achieve desirable loss. Sahiner et al. [48] and Ergen et al. [49] further explored the analysis of convex relaxations for self-attention, investigating

1224 potential optimization techniques and properties. In terms of expressive ability, Baldi and Vershynin
1225 [50] investigated the capacity of attention layers to capture complex patterns and information, while
1226 Dong et al. [51] provided additional insights into the expressive power of attention layers in various
1227 contexts. Likhoshesterov et al. [90] studied the model complexity for function approximation of the
1228 self-attention module, and Cordonnier et al. [91] provided sufficient and necessary conditions for
1229 multi-head self-attention structures to simulate convolution layers.

1230 Recent works have made progress in characterizing the optimization and generalization dynamics of
1231 attention. Jelassi et al. [52] studied gradient-based methods from random initialization and provided a
1232 theoretical analysis of the empirical finding that Vision Transformers learn position embeddings that
1233 recapitulate the spatial structure of the training data, even though this spatial structure is no longer
1234 explicitly represented after the image is split into patches. Li et al. [53] provided theoretical results on
1235 training three-layer ViTs for classification tasks. They quantified the importance of self-attention in
1236 terms of sample complexity for achieving zero generalization error, as well as the sparsity of attention
1237 maps when trained by stochastic gradient descent (SGD). In another related work, Nguyen et al.
1238 [95] proposed a primal-dual optimization framework that focuses on deriving attention as the dual
1239 expansion of a primal neural network layer. By solving a support vector regression problem, they
1240 gained a deeper understanding and explanation of various attention mechanisms. This framework
1241 also enables the creation of novel attention mechanisms, offering flexibility and customization in
1242 designing attention-based models. In another closely related work, Oymak et al. [17] analyzed the
1243 same attention model as ours, denoted by (ERM). However, it is important to note that all of these
1244 works make certain assumptions about the data. Specifically, they assume that tokens are tightly
1245 clusterable or can be clearly split into relevant and irrelevant sets. Additionally, Li et al. [53] require
1246 specific assumptions on the initialization of the model, while Jelassi et al. [52] consider a simplified
1247 attention structure where the attention matrix is not directly parameterized with respect to the input.

1248 In contrast, our work offers a comprehensive optimization-theoretic analysis of the attention model,
1249 establishing a formal connection to max-margin problems. This analysis allows us to gain a deeper
1250 understanding of the attention mechanism and its behavior during the training process. Notably, our
1251 work presents the first theoretical understanding of the implicit bias exhibited by gradient descent
1252 methods in the context of the attention model. By uncovering the underlying optimization principles,
1253 we provide valuable insights into the dynamics and generalization properties of attention-based
1254 models.