

---

# Optimizing Attention with Mirror Descent: Generalized Max-Margin Token Selection

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Attention mechanisms have revolutionized numerous domains of artificial intelligence, including natural language processing and computer vision, by enabling  
2 models to selectively focus on relevant parts of the input data. Building on recent  
3 results characterizing the optimization dynamics of gradient descent (GD) and the  
4 structural properties of its preferred solutions in attention-based models, this paper  
5 explores the convergence properties and implicit bias of a family of mirror descent  
6 (MD) algorithms designed for softmax attention mechanisms, with the potential  
7 function chosen as the  $p$ -th power of the  $\ell_p$ -norm. Specifically, we show the directional  
8 convergence of these algorithms to a generalized hard-margin SVM with  
9 an  $\ell_p$ -norm objective when applied to a classification problem using a one-layer  
10 softmax attention model. Our theoretical results demonstrate that these algorithms  
11 not only converge directionally to the generalized max-margin solutions but also  
12 do so at a rate comparable to that of traditional GD in simpler models, despite the  
13 highly nonlinear and nonconvex nature of the present problem. Additionally, we  
14 delve into the joint optimization dynamics of the key-query matrix and the decoder,  
15 establishing conditions under which this complex joint optimization converges to  
16 their respective hard-margin SVM solutions.  
17

## 18 1 Introduction

19 Attention mechanisms [4] have transformed natural language processing (NLP) and large language  
20 models (LLMs). Initially developed for encoder-decoder recurrent neural networks (RNNs), at-  
21 tention enables the decoder to focus on relevant input segments rather than relying solely on a  
22 fixed-length hidden state. This approach became fundamental in transformers [60], where attention  
23 layers—computing softmax similarities among input tokens—are the architecture’s backbone. Trans-  
24 formers have driven rapid advancements in NLP with models like BERT [19] and ChatGPT [42], and  
25 have become the preferred architecture for generative modeling [12, 46], computer vision [20, 45],  
26 and reinforcement learning [21, 11]. This has led to increased exploration of the mathematical  
27 foundations of attention’s optimization.

28 To understand the optimization dynamics of attention mechanisms, [53, 52] studied the *implicit*  
29 *bias* of gradient descent (GD) in binary classification with a fixed linear decoder. This bias reflects  
30 GD’s tendency to favor certain weight characteristics when multiple valid solutions exist. For  
31 instance, in linear logistic regression on separable data, GD aligns with the max-margin class  
32 separator [49, 31]. Similarly, [52, 53] propose a model akin to a hard-margin Support Vector Machine  
33 (SVM)—specifically,  $(\ell_p\text{-AttSVM})$  with  $p = 2$ —maximizing the margin between optimal and non-  
34 optimal tokens based on their softmax logits. These studies show that as training progresses, the  
35 key-query weights  $W(k)$  align with the locally optimal solution  $W_{\text{mm}}^\alpha$ , the minimizer of  $(\ell_p\text{-AttSVM})$ .  
36 Expanding on these insights, [58] explores global directional convergence and GD’s convergence  
37 rate under certain conditions. [48] extends this by relaxing assumptions about regularized paths  
38 for the  $(W_K, W_Q)$  parameterization, showing that gradient flow minimizes the nuclear norm of the  
39 key-query weight  $W = W_K W_Q^\top$ .

40 **Contributions.** While the above aforementioned works provide insights into the implicit bias and  
 41 token selection properties of attention mechanisms, their analyses are limited to GD. A broader  
 42 understanding of general descent algorithms, including the mirror descent (MD) family and their token  
 43 selection properties, is essential. We address this by examining a family of MD algorithms designed for  
 44 softmax attention, where the potential function is the  $p$ -th power of the  $\ell_p$ -norm, termed  $\ell_p$ -AttGD.  
 45 This generalizes both  $\ell_p$ -GD [2, 50, 51] and attention GD [53, 52], enabling the exploration of key  
 46 aspects of attention optimization via  $\ell_p$ -AttGD.

47 *Implicit bias of  $\ell_p$ -AttGD for attention optimization.* Building on [52, 58, 48], we examine a one-layer  
 48 attention model for binary classification, and extend the SVM formulation in [52] to ( $\ell_p$ -AttSVM),  
 49 defining a hard-margin SVM with the  $\ell_p$ -norm. The solution  $W_{\text{mm}}^\alpha$  separates locally optimal tokens  
 50  $(\alpha_i)_{i=1}^n$  with a generalized maximum margin. Theorem 3 shows sufficient conditions for  $\ell_p$ -AttGD to  
 51 converge directionally to  $W_{\text{mm}}^\alpha$ , while Theorem 2 demonstrates that  $\|W(k)\|_{p,p}$  diverges as  $k \rightarrow \infty$ .

52 *Convergence rate of  $\ell_p$ -AttGD to the solution of ( $\ell_p$ -AttSVM).* Theorem 4 shows that the iterates  
 53  $W(k)$  satisfy that  $D_\psi(W_{\text{mm}}^\alpha/\|W_{\text{mm}}^\alpha\|_{p,p}, W(k)/\|W(k)\|_{p,p})$  decreases at an inverse poly-log rate,  
 54 where  $D_\psi(\cdot, \cdot)$  denotes the Bregman divergence [9]. Despite optimizing a nonconvex softmax  
 55 function, the rate is similar to GD in linear binary classification [31, Theorem 1.1]. Though slower  
 56 than the  $O(k^{-3/4})$  rate in [58, Theorem 1], our result applies without assuming token orthogonality.

57 *Generalized Max-Margin Solutions and Joint Optimization of  $(v, W)$ .* We examine the joint problem  
 58 under logistic loss with  $\ell_p$ -norm regularization, solving (ERM) under relaxed  $\ell_p$ -norm constraints.  
 59 If the attention features  $\bar{X}_i = X_i^\top \sigma(X_i W z_i)$  are separable by labels  $y_i$ ,  $v$  acts as a generalized  
 60 max-margin classifier [3]. We show that under suitable geometric conditions,  $W$  and  $v$  converge to  
 61 their generalized max-margin solutions (Theorem 5).

62 We also provide experiments showing mirror descent improves generalization over GD, excelling in  
 63 optimal token selection and suppressing non-optimal tokens.

## 64 2 Preliminaries

65 **Notations.** Let  $N \geq 1$  and  $[N] = \{1, 2, \dots, N\}$ . Vectors are denoted by lowercase letters  
 66 (e.g.,  $a$ ), with components  $a_i$ , and matrices by uppercase letters (e.g.,  $A$ ). For a vector  $v \in \mathbb{R}^d$ ,  
 67 the  $p$ -norm is  $\|v\|_p = (\sum_{i=1}^d |v_i|^p)^{1/p}$ . For a matrix  $M \in \mathbb{R}^{d \times d}$ , the  $p, p$ -norm is  $\|M\|_{p,p} =$   
 68  $(\sum_{i=1}^d \sum_{j=1}^d |M_{ij}|^p)^{1/p}$ . For any two matrices  $X, Y$  of the same dimensions, we define  $\langle X, Y \rangle :=$   
 69  $\text{trace}(X^\top Y)$ . Asymptotic notations  $\mathcal{O}$  and  $\Omega$  hide constant factors, and all logarithms are natural.  
 70 For a differentiable function  $f : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}$ , we define  $D_f : \mathbb{R}^{d \times d} \times \mathbb{R}^{d \times d} \rightarrow \mathbb{R}$  as

$$D_f(W, V) := f(W) - f(V) - \langle \nabla f(V), W - V \rangle. \quad (1)$$

71 **Single-head attention model.** Given input sequences  $X, Z \in \mathbb{R}^{T \times d}$  with length  $T$  and embedding di-  
 72 mension  $d$ , the output of a single-head (cross)-attention layer is computed as:  $\sigma(XW_Q W_K^\top Z^\top) XW_V$ ,  
 73 where  $W_Q, W_K \in \mathbb{R}^{d \times d_1}$ ,  $W_V \in \mathbb{R}^{d \times d_2}$  are trainable key, query, value matrices, respectively;  
 74  $\sigma(XW_Q W_K^\top Z^\top)$  is the attention map; and  $\sigma(\cdot) : \mathbb{R}^{T \times T} \rightarrow \mathbb{R}^{T \times T}$  denotes the row-wise softmax  
 75 function applied row-wise on  $XW_Q W_K^\top Z^\top$ . Similar to [53, 52], we reparameterize the key-query  
 76 product matrix as  $W := W_Q W_K^\top \in \mathbb{R}^{d \times d}$ , and subsume the value weights  $W_V$  within the prediction  
 77 head  $v \in \mathbb{R}^d$ . Suppose the first token of  $Z$ , denoted by  $z$ , is used for prediction. Then, the attention  
 78 model can be formulated as

$$f(X, z) = v^\top X^\top \sigma(XWz). \quad (2)$$

79 **Attention-based empirical risk minimization.** We consider a one-layer attention model (2) for  
 80 binary classification. Consider the dataset  $(X_i, y_i, z_i)_{i=1}^n$ , where  $X_i \in \mathbb{R}^{T \times d}$  is the input with  $T$   
 81 tokens each of dimension  $d$ ,  $y_i \in \{\pm 1\}$  is the label, and  $z_i \in \mathbb{R}^d$  is the token used for comparison.  
 82 We use a smooth decreasing loss function  $l : \mathbb{R} \rightarrow \mathbb{R}$  and study empirical risk minimization (ERM):

$$\min_{v \in \mathbb{R}^d, W \in \mathbb{R}^{d \times d}} \mathcal{L}(v, W) := \frac{1}{n} \sum_{i=1}^n l(y_i v^\top X_i^\top \sigma(X_i W z_i)). \quad (\text{ERM})$$

83 Throughout, we will use  $\mathcal{L}(W)$  to denote the objective of (ERM) with fixed  $v$ .

84 Next, we provide an assumption on the loss function necessary to demonstrate the convergence of  
 85 MD for margin maximization within the attention mechanism.

86 **Assumption A.** Within any closed interval, the loss function  $l : \mathbb{R} \rightarrow \mathbb{R}$  is strictly decreasing and  
 87 differentiable, and its derivative  $l'$  is bounded and Lipschitz continuous.

88 Assumption A aligns with the assumptions on loss functions in [53, 52]. Commonly used loss  
 89 functions, such as  $l(x) = e^{-x}$ ,  $l(x) = -x$ , and  $l(x) = \log(1 + e^{-x})$ , satisfy this assumption.  
 90

91 **Preliminaries on mirror descent.** We review the mirror descent algorithm [7] for solving attention-  
 92 based (ERM). Mirror descent is defined using a *potential function*. We focus on differentiable and  
 93 strictly convex potentials  $\psi$  defined on the entire domain  $\mathbb{R}^{d \times d}$ . We call  $\nabla\psi$  the *mirror map*. The  
 94 natural “distance” associated with the potential  $\psi$  is given by the Bregman divergence [8].

95 **Definition 1** (Bregman Divergence). For a strictly convex function  $\psi : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}$ , the expression  
 96  $D_\psi(\cdot, \cdot)$  defined in (1) is called the Bregman divergence.

97 For more details, see [6]. MD with respect to the mirror map  $\psi$  is a generalization of GD where the  
 98 Bregman divergence is used as a measure of distance. Given a stepsize  $\eta > 0$ , the MD algorithm is as  
 99 follows:

$$W(k+1) \leftarrow \arg \min_{W \in \mathbb{R}^{d \times d}} \{ \eta^{-1} D_\psi(W, W(k)) + \langle \nabla \mathcal{L}(W(k)), W \rangle \}. \quad (\text{MD})$$

100 Equivalently, MD can be written as  $\nabla\psi(W(k+1)) = \nabla\psi(W(k)) - \eta \nabla \mathcal{L}(W(k))$ ; see [10, 34].

101 A useful fact about the Bregman divergence is that it is always non-negative and  $D_\psi(W, V) = 0$  if  
 102 and only if  $W = V$ . Using this notation, one property we will repeatedly use is the following [2]:

103 **Lemma 1.** For any  $W \in \mathbb{R}^{d \times d}$ , the following identities hold for MD:

$$\begin{aligned} D_\psi(W, W(k)) &= D_\psi(W, W(k+1)) + D_{\psi - \eta \mathcal{L}}(W(k+1), W(k)) \\ &\quad - \eta \langle \nabla \mathcal{L}(W(k)), W - W(k) \rangle - \eta \mathcal{L}(W(k)) + \eta \mathcal{L}(W(k+1)). \end{aligned} \quad (3)$$

104 **Preliminaries on attention SVM.** Following [53, 52], we use the following definition of token  
 105 scores.

106 **Definition 2** (Token Score). For prediction head  $v \in \mathbb{R}^d$ , the score of token  $X_{it}$  is  $\gamma_{it} = y_i v^\top X_{it}$ .

107 It is important to highlight that the score is determined solely based on the *value embeddings*  $v^\top X_{it}$   
 108 of the tokens. The softmax function  $\sigma(\cdot)$  minimizes (ERM) by selecting the token with the highest  
 109 score [52, Lemma 2]. Using (2), [52] defines globally optimal tokens  $(\text{opt}_i)_{i=1}^n$ , with each  $\text{opt}_i$   
 110 maximizing the score for  $X_{i\text{opt}_i}$ . For our MD analysis, we primarily consider locally optimal tokens,  
 111 as they are more general than globally optimal ones. Locally optimal tokens are characterized  
 112 by having scores that surpass those of nearby tokens. Intuitively, these are the tokens that locally  
 113 minimize (ERM) upon selection and can be defined based on *support tokens*. Before presenting the  
 114 mathematical notion of locally optimal tokens, we provide the formulation of the attention SVM  
 115 problem. Given a set of (locally) optimal token indices  $(\alpha_i)_{i=1}^n$ , [52] defines the following hard-  
 116 margin attention SVM problem, which aims to separate, with maximal margin, (locally) optimal  
 117 tokens from the rest of the tokens for every input sequence:

$$\begin{aligned} W_{\text{mm}}^\alpha &:= \arg \min_{W \in \mathbb{R}^{d \times d}} \|W\|_F \\ \text{subj. to} &\quad (X_{i\alpha_i} - X_{it})^\top W z_i \geq 1, \quad \text{for all } t \in [T] - \{\alpha_i\}, i \in [n]. \end{aligned} \quad (4)$$

118 The constraint  $(X_{i\alpha_i} - X_{it})^\top W z_i \geq 1$  indicates that in the softmax probability vector  $\sigma(X_i W z_i)$ ,  
 119 the  $\alpha_i$  component has a significantly higher probability compared to the rest, and so these problems  
 120 solve for a sort of probability separator that has the lowest norm.

121 **Definition 3** (Globally and Locally Optimal Tokens). Consider the dataset  $(X_i, y_i, z_i)_{i=1}^n$ .

122 **1.** The tokens with indices  $\text{opt} = (\text{opt}_i)_{i=1}^n$  are called globally optimal if they have the highest  
 123 scores, given by  $\text{opt}_i \in \arg \max_{t \in [T]} \gamma_{it}$ .

124 **2.** Fix token indices  $(\alpha_i)_{i=1}^n$  for which (4) is feasible to obtain  $W_{\text{mm}}^\alpha$ . Let the support tokens  $\mathcal{T}_i$  for  
 125 the  $i^{\text{th}}$  data be the set of tokens  $\tau$  such that  $(X_{i\alpha_i} - X_{i\tau})^\top W_{\text{mm}}^\alpha z_i = 1$ . The tokens with indices  
 126  $(\alpha_i)_{i=1}^n$  are called locally optimal if, for all  $i \in [n]$  and  $\tau \in \mathcal{T}_i$ , the scores per Def. 2 obey  $\gamma_{i\alpha_i} > \gamma_{i\tau}$ .

127 It is worth noting that token scoring and optimal token identification can help us understand the  
 128 importance of individual tokens and their impact on the overall objective. A token score measures  
 129 how much a token contributes to a prediction or classification task, while an optimal token is defined  
 130 as the token with the highest relevance in the corresponding input sequence [53, 52].

### 131 3 Implicit Bias of Mirror Descent for Optimizing Attention

#### 132 3.1 Optimizing Attention with Fixed Head $v$

133 In this section, we assume that the prediction head is fixed, allowing us to delve into the dynamics  
 134 of the token selection mechanism driven by the training of the key-query weight matrix  $W$ . The  
 135 analysis will later be expanded in Section 3.2 to include the joint optimization of both  $v$  and  $W$ .

136 We investigate the theoretical properties of the main algorithm of interest, namely **MD** with  $\psi(\cdot) =$   
 137  $\frac{1}{p} \|\cdot\|_{p,p}^p$  for  $p > 1$  for training (**ERM**) with fixed  $v$ . For conciseness, we will refer to this algorithm  
 138 by the shorthand  $\ell_p$ -**AttGD**. As noted by [3], this choice of mirror potential is particularly of practical  
 139 interest because the mirror map  $\nabla\psi$  updates become *separable* in coordinates and thus can be  
 140 implemented *coordinate-wise* independently of other coordinates.

$$\forall i, j \in [d], \quad \begin{cases} [W(k+1)]_{ij} \leftarrow |[W(k)]_{ij}^+|^{\frac{1}{p-1}} \cdot \text{sign}([W(k)]_{ij}^+), \\ [W(k)]_{ij}^+ := |[W(k)]_{ij}|^{p-1} \text{sign}([W(k)]_{ij}) - \eta[\nabla L(W(k))]_{ij}. \end{cases} \quad (\ell_p\text{-AttGD})$$

141 In the following, we first identify the conditions that guarantee the convergence of  $\ell_p$ -**AttGD**. The  
 142 intuition is that, for attention to exhibit implicit bias, the softmax nonlinearity should select the locally  
 143 optimal token within each input sequence. [52] shows that under certain assumptions, training an  
 144 attention model using GD causes its parameters' direction to converge.

145 This direction can be found by solving a simpler optimization problem, such as attention SVM (4),  
 146 which selects the locally optimal token. Here, we generalize (4) using the  $\ell_p$ -norm as follows:

147 **Definition 4** (Attention SVM with  $\ell_p$ -norm Objective). *For a dataset  $\{(X_i, y_i, z_i)\}_{i=1}^n$  with  $y_i \in$*   
 148  *$\{\pm 1\}$ ,  $X_i \in \mathbb{R}^{T \times d}$ , and token indices  $(\alpha_i)_{i=1}^n$ ,  $\ell_p$ -based attention SVM is defined as*

$$\begin{aligned} W_{\text{mm}}^\alpha &:= \arg \min_{W \in \mathbb{R}^{d \times d}} \|W\|_{p,p} \\ \text{subj. to } &(X_{i\alpha_i} - X_{it})^\top W z_i \geq 1, \text{ for all } t \in [T] - \{\alpha_i\}, i \in [n]. \end{aligned} \quad (\ell_p\text{-AttSVM})$$

149 Problem ( $\ell_p$ -**AttSVM**) is strictly convex, so it has unique solutions when feasible. Furthermore,  
 150 under mild overparameterization,  $d \geq \max\{T-1, n\}$ , the problem is almost always feasible [52,  
 151 Theorem 1]. We assert that the solution to the ( $\ell_p$ -**AttSVM**) problems determines the direction that  
 152 the attention model parameters approach as the training progresses.

153 **Theorem 1** ( $\ell_p$ -norm Regularization Path). *Suppose Assumption A on the loss function holds.*  
 154 *Consider the ridge-constrained solutions  $W^{(R)}$  of (ERM) defined as*

$$W^{(R)} := \arg \min_{W \in \mathbb{R}^{d \times d}} \mathcal{L}(W) \quad \text{subj. to } \|W\|_{p,p} \leq R. \quad (\ell_p\text{-AttRP})$$

155 *Then,  $\lim_{R \rightarrow \infty} W^{(R)} / R = W_{\text{mm}}^{\text{opt}} / \|W_{\text{mm}}^{\text{opt}}\|_{p,p}$ , where  $W_{\text{mm}}^{\text{opt}}$  is the solution of ( $\ell_p$ -**AttSVM**), with  $\alpha_i$*   
 156 *replaced by  $\text{opt}_i$ .*

157 Theorem 1 shows that as the regularization strength  $R$  increases, the optimal direction  $W^{(R)}$  aligns  
 158 more closely with the max-margin solution  $W_{\text{mm}}^\alpha$ . This theorem, which allows for globally optimal  
 159 tokens (see Definition 3), does not require any specific initialization for the  $\ell_p$ -**AttRP** algorithm and  
 160 demonstrates that max-margin token separation is an essential feature of the attention mechanism.

161 Next, we provide the convergence of **MD** applied to (**ERM**). We found that under certain initializations,  
 162 the parameter's  $\ell_p$ -norm increases to infinity as training progresses, and its direction approaches that  
 163 of the ( $\ell_p$ -**AttSVM**) solution. To describe the initialization that allows for these, we define the notion  
 164 of cone sets.

165 **Definition 5.** *Given a square matrix  $W \in \mathbb{R}^{d \times d}$ ,  $\mu \in (0, 1)$ , and some  $R > 0$ ,*

$$S_{p,\mu}(W) := \left\{ W' \in \mathbb{R}^{d \times d} \mid D_\psi \left( \frac{W}{\|W\|_{p,p}}, \frac{W'}{\|W'\|_{p,p}} \right) \leq \mu \right\}, \quad (5a)$$

$$C_{p,\mu,R}(W) := S_\mu(W) \cap \{W' \mid \|W'\|_{p,p} \geq R\}. \quad (5b)$$

166 These sets contain matrices with a similar direction to a reference matrix  $W$ , as captured by the inner  
 167 product in  $S_\mu(W)$ . For  $C_{p,\mu,R}(W)$ , there is an additional constraint that the matrices must have a  
 168 sufficiently high norm. We note that  $S_{p,\mu}(W)$  and  $C_{p,\mu,R}(W)$  reduce to their Euclidean variants  
 169 as described in [53, 52]. With this definition, we present our first theorem about the norm of the  
 170 parameter increasing during training.

171 **Theorem 2.** Suppose Assumption A holds. Let  $(\alpha_i)_{i=1}^n$  be locally optimal tokens as per Defi-  
 172 nition 3. Consider the sequence  $W(k)$  generated by Algorithm  $\ell_p$ -AttGD. For a small enough  
 173 stepsize  $\eta$ , if  $W(0) \in C_{p,\mu,R}(W_{\text{mm}}^\alpha)$  for some dataset-dependent constants  $\mu, R > 0$ , then we have  
 174  $\lim_{k \rightarrow \infty} \|W(k)\|_{p,p} = \infty$ .

175 **Remark 1.** The condition on the stepsize  $\eta$  is that it must be sufficiently small so that  $\psi(\cdot) - \eta\mathcal{L}(\cdot)$   
 176 remains convex for the matrices  $W$  along the path traced by the iterates  $W(k)$ . Specifically, there  
 177 exists an index  $k$  and a real number  $r \in [0, 1]$  such that  $W = rW(k) + (1 - r)W(k + 1)$ . This  
 178 restriction applies to all theorems in this paper that require a sufficiently small stepsize  $\eta$ .

179 This theorem implies that the parameters will increase and diverge to infinity, justifying the need to  
 180 characterize the convergence of their direction.

181 **Theorem 3** (Convergence of  $\ell_p$ -AttGD). Suppose Assumption A holds. Let  $(\alpha_i)_{i=1}^n$  be locally  
 182 optimal tokens as per Definition 3. Consider the sequence  $W(k)$  generated by Algorithm  $\ell_p$ -AttGD.  
 183 For a small enough  $\eta$ , if  $W(0) \in C_{p,\mu,R}(W_{\text{mm}}^\alpha)$  for some constants  $\mu > 0, R > \exp(2)$ , then

$$\lim_{k \rightarrow \infty} \frac{W(k)}{\|W(k)\|_{p,p}} = \frac{W_{\text{mm}}^\alpha}{\|W_{\text{mm}}^\alpha\|_{p,p}}.$$

184 These theorems show that as the parameters grow large enough and approach a locally optimal  
 185 direction, they will keep moving toward that direction.

186 **Theorem 4** (Convergence Rate of  $\ell_p$ -AttGD). Suppose Assumption A holds. Let  $(\alpha_i)_{i=1}^n$  be locally  
 187 optimal tokens as per Definition 3. Consider the sequence  $W(k)$  generated by Algorithm  $\ell_p$ -AttGD.  
 188 For a small enough  $\eta$ , if  $W(0) \in C_{p,\mu,R}(W_{\text{mm}}^\alpha)$  for some constants  $\mu > 0, R > \exp(2)$ , then

$$D_\psi \left( \frac{W_{\text{mm}}^\alpha}{\|W_{\text{mm}}^\alpha\|_{p,p}}, \frac{W(k)}{\|W(k)\|_{p,p}} \right) = \mathcal{O} \left( \begin{cases} \frac{\log \log k}{\log k} & \text{if } p > 2, \\ \frac{(\log \log k)^2}{\log k} & \text{if } p = 2, \\ \frac{1}{(\log k)^{p-1}} & \text{otherwise.} \end{cases} \right). \quad (6)$$

189 Despite optimizing a highly nonlinear, nonconvex softmax function, we achieve a convergence rate  
 190 similar to that of GD in linear binary classification [31, Theorem 1.1] (up to a  $\log \log k$  factor).

### 191 3.2 Training Dynamics of Mirror Descent for Joint Optimization of $W$ and $v$

192 This section delves into the training dynamics of simultaneously optimizing the prediction head  $v$   
 193 and the attention weights  $W$ . Unlike Section 3.1, the main challenge here is the evolving token  
 194 scores  $\gamma$  influenced by the changing nature of  $v$ . This requires additional technical considerations  
 195 beyond those in Section 3.1, which are also addressed in this section. Given stepsizes  $\eta_W, \eta_v > 0$ , we  
 196 consider the following *joint* updates for  $W$  and  $v$  applied to (ERM), respectively: For all  $i, j \in [d]$ :

$$\begin{cases} [W(k+1)]_{ij} \leftarrow |[W(k)]_{ij}^+|^{\frac{1}{p-1}} \cdot \text{sign}([W(k)]_{ij}^+), \\ [W(k)]_{ij}^+ := |[W(k)]_{ij}|^{p-1} \text{sign}([W(k)]_{ij}) - \eta_W [\nabla_W L(W(k), v(k))]_{ij}, \\ [v(k+1)]_i \leftarrow |[v(k)]_i^+|^{\frac{1}{p-1}} \cdot \text{sign}([v(k)]_i^+), \\ [v(k)]_i^+ := |[v(k)]_i|^{p-1} \text{sign}([v(k)]_i) - \eta_v [\nabla_v L(W(k), v(k))]_i. \end{cases} \quad (\ell_p\text{-JointGD})$$

197 We discuss the implicit bias and convergence for  $v(k)$  below. From previous results [3], one can expect  
 198  $v(k)$  to converge to the  $\ell_p$ -SVM solution, i.e., the max-margin classifier separating the set of samples  
 199  $\{(X_{i\alpha_i}, y_i)\}_{i=1}^n$ , where  $X_{i\alpha_i}$  denote the (locally) optimal token for each  $i \in [n]$ . Consequently, we  
 200 consider the following hard-margin SVM problem,

$$v_{\text{mm}} = \arg \min_{v \in \mathbb{R}^d} \|v\|_p \quad \text{subj. to} \quad y_i X_{i\alpha_i}^\top v \geq 1 \quad \text{for all } i \in [n]. \quad (\ell_p\text{-SVM})$$

201 In ( $\ell_p$ -SVM), define the *label margin* as  $1/\|v_{\text{mm}}\|_p$ . The label margin quantifies the distance between  
 202 the separating hyperplane and the nearest data point in the feature space. A larger label margin  
 203 indicates better generalization performance of the classifier, as it suggests that the classifier has a  
 204 greater separation between classes. From ( $\ell_p$ -SVM) and Definitions 2 and 3, an additional intuition



Label	Optimal Token	$\ell_{1.1}$ -MD Token Selection	GD Token Selection	Better Selector
+	fantastic	the movie was <b>fantastic</b>	<b>the</b> movie was fantastic	1.1
-	hated	i <b>hated</b> the movie	i hated <b>the</b> movie	1.1
-	boring	<b>the plot was</b> boring	<b>the plot was</b> boring	2
+	love	i love this <b>movie</b>	i love this <b>movie</b>	2
-	terrible	<b>the plot was</b> terrible	<b>the plot was</b> terrible	1.1

Figure 1: The attention map generated by the resulting models that were trained using  $\ell_{1.1}$  mirror descent and GD for five sample sentences. For three out of five of the sample sentences, the model trained using  $\ell_{1.1}$  mirror descent selects the optimal token better than the model trained using GD.

205 by [53] behind optimal tokens is that they maximize the label margin when selected; see Figure 4 in  
 206 the appendix for a visualization. Selecting the locally optimal token indices  $\alpha = (\alpha_i)_{i=1}^n$  from each  
 207 input data sequence achieves the largest label margin, meaning that including other tokens will reduce  
 208 the label margin as defined in ( $\ell_p$ -SVM). In the Appendix G, we show that  $W$  and  $v$  generated by  
 209  $\ell_p$ -JointRP converge to their respective max-margin solutions under suitable geometric conditions  
 210 (Theorem 5 in the appendix).

## 211 4 Experimental Results

212 We validate our theorems through numerical simulations in Appendix I, and present real data  
 213 experiments here. Our results show that training an attention network with mirror descent improves  
 214 generalization and token selection compared to GD.

Algorithm	Model Size 3	Model Size 4	Model Size 6
$\ell_{1.1}$ -MD	<b>83.47 ± 0.09%</b>	<b>83.36 ± 0.13%</b>	<b>83.65 ± 0.13%</b>
$\ell_2$ -MD	81.66 ± 0.09%	81.05 ± 0.17%	82.22 ± 0.13%
$\ell_3$ -MD	82.57 ± 0.09%	82.40 ± 0.12%	81.97 ± 0.10%

Table 1: Test accuracies of transformer classification models trained with  $\ell_{1.1}$ ,  $\ell_2$ , and  $\ell_3$ -MD on the **Stanford Large Movie Review Dataset**. The model sizes refers to the number of layers in the transformer model and the number of attention heads per layer.  $\ell_{1.1}$ -MD provides superior generalization performance.

215 We trained a transformer classification model on the Stanford Large Movie Review Dataset [39] using  
 216 MD with  $\ell_{1.1}$ ,  $\ell_2$ , and  $\ell_3$  potentials. The models are similar to the one in [60], with the last layer being  
 217 a linear classification layer on the feature representation of the first [CLS] token. Table 1 summarizes  
 218 the resulting test accuracy of several variants of that model when trained with the three algorithms,  
 219 which shows that the  $\ell_{1.1}$  potential mirror descent outperforms the other mirror descent algorithms,  
 220 including the one with the  $\ell_2$  potential, which is equivalent to the GD.

221 We investigate how the model’s attention layers select pivotal tokens in simple GPT-4o-generated  
 222 reviews, focusing on those that determine whether the review is positive or negative. These pivotal  
 223 tokens were also identified by GPT-4o. We compare the model trained using  $\ell_{1.1}$  mirror descent to  
 224 one trained with GD, with full results in the Appendix ( Figure 1 shows five examples). The  $\ell_{1.1}$   
 225 mirror descent outperforms GD in token selection.

## 226 5 Conclusion

227 We studied the optimization dynamics of mirror descent algorithms for softmax attention, focusing on  
 228  $\ell_p$ -AttGD, which generalizes GD using the  $p$ -th power of the  $\ell_p$ -norm as the potential function. Our  
 229 analysis and experiments show that  $\ell_p$ -AttGD converges to the solution of a generalized hard-margin  
 230 SVM with an  $\ell_p$ -norm objective in classification tasks using a one-layer softmax attention model.  
 231 This generalized SVM separates optimal from non-optimal tokens via linear constraints on token  
 232 pairs. We also analyzed the joint problem under logistic loss with  $\ell_p$ -norm regularization and proved  
 233 convergence of  $W$  and  $v$  to their generalized max-margin solutions under appropriate conditions.  
 234 Numerical experiments on synthetic data support our theoretical results.

235 **References**

- 236 [1] Sanjeev Arora, Nadav Cohen, Wei Hu, and Yuping Luo. Implicit regularization in deep matrix  
237 factorization. *Advances in Neural Information Processing Systems*, 32, 2019.
- 238 [2] Navid Azizan and Babak Hassibi. Stochastic gradient/mirror descent: Minimax optimality and  
239 implicit regularization. In *International Conference on Learning Representations*, 2018.
- 240 [3] Navid Azizan, Sahin Lale, and Babak Hassibi. Stochastic mirror descent on overparameterized  
241 nonlinear models. *IEEE Transactions on Neural Networks and Learning Systems*, 33(12):7717–  
242 7727, 2021.
- 243 [4] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly  
244 learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- 245 [5] Han Bao, Ryuichiro Hataya, and Ryo Karakida. Self-attention networks localize when qk-  
246 eigenspectrum concentrates. *arXiv preprint arXiv:2402.02098*, 2024.
- 247 [6] Heinz H Bauschke, Jérôme Bolte, and Marc Teboulle. A descent lemma beyond lipschitz  
248 gradient continuity: first-order methods revisited and applications. *Mathematics of Operations  
249 Research*, 42(2):330–348, 2017.
- 250 [7] Charles Blair. Problem complexity and method efficiency in optimization (a. s. nemirovsky and  
251 d. b. yudin). *SIAM Review*, 27(2):264–265, 1985.
- 252 [8] Lev M Bregman. The relaxation method of finding the common point of convex sets and  
253 its application to the solution of problems in convex programming. *USSR computational  
254 mathematics and mathematical physics*, 7(3):200–217, 1967.
- 255 [9] L.M. Bregman. The relaxation method of finding the common point of convex sets and  
256 its application to the solution of problems in convex programming. *USSR Computational  
257 Mathematics and Mathematical Physics*, 7(3):200–217, 1967.
- 258 [10] Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and  
259 Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- 260 [11] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter  
261 Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning  
262 via sequence modeling. In *Advances in Neural Information Processing Systems*, volume 34,  
263 pages 15084–15097, 2021.
- 264 [12] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared  
265 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large  
266 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 267 [13] Sitan Chen and Yuanzhi Li. Provably learning a multi-head attention layer. *arXiv preprint  
268 arXiv:2402.04084*, 2024.
- 269 [14] Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Training dynamics of multi-  
270 head softmax attention for in-context learning: Emergence, convergence, and optimality. *arXiv  
271 preprint arXiv:2402.19442*, 2024.
- 272 [15] Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural  
273 networks trained with the logistic loss. In *Conference on learning theory*, pages 1305–1338.  
274 PMLR, 2020.
- 275 [16] Liam Collins, Advait Parulekar, Aryan Mokhtari, Sujay Sanghavi, and Sanjay Shakkottai. In-  
276 context learning with transformers: Softmax attention adapts to function lipschitzness. *arXiv  
277 preprint arXiv:2402.11639*, 2024.
- 278 [17] Yichuan Deng, Zhao Song, Shenghao Xie, and Chiwun Yang. Unmasking transformers: A  
279 theoretical approach to data recovery via attention weights. *arXiv preprint arXiv:2310.12462*,  
280 2023.
- 281 [18] Puneesh Deora, Rouzbeh Ghaderi, Hossein Taheri, and Christos Thrampoulidis. On the  
282 optimization and generalization of multi-head attention. *arXiv preprint arXiv:2310.12680*,  
283 2023.

- 284 [19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of  
285 deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Confer-*  
286 *ence of the North American Chapter of the Association for Computational Linguistics: Human*  
287 *Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis,  
288 Minnesota, June 2019. Association for Computational Linguistics.
- 289 [20] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
290 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,  
291 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image  
292 recognition at scale. In *International Conference on Learning Representations*, 2021.
- 293 [21] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter,  
294 Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied  
295 multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- 296 [22] Tolga Ergen, Behnam Neyshabur, and Harsh Mehta. Convexifying transformers: Improving  
297 optimization and understanding of transformer networks. *arXiv:2211.11052*, 2022.
- 298 [23] Spencer Frei, Gal Vardi, Peter L Bartlett, Nathan Srebro, and Wei Hu. Implicit bias in leaky  
299 relu networks trained on high-dimensional data. *arXiv preprint arXiv:2210.07082*, 2022.
- 300 [24] Deqing Fu, Tian-Qi Chen, Robin Jia, and Vatsal Sharan. Transformers learn higher-order  
301 optimization methods for in-context learning: A study with linear models. *arXiv preprint*  
302 *arXiv:2310.17086*, 2023.
- 303 [25] Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias  
304 in terms of optimization geometry. In *International Conference on Machine Learning*, pages  
305 1832–1841. PMLR, 2018.
- 306 [26] Yu Huang, Yuan Cheng, and Yingbin Liang. In-context convergence of transformers. *arXiv*  
307 *preprint arXiv:2310.05249*, 2023.
- 308 [27] M Emrullah Ildiz, Yixiao Huang, Yingcong Li, Ankit Singh Rawat, and Samet Oymak. From  
309 self-attention to markov models: Unveiling the dynamics of generative transformers. *arXiv*  
310 *preprint arXiv:2402.13512*, 2024.
- 311 [28] Samy Jelassi, Michael Eli Sander, and Yuanzhi Li. Vision transformers provably learn spatial  
312 structure. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors,  
313 *Advances in Neural Information Processing Systems*, 2022.
- 314 [29] Hong Jun Jeon, Jason D Lee, Qi Lei, and Benjamin Van Roy. An information-theoretic analysis  
315 of in-context learning. *arXiv preprint arXiv:2401.15530*, 2024.
- 316 [30] Ziwei Ji, Nathan Srebro, and Matus Telgarsky. Fast margin maximization via dual acceleration.  
317 In *International Conference on Machine Learning*, pages 4860–4869. PMLR, 2021.
- 318 [31] Ziwei Ji and Matus Telgarsky. Risk and parameter convergence of logistic regression. *arXiv*  
319 *preprint arXiv:1803.07300*, 2018.
- 320 [32] Ziwei Ji and Matus Telgarsky. Directional convergence and alignment in deep learning. In  
321 H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural*  
322 *Information Processing Systems*, volume 33, pages 17176–17186. Curran Associates, Inc.,  
323 2020.
- 324 [33] Ziwei Ji and Matus Telgarsky. Characterizing the implicit bias via a primal-dual analysis. In  
325 *Algorithmic Learning Theory*, pages 772–804. PMLR, 2021.
- 326 [34] Anatoli Juditsky and Arkadi Nemirovski. First-order methods for nonsmooth convex large-scale  
327 optimization, i: General purpose methods. 2011.
- 328 [35] Hongkang Li, Meng Wang, Sijia Liu, and Pin-Yu Chen. A theoretical understanding of  
329 shallow vision transformers: Learning, generalization, and sample complexity. *arXiv preprint*  
330 *arXiv:2302.06015*, 2023.
- 331 [36] Yingcong Li, Yixiao Huang, Muhammed E Ildiz, Ankit Singh Rawat, and Samet Oymak.  
332 Mechanics of next token prediction with self-attention. In *International Conference on Artificial*  
333 *Intelligence and Statistics*, pages 685–693. PMLR, 2024.



- 334 [37] Zhiyuan Li, Yuping Luo, and Kaifeng Lyu. Towards resolving the implicit bias of gradient  
335 descent for matrix factorization: Greedy low-rank learning. In *International Conference on*  
336 *Learning Representations*, 2020.
- 337 [38] Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural  
338 networks. *arXiv preprint arXiv:1906.05890*, 2019.
- 339 [39] Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher  
340 Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting*  
341 *of the Association for Computational Linguistics: Human Language Technologies*, pages  
342 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics.
- 343 [40] Ashok Vardhan Makkuva, Marco Bondaschi, Adway Girish, Alliot Nagle, Martin Jaggi, Hyeji  
344 Kim, and Michael Gastpar. Attention with markov: A framework for principled analysis of  
345 transformers via markov chains. *arXiv preprint arXiv:2402.04161*, 2024.
- 346 [41] Mor Shpigel Nacson, Jason Lee, Suriya Gunasekar, Pedro Henrique Pamplona Savarese, Nathan  
347 Srebro, and Daniel Soudry. Convergence of gradient descent on separable data. In *The 22nd*  
348 *International Conference on Artificial Intelligence and Statistics*, pages 3420–3428. PMLR,  
349 2019.
- 350 [42] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 351 [43] Samet Oymak, Ankit Singh Rawat, Mahdi Soltanolkotabi, and Christos Thrampoulidis. On the  
352 role of attention in prompt-tuning. In *International Conference on Machine Learning*, 2023.
- 353 [44] Mary Phuong and Christoph H Lampert. The inductive bias of relu networks on orthogonally  
354 separable data. In *International Conference on Learning Representations*, 2020.
- 355 [45] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
356 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
357 models from natural language supervision. In *International conference on machine learning*,  
358 pages 8748–8763. PMLR, 2021.
- 359 [46] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark  
360 Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on*  
361 *Machine Learning*, pages 8821–8831. PMLR, 2021.
- 362 [47] Arda Sahiner, Tolga Ergen, Batu Ozturkler, John Pauly, Morteza Mardani, and Mert Pilanci.  
363 Unraveling attention via convex duality: Analysis and interpretations of vision transformers. In  
364 *International Conference on Machine Learning*, pages 19050–19088. PMLR, 2022.
- 365 [48] Heejune Sheen, Siyu Chen, Tianhao Wang, and Harrison H Zhou. Implicit regularization of  
366 gradient flow on one-layer softmax attention. *arXiv preprint arXiv:2403.08699*, 2024.
- 367 [49] Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The  
368 implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*,  
369 19(1):2822–2878, 2018.
- 370 [50] Haoyuan Sun, Kwangjun Ahn, Christos Thrampoulidis, and Navid Azizan. Mirror descent max-  
371 imizes generalized margin and can be implemented efficiently. *Advances in Neural Information*  
372 *Processing Systems*, 35:31089–31101, 2022.
- 373 [51] Haoyuan Sun, Khashayar Gatmiry, Kwangjun Ahn, and Navid Azizan. A unified approach to  
374 controlling implicit regularization via mirror descent. *Journal of Machine Learning Research*,  
375 24(393):1–58, 2023.
- 376 [52] Davoud Ataee Tarzanagh, Yingcong Li, Christos Thrampoulidis, and Samet Oymak. Transform-  
377 ers as support vector machines. *arXiv preprint arXiv:2308.16898*, 2023.
- 378 [53] Davoud Ataee Tarzanagh, Yingcong Li, Xuechen Zhang, and Samet Oymak. Max-margin token  
379 selection in attention mechanism. *Advances in Neural Information Processing Systems*, 36,  
380 2024.
- 381 [54] Yuandong Tian, Yiping Wang, Beidi Chen, and Simon Du. Scan and snap: Understanding  
382 training dynamics and token composition in 1-layer transformer. *arXiv:2305.16380*, 2023.
- 383 [55] Yuandong Tian, Yiping Wang, Zhenyu Zhang, Beidi Chen, and Simon Du. Joma: De-  
384 mystifying multilayer transformers via joint dynamics of mlp and attention. *arXiv preprint*  
385 *arXiv:2310.00535*, 2023.

- 386 [56] Gal Vardi. On the implicit bias in deep-learning algorithms. *Communications of the ACM*,  
387 66(6):86–93, 2023.
- 388 [57] Gal Vardi and Ohad Shamir. Implicit regularization in relu networks with the square loss. In  
389 *Conference on Learning Theory*, pages 4224–4258. PMLR, 2021.
- 390 [58] Bhavya Vasudeva, Puneesh Deora, and Christos Thrampoulidis. Implicit bias and fast conver-  
391 gence rates for self-attention. *arXiv preprint arXiv:2402.05738*, 2024.
- 392 [59] Bhavya Vasudeva, Deqing Fu, Tianyi Zhou, Elliott Kau, Youqi Huang, and Vatsal Sharan. Sim-  
393 plicity bias of transformers to learn low sensitivity functions. *arXiv preprint arXiv:2403.06925*,  
394 2024.
- 395 [60] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
396 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*  
397 *processing systems*, 30, 2017.
- 398 [61] Junhao Zheng, Shengjie Qiu, and Qianli Ma. Learn or recall? revisiting incremental learning  
399 with pre-trained language models. *arXiv preprint arXiv:2312.07887*, 2023.

400	<b>Contents</b>	
401	<b>1 Introduction</b>	<b>1</b>
402	<b>2 Preliminaries</b>	<b>2</b>
403	<b>3 Implicit Bias of Mirror Descent for Optimizing Attention</b>	<b>4</b>
404	3.1 Optimizing Attention with Fixed Head $v$ . . . . .	4
405	3.2 Training Dynamics of Mirror Descent for Joint Optimization of $W$ and $v$ . . . . .	5
406	<b>4 Experimental Results</b>	<b>6</b>
407	<b>5 Conclusion</b>	<b>6</b>
408	<b>A Related Work</b>	<b>11</b>
409	<b>B Auxiliary lemmas</b>	<b>12</b>
410	B.1 Additional Notations . . . . .	12
411	B.2 Lemma for Analyzing The $\ell_p$ -Norm . . . . .	13
412	B.3 Lemma for Analyzing ERM Objective and Its Gradient . . . . .	16
413	B.4 Lemma for Analyzing $\ell_p$ -AttGD . . . . .	23
414	B.5 Lemma for Analyzing Rate of Convergence . . . . .	25
415	<b>C Proof of Theorem 1</b>	<b>28</b>
416	<b>D Proof of Theorem 2</b>	<b>28</b>
417	<b>E Proof of Theorem 3</b>	<b>28</b>
418	<b>F Proof of Theorem 4</b>	<b>28</b>
419	<b>G On the Convergence of the <math>\ell_p</math> Regularization Path for Joint <math>W</math> and <math>v</math></b>	<b>29</b>
420	<b>H Proof of Theorem 5</b>	<b>29</b>
421	<b>I Implementation Details</b>	<b>31</b>
422	I.1 Illustrating Optimal Tokens . . . . .	31
423	I.2 Synthetic Data Experiment . . . . .	32
424	I.3 Additional Real Experiments . . . . .	34

425 **A Related Work**

426 **Transformers Optimization.** Recently, the study of optimization dynamics of attention mechanisms  
427 has garnered significant attention [18, 26, 55, 24, 36, 53, 52, 58, 48, 17, 40, 29, 61, 16, 13, 35, 48,  
428 27, 59, 5, 14]. We discuss the works most closely related to this paper. Studies such as [47, 22]  
429 investigate the optimization of attention models through convex relaxations. [28] demonstrate that  
430 Vision Transformers (ViTs) identify spatial patterns in binary classification via gradient methods.  
431 [35] provide sample complexity bounds and discuss attention sparsity in SGD for ViTs. [43] and [18]  
432 explore optimization dynamics in prompt-attention and multi-head attention models, respectively.  
433 [54, 55] study SGD dynamics and multi-layer transformer training. [53, 52] explored GD’s implicit  
434 bias in a binary classification setting with a fixed linear decoder. [58] discusses the global directional  
435 convergence and convergence rate of GD under specific data conditions. [48] notes that gradient flow

not only achieves minimal loss but also minimizes the nuclear norm of the key-query weight  $W$ . Our work extends these findings and those of [53, 52], focusing on the implicit bias of the general class of MD algorithms for attention training.

**Implicit Bias of First Order Methods.** In recent years, significant progress has been made in understanding the implicit bias of gradient descent on separable data, particularly highlighted by the works of [49, 31]. For linear predictors, [41, 33, 30] demonstrated that gradient descent methods rapidly converge to the max-margin predictor. Extending these insights to MLPs, [32, 38, 15] have examined the implicit bias of GD and gradient flow using exponentially-tailed classification losses, and show convergence to the Karush-Kuhn-Tucker (KKT) points of the corresponding max-margin problem, both in finite [32, 38] and infinite width scenarios [15]. Further, the implicit bias of GD for training ReLU and Leaky-ReLU networks has been investigated, particularly on orthogonal data [44, 23]. Additionally, the implicit bias towards rank minimization in regression settings with square loss has been explored in [57, 1, 37].

Our work is closely related to the implicit bias of MD [25, 2] for regression and classification, respectively. Specifically, [50] extended the findings of [25, 2] to classification problems, and developed a class of algorithms exhibiting an implicit bias towards a generalized SVM with  $\ell_p$  norms that effectively separates samples based on their labels; for a survey, we refer to [56].

## B Auxiliary lemmas

### B.1 Additional Notations

We denote the minimum and maximum of scalars  $a$  and  $b$  as  $a \wedge b$  and  $a \vee b$ , respectively. Consider the following constants for the proofs, depending on the dataset  $(X_i, Y_i, z_i)_{i=1}^n$ , the parameter  $v$ , and the locally optimal token  $(\alpha_i)_{i=1}^n$ :

$$\begin{aligned} \delta' &:= \frac{1}{2} \min_{i \in [n]} \min_{\tau \in \bar{\mathcal{T}}_i} ((X_{i\alpha_i} - X_{i\tau})^\top W_{\text{mm}}^\alpha z_i - 1) \\ &\leq \frac{1}{2} \min_{i \in [n]} \min_{t \in \mathcal{T}_i, \tau \in \bar{\mathcal{T}}_i} ((X_{it} - X_{i\tau})^\top W_{\text{mm}}^\alpha z_i); \end{aligned} \quad (7a)$$

$$\delta := \min\{0.25, \delta'\}. \quad (7b)$$

When  $\bar{\mathcal{T}}_i = \emptyset$  for all  $i \in [n]$  (i.e. globally-optimal indices), we set  $\delta' = \infty$  as all non-neighbor related terms will disappear. Further, recalling Definition 4 and using  $W_{\text{mm}}^\alpha$ —i.e., the minimizer of ( $\ell_p$ -AttSVM), we set

$$\begin{aligned} A' &:= \|W_{\text{mm}}^\alpha\|_{p,p} \max_{i \in [n], t \in [T]} \|X_{it} z_i^\top\|_{\frac{p}{p-1}, \frac{p}{p-1}}; \\ A &:= \max\{1, A'\}. \end{aligned} \quad (8)$$

Recalling Definition 5, we provide the following initial radius  $\mu = \mu_0$  which will be used later in Lemma 10:

$$\mu_0 := \begin{cases} \frac{1}{p} \left(\frac{\delta}{8A}\right)^p & \text{if } p \geq 2, \\ \frac{1}{p} \left(\frac{\delta(p-1)}{4Ad^{\frac{2}{p}-1}}\right)^2 & \text{otherwise.} \end{cases} \quad (9)$$

Furthermore, define the following sums for  $W$ :

$$S_i(W) := \sum_{t \in \mathcal{T}_i} [\sigma(X_i W z_i)]_t, \quad \text{and} \quad Q_i(W) := \sum_{t \in \bar{\mathcal{T}}_i} [\sigma(X_i W z_i)]_t.$$

For the samples  $i$  with non-empty supports  $\mathcal{T}_i$ , let

$$\gamma_i^{\text{gap}} := \gamma_{i\alpha_i} - \max_{t \in \mathcal{T}_i} \gamma_{it}, \quad \text{and} \quad \bar{\gamma}_i^{\text{gap}} := \gamma_{i\alpha_i} - \min_{t \in \bar{\mathcal{T}}_i} \gamma_{it}. \quad (10)$$

Furthermore, we define the *global score gap* as

$$\Gamma := \sup_{i \in [n], t, \tau \in [T]} |\gamma_{it} - \gamma_{i\tau}|. \quad (11)$$

467 **B.2 Lemma for Analyzing The  $\ell_p$ -Norm**

468 In this section of the Appendix, we provide some analysis on comparing the  $\ell_p$ -norm, the  $\ell_p$  Bregman  
 469 divergence, and the  $\ell_2$ -norm of matrices. Since the  $\ell_2$ -norm of matrices are much easier to analyze  
 470 and use, like in the inner product Cauchy-Schwarz inequality, having this comparison is valuable  
 471 when analyzing the  $\ell_p$ -AttGD.

472 **Lemma 2.** For any  $d \times d$  matrix  $W$ , let  $w$  denote its vectorization. Then,

$$\|w\|_p \in \left[ d^{\frac{2}{p}-1} \|w\|_2, \|w\|_2 \right]$$

473 for  $p \geq 2$ , and for  $1 < p \leq 2$ ,  $\|w\|_p$  is in a similar interval, with the two ends switched.

474 *Proof.* Let  $w_1, w_2, \dots, w_{d^2}$  be the entries of  $w$ . Therefore, for  $p \geq 2$ ,

$$\begin{aligned} \|w\|_p &= \sqrt[p]{\sum_{i=1}^{d^2} |w_i|^p} \\ &= \sqrt[p]{\sum_{i=1}^{d^2} (|w_i|^2)^{p/2}}, \end{aligned}$$

475 and because  $\frac{p}{2} \geq 1$ , we would have

$$\begin{aligned} \sqrt[p]{\sum_{i=1}^{d^2} (|w_i|^2)^{p/2}} &\leq \sqrt[p]{\left( \sum_{i=1}^{d^2} |w_i|^2 \right)^{p/2}} \\ &= \sqrt[p]{\|w\|_2^p} = \|w\|_2. \end{aligned}$$

476 Therefore,  $\|w\|_p \leq \|w\|_2$  whenever  $p \geq 2$ . A similar argument will get us  $\|w\|_p \geq \|w\|_2$  whenever  
 477  $1 < p \leq 2$ , so one end of the interval is solved for each case, now for the other end.

478 Using the power-mean inequality, we can get that whenever  $p \geq 2$ ,

$$\sqrt[p]{\frac{1}{d^2} \sum_{i=1}^{d^2} |w_i|^p} \geq \sqrt{\frac{1}{d^2} \sum_{i=1}^{d^2} |w_i|^2},$$

479

$$d^{-\frac{2}{p}} \|w\|_p \geq d^{-1} \|w\|_2,$$

480

$$\|w\|_p \geq d^{\frac{2}{p}-1} \|w\|_2.$$

481 Similarly, for  $1 < p \leq 2$ ,

$$\|w\|_p \leq d^{\frac{2}{p}-1} \|w\|_2.$$

482

□

483 **Lemma 3.** Let  $W_1, W_2 \in \mathbb{R}^{d \times d}$  be two matrices such that  $\|W_1\|_{p,p} = \|W_2\|_{p,p} = 1$ . Then, the  
 484 following inequalities hold:

485 L1. For  $p \geq 2$ ,

$$D_\psi(W_1, W_2) \geq \frac{1}{p \times 2^p} \|W_1 - W_2\|_{p,p}^p,$$

486 L2. For  $p \in (1, 2)$ ,

$$D_\psi(W_1, W_2) \geq \frac{(p-1)^2}{p} \|W_1 - W_2\|_{2,2}^2.$$

487 Here,  $D_\psi(\cdot, \cdot)$  denotes the Bregman divergence given in Definition 1.

488 *Proof.* Let  $W_1 = (x_{ij})_{i,j \in [d]}$  and  $W_2 = (y_{ij})_{i,j \in [d]}$ , then from Definition 1, we have

$$\begin{aligned} D_\psi(W_1, W_2) &= \frac{1}{p} \sum_{i,j \in [d]} |x_{ij}|^p - \frac{1}{p} \sum_{i,j \in [d]} |y_{ij}|^p - \sum_{i,j \in [d]} |y_{ij}|^{p-1} (x_{ij} - y_{ij}) \operatorname{sign}(y_{ij}) \\ &= \sum_{i,j \in [d]} \left( \frac{1}{p} |x_{ij}|^p + \frac{p-1}{p} |y_{ij}|^p - |y_{ij}|^{p-1} |x_{ij}| \operatorname{sign}(x_{ij} y_{ij}) \right). \end{aligned}$$

489 Therefore, it is enough to prove that whenever  $x, y \in [-1, 1]$ , the expression

$$\frac{1}{p} |x|^p + \frac{p-1}{p} |y|^p - |x||y|^{p-1} \operatorname{sign}(xy) \quad (12)$$

490 is at least  $\frac{1}{p2^p} |x - y|^p$  if  $p \geq 2$ , or is at least  $\frac{(p-1)^2}{p} |x - y|^2$  if  $p \in (1, 2)$ . We split the argument into  
491 two cases, the first is when the signs of  $x$  and  $y$  are the same, and the second for when they are not.

492 **Case 1:**  $\operatorname{sign}(xy) = 1$ , so both  $x$  and  $y$  have the same sign, WLOG both are non-negative. Let us fix  
493 the value  $\Delta \in [-1, 1]$  and find the minimum value of (12) when we constraint  $x$  and  $y$  to be positive  
494 and  $x - y = \Delta$ . Therefore, that expression can be written as

$$\frac{(y + \Delta)^p + (p-1)y^p}{p} - (y + \Delta)y^{p-1},$$

495 the first derivative with respect to  $y$  is

$$\begin{aligned} (y + \Delta)^{p-1} + (p-1)y^{p-1} - y^{p-1} - (p-1)(y + \Delta)y^{p-2} \\ = (y + \Delta)^{p-1} - y^{p-1} - (p-1)\Delta y^{p-2}. \end{aligned}$$

496 Since the function  $t \mapsto t^{p-1}$  is convex for  $p \geq 2$ , and concave for  $p \in (1, 2)$ , then that derivative is  
497 always non-negative when  $p \geq 2$  and always negative when  $p \in (1, 2)$ .

498 **Sub-Case 1.1:**  $p \geq 2$ . In this subcase, (12) reaches its minimum when  $(x, y) = (\Delta, 0)$  or  $(0, -\Delta)$ ,  
499 depending on the sign of  $\Delta$ , plugging them in gets us the minimum, which is  $\frac{1}{p} |\Delta|^p$  when  $\Delta \geq 0$  or  
500  $\frac{p-1}{p} |\Delta|^p$  otherwise.

501 **Sub-Case 1.2:**  $p \in (1, 2)$ . In this subcase, (12) reaches its minimum when  $(x, y) = (1, 1 - \Delta)$  if  $\Delta$   
502 is non-negative or  $(1 + \Delta, 1)$  otherwise. When  $\Delta$  is non-negative, the desired minimum is

$$\begin{aligned} \frac{1 + (p-1)(1 - \Delta)^p}{p} - (1 - \Delta)^{p-1} &= \frac{1}{p} (1 - (1 - \Delta)^{p-1} - (p-1)\Delta(1 - \Delta)^{p-1}) \\ &\geq \frac{1}{p} ((p-1)\Delta - (p-1)\Delta(1 - \Delta)^{p-1}) \\ &= \frac{(p-1)\Delta}{p} (1 - (1 - \Delta)^{p-1}) \geq \frac{(p-1)^2}{p} \Delta^2. \end{aligned}$$

503 Combining the results from the subcases, we get that the expression in (12) is lower-bounded by  
504  $\frac{1}{p} |x - y|^p$  when  $p \geq 2$ , or  $\frac{(p-1)^2}{p} |x - y|^2$  otherwise, which sufficiently satisfies the desired bounds  
505 for case 1.

506

507 **Case 2:**  $\operatorname{sign}(xy) = -1$ , so  $x$  and  $y$  has opposite sign. The expression in (12) can be simplified to

$$\frac{1}{p} |x|^p + \frac{p-1}{p} |y|^p + |x||y|^{p-1},$$

508 and we want to prove that it is at least  $\frac{1}{p2^p} (|x| + |y|)^p$  when  $p \geq 2$ , or is at least  $\frac{(p-1)^2}{p} (|x| + |y|)^2$   
509 when  $p \in (1, 2)$ . In the case that  $p \geq 2$ , one of  $|x|$  or  $|y|$  is at least  $\frac{|x| + |y|}{2}$ , so the above is at least



510  $\frac{1}{p} \left( \frac{|x|+|y|}{2} \right)^p = \frac{1}{p2^p} (|x| + |y|)^p$ . Otherwise,

$$\begin{aligned} \frac{1}{p}|x|^p + \frac{p-1}{p}|y|^p + |x|y|^{p-1} &= \frac{|x|(|x|^{p-1} + |y|^{p-1}) + (p-1)|y|^{p-1}(|x| + |y|)}{p} \\ &\geq \frac{(|x| + |y|)(|x| + (p-1)|y|^{p-1})}{p} \\ &\geq \frac{(|x| + |y|)((p-1)|x| + (p-1)|y|)}{p} \\ &= \frac{p-1}{p}(|x| + |y|)^2 \geq \frac{(p-1)^2}{p}(|x| + |y|)^2. \end{aligned}$$

511 Therefore, we have proven the bound for this case. □

512 **Lemma 4.** For any  $x \geq y \geq 0$ , we we have

$$\frac{p-1}{p}x^p - \frac{p-1}{p}y^p \geq y(x^{p-1} - y^{p-1}).$$

*Proof.*

$$\frac{d}{dx} \left( \frac{p-1}{p}x^p - \frac{p-1}{p}y^p \right) = (p-1)x^{p-1},$$

513

$$\frac{d}{dx} y(x^{p-1} - y^{p-1}) = (p-1)x^{p-2}y \leq (p-1)x^{p-1},$$

514 so as we increase  $x$ , the left side grows faster than the right side, so we simply need to prove that the  
515 inequality holds at  $x = y$ , which is trivially true. □

516 **Lemma 5.** For any  $x \geq y \geq 0$ , we we have that if  $q \geq 1$

$$x^q - y^q \leq qx^{q-1}(x - y),$$

517 and if  $0 < q < 1$ ,

$$x^q - y^q \leq qy^{q-1}(x - y)$$

*Proof.*

$$\frac{d}{dx} (x^q - y^q) = qx^{q-1},$$

518

$$\frac{d}{dx} qx^{q-1}(x - y) = q(q-1)x^{q-2}(x - y) + qx^{q-1}, \text{ and}$$

519

$$\frac{d}{dx} qy^{q-1}(x - y) = qy^{q-1}.$$

520 When  $q \geq 1$ ,

$$\frac{d}{dx} (x^q - y^q) \geq \frac{d}{dx} qx^{q-1}(x - y),$$

521 so because we have

$$x^q - y^q = qx^{q-1}(x - y) = 0$$

522 when  $x = y$ , then

$$x^q - y^q \geq qy^{q-1}(x - y)$$

523 when  $x \geq y \geq 0$  if  $q \geq 1$ . We can use a similar argument for the  $0 < q < 1$  case. □

524 **B.3 Lemma for Analyzing ERM Objective and Its Gradient**

525 In this section of the Appendix, we analyze the objective function. We especially want to know about  
 526 its gradient and the inner product of this gradient with the matrices of the cone set, as was mentioned  
 527 before in the main body of the paper. The first one bounds the loss objective,

528 **Lemma 6.** *Under Assumption A,  $\mathcal{L}(W)$  is bounded from above by  $\mathcal{L}_{max}$  and below by  $\mathcal{L}_{min}$  for  
 529 some dataset-dependent constants  $\mathcal{L}_{max}$  and  $\mathcal{L}_{min}$  that are finite.*

530 *Proof.* It is enough to show the same thing for each of the loss contributions of each sample,  
 531  $l_i(y_i v^\top X_i^\top \sigma(X_i W z_i))$ . By Assumption A, we simply need to show that  $y_i v^\top X_i^\top \sigma(X_i W z_i)$  is  
 532 bounded by dataset-dependent bounds. However,  $W$  only affects the softmax, so the above expression  
 533 is bounded above by  $\max_{t \in [T]} \gamma_{it}$  and bounded below by  $\min_{t \in [T]} \gamma_{it}$ , which are dataset dependent.  
 534  $\square$

535 **Lemma 7.** *If we denote  $h_i := X_i W z_i$  and  $l'_i := l'(\gamma_i^\top \sigma(h_i))$ , then*

$$\nabla \mathcal{L}(W) = \frac{1}{n} \sum_{i=1}^n l'_i X_i^\top (\text{diag}(\sigma(h_i)) - \sigma(h_i) \sigma(h_i)^\top) \gamma_i z_i^\top,$$

536 where  $\mathcal{L}(W)$  denotes the objective of (ERM) with fixed  $v$ .

537 *Proof.* We first calculate the derivatives of each term in the sum of  $\mathcal{L}(W)$ . The derivative of the  $i$ -th  
 538 term for the  $W_{j_1 j_2}$  component is

$$\begin{aligned} \frac{\partial}{\partial W_{j_1 j_2}} l(y_i v^\top X_i^\top \sigma(X_i W z_i)) &= l'_i \gamma_i^\top \frac{\partial}{\partial W_{j_1 j_2}} \sigma(X_i W z_i) \\ &= l'_i \gamma_i^\top \nabla \sigma(h_i) X_{i, :, j_1}^\top z_{i j_2} \\ &= l'_i X_{i, :, j_1}^\top \nabla \sigma(h_i)^\top \gamma_i z_{i j_2}. \end{aligned}$$

539 Therefore, the derivative for the  $j_2$ -th row of  $W$  is

$$l'_i X_i^\top \nabla \sigma(h_i)^\top \gamma_i z_{i j_2}.$$

540 Next, the full gradient for the  $i$ -th term equals

$$l'_i X_i^\top \nabla \sigma(h_i)^\top \gamma_i z_i^\top.$$

541 To finish the proof, we calculate the derivative of  $\sigma(h_i)$ . The derivative of the  $j_1$ -th component of  
 542  $\sigma(h_i)$  with respect to  $h_{i j_2}$  is

$$\begin{aligned} \frac{\partial}{\partial h_{i j_2}} \left( \frac{e^{h_{i j_1}}}{\sum_{l=1}^T e^{h_{i l}}} \right) &= \frac{e^{h_{i j_1}} \mathbf{1}_{j_1=j_2}}{\sum_{l=1}^T e^{h_{i l}}} - \frac{e^{h_{i j_1}} e^{h_{i j_2}}}{\left( \sum_{l=1}^T e^{h_{i l}} \right)^2} \\ &= \sigma(h_i)_{j_1} \mathbf{1}_{j_1=j_2} - \sigma(h_i)_{j_1} \sigma(h_i)_{j_2}. \end{aligned}$$

543 Thus, the derivative of  $\sigma(h_i)$  is a matrix in  $\mathbb{R}^{T \times T}$  defined as

$$\text{diag}(\sigma(h_i)) - \sigma(h_i) \sigma(h_i)^\top.$$

544 Therefore, the full gradient is

$$\frac{1}{n} \sum_{i=1}^n l'_i X_i^\top (\text{diag}(\sigma(h_i)) - \sigma(h_i) \sigma(h_i)^\top) \gamma_i z_i^\top.$$

545  $\square$

546 **Lemma 8.** *Under Assumption A,  $\|\nabla \mathcal{L}(W)\|_{p,p}$  is bounded by a dataset-dependent constant  $L$ .*

547 *Proof.* Using the expression in Lemma 7, since  $l'$  is bounded and the entries in  $\sigma(h_i)$  is always  
 548 between 0 and 1, then the entries of  $\nabla \mathcal{L}(W)$  is bounded by a dataset-dependent bounded, which  
 549 directly implies this lemma statement.  $\square$

550 In the following lemma, we analyze the behaviors of the ( $\ell_p$ -AttSVM) constraint  $(X_{it} - X_{i\tau})^\top W z_i$   
 551 for all  $W \in S_{p,\mu_0}(W_{\text{mm}}^\alpha)$  satisfying  $\|W\|_{p,p} = \|W_{\text{mm}}^\alpha\|_{p,p}$ , the result of which is a generalization of  
 552 [52, Equation 64] for a general  $\ell_p$  norm.

553 **Lemma 9.** Let  $\alpha = (\alpha_i)_{i=1}^n$  be locally optimal tokens as per Definition 3, and let  $W_{\text{mm}}^\alpha$  be the  
 554 ( $\ell_p$ -AttSVM) solution. Let  $(\mathcal{T}_i)_{i=1}^n$  be the index set of all support tokens per Definition 3. Let  
 555  $\bar{\mathcal{T}}_i = [T] - \mathcal{T}_i - \{\alpha_i\}$ . For any  $W \in S_{p,\mu_0}(W_{\text{mm}}^\alpha)$  with  $\mu_0$  defined in (9) and  $\|W\|_{p,p} = \|W_{\text{mm}}^\alpha\|_{p,p}$ ,  
 556 we have

$$(X_{it} - X_{i\tau})^\top W z_i \geq \frac{3}{2}\delta > 0, \quad (13a)$$

$$(X_{i\alpha_i} - X_{i\tau})^\top W z_i \geq 1 + \frac{3}{2}\delta, \quad (13b)$$

$$1 + \frac{1}{2}\delta \geq (X_{i\alpha_i} - X_{it})^\top W z_i \geq 1 - \frac{1}{2}\delta, \quad (13c)$$

557 for all  $t \in \mathcal{T}_i$  and  $\tau \in \bar{\mathcal{T}}_i$

558 *Proof.* Let

$$\bar{W} := \frac{W}{\|W\|_{p,p}} \quad \text{and} \quad \bar{W}_{\text{mm}}^\alpha := \frac{W_{\text{mm}}^\alpha}{\|W_{\text{mm}}^\alpha\|_{p,p}}.$$

559 Using Lemma 3 and the definition of  $S_{p,\mu_0}(W_{\text{mm}}^\alpha)$  in (5a), when  $p \geq 2$ ,

$$\begin{aligned} \|\bar{W} - \bar{W}_{\text{mm}}^\alpha\|_{p,p}^p &\leq 2^p p D_\psi(\bar{W}_{\text{mm}}^\alpha, \bar{W}) \\ &\leq 2^p p \mu_0 \\ &= \left(\frac{\delta}{4A}\right)^p, \end{aligned}$$

560 which implies that

$$\|\bar{W} - \bar{W}_{\text{mm}}^\alpha\|_{p,p} \leq \frac{\delta}{4A}.$$

561 When  $p \in (1, 2)$ , we can also use Lemmas 2 and 3 to obtain

$$\begin{aligned} \|\bar{W} - \bar{W}_{\text{mm}}^\alpha\|_{p,p} &\leq d^{\frac{2}{p}-1} \|\bar{W} - \bar{W}_{\text{mm}}^\alpha\|_{2,2} \\ &\leq d^{\frac{2}{p}-1} \frac{\sqrt{p}}{p-1} \sqrt{D_\psi(\bar{W}_{\text{mm}}^\alpha, \bar{W})} \\ &\leq d^{\frac{2}{p}-1} \frac{\sqrt{p}}{p-1} \sqrt{\mu_0} = \frac{\delta}{4A}, \end{aligned}$$

562 where the last inequality uses the definition of  $S_{p,\mu_0}(W_{\text{mm}}^\alpha)$  in (5a).

563 Therefore, either way, we have

$$\|W - W_{\text{mm}}^\alpha\|_{p,p} \leq \frac{\delta}{4A} \|W_{\text{mm}}^\alpha\|_{p,p}.$$

564 We will proceed to show a bound on  $(X_{it_1} - X_{it_2})^\top (W - W_{\text{mm}}^\alpha) z_i$  for any  $i \in [n]$  and any token  
 565 indices  $t_1, t_2 \in [T]$ . To do that, let us focus on the term  $X_{it_1}^\top (W - W_{\text{mm}}^\alpha) z_i$  first,

$$\begin{aligned} |X_{it_1}^\top (W - W_{\text{mm}}^\alpha) z_i| &= |\langle W - W_{\text{mm}}^\alpha, X_{it_1} z_i^\top \rangle| \\ &\leq \|W - W_{\text{mm}}^\alpha\|_{p,p} \cdot \|X_{it_1} z_i^\top\|_{\frac{p}{p-1}, \frac{p}{p-1}} \\ &\leq \frac{\delta}{4A} \|W_{\text{mm}}^\alpha\|_{p,p} \cdot \|X_{it_1} z_i^\top\|_{\frac{p}{p-1}, \frac{p}{p-1}} \\ &\leq \frac{\delta}{4A} \cdot A \\ &= \frac{\delta}{4}. \end{aligned}$$

566 The first inequality above uses Hölder's Inequality. We now have

$$|(X_{it_1} - X_{it_2})^\top (W - W_{\text{mm}}^\alpha) z_i| \leq \frac{1}{2} \delta.$$

567 To obtain the first inequality of the lemma in (13a), for all  $t \in \mathcal{T}_i$  and  $\tau \in \bar{\mathcal{T}}_i$ , we have

$$\begin{aligned} (X_{it} - X_{i\tau})^\top W z_i &\geq (X_{it} - X_{i\tau})^\top W_{\text{mm}}^\alpha z_i + (X_{it} - X_{i\tau})^\top (W - W_{\text{mm}}^\alpha) z_i \\ &\geq 2\delta' - \frac{1}{2} \delta \geq \frac{3}{2} \delta. \end{aligned}$$

568 To get the second inequality in (13b), for all  $\tau \in \bar{\mathcal{T}}_i$ , we have

$$\begin{aligned} (X_{i\alpha_i} - X_{i\tau})^\top W z_i &\geq (X_{i\alpha_i} - X_{i\tau})^\top W_{\text{mm}}^\alpha z_i + (X_{i\alpha_i} - X_{i\tau})^\top (W - W_{\text{mm}}^\alpha) z_i \\ &\geq 1 + 2\delta' - \frac{1}{2} \delta \geq 1 + \frac{3}{2} \delta. \end{aligned}$$

569 Finally, to get the last inequality in (13c), for all  $t \in \mathcal{T}_i$ , we have

$$\begin{aligned} |(X_{i\alpha_i} - X_{it})^\top W z_i - 1| &= |(X_{i\alpha_i} - X_{it})^\top W_{\text{mm}}^\alpha z_i + (X_{i\alpha_i} - X_{it})^\top (W - W_{\text{mm}}^\alpha) z_i - 1| \\ &= |(X_{i\alpha_i} - X_{it})^\top (W - W_{\text{mm}}^\alpha) z_i| \leq \frac{1}{2} \delta, \end{aligned}$$

570 which implies that

$$1 + \frac{1}{2} \delta \geq (X_{i\alpha_i} - X_{it})^\top W z_i \geq 1 - \frac{1}{2} \delta.$$

571

□

572 The following two lemmas aim at bounding the correlation between the gradient and the attention  
573 matrix parameter, each of which is a generalization of [52, Lemmas 13 and 14] for the generalized  $\ell_p$   
574 norm.

575 **Lemma 10.** *Suppose Assumption A holds. Let  $\alpha = (\alpha_i)_{i=1}^n$  be locally optimal tokens as per  
576 Definition 3, and let  $W_{\text{mm}}^\alpha$  be the solution to ( $\ell_p$ -AttSVM). There exists a dataset-dependent constant  
577  $R_\delta = \mathcal{O}(1/\delta)$  such that for all  $W, V \in C_{p, \mu_0, R_\delta}(W_{\text{mm}}^\alpha)$  with  $\|V\|_{p,p} = \|W_{\text{mm}}^\alpha\|_{p,p}$ ,  $\delta$  and  $\mu_0$   
578 defined in (7) and (9), respectively,*

$$-\langle \nabla \mathcal{L}(W), V \rangle = \Omega \left( e^{-\frac{\|W\|_{p,p}}{\|W_{\text{mm}}^\alpha\|_{p,p}} (1 + \frac{1}{2} \delta)} \right) > 0.$$

579 *Proof.* Let

$$h_i := X_i W z_i, \quad \tilde{h}_i := X_i V z_i, \quad l'_i := l'(\gamma_i^\top \sigma(h_i)), \quad \text{and} \quad s_i = \sigma(h_i).$$

580 Therefore,

$$\begin{aligned} \langle \nabla \mathcal{L}(W), V \rangle &= \frac{1}{n} \sum_{i=1}^n l'_i \langle X_i^\top (\text{diag}(s_i) - s_i s_i^\top) \gamma_i z_i^\top, V \rangle \\ &= \frac{1}{n} \sum_{i=1}^n l'_i \langle (\text{diag}(s_i) - s_i s_i^\top) \gamma_i, X_i V z_i \rangle \\ &= \frac{1}{n} \sum_{i=1}^n l'_i \langle (\text{diag}(s_i) - s_i s_i^\top) \gamma_i, \tilde{h}_i \rangle \\ &= \frac{1}{n} \sum_{i=1}^n l'_i \tilde{h}_i^\top (\text{diag}(s_i) - s_i s_i^\top) \gamma_i, \end{aligned}$$

581

$$-\langle \nabla \mathcal{L}(W), V \rangle = \frac{1}{n} \sum_{i=1}^n (-l'_i) \tilde{h}_i^\top (\text{diag}(s_i) - s_i s_i^\top) \gamma_i. \quad (14)$$

582 The value  $\gamma_i^\top \sigma(h_i)$  for any  $i \in [n]$  must be bounded, and the bound is only dataset-dependent, so by  
 583 Assumption A,  $l'_i$  is bounded for any  $i \in [n]$  by some bound that is dataset-dependent. Furthermore,  
 584 because  $l$  is decreasing,  $-l'$  is always non-negative, so an easier approach is to lower-bound the  
 585 following for each  $i \in [n]$ ,

$$\tilde{h}_i^\top s_i s_i^\top \gamma_i - \tilde{h}_i^\top \text{diag}(s_i) \gamma_i.$$

586 Next, we can get for all  $i \in [n]$  and  $t \in [T]$  that

$$\begin{aligned} \tilde{h}_{it} &= X_{it}^\top V z_i = \langle X_{it} z_i^\top, V \rangle \\ &\leq \|V\|_{p,p} \|X_{it} z_i^\top\|_{\frac{p}{p-1}} \\ &\leq A, \end{aligned}$$

587 where  $A$  is defined in (8).

588 Therefore, if we drop the  $i$  notation and let  $\alpha_i = 1$ , and use [52, Lemma 7],

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

589 Let us attempt to remove the non-support tokens from the sum above by bounding the sum of the  
 590 term for the non-supports,

$$\left| \sum_{t \in \mathcal{T}} (\tilde{h}_1 - \tilde{h}_t) s_t (\gamma_1 - \gamma_t) \right| \leq 2 \max_{t \in [T]} \{|\tilde{h}_t|\} Q(W) \Gamma \leq 2A Q(W) \Gamma.$$

591 Therefore,

$$\left| \tilde{h}^\top s s^\top \gamma - \tilde{h}^\top \text{diag}(s) \gamma - \sum_{t \in \mathcal{T}} (\tilde{h}_1 - \tilde{h}_t) s_t (\gamma_1 - \gamma_t) \right| \leq 2\Gamma A ((1 - s_1)^2 + Q(W)),$$

592 which implies that

$$\tilde{h}^\top s s^\top \gamma - \tilde{h}^\top \text{diag}(s) \gamma \geq \sum_{t \in \mathcal{T}} (\tilde{h}_1 - \tilde{h}_t) s_t (\gamma_1 - \gamma_t) - 2\Gamma A ((1 - s_1)^2 + Q(W)).$$

593 Using Lemma 9, we have

$$\tilde{h}^\top s s^\top \gamma - \tilde{h}^\top \text{diag}(s) \gamma \geq \left(1 - \frac{1}{2}\delta\right) \sum_{t \in \mathcal{T}} s_t (\gamma_1 - \gamma_t) - 2\Gamma A ((1 - s_1)^2 + Q(W)). \quad (15)$$

594 To proceed, we can upper-bound  $1 - s_1$  and  $Q(W)$ . For bounding  $1 - s_1$ , let  $\tau > 1$  be some index  
 595 that maximizes  $X_\tau^\top W z$ , so

$$\begin{aligned} 1 - s_1 &= \frac{\sum_{t=2}^T e^{X_t^\top W z}}{\sum_{t=1}^T e^{X_t^\top W z}} \leq \frac{(T-1)e^{X_\tau^\top W z}}{(T-1)e^{X_\tau^\top W z} + e^{X_1^\top W z}} \\ &\leq \frac{T}{T + e^{(X_1 - X_\tau)^\top W z}} \\ &\leq \frac{T}{T + e^{\frac{\|W\|_{p,p}}{\|W_{\text{mm}}^\alpha\|_{p,p}} (1 - \frac{1}{2}\delta)}} \\ &\leq \frac{T}{e^{\frac{\|W\|_{p,p}}{\|W_{\text{mm}}^\alpha\|_{p,p}} (1 - \frac{1}{2}\delta)}}, \end{aligned}$$

596 with the last inequality using the third inequality Lemma 9.

597 For ease of notation, denote

$$R' := \frac{\|W\|_{p,p}}{\|W_{\text{mm}}^\alpha\|_{p,p}}. \quad (16)$$

598 To upper bound  $Q(W)$ , we use a method similar to that for  $1 - s_1$ , but we utilize the second inequality  
 599 of Lemma 9 instead of the first. This gives:

$$Q(W) \leq \frac{T}{T + e^{(1+\frac{3}{2}\delta)R'}} \leq \frac{T}{e^{(1+\frac{3}{2}\delta)R'}}.$$

600 Therefore, we have

$$\begin{aligned} 2\Gamma A((1 - s_1)^2 + Q(W)) &\leq 2\Gamma A\left(\frac{T^2}{e^{(2-\delta)R'}} + \frac{T}{e^{(1+\frac{3}{2}\delta)R'}}\right) \\ &\leq \frac{2\Gamma AT(T+1)}{e^{(1+\frac{3}{2}\delta)R'}}. \end{aligned} \quad (17)$$

601 Now it is time to lower-bound the sum on the right side of Equation (15). When the set of supports is  
 602 empty, that sum is zero. However, if it is not empty,

$$\sum_{t \in \mathcal{T}} s_t(\gamma_1 - \gamma_t) \geq S(W)\gamma^{\text{gap}}.$$

603 If we let  $\tau \in \mathcal{T}$  be the support index that minimizes  $X_\tau^\top Wz$ , then

$$\begin{aligned} S(W) &= \frac{\sum_{t \in \mathcal{T}} e^{X_t^\top Wz}}{\sum_{t=1}^T e^{X_t^\top Wz}} \geq \frac{e^{X_\tau^\top Wz}}{T e^{X_1^\top Wz}} = \frac{1}{T e^{(X_1 - X_\tau)^\top Wz}} \\ &\geq \frac{1}{T e^{(1+\frac{1}{2}\delta)R'}}, \end{aligned}$$

604 with the last inequality coming from the third inequality of Lemma 9.

605 Therefore,

$$\sum_{t \in \mathcal{T}} s_t(\gamma_1 - \gamma_t) \geq \frac{\gamma^{\text{gap}}}{T e^{(1+\frac{1}{2}\delta)R'}} > 0.$$

606 Using Equation (15), we get that if the support index set is empty,

$$\tilde{h}^\top s s^\top \gamma - \tilde{h}^\top \text{diag}(s)\gamma \geq -\frac{2\Gamma AT(T+1)}{e^{(1+\frac{3}{2}\delta)R'}},$$

607 otherwise,

$$\tilde{h}^\top s s^\top \gamma - \tilde{h}^\top \text{diag}(s)\gamma \geq \frac{\gamma^{\text{gap}}}{T e^{(1+\frac{1}{2}\delta)R'}} \left(1 - \frac{1}{2}\delta\right) - \frac{2\Gamma AT(T+1)}{e^{(1+\frac{3}{2}\delta)R'}}.$$

608 Plugging everything back into Equation (14), and considering that some samples will have non-empty  
 609 support index sets, we have:

$$\begin{aligned} -\langle \mathcal{L}(W), V \rangle &\geq -\frac{\min_{i \in \mathcal{T}_i} \{\gamma_i^{\text{gap}}\}}{nT e^{(1+\frac{1}{2}\delta)R'}} \left(1 - \frac{1}{2}\delta\right) \max_{i=1}^n \{l'_i\} \\ &\quad + \frac{2\Gamma AT(T+1)}{e^{(1+\frac{3}{2}\delta)R'}} \sum_{i=1}^n l'_i = \Omega\left(e^{-(1+\frac{1}{2}\delta)R'}\right). \end{aligned} \quad (18)$$

610 Let

$$\bar{L} := \frac{\sum_{i=1}^n l'_i}{\max_{i=1}^n \{l'_i\}}. \quad (19)$$

611 Note that using Assumption A,  $\bar{L}$  is positive. Hence, using (19) and (18), the term  $-\langle \mathcal{L}(W), V \rangle$  is  
 612 positive when

$$R' \geq \frac{1}{\delta} \log \left( \frac{2\Gamma \bar{L} AT^2(T+1)n}{\min_{i \in \mathcal{T}_i} \{\gamma_i^{\text{gap}}\} (1 - \frac{1}{2}\delta)} \right),$$

613 or equivalently, from (16), we have

$$\|W\|_{p,p} \geq \frac{\|W_{\text{mm}}^\alpha\|_{p,p}}{\delta} \log \left( \frac{2\Gamma \bar{L} AT^2(T+1)}{\min_{i \in \mathcal{T}_i} \{\gamma_i^{\text{gap}}\} (1 - \frac{1}{2}\delta)} \right).$$

614

□



615 Finally, we introduce the following lemma to help understand the correlation between the gradient of  
616 the objective and the parameter.

617 **Lemma 11.** *Suppose Assumption A holds. Let  $\alpha = (\alpha_i)_{i=1}^n$  be locally optimal tokens as per  
618 Definition 3, let  $W_{\text{mm}}^\alpha$  be the ( $\ell_p$ -AttSVM) solution, and let  $R_\delta$  be the constant from Lemma 10. For  
619 any choice of  $\pi \in (0, 1)$ , there exists  $R_\pi$  that depends on  $\pi$  defined as*

$$R_\pi := \max \left\{ R_\delta, \mathcal{O} \left( \frac{1}{\pi\delta} \log \frac{\delta}{\pi} \right) \right\},$$

620 such that for all  $W \in C_{p, \mu_0, R_\pi}(W_{\text{mm}}^\alpha)$ ,

$$\left\langle \nabla \mathcal{L}(W), \frac{W}{\|W\|_{p,p}} \right\rangle \geq (1 + \pi) \left\langle \nabla \mathcal{L}(W), \frac{W_{\text{mm}}^\alpha}{\|W_{\text{mm}}^\alpha\|_{p,p}} \right\rangle.$$

621 *Proof.* Let

$$\begin{aligned} h_i &:= X_i W z_i, & \tilde{h}_i &:= X_i W_{\text{mm}}^\alpha z_i, & l'_i &:= l'(\gamma_i^\top \sigma(h_i)), \\ s_i &:= \sigma(h_i), & \bar{W} &:= \frac{\|W_{\text{mm}}^\alpha\|_{p,p} W}{\|W\|_{p,p}}, & \text{and } \bar{h}_i &:= X_i \bar{W} z_i. \end{aligned} \quad (20)$$

622 By decomposing  $\mathcal{L}(W)$  into its sum and using Lemma 7, the main inequality is equivalent to the  
623 following,

$$\begin{aligned} & \sum_{i=1}^n (-l'_i) \langle X_i^\top (\text{diag}(s_i) - s_i s_i^\top) \gamma_i z_i^\top, \bar{W} \rangle \\ & \leq (1 + \pi) \sum_{i=1}^n (-l'_i) \langle X_i^\top (\text{diag}(s_i) - s_i s_i^\top) \gamma_i z_i^\top, W_{\text{mm}}^\alpha \rangle, \end{aligned}$$

624 which implies that

$$\begin{aligned} & \sum_{i=1}^n (-l'_i) \langle (\text{diag}(s_i) - s_i s_i^\top) \gamma_i, X_i \bar{W} z_i \rangle \\ & \leq (1 + \pi) \sum_{i=1}^n (-l'_i) \langle (\text{diag}(s_i) - s_i s_i^\top) \gamma_i, X_i W_{\text{mm}}^\alpha z_i \rangle. \end{aligned}$$

625 Using (20), we get

$$\sum_{i=1}^n (-l'_i) \langle (\text{diag}(s_i) - s_i s_i^\top) \gamma_i, \bar{h}_i \rangle \leq (1 + \pi) \sum_{i=1}^n (-l'_i) \langle (\text{diag}(s_i) - s_i s_i^\top) \gamma_i, \tilde{h}_i \rangle,$$

626 which gives

$$\sum_{i=1}^n (-l'_i) \bar{h}_i^\top (\text{diag}(s_i) - s_i s_i^\top) \gamma_i \leq (1 + \pi) \sum_{i=1}^n (-l'_i) \tilde{h}_i^\top (\text{diag}(s_i) - s_i s_i^\top) \gamma_i.$$

627 Hence,

$$\sum_{i=1}^n (-l'_i) \left[ (1 + \pi) \left( \tilde{h}_i^\top \text{diag}(s_i) \gamma_i - \tilde{h}_i^\top s_i s_i^\top \gamma_i \right) - \left( \bar{h}_i^\top \text{diag}(s_i) \gamma_i - \bar{h}_i^\top s_i s_i^\top \gamma_i \right) \right] \geq 0.$$

628 Using a similar technique as the one we used to prove Lemma 10,

$$\begin{aligned} & \left| \tilde{h}_i^\top \text{diag}(s_i) \gamma_i - \tilde{h}_i^\top s_i s_i^\top \gamma_i - \sum_{t \in \mathcal{T}_i} (\tilde{h}_{i\alpha_i} - \tilde{h}_{it}) s_{it} (\gamma_{i\alpha_i} - \gamma_{it}) \right| \\ & \leq 2\Gamma A((1 - s_{i\alpha_i})^2 + Q_i(W)). \end{aligned}$$

629 Similarly,

$$\begin{aligned} & \left| \bar{h}_i^\top \text{diag}(s_i) \gamma_i - \bar{h}_i^\top s_i s_i^\top \gamma_i - \sum_{t \in \mathcal{T}_i} (\bar{h}_{i\alpha_i} - \bar{h}_{it}) s_{it} (\gamma_{i\alpha_i} - \gamma_{it}) \right| \\ & \leq 2\Gamma A((1 - s_{i\alpha_i})^2 + Q_i(W)). \end{aligned}$$

630 Therefore, it is enough to prove that

$$\begin{aligned} & \sum_{i=1}^n (-l'_i) \left( (1 + \pi) \left( \sum_{t \in \mathcal{T}_i} (\tilde{h}_{i\alpha_i} - \tilde{h}_{it}) s_{it} (\gamma_{i\alpha_i} - \gamma_{it}) - 2\Gamma A((1 - s_{i\alpha_i})^2 + Q_i(W)) \right) \right. \\ & \quad \left. - \left( \sum_{t \in \mathcal{T}_i} (\bar{h}_{i\alpha_i} - \bar{h}_{it}) s_{it} (\gamma_{i\alpha_i} - \gamma_{it}) + 2\Gamma A((1 - s_{i\alpha_i})^2 + Q_i(W)) \right) \right), \end{aligned} \quad (21)$$

631 Using the fact that  $\pi < 1$  and using Equation (17), we get another lower-bound

$$\sum_{i=1}^n \sum_{t \in \mathcal{T}_i} (-l'_i) (1 + \pi - (\bar{h}_{i\alpha_i} - \bar{h}_{it})) s_{it} (\gamma_{i\alpha_i} - \gamma_{it}) + \frac{6\Gamma A T (T + 1)}{e^{(1 + \frac{3}{2}\delta)R'}} \sum_{i=1}^n l'_i, \quad (22)$$

632 with  $R'$  again being  $\frac{\|W\|_{p,p}}{\|W_{\text{mm}}^\alpha\|_{p,p}}$ . Next, we analyze the softmax probability  $s_{it}$ , and lower and upper-  
633 bound them in terms of  $R'$  and  $\bar{h}_{i\alpha_i} - \bar{h}_{it}$ . For the lower-bound,

$$\begin{aligned} s_{it} &= \frac{e^{\bar{h}_{it}R'}}{\sum_{\tau \in [T]} e^{\bar{h}_{i\tau}R'}} \geq \frac{e^{\bar{h}_{it}R'}}{T e^{\bar{h}_{i\alpha_i}R'}} \\ &= \frac{1}{T} e^{-(\bar{h}_{i\alpha_i} - \bar{h}_{it})R'}. \end{aligned}$$

634 For the upper-bound,

$$\begin{aligned} s_{it} &= \frac{e^{\bar{h}_{it}R'}}{\sum_{\tau \in [T]} e^{\bar{h}_{i\tau}R'}} \leq \frac{e^{\bar{h}_{it}R'}}{e^{\bar{h}_{i\alpha_i}R'}} \\ &= e^{-(\bar{h}_{i\alpha_i} - \bar{h}_{it})R'}. \end{aligned}$$

635 In both bounds, the main inequality derivation stems from the fact that  $\bar{h}_{i\alpha_i} > \bar{h}_{i\tau}$  for all  $\tau \in [T]$ ,  
636 which we obtain from Lemma 9. Now, we analyze the left double-summation in Equation (22). To  
637 analyze the sum, let  $\mathcal{I}$  be the subset of  $[n] \times [T]$  that contains all  $(i, t)$  such that  $t \in \mathcal{T}_i$ . Furthermore,  
638 let

$$\begin{aligned} \mathcal{I}_1 &:= \{(i, t) \in \mathcal{I} \mid \bar{h}_{i\alpha_i} - \bar{h}_{it} \leq 1\}, \\ \mathcal{I}_2 &:= \{(i, t) \in \mathcal{I} \mid 1 < \bar{h}_{i\alpha_i} - \bar{h}_{it} \leq 1 + \pi\}, \\ \mathcal{I}_3 &:= \{(i, t) \in \mathcal{I} \mid \bar{h}_{i\alpha_i} - \bar{h}_{it} > 1 + \pi\}. \end{aligned}$$

639 Therefore, we can split the sum above into the sum over  $\mathcal{I}_1, \mathcal{I}_2$ , and  $\mathcal{I}_3$ . The set  $\mathcal{I}_1$  in particular must  
640 be non-empty because  $\|\bar{W}\|_{p,p} = \|W_{\text{mm}}^\alpha\|_{p,p}$ , meaning that one of the constraints in the  $\ell_p$ -AttSVM  
641 problem must either be fulfilled exactly or violated.

642 The sum over  $\mathcal{I}_1$  must be positive and is at least

$$-\frac{\pi}{T} \min_{i \in \mathcal{I}_1} \{\gamma_i^{gap}\} e^{-R'} \max_{i=1}^n \{l'_i\}.$$

643 The sum over  $\mathcal{I}_2$  must be non-negative, and the sum over  $\mathcal{I}_3$  is negative can be bounded from below  
644 using Lemma 9

$$\frac{1}{2} \delta \max_{i \in \mathcal{I}_3} \{\bar{\gamma}_i^{gap}\} T e^{-(1+\pi)R'} \sum_{i=1}^n l'_i.$$

645 Putting things together into Equation (22), we get that we want the following to be non-negative

$$-\frac{\pi}{T} \min_{i \in \mathcal{I}_1} \{\gamma_i^{gap}\} e^{-R'} \max_{i=1}^n \{l'_i\} + \frac{1}{2} \delta \max_{i \in \mathcal{I}_3} \{\bar{\gamma}_i^{gap}\} T e^{-(1+\pi)R'} \sum_{i=1}^n l'_i \\ + 6\Gamma AT(T+1) e^{-(1+\frac{3}{2}\delta)R'} \sum_{i=1}^n l'_i.$$

646 This can be achieved when

$$R' \geq \frac{1}{\min\{\pi, \frac{3}{2}\delta\}} \log \left( \frac{\frac{1}{2} \delta \max_{i \in \mathcal{I}_3} \{\bar{\gamma}_i^{gap}\} T^2 + 6\Gamma AT^2(T+1) \sum_{i=1}^n l'_i}{\pi \min_{i \in \mathcal{I}_1} \{\gamma_i^{gap}\} \max_{i=1}^n \{l'_i\}} \right),$$

647 or equivalently,

$$\|W\|_{p,p} \geq \frac{\|W_{mm}^\alpha\|_{p,p}}{\min\{\pi, \frac{3}{2}\delta\}} \log \left( \frac{\frac{1}{2} \delta \max_{i \in \mathcal{I}_3} \{\bar{\gamma}_i^{gap}\} T^2 + 6\Gamma AT^2(T+1) \sum_{i=1}^n l'_i}{\pi \min_{i \in \mathcal{I}_1} \{\gamma_i^{gap}\} \max_{i=1}^n \{l'_i\}} \right),$$

648 which means that such dataset dependent  $R_\pi$  exists.  $\square$

#### 649 **B.4 Lemma for Analyzing $\ell_p$ -AttGD**

650 We introduce the lemmas for analyzing  $\ell_p$ -AttGD. The first we prove is Lemma 12, which describes  
651 the lower bound of the  $W$  parameter at every iterate.

652 **Lemma 12.** *Suppose Assumption A holds. For the sequence  $\{W(k)\}_{k \geq 0}$  generated by  $\ell_p$ -AttGD,*  
653 *we have*

$$\|W(k+1)\|_{p,p}^{p-1} \geq \|W(k)\|_{p,p}^{p-1} + \frac{\eta}{\|W(k)\|_{p,p}} \langle -\nabla \mathcal{L}(W(k)), W(k) \rangle.$$

654 *Proof.* With  $\psi(W) = \frac{1}{p} \|W\|_{p,p}$ , the derivative  $\nabla \psi(\cdot)$  is computed as follows:

$$\nabla \psi(W) = (\text{sign}(W_{ij}) |W_{ij}|^{p-1})_{1 \leq i,j \leq d}.$$

655 Thus, we have

$$\langle \nabla \psi(W), W \rangle = \sum_{i,j} \text{sign}(W_{ij}) |W_{ij}|^{p-1} W_{ij} = \|W\|_{p,p}^p.$$

656 Using this fact, we take the inner product of both sides of (3) with  $W(k)$ :

$$\langle \nabla \psi(W(k+1)), W(k) \rangle = \langle \nabla \psi(W(k)), W(k) \rangle + \eta \langle -\nabla \mathcal{L}(W(k)), W(k) \rangle,$$

657

$$\langle \nabla \psi(W(k+1)), W(k) \rangle = \|W(k)\|_{p,p}^p + \eta \langle -\nabla \mathcal{L}(W(k)), W(k) \rangle. \quad (23)$$

658 The left side of the above equation is upper-bounded by

$$\sum_{i,j} \text{sign}(W_{ij}(k+1)) |W_{ij}(k+1)|^{p-1} W_{ij}(k) \leq \sum_{i,j} |W_{ij}(k+1)|^{p-1} |W_{ij}(k)|.$$

659 Using Hölder's inequality:

$$\sum_{i,j} |W_{ij}(k+1)|^{p-1} |W_{ij}(k)| \leq \left( \sum_{i,j} (|W_{ij}(k+1)|^{p-1})^{\frac{p}{p-1}} \right)^{\frac{p-1}{p}} \left( \sum_{i,j} |W_{ij}(k)|^p \right)^{\frac{1}{p}} \\ = \|W(k+1)\|_{p,p}^{p-1} \|W(k)\|_{p,p}.$$

660 Combining this result with (23), we get:

$$\|W(k+1)\|_{p,p}^{p-1} \geq \|W(k)\|_{p,p}^{p-1} + \frac{\eta}{\|W(k)\|_{p,p}} \langle -\nabla \mathcal{L}(W(k)), W(k) \rangle.$$

661  $\square$

662 Next, we show several tools for analyzing the algorithm further and for analyzing the Bregman  
 663 divergence. The following two specifically are from [50, Lemma 18, 3], and so the proofs are  
 664 omitted.

665 **Lemma 13.** *Suppose Assumptions A hold and  $\eta$  is small enough. For the sequence  $\{W(k)\}_{k \geq 0}$   
 666 generated by  $\ell_p$ -AttGD, we have*

$$\begin{aligned} \frac{p-1}{p} \|W(k+1)\|_{p,p}^p - \frac{p-1}{p} \|W(k)\|_{p,p}^p + \eta \mathcal{L}(W(k+1)) - \eta \mathcal{L}(W(k)) \\ \leq \langle -\eta \nabla \mathcal{L}(W(k)), W(k) \rangle. \end{aligned} \quad (24)$$

667 **Lemma 14.** *Suppose Assumptions A hold. Consider the sequence  $W(k)$  generated by Algorithm  
 668  $\ell_p$ -AttGD. Given that the step size  $\eta$  is sufficiently small, then the ERM objective  $\mathcal{L}(W(k))$  is  
 669 decreasing in  $k$ .*

670 This following is a well-known lemma, so the proof is omitted.

671 **Lemma 15** (Bregman Divergences Cosine Law). *For any  $w, w', w''$  that are all vectors or matrices  
 672 with the same dimensionalities, we have*

$$D_\psi(w, w') = D_\psi(w, w'') + D_\psi(w'', w') - \langle \nabla \psi(w') - \nabla \psi(w''), w - w'' \rangle.$$

673 The following is adapted from [50, Equation 12] for the case of our attention model. Our proof is  
 674 quite similar, except that we use our version of the gradient correlation lemma.

675 **Lemma 16.** *Suppose Assumptions A hold. Consider the sequence  $W(k)$  generated by Algorithm  
 676  $\ell_p$ -AttGD. For any  $\pi \in (0, 1)$ , if  $W(k) \in C_{p, \mu_0, R_\pi}(W_{\text{mm}}^\alpha)$ , with  $R_\pi$  being the constant from Lemma  
 677 11, then for a small enough step size  $\eta$ ,*

$$\begin{aligned} \langle \nabla \psi(W(k+1)) - \nabla \psi(W(k)), \bar{W}_{\text{mm}}^\alpha \rangle \geq \frac{1}{1+\pi} (\|W(k+1)\|_{p,p}^{p-1} - \|W(k)\|_{p,p}^{p-1}) \\ + \frac{\eta}{\|W(k)\|_{p,p}} (\mathcal{L}(W(k+1)) - \mathcal{L}(W(k))). \end{aligned} \quad (25)$$

678 *Proof.* Let  $\bar{W}_{\text{mm}}^\alpha = \frac{W_{\text{mm}}^\alpha}{\|W_{\text{mm}}^\alpha\|_{p,p}}$ . Using the  $\ell_p$ -AttGD algorithm equation,

$$\langle \nabla \psi(W(k+1)) - \nabla \psi(W(k)), \bar{W}_{\text{mm}}^\alpha \rangle = \langle -\eta \nabla \mathcal{L}(W(k)), \bar{W}_{\text{mm}}^\alpha \rangle.$$

679 Then, using Lemma 11, we get that

$$\langle -\eta \nabla \mathcal{L}(W(k)), \bar{W}_{\text{mm}}^\alpha \rangle \geq \frac{1}{(1+\pi)\|W(k)\|_{p,p}} \langle -\eta \nabla \mathcal{L}(W(k)), W(k) \rangle,$$

680 and using Lemma 13, we get that this is lower-bounded by

$$\frac{p-1}{p(1+\pi)\|W(k)\|_{p,p}} (\|W(k+1)\|_{p,p}^p - \|W(k)\|_{p,p}^p) + \frac{\eta}{(1+\pi)\|W(k)\|_{p,p}} (\mathcal{L}(W(k+1)) - \mathcal{L}(W(k))).$$

681 By Lemma 10,  $\langle -\eta \nabla \mathcal{L}(W(k)), W(k) \rangle > 0$ , so by Lemma 12,  $\|W(k+1)\|_{p,p} \geq \|W(k)\|_{p,p}$ .  
 682 Therefore, we can use Lemma 4 to get that the above is lower-bounded by

$$\frac{1}{1+\pi} (\|W(k+1)\|_{p,p}^{p-1} - \|W(k)\|_{p,p}^{p-1}) + \frac{\eta}{(1+\pi)\|W(k)\|_{p,p}} (\mathcal{L}(W(k+1)) - \mathcal{L}(W(k))).$$

683 From Lemma 14, we get that we can lower-bound the above further using the right hand side of  
 684 (25).  $\square$

685 With all these lemmas in hand, we provide the following  
 686 Lemma 17.

687 **Lemma 17.** *Suppose Assumptions A holds and that the  
 688 step size  $\eta$  is sufficiently small. For any  $\mu \in (0, \mu_0]$   
 689 and any locally optimal tokens  $(\alpha_i)_{i=1}^n$  as per Definition 3,  
 690 there exists constants  $R_\mu$  and  $\mu' \in (0, \mu]$  that  
 691 depends on the dataset and  $\mu$  such that if  $C_1$  is the  
 692 wider cone  $C_{p, \mu, R_\mu}(W_{\text{mm}}^\alpha)$  and  $C_2$  is the thinner cone  
 693  $C_{p, \mu', R_\mu}(W_{\text{mm}}^\alpha)$ , then if  $W(0) \in C_2$ , then  $W(k) \in C_1$   
 694 for all positive indices  $k$ .*

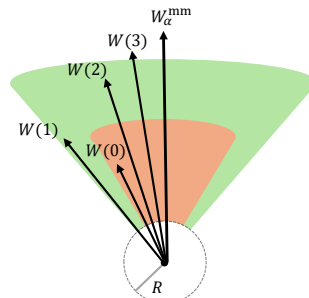


Figure 2: Illustration of Lemma 17.  $W(k)$  for all positive indices  $k$  are within the larger set.

695 *Proof.* Let  $\pi$  be some positive real number that we deter-  
 696 mine later, and let  $R_\pi$  be as described in Lemma 11.

697 For the proof, we use induction with the assumption that  
 698  $W(k) \in C_{p,\mu,R_\pi}(W_{\text{mm}}^\alpha)$  for all  $k = 0, \dots, K-1$ . We  
 699 aim to find the correct  $\mu'$  and  $R_\mu$  such that  $W(K) \in$   
 700  $C_{p,\mu,R_\pi}(W_{\text{mm}}^\alpha)$ .

701 Denote  $\bar{W}(k) := \frac{W(k)}{\|W(k)\|_{p,p}}$ , so

$$\begin{aligned} D_\psi(\bar{W}_{\text{mm}}^\alpha, \bar{W}(k)) &= \frac{1}{p} \|\bar{W}_{\text{mm}}^\alpha\|_{p,p} - \frac{1}{p} \|\bar{W}(k)\|_{p,p} - \langle \nabla\psi(\bar{W}(k)), \bar{W}_{\text{mm}}^\alpha - \bar{W}(k) \rangle \\ &= 1 - \langle \nabla\psi(\bar{W}(k)), \bar{W}_{\text{mm}}^\alpha \rangle. \end{aligned}$$

702 So now, let us analyze the term  $\langle \nabla\psi(\bar{W}(K)), \bar{W}_{\text{mm}}^\alpha \rangle$  using the inductive hypothesis on  $k =$   
 703  $0, 1, \dots, K-1$ . Lemma 16 tells us that

$$\begin{aligned} \langle \nabla\psi(W(k+1)) - \nabla\psi(W(k)), \bar{W}_{\text{mm}}^\alpha \rangle &\geq \frac{\|W(k+1)\|_{p,p}^{p-1} - \|W(k)\|_{p,p}^{p-1}}{(1+\pi)} \\ &\quad + \frac{\eta}{\|W(k)\|_{p,p}} (\mathcal{L}(W(k+1)) - \mathcal{L}(W(k))). \end{aligned} \quad (26)$$

704 Since this is true for all  $k = 0, 1, \dots, K-1$ , and since  $\|W(k)\|_{p,p}$  is increasing in  $k$ , we can sum all  
 705 the above inequalities and get the following,

$$\begin{aligned} \langle \nabla\psi(W(K)) - \nabla\psi(W(0)), \bar{W}_{\text{mm}}^\alpha \rangle &\geq \frac{\|W(K)\|_{p,p}^{p-1} - \|W(0)\|_{p,p}^{p-1}}{(1+\pi)} \\ &\quad + \frac{\eta}{\|W(0)\|_{p,p}} (\mathcal{L}(W(K)) - \mathcal{L}(W(0))). \end{aligned}$$

706 Rearranging this, we get

$$\begin{aligned} \|W(K)\|_{p,p}^{p-1} - \langle \nabla\psi(W(K)), \bar{W}_{\text{mm}}^\alpha \rangle &\leq \|W(0)\|_{p,p}^{p-1} - \langle \nabla\psi(W(0)), \bar{W}_{\text{mm}}^\alpha \rangle \\ &\quad + \frac{\pi}{1+\pi} (\|W(K)\|_{p,p}^{p-1} - \|W(0)\|_{p,p}^{p-1}) \\ &\quad + \frac{\eta}{\|W(0)\|_{p,p}} (\mathcal{L}(W(0)) - \mathcal{L}(W(K))). \end{aligned}$$

707 Dividing by  $\|W(K)\|_{p,p}^{p-1}$ , we get

$$\begin{aligned} D_\psi(\bar{W}_{\text{mm}}^\alpha, \bar{W}(K)) &\leq \frac{\|W(0)\|_{p,p}^{p-1}}{\|W(K)\|_{p,p}^{p-1}} D_\psi(\bar{W}_{\text{mm}}^\alpha, \bar{W}(0)) + \frac{\pi}{1+\pi} \left( 1 - \frac{\|W(0)\|_{p,p}^{p-1}}{\|W(K)\|_{p,p}^{p-1}} \right) \\ &\quad + \frac{\eta}{\|W(K)\|_{p,p}^{p-1} \|W(0)\|_{p,p}} (\mathcal{L}(W(0)) - \mathcal{L}(W(K))) \\ &\leq \mu' + \pi + \frac{\eta(\mathcal{L}(W(0)) - \mathcal{L}(W(K)))}{R_\mu^p}. \end{aligned} \quad (27)$$

708 Therefore, we can simply choose  $\mu' = \frac{1}{3}\mu$ ,  $\pi$  be any real number below  $\frac{1}{3}\mu$ , and have  $R_\mu$  big enough  
 709 so that  $\frac{\eta(\mathcal{L}(W(0)) - \mathcal{L}(W(K)))}{R_\mu^p} \leq \frac{1}{3}\mu$  and  $R_\mu \geq R_\pi$ , such  $R_\mu$  exists because  $\mathcal{L}$  is bounded.  $\square$

## 710 B.5 Lemma for Analyzing Rate of Convergence

711 **Lemma 18.** Suppose Assumptions A holds. Let  $R_\delta$  be from Lemma 10, let  $c$  be from Lemma 16, let  $\mu'$   
 712 and  $R_\mu$  be from Lemma 17 when  $\mu = \mu_0$ , and let  $R := \max\{R_\mu, R_\delta, e^{1/c}\}$ . If the initialization  $W(0)$   
 713 is in  $C_{p,\mu',R}(W_{\text{mm}}^\alpha)$ , then for a sufficiently small step size  $\eta$ , the sequence  $\{W(k)\}_{k \geq 0}$  generated by

714  $\ell_p$ -AttGD satisfies

$$D_\psi(\bar{W}_{\text{mm}}^\alpha, \bar{W}(k)) = \begin{cases} \mathcal{O}\left(\frac{\log \|W(k)\|_{p,p}}{\|W(k)\|_{p,p}}\right) & \text{if } p > 2, \\ \mathcal{O}\left(\frac{(\log \|W(k)\|_{p,p})^2}{\|W(k)\|_{p,p}}\right) & \text{if } p = 2, \\ \mathcal{O}\left(\frac{1}{\|W(k)\|_{p,p}^{p-1}}\right) & \text{otherwise.} \end{cases} \quad (28)$$

715 *Proof.* Using Lemma 11, setting  $c$  as the dataset dependent constant hidden by the  $\mathcal{O}$  notation for  
 716  $R_\pi$ , we can get that by setting  $\pi = \min\{\frac{c \log \|W(k)\|_{p,p}}{\delta \|W(k)\|_{p,p}}, 1\}$ , we can use the result of Lemma 16 on  $k$ ,  
 717 so rearranging that result, we get

$$\begin{aligned} \|W(k+1)\|_{p,p}^{p-1} - \langle \nabla \psi(W(k+1)), \bar{W}_{\text{mm}}^\alpha \rangle &\leq \|W(k)\|_{p,p}^{p-1} - \langle \nabla \psi(W(k)), \bar{W}_{\text{mm}}^\alpha \rangle \\ &\quad + \frac{\pi}{1+\pi} (\|W(k+1)\|_{p,p}^{p-1} - \|W(k)\|_{p,p}^{p-1}) \\ &\quad + \frac{\eta}{\|W(k)\|_{p,p}} (\mathcal{L}(W(k)) - \mathcal{L}(W(k+1))). \end{aligned}$$

718 From Lemma 10 and Lemma 12,  $\|W(k)\|_{p,p}$  is increasing, so focusing on the second line, we can  
 719 use Lemma 5 and get

$$\begin{aligned} \frac{\pi}{1+\pi} (\|W(k+1)\|_{p,p}^{p-1} - \|W(k)\|_{p,p}^{p-1}) &\leq \pi (\|W(k+1)\|_{p,p}^{p-1} - \|W(k)\|_{p,p}^{p-1}) \\ &\leq \frac{cp}{\delta \|W(k)\|_{p,p}} \max\{\|W(k)\|_{p,p}^{p-2}, \|W(k+1)\|_{p,p}^{p-2}\} \\ &\quad \times \log \|W(k)\|_{p,p} \\ &\quad \times (\|W(k+1)\|_{p,p} - \|W(k)\|_{p,p}). \end{aligned}$$

720 From Lemma 8, we know that for all index  $k$ ,

$$\|W(k+1)\|_{p,p} \leq \|W(k)\|_{p,p} + \eta L, \quad (29)$$

721 so we can use integral approximation when bounding the sums of  $\Delta(k)$ 's. Let

$$\begin{aligned} \Delta(k) &= \frac{cp}{\delta \|W(k)\|_{p,p}} \max\{\|W(k)\|_{p,p}^{p-2}, \|W(k+1)\|_{p,p}^{p-2}\} \log \|W(k)\|_{p,p} \\ &\quad \times (\|W(k+1)\|_{p,p} - \|W(k)\|_{p,p}), \end{aligned}$$

722 so we can get that

$$\begin{aligned} \|W(K)\|_{p,p}^{p-1} - \langle \nabla \psi(W(K)), \bar{W}_{\text{mm}}^\alpha \rangle &\leq \|W(0)\|_{p,p}^{p-1} - \langle \nabla \psi(W(0)), \bar{W}_{\text{mm}}^\alpha \rangle \\ &\quad + \sum_{k=0}^{K-1} \Delta(k) + \frac{\eta}{c} (\mathcal{L}(W(0)) - \mathcal{L}(W(K))), \end{aligned}$$

723

$$\begin{aligned} \|W(K)\|_{p,p}^{p-1} D_\psi(\bar{W}_{\text{mm}}^\alpha, \bar{W}(K)) &\leq \|W(0)\|_{p,p}^{p-1} D_\psi(\bar{W}_{\text{mm}}^\alpha, \bar{W}(0)) \\ &\quad + \sum_{k=0}^{K-1} \Delta(k) + \frac{\eta}{c} (\mathcal{L}(W(0)) - \mathcal{L}(W(K))). \end{aligned} \quad (30)$$

724 When  $p > 2$ , we have

$$\Delta(k) = \frac{cp}{\delta \|W(k)\|_{p,p}} (\|W(k)\|_{p,p} + \eta L)^{p-2} \log \|W(k)\|_{p,p} (\|W(k+1)\|_{p,p} - \|W(k)\|_{p,p}).$$

725 We can see that

$$\frac{d}{dx} (x + \eta L)^{p-2} (\log x - \log c) > \frac{p-2}{x} (x + \eta L)^{p-2} \log x$$



726 for all  $x > 0$ , so from Equation (29), we can get that

$$\sum_{k=0}^{K-1} \Delta(k) = O(\|W(K)\|^{p-2} \log \|W(K)\|_{p,p}).$$

727 When  $p = 2$ , we have

$$\Delta(k) = \frac{cp}{\|W(k)\|_{p,p}} \log \|W(k)\|_{p,p} (\|W(k+1)\|_{p,p} - \|W(k)\|_{p,p}).$$

728 We can see that

$$\frac{d}{dx} (\log x)^2 > \frac{2}{x} (\log x)$$

729 for all  $x \geq c$ , so from Equation (29), we can get that

$$\sum_{k=0}^{K-1} \Delta(k) = O((\log \|W(K)\|_{p,p})^2).$$

730 When  $p < 2$ , we have

$$\Delta(k) = cp \|W(k)\|_{p,p}^{p-3} \log \|W(k)\|_{p,p} (\|W(k+1)\|_{p,p} - \|W(k)\|_{p,p}).$$

731 From Equation (29), we can get that

$$\sum_{k=0}^{K-1} \Delta(k) = O(1).$$

732 Combining the above cases with Equation (30), we get that

$$\|W(K)\|_{p,p}^{p-1} D_\psi(\bar{W}_{\text{mm}}^\alpha, \bar{W}(K)) = \begin{cases} O(\|W(K)\|_{p,p}^{p-2} \log \|W(K)\|_{p,p}) & \text{if } p > 2, \\ O((\log \|W(K)\|_{p,p})^2) & \text{if } p = 2, \\ O(1) & \text{otherwise} \end{cases},$$

733 Dividing both sides by  $\|W(K)\|_{p,p}^{p-1}$  gives (28). □

734 **Lemma 19.** Suppose Assumptions A holds. Let  $\mu'$  be that from Lemma 17 if  $\mu = \mu_0$ , and let  $R$  the  
 735 maximum of the  $R_\mu$  from 17 and  $R_\delta$  10. Let  $\{W(k)\}_{k \geq 0}$  be the sequence generated by  $\ell_p$ -AttGD.  
 736 If the initialization  $W(0)$  is in  $C_{p,\mu',R}(W_{\text{mm}}^\alpha)$ , then with a small enough step size  $\eta$ , we have the  
 737 following for each  $k \geq 0$ ,

$$\|W(k)\|_{p,p} = \Omega(\log k).$$

738 *Proof.* For each  $k \geq 0$ , Lemma 12 gives

$$\|W(k+1)\|_{p,p}^{p-1} \geq \|W(k)\|_{p,p}^{p-1} + \frac{\eta}{\|W(k)\|_{p,p}} \langle -\nabla \mathcal{L}(W(k)), W(k) \rangle.$$

739 Lemma 17 gives us that  $W(k) \in C_{p,\mu,R}(W_{\text{mm}}^\alpha)$  for each  $k \geq 0$ , so by Lemma 10,

$$\frac{\eta}{\|W(k)\|_{p,p}} \langle -\nabla \mathcal{L}(W(k)), W(k) \rangle = \Omega \left( e^{-\frac{\|W(k)\|_{p,p}}{\|W_{\text{mm}}^\alpha\|_{p,p}} (1 + \frac{1}{2}\delta)} \right),$$

740 so there exists dataset dependent constants  $R_1, R_2 > 0$  such that

$$\frac{\eta}{\|W(k)\|_{p,p}} \langle -\nabla \mathcal{L}(W(k)), W(k) \rangle \geq R_1 e^{-R_2 \|W(k)\|_{p,p}},$$

741 so for each  $k \geq 0$ ,

$$\|W(k+1)\|_{p,p}^{p-1} \geq \|W(k)\|_{p,p}^{p-1} + R_1 e^{-R_2 \|W(k)\|_{p,p}}.$$

742 Set  $k_0 = 0$ , and let  $k_{i+1}$  be the lowest indices such that  $\|W(k_{i+1})\|_{p,p} \geq \|W(k_i)\|_{p,p} + 1$  for all  
 743 index  $i \geq 0$ . Therefore,

$$k_{i+1} - k_i \leq \frac{(\|W(k_i)\|_{p,p} + 1)^{p-1} - \|W(k_i)\|_{p,p}^{p-1}}{R_1 e^{-R_2 (\|W(k_i)\|_{p,p} + 1)}} = e^{O(\|W(k_i)\|_{p,p})}.$$

744 Therefore,

$$\|W(k)\|_{p,p} = \Omega(\log k).$$

745 □

746 **C Proof of Theorem 1**

747 *Proof.* The proof is similar to the proof of [53, Theorem 1]. Specifically, we need to show that  
 748  $f(X) = v^\top X^\top \sigma(XW)$  satisfies the assumptions of [53, Lemma 14], where the nonlinear head is  
 749 replaced by the linear term  $v$ . This holds independently of the choice of algorithm or the attention  
 750 SVM solution. Thus, we omit the details and refer to the proof of [53, Theorem 1].  $\square$

751 **D Proof of Theorem 2**

752 *Proof.* It is enough to show the existence of such constants  $\mu, R > 0$  such that if  $W(0)$  is in  
 753  $C_{p,\mu,R}(W_{\text{mm}}^\alpha)$ , then the norm diverges to infinity. As discussed in Lemma 12, for any timestep  $k$ ,

$$\|W(k+1)\|_p^{p-1} \geq \|W(k)\|_p^{p-1} - \frac{\eta}{\|W(k)\|_p} \langle \nabla \mathcal{L}(W(k)), W(k) \rangle. \quad (31)$$

754 Let  $R_1$  be the  $R$  from Lemma 10, set  $\mu$  and  $R_2$  to be the  $\mu'$  and  $R$  for  $\mu = \mu_0$  of Lemma 17, and set  
 755  $R := \max\{R_1, R_2\}$ . From Lemma 17, we know that  $W(k) \in C_{p,\mu_0,R}(W_{\text{mm}}^\alpha)$  for any timestep  $k$ ,  
 756 so from Lemma 10,

$$\langle \nabla \mathcal{L}(W(k)), W(k) \rangle < 0,$$

757 for all timesteps  $k$ .

758 Therefore, the  $l_p$ -norm is always increasing, but this does not immediately imply that the  $l_p$ -norm  
 759 will approach infinity; it could converge to a finite value. However, if  $\|W(k)\|_p$  converges to a finite  
 760 value, then again by Lemma 10, we get a lower bound for  $-\frac{\eta}{\|W(k)\|_p} \langle \nabla \mathcal{L}(W(k)), W(k) \rangle$  at any  
 761 timestep  $k$ . Therefore, by Equation (31),

$$\lim_{k \rightarrow \infty} \|W(k)\|_p^{p-1} = \infty,$$

762 a contradiction, so  $\|W(k)\|_p$  converges to infinity.  $\square$

763 **E Proof of Theorem 3**

764 *Proof.* This is a direct consequence of Theorem 4.  $\square$

765 **F Proof of Theorem 4**

766 *Proof.* Let  $R$  be the one from Lemma 18. Given  $W(0) \in C_{p,\mu,R}(W_{\text{mm}}^\alpha)$ , by Lemma 18, we have

$$D_\psi(\bar{W}_{\text{mm}}^\alpha, \bar{W}(k)) = \begin{cases} \mathcal{O}\left(\frac{\log \|W(k)\|_{p,p}}{\|W(k)\|_{p,p}}\right) & \text{if } p > 2, \\ \mathcal{O}\left(\frac{(\log \|W(k)\|_{p,p})^2}{\|W(k)\|_{p,p}}\right) & \text{if } p = 2, \\ \mathcal{O}\left(\frac{1}{\|W(k)\|_{p,p}^{p-1}}\right) & \text{otherwise.} \end{cases}$$

767 From Lemma 19, we know that

$$\|W(k)\|_{p,p} = \Omega(\log k).$$

768 The derivative  $\frac{d}{dx} \left(\frac{\log x}{x}\right) = \frac{1-\log x}{x^2}$  is negative when  $x > e$ , so  $\frac{\log x}{x}$  is decreasing when  $x > e$ .

769 Similarly,  $\frac{(\log x)^2}{x}$  is decreasing when  $x > e^2$ .

770 Thus when  $p > 2$ , for a large enough  $k$ ,

$$D_\psi(\bar{W}_{\text{mm}}^\alpha, \bar{W}(k)) = O\left(\frac{\log \log k}{\log k}\right). \quad (32a)$$

771 Similarly, when  $p = 2$ , for a large enough  $k$ ,

$$D_\psi(\bar{W}_{\text{mm}}^\alpha, \bar{W}(k)) = O\left(\frac{(\log \log k)^2}{\log k}\right). \quad (32b)$$

772 Finally, when  $1 < p < 2$ ,

$$D_\psi(\bar{W}_{\text{mm}}^\alpha, \bar{W}(k)) = O\left(\frac{1}{(\log k)^{p-1}}\right). \quad (32c)$$

773

□

## 774 G On the Convergence of the $\ell_p$ Regularization Path for Joint $W$ and $v$

775 **Assumption B.** Let  $\Gamma, \Gamma' > 0$  denote the label margins when solving ( $\ell_p$ -SVM) with  $X_{i\alpha_i}$  and its  
776 replacement with  $X_i^\top \sigma(X_i W z_i)$ , for all  $i \in [n]$ , respectively. There exists  $\nu > 0$  such that for all  
777  $i \in [n]$  and  $W \in \mathbb{R}^{d \times d}$ ,

$$\Gamma - \Gamma' \geq \nu \cdot (1 - s_{i\alpha_i}), \quad \text{where } s_{i\alpha_i} = [\sigma(X_i W z_i)]_{\alpha_i}.$$

778 Assumption B is similar to [53] and highlights that selecting optimal tokens is key to maximizing  
779 the classifier's label margin. When attention features, a weighted combination of all tokens, are  
780 used, the label margin shrinks based on how much attention is given to the optimal tokens. The term  
781  $\nu \cdot (1 - s_{i\alpha_i})$  quantifies this minimum shrinkage. If the attention mechanism fails to focus on these  
782 tokens (i.e., low  $s_{i\alpha_i}$ ), the margin decreases, reducing generalization. This assumption implies that  
783 optimal performance is achieved when attention converges on the most important tokens, aligning  
784 with the max-margin attention SVM solution.

785 Similar to how we provided the characterization of convergence for the regularization path of  
786  $\ell_p$ -AttGD, we offer a similar characterization here for  $\ell_p$ -JointGD.

787 **Theorem 5** (Joint  $\ell_p$ -norm Regularization Path). Consider (ERM) with a logistic loss  $l(x) =$   
788  $\log(1 + e^{-x})$ , and define

$$(v^{(r)}, W^{(R)}) := \arg \min_{(v, W)} \mathcal{L}(v, W) \quad \text{subj. to } \|W\|_{p,p} \leq R \text{ and } \|v\|_p \leq r. \quad (\ell_p\text{-JointRP})$$

789 Suppose there are token indices  $\alpha = (\alpha_i)_{i=1}^m$  for which  $W_{\text{mm}}^\alpha$  of ( $\ell_p$ -AttSVM) exists and Assump-  
790 tion B holds for some  $\Gamma, \nu > 0$ . Then,

$$\lim_{(r,R) \rightarrow (\infty, \infty)} \left( \frac{v^{(r)}}{r}, \frac{W^{(R)}}{R} \right) = \left( \frac{v_{\text{mm}}}{\|v_{\text{mm}}\|_p}, \frac{W_{\text{mm}}^\alpha}{\|W_{\text{mm}}^\alpha\|_{p,p}} \right). \quad (33)$$

791 Here,  $v_{\text{mm}}$  and  $W_{\text{mm}}^\alpha$  are the solution of ( $\ell_p$ -SVM) and ( $\ell_p$ -AttSVM), respectively.

792 Theorem 5 extends the results of Theorem 1 to the case of joint optimization of head  $v$  and attention  
793 weights  $W$  using a logistic loss function.

## 794 H Proof of Theorem 5

795 *Proof.* The proof is similar to the proof of [53, Theorem 5]. We provide the revised version for the  
796 generalized attention SVM, tracking the required changes. Without loss of generality, we set  $\alpha_i = 1$   
797 for all  $i \in [n]$ , and we use  $W_{\text{mm}}$  instead of  $W_{\text{mm}}^\alpha$ . Suppose the claim is incorrect, meaning either  
798  $W^{(R)}/R$  or  $v^{(r)}/r$  fails to converge as  $R$  and  $r$  grow. Define

$$\Xi = \frac{1}{\|\bar{W}_{\text{mm}}\|_{p,p}}, \quad \Gamma = \frac{1}{\|v_{\text{mm}}\|_p},$$

$$\bar{W}_{\text{mm}} := R\Xi W_{\text{mm}}, \quad \bar{v}_{\text{mm}} := r\Gamma v_{\text{mm}} \quad (34)$$

799 Our strategy is to show that  $(\bar{v}_{\text{mm}}, \bar{W}_{\text{mm}})$  is a strictly better solution compared to  $(v^{(r)}, W^{(R)})$  for  
800 large  $R$  and  $r$ , leading to a contradiction.

801 • **Case 1:**  $W^{(R)}/R$  does not converge to  $\bar{W}_{\text{mm}}/R$ . In this case, there exists  $\delta, \gamma = \gamma(\delta) > 0$  such  
 802 that we can find arbitrarily large  $R$  with

$$\|W^{(R)}/R - \bar{W}_{\text{mm}}/R\| \geq \delta$$

803 and the margin induced by  $W^{(R)}/R$  is at most  $\Xi(1 - \gamma)$ .

804 We bound the amount of non-optimality  $q_i^*$  of  $\bar{W}_{\text{mm}}$ :

$$\begin{aligned} q_i^* &:= \frac{\sum_{t \neq \alpha_i} \exp(X_{it}^\top \bar{W}_{\text{mm}} z_i)}{\sum_{t \in [T]} \exp(X_{it}^\top \bar{W}_{\text{mm}} z_i)} \leq \frac{\sum_{t \neq \alpha_i} \exp(X_{it}^\top \bar{W}_{\text{mm}} z_i)}{\exp(X_{i\alpha_i}^\top \bar{W}_{\text{mm}} z_i)} \\ &\leq T \exp(-\Xi R). \end{aligned}$$

805 Thus,

$$q_{\max}^* := \max_{i \in [n]} q_i^* \leq T \exp(-\Xi R). \quad (35a)$$

806 Next, assume without loss of generality that the first margin constraint is  $\gamma$ -violated by  $W^{(R)}$ ,  
 807 meaning

$$\min_{t \neq \alpha_1} (X_{1\alpha_1} - X_{1t})^\top W^{(R)} z_1 \leq \Xi R(1 - \gamma).$$

808 Denoting the amount of non-optimality of the first input of  $W^{(R)}$  as  $\hat{q}_1$ , we find

$$\begin{aligned} \hat{q}_1 &:= \frac{\sum_{t \neq \alpha_1} \exp(X_{1t}^\top W^{(R)} z_1)}{\sum_{t \in [T]} \exp(X_{1t}^\top W^{(R)} z_1)} \geq \frac{1}{T} \frac{\sum_{t \neq \alpha_1} \exp(X_{1t}^\top W^{(R)} z_1)}{\exp(X_{1\alpha_1}^\top W^{(R)} z_1)} \\ &\geq T^{-1} \exp(-(1 - \gamma)R\Xi). \end{aligned}$$

809 This implies that

$$\hat{q}_{\max} := \max_{i \in [n]} \hat{q}_i^* \geq T^{-1} \exp(-\Xi R(1 - \gamma)). \quad (35b)$$

810 We similarly have

$$q_{\max}^* \geq T^{-1} \exp(-\Xi R). \quad (35c)$$

811 Thus, (35) gives the following relationship between the upper and lower bounds on non-optimality:

$$\begin{aligned} -(1 - \gamma)\Xi R - \log T &\leq \log(\hat{q}_{\max}), \\ -\Xi R - \log T &\leq \log(q_{\max}^*) \leq -\Xi R + \log T. \end{aligned} \quad (36)$$

812 In other words,  $\bar{W}_{\text{mm}}$  has exponentially less non-optimality compared to  $W^{(R)}$  as  $R$  grows. To  
 813 proceed, we need to upper and lower bound the logistic loss of  $(\bar{v}_{\text{mm}}, \bar{W}_{\text{mm}})$  and  $(v^{(r)}, W^{(R)})$   
 814 respectively, to establish a contradiction.

815 • **Sub-Case 1.1: Upper bound for  $\mathcal{L}(\bar{v}_{\text{mm}}, \bar{W}_{\text{mm}})$ .** We now bound the logistic loss for the limiting  
 816 solution. Set  $\bar{r}_i = X_i^\top \sigma(X_i \bar{W}_{\text{mm}} z_i)$ . If  $\|\bar{r}_i - X_{i1}\|_p \leq \epsilon_i$ , then  $v_{\text{mm}}$  satisfies the SVM constraints  
 817 on  $\bar{r}_i$  with  $Y_i \cdot \bar{r}_i^\top v_{\text{mm}} \geq 1 - \epsilon_i/\Gamma$ . Setting  $\epsilon_{\max} = \sup_{i \in [n]} \epsilon_i$ ,  $v_{\text{mm}}$  achieves a label-margin of  
 818  $\Gamma - \epsilon_{\max}$  on the dataset  $(Y_i, \bar{r}_i)_{i \in [n]}$ . Let  $M = \sup_{i \in [n], t, \tau \in [T]} \|X_{it} - X_{i\tau}\|_p$ . Recalling (36), the  
 819 worst-case perturbation is

$$\epsilon_{\max} \leq M \exp(-\Xi R + \log T) = MT \exp(-\Xi R).$$

820 This implies the upper bound on the logistic loss:

$$\begin{aligned} \mathcal{L}(\bar{v}_{\text{mm}}, \bar{W}_{\text{mm}}) &\leq \max_{i \in [n]} \log(1 + \exp(-Y_i \bar{r}_i^\top \bar{v}_{\text{mm}})) \\ &\leq \max_{i \in [n]} \exp(-Y_i \bar{r}_i^\top \bar{v}_{\text{mm}}) \\ &\leq \exp(-r\Gamma + r\epsilon_{\max}) \\ &\leq e^{rMT \exp(-\Xi R)} e^{-r\Gamma}. \end{aligned} \quad (37)$$

821 • **Sub-Case 1.2: Lower bound for  $\mathcal{L}(v^{(r)}, W^{(R)})$ .** We now bound the logistic loss for the finite  
 822 solution. Set  $\bar{r}_i = X_i^\top \sigma(X_i W^{(R)} z_i)$ . Using Assumption B, solving ( $\ell_p$ -SVM) on  $(y_i, \bar{r}_i)_{i \in [n]}$   
 823 achieves at most  $\Gamma - \nu e^{-(1-\gamma)\Xi R}/T$  margin. Consequently, we have:

$$\begin{aligned}
\mathcal{L}(v^{(r)}, W^{(R)}) &\geq \frac{1}{n} \max_{i \in [n]} \log(1 + \exp(-Y_i \bar{r}_i^\top v^{(r)})) \\
&\geq \left( \frac{1}{2n} \max_{i \in [n]} \exp(-Y_i \bar{r}_i^\top v^{(r)}) \right) \wedge \log 2 \\
&\geq \left( \frac{1}{2n} \exp(-r(\Gamma - \nu e^{-(1-\gamma)\Xi R}/T)) \right) \wedge \log 2 \\
&\geq \left( \frac{1}{2n} e^{r(\nu/T) \exp(-(1-\gamma)\Xi R)} e^{-r\Gamma} \right) \wedge \log 2.
\end{aligned}$$

824 Observe that this lower bound dominates the upper bound from (37) when  $R$  is large, specifically  
825 when (ignoring the multiplier  $1/2n$  for simplicity):

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

826 Thus, we obtain the desired contradiction since such a large  $R$  is guaranteed to exist when  $W^{(R)}/R \not\rightarrow$   
827  $\bar{W}_{\text{mm}}$ . Therefore,  $W^{(R)}/R$  must converge to  $\bar{W}_{\text{mm}}/R$ .

828 • **Case 2: Suppose  $v^{(r)}/r$  does not converge.** In this case, there exists  $\delta > 0$  such that we  
829 can find arbitrarily large  $r$  obeying  $\text{dist}(v^{(r)}/r, \bar{v}_{\text{mm}}/r) \geq \delta$ . If  $\text{dist}(W^{(R)}/R, \Xi W_{\text{mm}}) \not\rightarrow 0$ ,  
830 then "Case 1" applies. Otherwise, we have  $\text{dist}(W^{(R)}/R, \Xi W_{\text{mm}}) \rightarrow 0$ , thus we can assume  
831  $\text{dist}(W^{(R)}/R, \Xi W_{\text{mm}}) \leq \epsilon$  for an arbitrary choice of  $\epsilon > 0$ .

832 On the other hand, due to the strong convexity of ( $\ell_p$ -AttSVM), for some  $\gamma := \gamma(\delta) > 0$ ,  $v^{(r)}$   
833 achieves a margin of at most  $(1 - \gamma)\Gamma r$  on the dataset  $(Y_i, X_{i1})_{i \in [n]}$ , where  $X_{i1}$  denotes the optimal  
834 token for each  $i \in [n]$ . Additionally, since  $\text{dist}(W^{(R)}/R, \Xi W_{\text{mm}}) \leq \epsilon$ ,  $W^{(R)}$  strictly separates  
835 all optimal tokens (for small enough  $\epsilon > 0$ ) and  $\hat{q}_{\text{max}} := f(\epsilon) \rightarrow 0$  as  $R \rightarrow \infty$ . Note that  $f(\epsilon)$   
836 quantifies the non-optimality of  $W^{(R)}$  compared to  $W_{\text{mm}}$ ; as  $\epsilon \rightarrow 0$ , meaning  $W^{(R)}/R$  converges  
837 to  $\Xi W_{\text{mm}}/R$ ,  $f(\epsilon) \rightarrow 0$ . Consequently, setting  $r_i = X_i^\top \sigma(X_i W^{(R)} z_i)$ , for sufficiently large  $R > 0$   
838 and setting  $M = \sup_{i \in [n], t \in [T]} \|X_{it}\|$ , we have that

$$\begin{aligned}
\min_{i \in [n]} y_i (v^{(r)})^\top r_i &\leq \min_{i \in [n]} y_i (v^{(r)})^\top X_{i1} + \sup_{i \in [n]} |(v^{(r)})^\top (X_{it} - X_{i1})| \\
&\leq (1 - \gamma)\Gamma r + M f(\epsilon) r \\
&\leq (1 - \gamma/2)\Gamma r.
\end{aligned} \tag{38}$$

839 This in turn implies that logistic loss is lower bounded by

$$\mathcal{L}(v^{(r)}, W^{(R)}) \geq \left( \frac{1}{2n} e^{\gamma\Gamma r/2} e^{-\Gamma r} \right) \wedge \log 2.$$

840 Going back to (37), this exponentially dominates the upper bound of  $(\bar{W}_{\text{mm}}, \bar{v}_{\text{mm}})$  whenever  
841  $rMT \exp(-\Xi R) < r\gamma\Gamma/2$  (that is, whenever  $R, r$  are sufficiently large), again concluding the  
842 proof.

843 □

## 844 I Implementation Details

845 The experiments were run on an Intel i7 core and a single V100 GPU using the pytorch and  
846 huggingface libraries. They should be runnable on any generic laptop.

### 847 I.1 Illustrating Optimal Tokens

848 **Example 1.** Consider the matrices  $X_1 = [5, 0; 0, 1]$  and  $X_2 = [-5, 0; 0, -1]$  with  $y_1 = -y_2 = 1$ .  
849 Let  $x_{i1}$  be the optimal token and  $x_{it}$  be the others. Problem ( $\ell_p$ -AttSVM) with  $p = 3$  and  $z_i = X_{i1}$   
850 yields  $W_{\text{mm}}^\alpha = W_{\text{mm}} = [0.03846, 0; -0.00769, 0]$ . Figure 3 illustrates how the optimal tokens  
851  $X_{11}$  and  $X_{21}$  are separated by the dashed lines (orthogonal to  $W_{\text{mm}} z_i$ ) for each sequence.

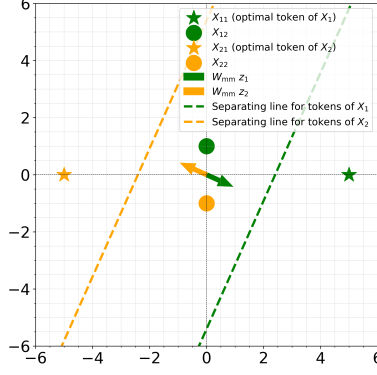


Figure 3: Visualization of Problem ( $\ell_p$ -AttSVM) with  $p = 3$ .

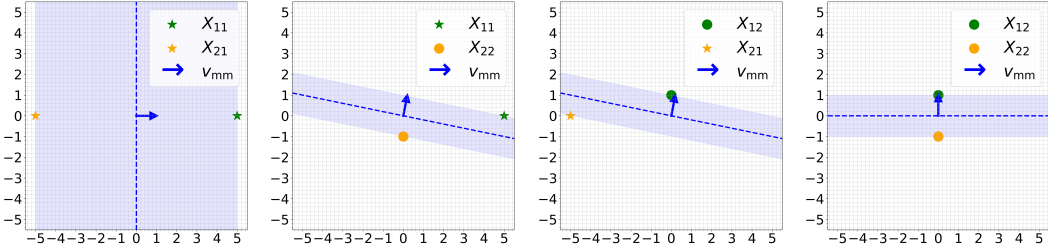


Figure 4: Visualizing the effect of token selection on margin size in ( $\ell_p$ -SVM) for Example 1. The first plot illustrates the largest class margin, indicating the optimality of tokens  $X_{11}$  and  $X_{21}$ . In subsequent plots, as different tokens are used, the class margin (light blue shaded area) decreases, reflecting suboptimal class separation.

## 852 I.2 Synthetic Data Experiment

853 We describe the setup of the experiments for  $\ell_p$ -AttGD and  $\ell_p$ -JointGD and their results.

854  **$\ell_p$ -AttGD Experiment.** To measure the directional distance between  $W_\alpha^{\text{mm}}$  ( $\ell_p$ -AttSVM) solution)  
 855 and  $W(k)$  (output of  $\ell_p$ -AttGD), we use a directional Bregman divergence, defined for  $W, V \in \mathbb{R}^{d \times d}$   
 856 as  $D_\psi(W/\|W\|_{p,p}, V/\|V\|_{p,p})$ . We compare the ( $\ell_p$ -AttSVM) solution with the  $\ell_q$  optimization  
 857 path for all  $p, q \in \{1.75, 2, 3\}$  for synthetically generated data. The experiment is repeated 100 times,  
 858 and the average directional Bregman divergence is reported. A closer look at one sample trial is also  
 859 provided.

860 The dataset  $(X_i, Y_i, z_i)_{i=1}^n$  used for the experiment is generated randomly:  $X_i$  and  $z_i$  are sampled  
 861 from the unit sphere, and  $Y_i$  is uniformly sampled from  $\{\pm 1\}$ . Additionally,  $v$  is randomly selected  
 862 from the unit sphere. We use  $n = 6$  samples,  $T = 8$  tokens per sample, and  $d = 10$  dimensions per  
 863 token, fulfilling the overparameterized condition for the  $\ell_p$ -AttSVM problem to be almost always  
 864 feasible.

865 The model parameter is initialized near the origin, and it is trained using Algorithms  $\ell_p$ -AttGD  
 866 with  $p = 1.75, 2$ , and  $3$ , and a learning rate of  $0.1$ . Training lasted for  $1,500$  epochs for  $p = 1.75$ ,  
 867  $2,000$  epochs for  $p = 2$ , and  $20,000$  epochs for  $p = 3$ . Gradients are normalized to accelerate  
 868 convergence without altering results significantly. We refer to the parameter histories as the  $\ell_{1.75}, \ell_2$ ,  
 869 and  $\ell_3$  optimization paths. We compute the chosen tokens  $(\alpha_i)_{i=1}^n$  for the ( $\ell_p$ -AttSVM) problem by  
 870 selecting the token with the highest softmax probability for each sample. This process is repeated for  
 871  $p = 1.75, 2$ , and  $3$ .

872 Figure 5 shows the directional Bregman divergence between the ( $\ell_p$ -AttSVM) solution and the  $\ell_q$   
 873 optimization path for each pair  $p, q \in \{1.75, 2, 3\}$ . In Figure 5a, the divergence converges to  $0$   
 874 only for the ( $\ell_p$ -AttSVM) ( $p = 1.75$ ) solution, indicating that the  $\ell_{1.75}$  path does not converge to  
 875 the  $p = 2$  or  $3$  solutions. The shrinking standard deviation shows consistent behavior. Similarly,  
 876 Figures 5b and 5c show the divergence converging to  $0$  for the corresponding ( $\ell_p$ -AttSVM) solution,  
 877 demonstrating that the  $\ell_p$  optimization path converges to the ( $\ell_p$ -AttSVM) solution, with the direction  
 878 of convergence changing with  $p$ .

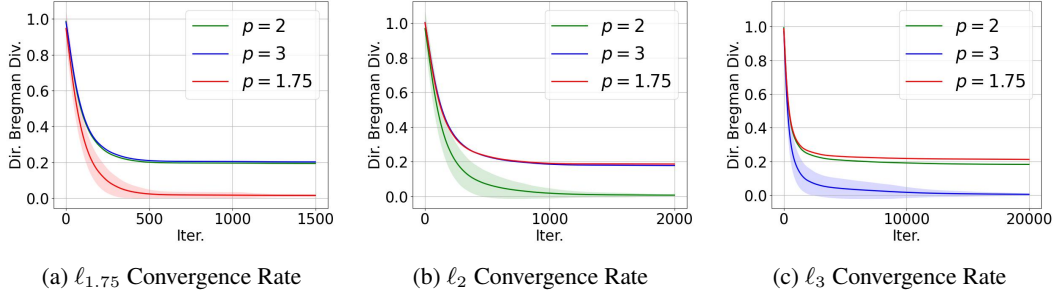


Figure 5: Average directional Bregman divergence between the (a)  $\ell_{1.75}$ , (b)  $\ell_2$ , and (c)  $\ell_3$  optimization paths and the ( $\ell_p$ -AttSVM) solutions for  $p = 1.75, 2$ , and  $3$  at each training iteration from 100 trials. The shaded area represents the standard deviation of the directional Bregman divergence.

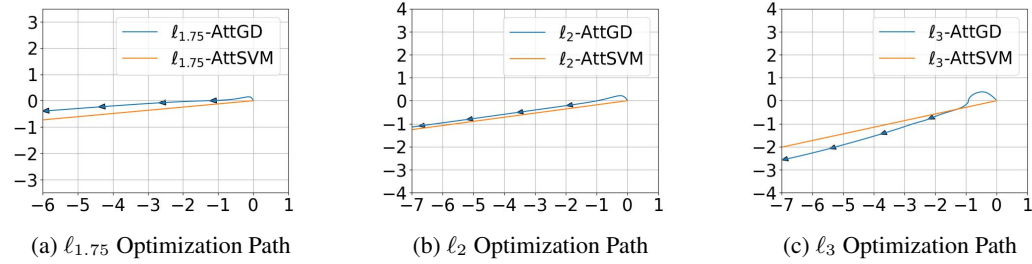


Figure 6: Direction of change of two entries of  $W$  updated by  $\ell_p$ -AttGD with  $p = 1.75, p = 2$ , and  $p = 3$  for one trial, shown in (a), (b), and (c). Each axis represents a different entry. The orange line shows the direction of ( $\ell_p$ -AttSVM).

879 Using this same synthetic data, we can also observe the convergence in direction for one of the  
880 trials directly by plotting how two of the entries of  $W$  change during training simultaneously and  
881 plotting it on a Cartesian graph, then showing that the path it follows converges to the direction of  
882 the ( $\ell_p$ -AttSVM) solution. As we can see in Figure 6, each of the  $\ell_p$  optimization paths follows the  
883 direction of the corresponding ( $\ell_p$ -AttSVM) solution.

884  **$\ell_p$ -JointGD Experiment.** We use the data from the following to train a model using  $\ell_p$ -JointGD  
885 for  $p = 1.75, 2$ , and  $3$ .

886 **Example 2.** Let  $n = 2, T = 3, d = 2$ . Let  $y_1 = 1, y_2 = -1$ . Let:

$$X_1 = \begin{pmatrix} X_{11} \\ X_{12} \\ X_{13} \end{pmatrix} = \begin{pmatrix} -5.4 & 2.4 \\ 2.8 & 4.2 \\ 2.6 & -0.2 \end{pmatrix}, \quad \text{and} \quad X_2 = \begin{pmatrix} X_{21} \\ X_{22} \\ X_{23} \end{pmatrix} = \begin{pmatrix} 0.8 & -4.4 \\ -2.2 & -0.8 \\ 1.8 & 0.2 \end{pmatrix}. \quad (39)$$

887 Let  $z_1 = X_{11}, z_2 = X_{21}$ .

888 We use learning rates 0.1 and we trained the model for 1,500 epochs for when  $p = 1.75$ , 2,000  
889 epochs for  $p = 2$ , and 20,000 epochs for  $p = 3$ . As it was done in the previous experiment, we  
890 obtain the parameter history for each  $p$ , and compute the optimal token for the ( $\ell_p$ -AttSVM) and  
891  $\ell_p$ -SVM problems.

892 The comparison between the iterates and the SVM solutions in Figure 7 shows that the iterates of  $W$   
893 and  $v$  converge to the  $\ell_p$ -AttSVM and  $\ell_p$ -SVM directions, respectively, for each of  $p = 1.75, 2$ , and  
894  $3$ . These convergence are similar to Theorem 5, as in both this experiment and that theorem, we get  
895 that the iterates converge to the SVM problem solutions. In addition to these iterates, we record  
896 the evolution of the average softmax probability of the optimal token, along with the average logistic  
897 probability of the model, which we define to be  $1/n \sum_{i=1}^n 1/(1 + e^{-\gamma i \alpha_i})$ .

898 As we can see in Figure 8, each of the average softmax probability converges to 1, indicating that the  
899 attention mechanism produces a softmax probability vector that converges to a one-hot vector for the  
900 different  $\ell_p$ -JointGD training. Furthermore, the average logistic probability also converges to 1,  
901 indicating that the model's prediction converges to a 100% accuracy.



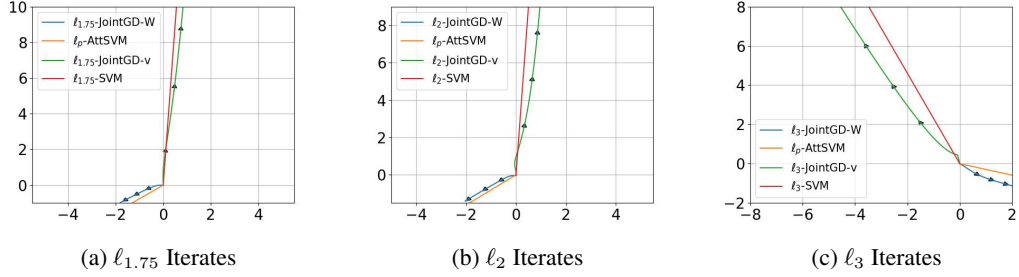


Figure 7: Iterates of the  $W$  and  $v$  parameters of the model as they are trained using  $\ell_p$ -JointGD for  $p = 1.75, 2$ , and  $3$ , along with the corresponding  $\ell_p$ -AttSVM and  $\ell_p$ -SVM directions.

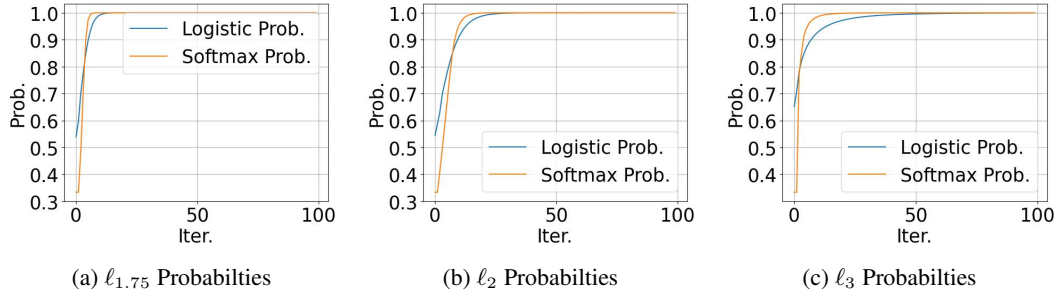


Figure 8: Softmax probability evolution of the optimal token and logistic probability evolution for  $p = 1.75, 2$ , and  $3$ .

### 902 I.3 Additional Real Experiments

903 We collect the training weights from the resulting models trained by  $\ell_{1.1}$  mirror descent and the  
 904 gradient descent and plot a histogram of their absolute values in Figure 9. Specifically, we take the  
 905 histogram of the weights responsible for determining the softmax within the model and the value  
 906 matrices. The figures shows us that the resulting model that was trained using  $\ell_{1.1}$  mirror descent is  
 907 sparser than the one trained using gradient descent, which could hint at a potential explanation as to  
 908 why  $\ell_{1.1}$  mirror descent can outperform the standard gradient descent algorithm when it is used to  
 909 train attention-based models.

910 Finally, the following figures show the full attention maps.

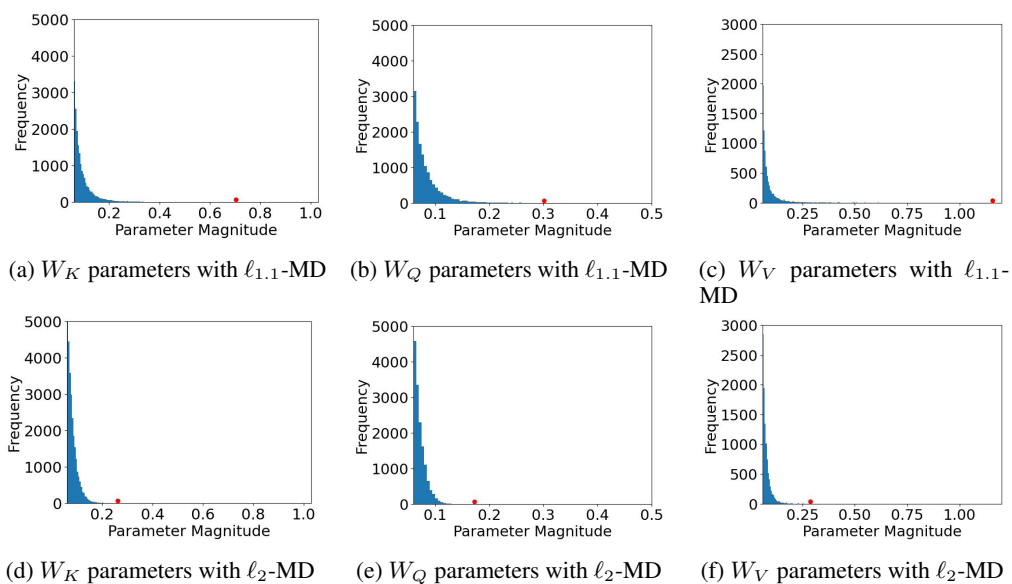


Figure 9: Histogram of the absolute values of the  $W_K$ ,  $W_Q$ , and  $W_V$  components of transformer models trained with  $\ell_{1.1}$  and  $\ell_2$ -MD on the **Stanford Large Movie Review Dataset**. Only large parameters ( $\geq 0.06$ ) are shown, with the maximum magnitude component marked by a red dot. The  $\ell_{1.1}$ -MD model has 18,206 components in  $W_K$ , 13,964 in  $W_Q$ , and 7,643 in  $W_V$  with magnitudes  $\geq 0.06$ , while the  $\ell_2$ -MD model has 27,224 in  $W_K$ , 14,654 in  $W_Q$ , and 10,127 in  $W_V$  with such magnitudes. These results imply that the  $\ell_{1.1}$ -MD algorithm yields sparser parameters and that it would have a stronger token selection ability.

Label	Optimal Token	$\ell_{1,1}$ -MD Token Selection	GD Token Selection	Better Selector
+	fantastic	the movie was fantastic	the movie was fantastic	1.1
-	hated	i hated the movie	i hated the movie	1.1
-	boring	the plot was boring	the plot was boring	2
+	love	i love this movie	i love this movie	2
-	terrible	the plot was terrible	the plot was terrible	1.1
+	great	this movie is great	this movie is great	1.1
-	dirty	the scenes were dirty	the scenes were dirty	2
+	satisfied	i m satisfied with movie	i ' m satisfied with movie	2
-	late	the dvd arrived late	the dvd arrived late	1.1
+	perfectly	the sub ##titles work perfectly	the sub ##titles work perfectly	1.1
-	disappointing	the movie was disappointing	the movie was disappointing	1.1
-	unreliable	the pacing is unreliable	the pacing is unreliable	1.1
+	friendly	the cast were friendly	the cast were friendly	2
-	slow	the script is slow	the script is slow	1.1
+	great	the movie was great	the movie was great	1.1
-	poor	the dvd was poor	the dvd was poor	1.1
+	fascinating	the plot was fascinating	the plot was fascinating	1.1
+	sturdy	the set was sturdy	the set was sturdy	2
-	ruined	the cinematography was ruined	the cinematography was ruined	1.1
+	engaging	the documentary was engaging	the documentary was engaging	1.1
-	crashes	the dvd crashes often	the dvd crashes often	1.1
+	delicious	the scenes were delicious	the scenes were delicious	1.1
-	broke	the dvd broke down	the dvd broke down	2
+	prompt	the service was prompt	the service was prompt	2
-	predictable	the plot was predictable	the plot was predictable	1.1
+	excellent	the service was excellent	the service was excellent	2
+	scenic	the theater is scenic	the theater is scenic	2
-	stopped	the project ##or stopped	the project ##or stopped	1.1
+	vibrant	the festival was vibrant	the festival was vibrant	1.1
+	fun	the movie was fun	the movie was fun	1.1
-	delayed	the screening was delayed	the screening was delayed	1.1

+	fun	the movie was <b>fun</b>	the movie was fun	1.1
-	delayed	the <b>screening</b> was delayed	the screening was delayed	1.1
+	pleasant	the impact was pleasant	the impact <b>was</b> pleasant	2
-	unstable	<b>the</b> streaming is unstable	the streaming <b>is</b> unstable	=
+	fresh	the snacks <b>are</b> fresh	the snacks <b>are</b> fresh	2
-	cracked	the <b>dvd</b> cracked	the <b>dvd</b> cracked	2
+	selection	<b>the</b> theater has selection	the theater <b>has</b> selection	=
-	difficult	the interface is <b>difficult</b>	the interface <b>is</b> difficult	1.1
+	spacious	<b>the</b> cinema is spacious	the cinema <b>is</b> spacious	2
-	broke	<b>the</b> equipment broke	the equipment <b>broke</b>	2
+	friendly	the staff <b>are</b> friendly	the staff are <b>friendly</b>	2

Figure 10: The attention map generated by the resulting models that were trained using  $\ell_{1,1}$  mirror descent and gradient descent for five sample sentences. For three out of five of the sample sentences, the model trained using  $\ell_{1,1}$  mirror descent selects the optimal token better than the model trained using gradient descent.