

---

# On the Role of Attention in Prompt-tuning

---

Samet Oymak<sup>\*1</sup> Ankit Singh Rawat<sup>\*2</sup> Mahdi Soltanolkotabi<sup>\*3</sup> Christos Thrampoulidis<sup>\*4</sup>

## Abstract

*Prompt-tuning* is an emerging strategy to adapt large language models (LLM) to downstream tasks by learning a (soft-)prompt parameter from data. Despite its success in LLMs, there is limited theoretical understanding of the power of prompt-tuning and the role of the attention mechanism in prompting. In this work, we explore prompt-tuning for one-layer attention architectures and study contextual mixture-models where each input token belongs to a context-relevant or -irrelevant set. We isolate the role of prompt-tuning through a self-contained *prompt-attention* model. Our contributions are as follows: (1) We show that softmax-prompt-attention is provably more expressive than softmax-self-attention and linear-prompt-attention under our contextual data model. (2) We analyze the initial trajectory of gradient descent and show that it learns the prompt and prediction head with near-optimal sample complexity and demonstrate how the prompt can provably attend to sparse context-relevant tokens. (3) Assuming a known prompt but an unknown prediction head, we characterize the exact finite sample performance of prompt-attention which reveals the fundamental performance limits and the precise benefit of the context information. We also provide experiments that verify our theoretical insights on real datasets and demonstrate how prompt-tuning enables the model to attend to context-relevant information.

## 1. Introduction

Transformer models have achieved remarkable success in a wide array of machine learning domains spanning language

<sup>\*</sup>In alphabetical order <sup>1</sup>University of Michigan & UC Riverside, USA <sup>2</sup>Google Research NYC, USA <sup>3</sup>University of Southern California, USA <sup>4</sup>University of British Columbia, Canada. Correspondence to: Christos Thrampoulidis <cthrampo@ece.ubc.ca>.

modeling, vision, and decision making. Recently, one of the key techniques that has helped pave the way for the deployment of transformers to ever increasing application areas is their ability to adapt to multiple unseen tasks by conditioning their predictions through their inputs – a technique known as prompt-tuning (Lester et al., 2021; Li & Liang, 2021). Concretely, prompt-tuning provides a more efficient (cheaper/faster) alternative to fine-tuning the entire weights of the transformer by instead training (fewer) so-called prompt parameters that are appended on the input and can be thought of as an input interface. In fact, several recent works have demonstrated experimentally that prompt-tuning is not only more efficient, but often even becomes competitive to fine-tuning in terms of accuracy (Lester et al., 2021; Liu et al., 2023). However, there is currently limited formal justification of such observations. This motivates the first question of this paper:

*How does prompt-tuning compare to fine-tuning in terms of expressive power? Are there scenarios prompt-tuning outperforms fine-tuning in that regard?*

The core constituent of a transformer, and thus of prompt-tuning, is the attention mechanism. Through the attention layer, prompts get to interact with other input features, create/modify attention weights, and facilitate the model to attend on latent task-specific information. The standard attention layer relies on softmax nonlinearities. Operationally, the softmax function allows a model to selectively focus on certain parts of the input tokens when generating attention outputs. However, there is little rigorous understanding of attention-based prompt-tuning. Concretely,

*What is the role of the softmax-attention in prompt-tuning in terms of optimization and generalization? How does it locate and extract relevant contextual information?*

**Contributions.** Our contributions are as follows:

- We show that a particular form of attention which we refer to as *prompt-attention* naturally arises from the self-attention model with prompt-tuning. We identify provable scenarios where it is more expressive than self-attention and linear-prompt-attention. <sup>1</sup> This separation result reveals insightful data models where prompt-tuning can be superior to fine-tuning with self-attention.

<sup>1</sup>Our emphasis is on the role of attention (whether prompt- or self-). However, we analyze the general problem where attention weights are optimized jointly with the linear classifier head.

- We develop new statistical foundations for gradient-based prompt-tuning: we study the optimization and generalization dynamics of the initial trajectory of gradient descent for optimizing prompt-attention. Concretely, we show the first few iterations learn the prompt and prediction head with near-optimal sample complexity while achieving high accuracy.
- Our results provide insights into the critical role of softmax in facilitating attention: we show how the initial trajectory of gradient descent utilizes softmax to provably attend to sparse context-relevant tokens, ignoring noisy/nuisance tokens.
- We also characterize the exact finite sample performance of prompt-attention assuming known prompt but unknown prediction head. This reveals the fundamental performance limits and precisely quantifies the benefits of context information.
- Our results highlight various trade-offs among different model parameters: (i) the role of sparsity, i.e., the fraction of context-relevant tokens, and (ii) the relative effects of the different constituents of context-relevant tokens.
- Finally, we empirically validate our theoretical insights on both synthetic contextual-mixture datasets and image-classification datasets. Specifically, we compare multiple variants of prompt-tuning against standard fine-tuning on the latter. Our results highlight the role of prompt-attention in selecting relevant tokens in the image classification setting.

**Related works.** Attention, specifically the so-called self-attention, is the backbone mechanism of transformers (Vaswani et al., 2017). It differs from conventional models (e.g., multi-layer perceptrons and convolutional neural networks) in that it computes feature representations by globally modeling interactions between different parts of an input sequence. Despite tremendous empirical success (see, e.g., Vaswani et al., 2017; Brown et al., 2020; Saharia et al., 2022; Ramesh et al., 2022; Cha; Narayanan et al., 2021; Reed et al., 2022, and references therein), the underlying mechanisms of the attention layer remain largely unknown: How does it learn? What makes it better (and when) compared to conventional architectures? Yun et al. (2020) show that self-attention based transformers with large enough depth are universal approximators of seq2seq functions. Focusing on the self-attention component, Edelman et al. (2021) show that self-attention can efficiently represent sparse functions of its input space, while Sahiner et al. (2022); Ergen et al. (2022) analyze convex-relaxations of Self-attention, and Baldi & Vershynin (2022); Dong et al. (2021) study the expressive ability of attention layers. However, these works do *not* characterize the optimization and generalization dynamics of attention. To the best of our knowledge, the only prior works attempting this are Jelassi

et al. (2022) and Li et al. (2023). Jelassi et al. (2022) assume a simplified attention structure in which the attention matrix is *not* directly parameterized in terms of the input sequence. Our paper also distinguishes itself from contemporaneous work by Li et al. (2023) in several ways: (1) Unlike their data model, ours incorporates a context vector and employs a sub-Gaussian noise model instead of assuming bounded noise. (2) We provide a precise asymptotic analysis that elucidates the role of various problem parameters. (3) While Li et al. (2023) primarily focuses on vanilla self-attention, our study centers on understanding the potential benefits of prompt-tuning through prompt-attention.

## 2. Problem setting

### 2.1. Motivation: Prompt-tuning

Consider a single-head self-attention layer

$$\mathbf{O}_{\text{pre}} = \phi(\mathbf{X}\mathbf{W}_Q\mathbf{W}_K^\top\mathbf{X}^\top)\mathbf{X}\mathbf{W}_V, \quad (1)$$

with input  $\mathbf{X} \in \mathbb{R}^{T \times d}$  consisting of  $T$  tokens of dimension  $d$  each, trainable parameters  $(\mathbf{W}_K, \mathbf{W}_Q, \mathbf{W}_V)$  and a softmax nonlinearity  $\phi: \mathbb{R}^T \mapsto \mathbb{R}^T$ ,  $[\phi(\mathbf{v})]_t = e^{v_t} / \sum_{t' \in [T]} e^{v_{t'}}$  that acts row-wise when its argument is a  $T \times T$  matrix. We scalarize the output of the self-attention layer with a trainable linear head  $\bar{\mathbf{U}}$  which yields

$$y_{\text{pre}} = \langle \bar{\mathbf{U}}, \mathbf{O}_{\text{pre}} \rangle = \langle \mathbf{U}, \phi(\mathbf{X}\mathbf{W}_Q\mathbf{W}_K^\top\mathbf{X}^\top)\mathbf{X} \rangle. \quad (2)$$

Note here that we implicitly subsume the value matrix  $\mathbf{W}_V$  in the linear head via  $\mathbf{U} := \bar{\mathbf{U}}\mathbf{W}_V^\top$ .

We assume that the model above is pre-trained so that  $\mathbf{W}_K, \mathbf{W}_Q, \mathbf{U}$  are fixed. Our goal is to use the pretrained transformer on (potentially) new classification tasks. Towards this goal, we explore the use of prompt-tuning, introduced in Li & Liang (2021); Lester et al. (2021) as an alternative to fine-tuning the existing transformer weights.

Prompt-tuning appends a trainable prompt  $\mathbf{P} \in \mathbb{R}^{m \times d}$  to the input features  $\mathbf{X} \in \mathbb{R}^{T \times d}$  with the goal of conditioning the transformer to solve the new classification task. Let  $\mathbf{X}_{\mathbf{P}} := \begin{bmatrix} \mathbf{P} \\ \mathbf{X} \end{bmatrix} \in \mathbb{R}^{(T+m) \times d}$  be the new transformer input. The output of the attention-layer is thus is of the form

$$\mathbf{O} = \phi(\mathbf{X}_{\mathbf{P}}\mathbf{W}_Q\mathbf{W}_K^\top)\mathbf{X}.$$

Note that this is slightly different from (1) in that now the layer computes a cross-attention between the augmented inputs  $\mathbf{X}_{\mathbf{P}}$  and the original inputs  $\mathbf{X}$ . This is also equivalent to self-attention on  $\mathbf{X}_{\mathbf{P}}$  after masking the prompt  $\mathbf{P}$  from keys. This masking is used to cleanly isolate the residual contribution of the prompt without impacting the frozen attention output. Concretely, let  $\mathbf{W}_{\text{head}}$  be the prediction head associated with the prompt tokens. As before, we

scalarize the output by projecting with a linear head of size  $(T + m) \times d$  as follows:

$$y = \langle [\mathbf{W}_{\text{head}}^\top \mathbf{U}^\top]^\top, \phi(\mathbf{X}_P \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{X}^\top) \mathbf{X} \rangle \quad (3)$$

$$= \underbrace{\langle \mathbf{W}_{\text{head}}, \phi(\mathbf{P} \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{X}^\top) \mathbf{X} \rangle}_{\text{prompt-attention } y_{\text{new}}} + \underbrace{\langle \mathbf{U}, \phi(\mathbf{X} \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{X}^\top) \mathbf{X} \rangle}_{\text{frozen self-attention } y_{\text{pre}}}.$$

Here,  $y_{\text{new}}$  captures the additive impact of prompt-tuning on the prediction. We denote the trainable parameters in the model above as  $\theta := (\mathbf{W}_{\text{head}}, \mathbf{P})$ <sup>2</sup> Since the  $y_{\text{new}}$  term becomes a self-contained module and the features attend directly to the prompt vector, we will refer to it as *prompt-attention*.

Our goal is understanding the expressivity, training dynamics, and generalization properties of the above model. To simplify our analysis, we consider the following setting.

1. We focus our attention on the novel component  $y_{\text{new}}$  of the model output in (3) so as to isolate and fully understand the capabilities of prompt-attention.
2. We assume  $\mathbf{W}_K, \mathbf{W}_Q \in \mathbb{R}^{d \times d}$  are full-rank.
3. We assume a single trainable prompt  $\mathbf{q} \in \mathbb{R}^d$  i.e.,  $m = 1$ .

**Prompt-attention model.** Using these assumptions and setting  $\mathbf{q} := \mathbf{W}_K \mathbf{W}_Q^\top \mathbf{P}^\top \in \mathbb{R}^d$  and  $\mathbf{w} = \mathbf{W}_{\text{head}}^\top \in \mathbb{R}^d$ , we arrive at our core *prompt-attention* model  $f_{\theta}^{\text{ATT}}$  (or simply  $f^{\text{ATT}}$ ):

$$f_{\theta}^{\text{ATT}}(\mathbf{X}) = \langle \mathbf{w}, \mathbf{X}^\top \phi(\mathbf{X} \mathbf{q}) \rangle, \quad \theta = (\mathbf{w}, \mathbf{q}). \quad (4)$$

We shall see that this model exhibits interesting properties to learn rich contextual relationships within the data and can even be more expressive than a single self-attention layer.

We remark that the model above is of interest even beyond the context of prompting: the prompt-attention model in (4) is reminiscent of the simplified model proposed in earlier seq2seq architectures (Bahdanau et al., 2014; Xu et al., 2015; Chan et al., 2015) preceding self-attention and Transformers (Vaswani et al., 2017). Indeed, in the simplified attention mechanism of (Bahdanau et al., 2014; Xu et al., 2015; Chan et al., 2015), the tokens' *relevance scores* and corresponding *attention weights* are determined by  $\mathbf{a} = \phi(\mathbf{X} \mathbf{q})$  in which  $\mathbf{q}$  is a trainable vector and  $\phi$  is the softmax-score transformation. Note here that the trainable parameter  $\mathbf{q}$  corresponds exactly to the trainable prompt vector in (4).

<sup>2</sup>In our model, we train the classifier head  $\mathbf{W}_{\text{head}}$  in addition to the prompt vectors  $\mathbf{P}$ . Despite the additional training for the classifier head, the computational overhead remains minimal, and the overall scheme remains significantly more efficient compared to updating the entire model  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{U}$ .

## 2.2. Contextual data model

Consider supervised classification on IID data  $(\mathbf{X}, y) \sim \mathcal{D}$  with features  $\mathbf{X} \in \mathbb{R}^{T \times d}$  and binary label  $y \in \{\pm 1\}$ .

**Dataset model.** We assume the following about an example  $(\mathbf{X}, y)$  drawn from  $\mathcal{D}$ : The labels  $y$  are distributed as  $\mathbb{P}(y = 1) = 1 - \mathbb{P}(y = -1) = \pi$ ; for simplicity, we set  $\pi = 1/2$  so that  $\mathbb{E}[y] = 0$ . The tokens  $\mathbf{x}_t, t \in [T]$  of input example  $\mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_T]$  are split into a *context-relevant* set  $\mathcal{R} \subset [T]$  and *context-irrelevant* set  $\mathcal{R}^c := [T] - \mathcal{R}$ . Specifically, conditioned on the labels and relevance set  $\mathcal{R}$ , tokens  $\mathbf{x}_t, t \in [T]$  within  $\mathbf{X}$  are i.i.d. as follows

$$\mathbf{x}_t | y = \begin{cases} \mathbf{q}_* + y \mathbf{w}_*, & t \in \mathcal{R} \quad (\text{relevant token}) \\ -\delta^q \mathbf{q}_* - \delta^w y \mathbf{w}_* + \mathbf{z}_t, & t \notin \mathcal{R} \quad (\text{irrelevant token}). \end{cases}$$

(DATA)

In the above,  $\mathbf{q}_*$  is a context-vector that indicates token relevance and  $\mathbf{w}_*$  is a regressor vector.  $y, \delta := (\delta^q, \delta^w), (\mathbf{z}_t)_{t=1}^T$  are independent variables as follows:

- $\delta = (\delta^q, \delta^w) \in \mathbb{R}^2$  is a bounded random variable that obeys  $\delta^q \geq 0$ . Thus,  $\delta$  reflects *out-of-context* information within irrelevant tokens. However,  $\delta$  is allowed to expose label information through  $\delta^w$ . When  $\delta = (0, 0)$  almost surely, we call the resulting distribution **core dataset model**.
- $\mathbf{z}_t$  are independent centered subgaussian and random variables with covariance  $\Sigma$  (see Ass. 1.a). When  $\Sigma = 0$ , we call the resulting distribution **discrete dataset model**.
- We allow the relevance set  $\mathcal{R}$  to be non-stochastic. This includes  $\mathcal{R}$  being adversarial to the classification model.
- We assume constant fraction  $\zeta = |\mathcal{R}|/T \in (0, 1)$  of label-relevant tokens for each input example  $\mathbf{X}$  drawn from  $\mathcal{D}$ . Thus,  $\zeta$  represents the sparsity of relevant signal tokens.

**Training dataset**  $\mathcal{S} := (\mathbf{X}_i, y_i)_{i=1}^n$ . We draw  $n$  i.i.d. samples from  $\mathcal{D}$  to form our training dataset  $\mathcal{S} := (\mathbf{X}_i, y_i)_{i=1}^n$ . For  $i$ 'th example  $(\mathbf{X}_i, y_i)$ , we denote the tokens by  $(\mathbf{x}_{i,t})_{t=1}^T$ , noise by  $(\mathbf{z}_{i,t})_{t=1}^T$ , relevance set by  $\mathcal{R}_i$ , and out-of-context variable by  $\delta_i = (\delta_i^q, \delta_i^w)$ .

Ideally, for  $i$ 'th example, we would like to identify its context-relevant set  $\mathcal{R}_i$  and discard the rest. This would especially help when the signal-to-noise-ratio is small, i.e.  $\zeta = |\mathcal{R}_i|/T \ll 1$ . This is precisely the role of the context-vector  $\mathbf{q}_*$ : Observe that, per our construction, relevant tokens have positive correlation with  $\mathbf{q}_*$  whereas irrelevant tokens have non-positive correlation with  $\mathbf{q}_*$  in expectation. Thus, by focusing attention onto tokens based on their  $\mathbf{q}_*$  correlation, we can potentially select the relevant set.

**Remark 1** (Model interpretation). (DATA) can be thought of as a simplified model for binary image classification with tokens being image patches of two types: ones revealing

information about the label (set  $\mathcal{R}$ ) and uninformative ones containing noise. Tokens in  $\mathcal{R}$  contain information indicating: (i) class-membership via signed-regressor  $y\mathbf{w}_*$  and (ii) context-relevance via context-vector  $\mathbf{q}_*$ . The signed-regressor differs across tokens of examples belonging to different classes  $y \in \{\pm 1\}$ , while the context-vector is common for all context-relevant tokens across classes. For a concrete example, consider images each depicting multiple, say  $|\mathcal{R}|$ , birds of one type surrounded by label-irrelevant/noisy background. The goal is to classify images according to one of two types of birds. Here, think of “context” as feature-information indicating corresponding pixels belong to “bird” (of either type) rather than “background,” while the “regressor” represents feature information useful to distinguish between two bird types. Alternatively, (DATA) may be modeling deep representations (rather than raw pixels) of the original images. Overall, simplified models similar to (DATA) have been used previously to analyze optimization and generalization dynamics of training fully-connected (Frei et al., 2022) and convolutional models (Cao et al., 2022). Specifically, (DATA) is an extension of the commonly used (sub)-Gaussian mixture model customized to the nature of attention: each example is tokenized and context-relevant information is described in terms of both a regressor (differing between classes) and a context (common across classes).

### 2.3. Baseline Models

We compare performance of the prompt-attention model in (4) with the following three baseline models.

**The linear model** is parameterized by  $\theta = \mathbf{w}$  and outputs

$$f^{\text{LIN}}(\mathbf{w}) = \frac{1}{T} \mathbf{w}^\top \mathbf{X}^\top \mathbf{1}_T = \frac{1}{T} \sum_{t \in [T]} \mathbf{w}^\top \mathbf{x}_t. \quad (5)$$

Note this corresponds to a prompt-attention model with uniform attention weights  $[\mathbf{a}]_t = 1/T, t \in [T]$ .

**The self-attention model** is a strict generalization of the linear model. Recalling (2), let us merge the key-query weights  $\mathbf{W} := \mathbf{W}_Q \mathbf{W}_K^\top$  (without losing generality) and gather weights into  $\theta = (\mathbf{U}, \mathbf{W})$ ; We then write it as

$$f^{\text{SATT}}(\mathbf{U}, \mathbf{W}) = \frac{1}{T} \langle \mathbf{U}, \phi(\mathbf{X} \mathbf{W} \mathbf{X}^\top) \mathbf{X} \rangle. \quad (6)$$

Rather than using a  $Td$  dimensional  $\mathbf{U}$ , we will also consider the simpler token-pooling via  $\mathbf{U} = \mathbf{1}_T \mathbf{u}^\top$  for  $\mathbf{u} \in \mathbb{R}^d$ .

**The linear-attention model** parameterized by  $\theta = (\mathbf{w}, \mathbf{q})$  replaces the softmax score transformation in (4) with a linear function and outputs

$$f^{\text{LIN-ATT}}(\mathbf{w}, \mathbf{q}) = \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{q} / T. \quad (7)$$

### 2.4. Training

Given  $\mathcal{S} = (\mathbf{X}_i, y_i)_{i=1}^n$  drawn i.i.d. from  $\mathcal{D}$ , we solve square-loss empirical risk minimization to obtain  $\hat{\theta} = (\hat{\mathbf{w}}, \hat{\mathbf{q}})$

$$\hat{\theta} = \arg \min_{\theta} \widehat{\mathcal{L}}_{\mathcal{S}}(\theta) := \frac{1}{2n} \sum_{i=1}^n (y_i - f_{\theta}(\mathbf{X}_i))^2. \quad (8)$$

Within our theoretical investigation, we are interested in the following performance criteria for models  $f \in \{f^{\text{ATT}}, f^{\text{LIN-ATT}}, f^{\text{SATT}}\}$ :

- **Classification error:** For a model  $f_{\hat{\theta}}$  this is defined as  $\text{ERR}(\hat{\theta}) := \mathbb{P}(y f_{\hat{\theta}}(\mathbf{X}) < 0)$ .
- **Test risk:**  $\mathcal{L}(\hat{\theta}) = \mathbb{E}_{(y, \mathbf{X}) \sim \mathcal{D}}[(y - f_{\hat{\theta}}(\mathbf{X}))^2]$ .

### 2.5. Assumptions and notations

First, we formally state our assumptions on the noisy tokens. The more general condition is that noise is subgaussian and satisfies a mild zero third-moment condition.

**Assumption 1.a.** *The noise vector  $\mathbf{z} \sim \mathcal{SN}(\sigma)$  is centered  $\sigma$ -subGaussian, i.e.  $\|\mathbf{z}\|_{\psi_2} = \sigma$ . Moreover, its distribution is symmetric and has zero-third moment, i.e.  $\mathbb{E}[\mathbf{z} \otimes (\mathbf{z}^\top \mathbf{z})] = 0$ . Let  $\Sigma := \mathbb{E}[\mathbf{z} \mathbf{z}^\top]$  denote the noise covariance.*

For some of our results it will be convenient to further assume that noise is Gaussian since this leads to precise formulas that are easily interpretable.

**Assumption 1.b.** *The noise vector  $\mathbf{z} \sim \mathcal{N}(0, \sigma^2 \mathbb{I})$  is isotropic Gaussian with variance  $\sigma^2$ .*

Second, we require a mild assumption on the correlation between the context  $\mathbf{q}_*$  and classifier  $\mathbf{w}_*$  to guarantee that pure signal tokens  $\mathbf{q}_* + y\mathbf{w}_*$  are correctly classified by the true regressor  $\mathbf{w}_*$ , i.e.  $y\mathbf{w}_*^\top(\mathbf{q}_* + y\mathbf{w}_*) > 0$ . For convenience, we denote

$$W := \|\mathbf{w}_*\|, \quad Q := \|\mathbf{q}_*\|, \quad \rho := \mathbf{q}_*^\top \mathbf{w}_* / (\|\mathbf{q}_*\| \|\mathbf{w}_*\|).$$

**Assumption 3.a.** *Correlation satisfies  $|\rho| < W/Q$ .*

We will also often consider the special case of zero correlation  $\rho$  and thus state it as separate assumption below. This orthogonality assumption, is useful for more tractable analysis as it helps decouple feature selection and prediction.

**Assumption 3.b.** *The context and classifier vectors are orthogonal, i.e.  $\mathbf{q}_* \perp \mathbf{w}_*$ .*

**Notation.** We use boldface letters for vectors and matrices.  $\mathbf{1}_m$  represents an  $m$ -dimensional all-ones vector. For a vector  $\mathbf{v}$ ,  $\|\mathbf{v}\|$  denotes its Euclidean norm and  $\mathbf{v}/\|\mathbf{v}\|$  its normalization.  $\phi(\cdot)$  denotes the softmax transformation.  $\mathcal{Q}(\cdot)$  denotes the tail function of the standard normal distribution.  $\wedge$  and  $\vee$  denote the minimum and maximum of two numbers, respectively.  $\tilde{O}()$  and  $\gtrsim$  notations suppress logarithmic dependencies. Finally,  $\propto$  denotes proportionality.

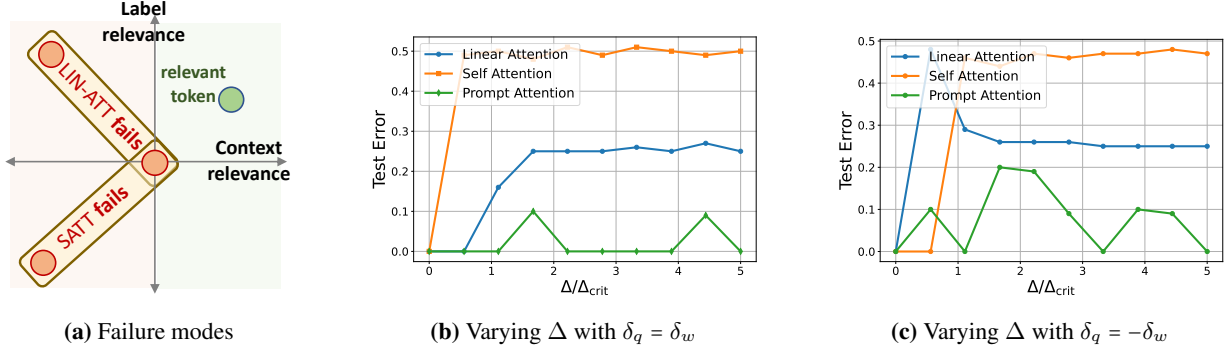


Figure 1. This figure summarizes and verifies the outcomes of Theorem 1. **Fig (a)** depicts the outcome of our Theorem 1. Relevant token is at position  $y\mathbf{w}_* + \mathbf{q}_*$  whereas red tokens (irrelevant) are in positions  $-\delta(\mathbf{q}_* \pm y\mathbf{w}_*)$  with  $\delta \in \{0, \Delta\}$ . **Figures (b) & (c)** plot the performance of our attention models under the contextual dataset with  $\delta$  equally-likely over  $\{0, \Delta\}$  for a synthetic setup (cf. Section 5.1). We set  $n = 100$ ,  $d = T = 10$ ,  $\zeta = 0.4$  and train with 100 SGD epochs. We report the median test accuracy over 20 trials. Fig (b) sets  $\delta = \delta_q = \delta_w$  and verifies self-attention has  $\geq 50\%$  error (for  $\Delta \geq \Delta_{\text{crit}} = (1 - \zeta)^{-2}$ ). Fig (c) sets  $\delta = \delta_q = -\delta_w$  and verifies linear-prompt-attention has  $\geq 25\%$  error (when  $\delta \geq \Delta_{\text{crit}} = \sqrt{\zeta}/(1 - \zeta)$  in this case).

### 3. Contrasting prompt-attention to baselines

In this section, we establish separation results between prompt-attention (cf. (4)) and the baselines of self-attention (cf. (6)) and linear attention (cf. (7)). For this, we focus on the discrete dataset model with noiseless irrelevant tokens ( $\mathbf{z}_t = 0$ ,  $t \in [T]$ ).

We first observe that if  $\delta = (\delta^q, \delta^w)$  admits a single value, even a linear model can solve the discrete dataset model.

**Observation 1** (Linear model solves singleton). *Suppose  $(\delta^q, \delta^w) = (\Delta^q, \Delta^w)$  almost surely for  $\Delta^q, \Delta^w \in \mathbb{R}$ . Set  $\mathbf{w}'_* = (\mathbf{I} - \bar{\mathbf{q}}_* \bar{\mathbf{q}}_*^\top) \mathbf{w}_*$ . As long as  $\Delta^w \neq \zeta/(1 - \zeta)$  and  $\mathbf{w}'_* \neq 0$ ,  $f^{\text{LIN}}(\mathbf{w}'_*)$  or  $f^{\text{LIN}}(-\mathbf{w}'_*)$  achieves perfect accuracy.*

Thus, to investigate the expressivity of  $f^{\text{ATT}}, f^{\text{SATT}}, f^{\text{LIN}}$ ,  $(\delta^q, \delta^w)$  would need to admit two or more values. Perhaps surprisingly, we prove that, as soon as,  $(\delta^q, \delta^w)$  comes from a binary distribution, then both  $f^{\text{SATT}}$  and  $f^{\text{LIN}}$  can indeed provably fail. Importantly, this happens in the regime  $\delta^q \geq 0$  where prompt-attention thrives.

**Theorem 1** (Separation of population accuracies). *Consider the discrete dataset model where we set  $\Sigma = 0$  in (DATA). The following statements hold:*

1. **Prompt-attention:** *Suppose  $\rho^2 < 1$ ,  $\delta^q \geq 0$ , and  $|\delta^w| \leq C$  almost surely. Define  $\mathbf{q}'_* = (\mathbf{I} - \bar{\mathbf{w}}_* \bar{\mathbf{w}}_*^\top) \mathbf{q}_*$ ,  $\mathbf{w}'_* = (\mathbf{I} - \bar{\mathbf{q}}_* \bar{\mathbf{q}}_*^\top) \mathbf{w}_*$ . For  $\Gamma > \frac{\log(C/(1-\zeta))}{Q^2(1-\rho^2)}$ , choosing  $\theta = (\mathbf{w}'_*, \Gamma \mathbf{q}'_*)$ ,  $f_\theta^{\text{ATT}}$  achieves perfect classification accuracy on (DATA).*

2. **Self-attention:** *In (DATA), choose  $(\delta^q, \delta^w)$  to be  $(0, 0)$  or  $(\Delta, \Delta)$  equally-likely with  $\Delta > (1 - \zeta)^{-2}$ .*

- For any choice of  $(\mathbf{U} = \mathbb{1}\mathbf{u}^\top, \mathbf{W})$ ,  $f^{\text{SATT}}(\mathbb{1}\mathbf{u}^\top, \mathbf{W})$  achieves 50% accuracy (i.e. random guess).
- For any choice of  $(\mathbf{U}, \mathbf{W})$ , there exists a (DATA) distri-

bution with adversarial relevance set choices such that  $f^{\text{SATT}}(\mathbf{U}, \mathbf{W})$  achieves 50% accuracy.

3. **Linear-attention:** *In (DATA), choose  $(\delta^q, \delta^w)$  to be  $(0, 0)$  or  $(\Delta, -\Delta)$  equally-likely with  $\Delta > \sqrt{\zeta}/(1 - \zeta)$ . For any choice of  $(\mathbf{w}, \mathbf{q})$ ,  $f^{\text{LIN-ATT}}(\mathbf{w}, \mathbf{q})$  achieves at most 75% accuracy.*

See Fig. 1 for an illustration of the main takeaways from Thm. 1 and a numerical validation of its conclusion on synthetic data. While surprising, the reason prompt-attention can provably beat self-attention is because it is optimized for context-retrieval and can *attend* perfectly on the relevant contextual information. In contrast, self-attention scores are fully feature-based; thus, context information is mixed with other features and can be lost during aggregation of the output. Also note that all results, with the exception of self-attention for general  $\mathbf{U}$ , hold for arbitrary choices of the relevance sets (including adversarial ones). The reason is that tokens are pooled and the particular choice of  $\mathcal{R}$  does not matter. Only for  $f^{\text{SATT}}(\mathbf{U}, \mathbf{W})$  we need to adapt the relevance set  $\mathcal{R}$  to the output layer  $\mathbf{U}$  (as well as  $(y, \delta)$  variables) to promote misclassification. Otherwise, with the hindsight knowledge of the relevance set,  $\mathbf{U}$  can intelligently process individual tokens of the self-attention output to filter out “confusing” tokens. In fact, for the same failure dataset model, self-attention can achieve perfect accuracy by choosing  $\mathbf{U} = \mathbb{1}_{\mathcal{R}} \mathbf{w}'_*{}^\top$  where  $\mathbb{1}_{\mathcal{R}}$  is the vector of ones over the (known!) relevance set  $\mathcal{R}$  (see Lemma 17). However, this is of course only known in hindsight.

## 4. Gradient-based analysis of prompt-attention

This section investigates how gradient-descent optimization of the prompt-attention model learns (DATA). Concretely, it shows that a few gradient steps can provably attend to the context-relevant tokens leading to high-classification accuracy. Our results capture requirements on sample complexity in terms of all problem parameters, i.e. dimension  $d$ , correlation  $\rho$ , context / signal energies  $Q / W$ , number of tokens  $T$ , and sparsity  $\zeta$ . This allows studying tradeoffs in different regimes.

Our analysis in this section concerns the prompt-attention model  $f_{\theta}^{\text{ATT}}$ , so we simply write  $f_{\theta}$ . Also, without any further explicit reference, we focus on the *core dataset model*, i.e. (DATA) with  $\delta = (0, 0)$ . All our results here hold under the mild noise and correlation assumptions: Assumption 1.a and Assumption 3.a (we will not further state these). Finally, for simplicity of presentation we assume here isotropic noise  $\Sigma = \sigma^2 \mathbb{I}$  and handle the general case in the appendix.

### 4.1. Gradient-based algorithm

For data generated from (DATA), we show the three-step gradient-based algorithm described below achieves high test accuracy. Our analysis also explains why three appropriately chosen steps suffice.

**Algorithm:** We split the train set in three separate subsets  $\mathcal{S}_1, \mathcal{S}_2, \mathcal{S}_3$  of size  $n$  each. Starting from  $\mathbf{w}_0 = 0, \mathbf{q}_0 = 0$ , the algorithm proceeds in three gradient steps for step sizes  $\eta > 0$  and  $\gamma > 0$  and a final debiasing step as follows:

$$\widehat{\mathbf{w}}_1 := -\eta \nabla_{\mathbf{w}} \widehat{\mathcal{L}}_{\mathcal{S}_1}(0, 0), \quad (9a)$$

$$\widehat{\mathbf{q}}_1 := -\gamma \nabla_{\mathbf{q}} \widehat{\mathcal{L}}_{\mathcal{S}_2}(0, \widehat{\mathbf{w}}_1), \quad (9b)$$

$$\widehat{\mathbf{w}}_2 := -\eta \nabla_{\mathbf{w}} \widehat{\mathcal{L}}_{\mathcal{S}_3}(\widehat{\mathbf{q}}_1, \widehat{\mathbf{w}}_1), \quad (9c)$$

where  $\widehat{\mathcal{L}}_{\mathcal{S}_j}, j = 1, 2, 3$  is the loss in (8) evaluated on sets  $\mathcal{S}_j$ . The debiasing step is defined in Section 4.3.

### 4.2. Population analysis

To gain intuition we first present results on the population counterpart of the algorithm, i.e., substituting  $\widehat{\mathcal{L}}(\mathbf{w}, \mathbf{q})$  with its population version  $\mathcal{L}(\mathbf{w}, \mathbf{q}) = \mathbb{E}[\widehat{\mathcal{L}}(\mathbf{w}, \mathbf{q})]$  in all three steps in (9). It is convenient to introduce the following shorthand notation for the negative gradient steps  $\mathbf{G}_{\mathbf{q}}(\mathbf{q}, \mathbf{w}) := -\nabla_{\mathbf{q}} \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{X}, y) \sim \mathcal{D}}[(y - f_{\boldsymbol{\theta}}(\mathbf{X})) \nabla_{\mathbf{q}} f_{\boldsymbol{\theta}}(\mathbf{X})]$  and  $\mathbf{G}_{\mathbf{w}}(\mathbf{q}, \mathbf{w}) := -\nabla_{\mathbf{w}} \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{X}, y) \sim \mathcal{D}}[(y - f_{\boldsymbol{\theta}}(\mathbf{X})) \nabla_{\mathbf{w}} f_{\boldsymbol{\theta}}(\mathbf{X})]$ .

The first gradient step (cf. (9a)) is easy to calculate and returns a classifier estimate that is already in the direction of  $\mathbf{w}_*$ .

**Lemma 1** (Population first step). *The first population gradi-*

*ent step  $\mathbf{w}_1 = \eta \mathbf{G}_{\mathbf{w}}(0, 0)$  satisfies  $\mathbf{w}_1 = \eta \zeta \mathbf{w}_*$  since under (DATA),  $\mathbf{G}_{\mathbf{w}}(0, 0) = \mathbb{E}_{(\mathbf{X}, y)}[y \mathbf{X}^{\top} \mathbb{1}] / T = \zeta \mathbf{w}_*$ .*

The second gradient step  $\mathbf{q}_1 = \gamma \mathbf{G}_{\mathbf{q}}(\mathbf{w}_1, 0)$  is more intricate: unless  $\mathbf{q}_* \perp \mathbf{w}_*$ ,  $\mathbf{q}_1$  also has nonzero components in both directions  $\mathbf{q}_*$  and  $\mathbf{w}_*$ .

**Lemma 2** (Population second step). *The second population gradient step  $\mathbf{q}_1 = \gamma \mathbf{G}_{\mathbf{q}}(\mathbf{w}_1, 0)$  satisfies for  $\alpha := \eta \zeta$*

$$\mathbf{q}_1 = \gamma \alpha W^2 (\zeta - \zeta^2) \left(1 + \frac{\alpha \sigma^2}{T} - \alpha \zeta (W^2 + \rho^2 Q^2)\right) \mathbf{q}_* \quad (10) \\ + \gamma \alpha \rho Q W (\zeta - \zeta^2) \left(1 - 2\zeta \alpha W^2 - \left(1 + \frac{2}{T}\right) \alpha \sigma^2\right) \mathbf{w}_*.$$

*Proof.* Since this computation involves several terms, we defer complete proof to Appendix D.1. The above simplification is made possible by leveraging the assumption on the third-moment of noise (cf. Assumption 1.a).  $\square$

Lemma 2 highlights the following key aspects: (i) As mentioned,  $\mathbf{q}_1$  also picks up the  $\mathbf{w}_*$  direction unless  $\rho = 0$ . However, we can control the magnitude of this undesired term by choosing small step-size  $\eta$  (see Cor. 2). (ii) As  $\alpha W^2$  grows, the gradient component in the  $\mathbf{q}_*$  direction might end up pointing in the direction of  $-\mathbf{q}_*$ . This is because large signal along the  $\mathbf{w}_*$  direction might still allow to predict  $\pm 1$  label. However, this can always be avoided by choosing sufficiently small step-size  $\eta$  (see Cor. 2). (iii) Similarly, as the noise strength  $\sigma^2$  grows, gradient in the  $\mathbf{q}_*$  direction grows as well. This is because, going along  $\mathbf{q}_*$  direction attenuates the noise and cleans up the prediction. (iv) Finally, as  $\zeta \rightarrow 1$  and  $\zeta - \zeta^2 \rightarrow 0$  the magnitude of the gradient decays because all tokens contain signal information and there is no need for  $\mathbf{q}_*$ .

To see how  $\mathbf{q}_1$  selects good tokens, we investigate the relevance scores (normalized by the step size  $\gamma$ )  $r_t := \mathbf{x}_t^{\top} \mathbf{q}_1 / \gamma$  of relevant vs irrelevant tokens. Attending to context-relevant tokens requires their relevance scores to be larger than those of the noisy ones. Concretely, suppose we have

$$B := \min_{t \in \mathcal{R}} \{r_t = \frac{(\mathbf{q}_* + y \mathbf{w}_*^{\top}) \mathbf{q}_1}{\gamma}\} \geq 2 \max_{t \in \mathcal{R}^c} \{r_t = \frac{\mathbf{z}_t^{\top} \mathbf{q}_1}{\gamma}\}. \quad (11)$$

Note above that the relevance scores are the same for each  $t \in \mathcal{R}$ . Thus,  $|\mathcal{R}| e^{\gamma B} + |\mathcal{R}^c| e^{\gamma B/2} \geq S := \sum_{t' \in [T]} e^{\gamma r_{t'}} \geq |\mathcal{R}| e^{\gamma B}$ , which implies the following for the attention weights as step size increases  $\gamma \rightarrow \infty$ :

$$a_t = [\phi(\mathbf{X} \mathbf{q}_1)]_t = e^{\gamma r_t} / S \begin{cases} = \frac{e^{\gamma B}}{S} \rightarrow \frac{1}{\zeta T} & t \in \mathcal{R} \\ \leq \frac{e^{\gamma B/2}}{S} \rightarrow 0 & t \in \mathcal{R}^c \end{cases}. \quad (12)$$

Provided (11) holds, a large enough second gradient step (i.e. large  $\gamma$ ) finds  $\mathbf{q}_1$  that attends (nearly) perfectly to context-relevant tokens in  $\mathcal{R}$  and attenuates (almost) all irrelevant tokens in  $\mathcal{R}^c$ . The following theorem formalizes the above intuition. We defer the complete proof to Appendix D.3.

**Theorem 2** (Main theorem: Population). *Consider the model  $\theta^\gamma = (\mathbf{w}_2^\gamma, \mathbf{q}_1^\gamma)$  where  $\mathbf{q}_1^\gamma = \gamma \mathbf{G}_q(\mathbf{w}_1, 0)$ ,  $\mathbf{w}_2^\gamma = \mathbf{G}_w(0, \mathbf{q}_1^\gamma)$  and  $\mathbf{w}_1 = \eta \mathbf{G}_w(0, 0)$  for step-size  $\eta$  small enough (see Eq. (66) for details). Then, there exists an absolute constant  $c > 0$ , sufficiently large context strength  $Q$  and step-size  $\gamma > 0$  such that*

$$\text{ERR}(f_{\theta^\gamma}) \leq 2Te^{-c\frac{\alpha^2 Q^2}{\sigma^2}},$$

provided

$$\frac{(1 - \rho^2/2)Q - 2|\rho|W}{\sqrt{1 + 3\rho^2}} \geq \alpha Q. \quad (13)$$

Eq. (13) guarantees the desired condition (11) holds. When  $\rho = 0$  ( $\mathbf{q}_* \perp \mathbf{w}_*$ ),  $\alpha$  can be as large as 1 in (13) in which case the rate is  $2Te^{-cQ^2/\sigma^2}$ . For  $\rho \neq 0$ , (13) imposes  $|\rho| < \frac{Q}{2W}$ , in which the role of  $Q, W$  is reversed compared to  $|\rho| \leq \frac{W}{Q}$  in Assumption 3.a: the latter guarantees classifier energy is larger so that signal  $y\mathbf{w}_*$  dominates  $\mathbf{q}_*$ , while for prompt-attention to attend to relevant tokens it is favorable that energy of  $\mathbf{q}_*$  dominates  $\mathbf{w}_*$ . Finally, we compare the theorem's error to the error  $Q(\sqrt{\frac{\zeta^2 W^2 T}{1-\zeta}})$  of the linear model in Fact A.1. For concreteness, consider a setting of extreme sparsity  $\zeta = O(1/\sqrt{T})$  and  $W = O(1)$ . Then the error of linear model is  $O(1)$ , while the (population) algorithm in (9) for prompt-attention achieves an error of  $e^{-O(Q^2)}$ , which is exponentially decreasing in  $Q$ .

### 4.3. Finite-sample analysis

Here, we investigate the behavior of the algorithm in (9) with finite sample-size  $n$ . For convenience, we first introduce an additional de-biasing step after calculating the three gradients in (9). Specifically, for a sample  $\mathcal{S}_4$  of size  $n$  we compute a bias variable  $\widehat{b} := \frac{1}{n} \sum_{i=1}^n f(\widehat{q}_1, \widehat{w}_2)(\mathbf{X}_i)$ , and use it to de-bias the model's prediction by outputting  $f_{(\theta, b)}(\mathbf{X}) := f_\theta(\mathbf{X}) - b$ . While this extra step is not necessary, it simplifies the statement of our results. Intuitively,  $\widehat{b}$  helps with adjusting the decision boundary by removing contributions of the context vector in the final prediction (the context vector is useful only for token-selection rather than final prediction).

Below we provide a simplified version of our main result where noise variance  $\sigma \propto 1$  and  $\gtrsim$  subsumes constants. Refer to Theorem 6 in the appendix for precise details.

**Theorem 3** (Main theorem: Finite-sample). *Suppose  $Q, W$  and  $\rho$  are such that there exists  $\alpha \in (0, 3/16)$  for which*

$$(3/16 - \rho^2/8)Q - (9/4|\rho| + 1/16)W \geq \alpha \cdot Q \gtrsim \sqrt{\log(nT)}.$$

*Fix any  $\epsilon > 0$ . For sufficiently small step-size  $\eta \lesssim Q^{-2}$ , sufficiently large step-size  $\gamma = \gamma(\epsilon)$ , and*

$$n \gtrsim d(Q/\zeta W)^4 \log^5(nd),$$

*the following statements hold with high probability (see Eq. (63)) over the training set:*

**1. Prompt attends to relevant tokens.** *Concretely, for any fresh sample  $(\mathbf{X}, y)$ , with probability at least  $1 - 2Te^{-c\alpha^2 Q^2}$ , the attention coefficients  $a_t = [\phi(\mathbf{X}\widehat{\mathbf{q}}_1)]_t$  satisfy:*

$$a_t \begin{cases} \geq \frac{1-\epsilon}{\zeta T}, & t \in \mathcal{R}, \\ \leq \frac{\epsilon}{(1-\zeta)T}, & t \notin \mathcal{R}. \end{cases}$$

**2. Prompt-attention learns relevant features.** *Concretely, for some absolute constant  $c > 0$ ,*

$$\mathbb{P}_{(\mathbf{X}, y)}(\|\mathbf{X}^\top \phi(\mathbf{X}\widehat{\mathbf{q}}_1) - (\mathbf{q}_* + y\mathbf{w}_*)\| < \epsilon) \leq 2Te^{-c\alpha^2 Q^2}.$$

**3. The test error of the model  $f'_\theta$  similarly satisfies**

$$\text{ERR}(f'_\theta) \leq 2Te^{-c\alpha^2 Q^2}.$$

Assuming small correlation coefficient  $\rho$  and  $W < Q$ , we can set  $\alpha = O(1)$ . Then, similar to Theorem 2, prompt attention achieves a test error rate of  $e^{-O(Q^2)}$  which is a strict improvement over the linear baseline of Fact A.1 whenever  $Q^2 \gtrsim \zeta^2 W^2 T$ . Secondly, our bound achieves a sample complexity of  $n \gtrsim d(Q/\zeta W)^4$ . The linear growth in  $d$  is intuitive from counting degrees of freedom. Interestingly, large  $Q$  does improve the test error, however, it degrades sample complexity. This is because it makes the estimation of parameters more challenging. Finally, larger  $\zeta W$  improves sample complexity since  $\zeta W$  (combining sparsity level and magnitude) captures the strength of the label-relevant information within relevant tokens  $t \in \mathcal{R}$ .

**Sharp error rates:** Finally, in Appendix A we provide an exact analysis of the classification error when  $\mathbf{q}_*$  is known and only  $\mathbf{w}_*$  is estimated from data. This analysis exactly quantifies the value of context-information and how prompt-tuning retrieves it. Specifically, we prove a sharp asymptotic error rate of  $Q \left( \frac{e^{Q^2/4}}{\sqrt{1 + \text{ISNR}(n/d)}} \cdot \sqrt{\frac{\zeta^2 W^2 T}{1-\zeta}} \right)$

where  $\text{ISNR}(\alpha) := \alpha^{-1} \frac{(1-\zeta)e^{-Q^2/2}}{\text{rate}_{\text{LIN}}}$ , This uniformly (for all  $Q \geq 0$  values) improves the optimal rates for (context-free) Gaussian mixture models thanks to the context information.

## 5. Experiments

First, we verify the utility of prompt-attention via experiments on a synthetic setting that precisely follows the contextual data model from Section 2.2. Subsequently, we explore prompt-tuning on various image classification tasks that are motivated by the contextual data model and compare it with the standard fine-tuning method. Finally, we validate the utility of prompt vectors in distinguishing relevant tokens from irrelevant tokens via prompt-attention under an image classification setting.

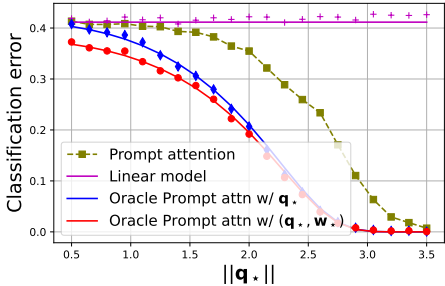


Figure 2. Performance of prompt-attention on the synthetic setting described in Section 5.1. For prompt-attention, we employ the algorithm in (9) to obtain  $\hat{q}$  and  $\hat{w}$ . For the baseline linear model and two oracle settings, we have closed-form expressions for their asymptotic test error (cf. Theorem 1), which are depicted by solid lines. On the other hand, markers show the finite sample performance of these three methods. All finite sample performances are obtained by averaging over 20 independent runs.

### 5.1. Synthetic experiments

Here, we generate a synthetic dataset according to the *core dataset model*, i.e. we have  $\delta = (\delta^q, \delta^w) = (0, 0)$  for all examples in the dataset. In particular, we consider a setting with  $T = 500$ ,  $d = 50$ , and  $\zeta = 0.1$ , i.e. each example has 500 tokens out of which 10% tokens are relevant. As for the noisy tokens, they consists of i.i.d.  $\mathcal{N}(0, I)$  vectors. Assuming that  $q_* \perp w_*$  and  $\sqrt{TW} = 3$ , we generate  $n = 10 \cdot d$  training examples from the core dataset model for varying  $Q$ . Fig. 2 showcases the performance of prompt-attention (cf. (4)) when combined with the estimates  $\hat{q}$  and  $\hat{w}$  produced by gradient-based algorithm in (9). We also showcase the performance of the linear model (cf. (5)) and two oracle methods where we assume access to true  $q_*$  and true  $(q_*, w_*)$ , respectively, while applying the prompt-attention. Note that prompt-attention achieves a vanishing classification test error in this setting whereas a natural baseline (linear model) can fail to achieve a good performance. On the other hand, the prompt-attention enabled by (9) successfully achieves a high accuracy as the context energy (defined by  $Q$ ) increases, validating the utility of prompt-attention as well as our gradient-based algorithm in (9). Finally, we also consider a stochastic  $\delta = (\delta^q, \delta^w) \neq (0, 0)$  to validate Theorem 1 as described in Fig. 1.

### 5.2. Image classification experiments

**Dataset.** Motivated by our contextual data model, we construct three datasets based on CIFAR-10 (Krizhevsky et al., 2009) to conduct our evaluation (see Fig. 3 for examples). Due to spaces constraints, we defer the additional details regarding datasets to Appendix H.1. We also refer the reader to Appendix H.1 for details regarding the model architecture and training procedure.

**Methods.** In our fine-tuning experiments, we update all

pre-trained model parameters. As for prompt-tuning, we only update newly introduced (prompt) variables and keep the pre-trained network frozen. We consider three variants of prompt-tuning: 1) PROMPT-TUNING-I (Lester et al., 2021), where we add trainable vector between CLS token embedding and first image (patch) embeddings only at the input; 2) PROMPT-TUNING-II (Li & Liang, 2021), where we add the *same* trainable vectors between the CLS embedding and the first image embedding at the input of every transformer layer; and 3) PROMPT-TUNING-III, where we add *different* trainable vectors between the CLS embedding and the first image patch embedding at the input of every transformer layer. Note that the number of trainable parameters in PROMPT-TUNING-I and PROMPT-TUNING-II do not scale with the number of layers whereas we have linear scaling with number of layers in PROMPT-TUNING-III. Interestingly, all three prompt-tuning variant are identical when the number of layers is 1, which corresponds to the setup we theoretically analyzed in the paper. However, they exhibit remarkably different behavior for a multi-layer transformer model, as we show next. (In Appendix H.2, we also compare prompt-tuning with fine-tuning only first layer self-attention parameters for the underlying ViT model as per the single-layer nature of our theoretical results.)

**Results.** Here, the main goal of our exploration is to highlight the different behavior of fine-tuning and prompt-tuning. We utilize a model trained on FULL-TILED dataset as the pre-trained model. This model achieves top-1 (in-domain) accuracy of 80.43 on FULL-TILED test set. In contrast, it achieves *zero-shot* top-1 test accuracy of 56.35 and 17.97 on PARTIAL-TILED and EMBED-IN-IMAGENET, respectively. This alludes to the fact that EMBED-IN-IMAGENET corresponds to a larger distribution shift from the pre-training distribution (FULL-TILED), as compared to PARTIAL-TILED.

Fig. 4 and Fig. 5 (cf. Appendix H.2) showcase the performance of fine-tuning and prompt-tuning approaches on EMBED-IN-IMAGENET and PARTIAL-TILED, respectively. Note that fine-tuning outperforms prompt-tuning in a data-rich setting (cf. Fig. 4a and 5a). This is due to fine-tuning having the ability to update a large number of model parameters (5.4M in our case). In contrast, with 2000 prompt vectors, PROMPT-TUNING-III (the most expensive prompt-tuning method out of all three) only updates 460.8K parameters.

Interestingly, in the data-limited regimes, prompt-tuning becomes more competitive. In fact, the best performing prompt-tuning method outperforms the fine-tuning (cf. Fig. 4b and 4c) on EMBED-IN-IMAGENET, where fine-tuning can easily overfit as it cannot leverage the benefits of the pre-training data due to a large distribution-shift between FULL-TILED and EMBED-IN-IMAGENET.

Part of the performance gap between PROMPT-TUNING-III and PROMPT-TUNING-II can be attributed to the larger num-



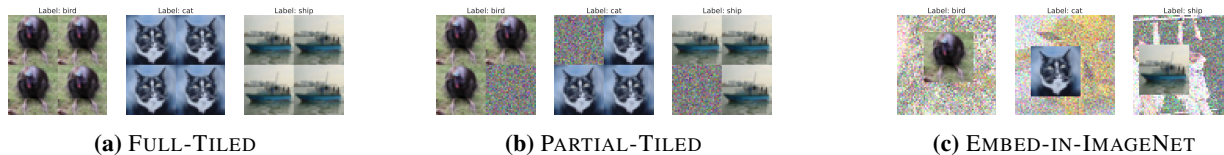


Figure 3. Illustration of different CIFAR-10 based datasets utilized in image classification experiments (cf. Section 5.2). Note that all three variants correspond to 10-way multiclass classification tasks corresponding to 10 original classes in CIFAR-10.

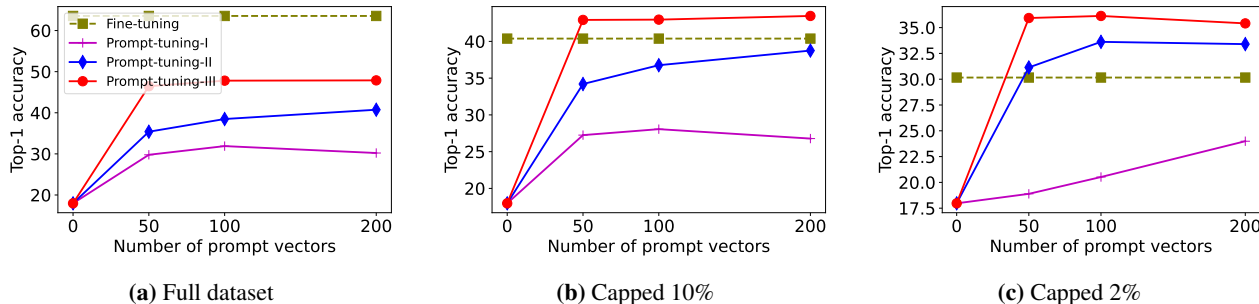


Figure 4. Performance of fine-tuning vs. prompt-tuning on 10-way classification tasks defined by EMBED-IN-IMAGENET dataset. Full dataset has 50K training examples. Capped 10% and 2% correspond to sub-sampled *train* sets where we select exactly 500 and 100 examples per class from the full dataset. Note that number of prompt vectors equal to 0 corresponds to *zero-shot* performance.

ber of trainable parameters available to PROMPT-TUNING-III. Note that PROMPT-TUNING-II consistently outperforms PROMPT-TUNING-I, even with the same number of trainable parameters. This alludes to the fact that optimization and architecture choices play a major role beyond just the number of trainable parameters. As mentioned earlier, our theoretical treatment for a single-layer model cannot distinguish among these different prompt-tuning approaches. As a result, we believe that our empirical observations point to multiple interesting avenues for future work.

### 5.3. Attention weights for prompt vectors

Finally, we explore what role prompt-attention, i.e. the attention weights with prompt vectors as keys and image patches/tokens as values, plays towards underlying task. In Fig. 6 (cf. Appendix H.2), we illustrate one representative example. The figure shows how average attention weights from prompt vectors to image tokens/patches evolve across transformer layers, when we employ PROMPT-TUNING-III. Indeed, the figure verifies that prompt-attention helps distinguish the relevant tokens/patches from the irrelevant patches, validating our starting hypothesis in Section 2.1 and 2.2.

## 6. Discussion

In light of remarkable success of attention architectures, we initiate a theoretical investigation of one of its core mechanisms: prompt-tuning. For a one-layer attention model, we motivate and analyze a simplified model for prompt-tuning, which we coin prompt-attention. Through this model, we developed new statistical foundations for gradient-based

prompt-tuning, characterized its optimization and generalization dynamics, and explored how it facilitates attending to context-relevant information. We also showed that under (DATA) one-layer softmax-prompt-tuning can provably be superior to alternatives including one layer self-attention. Thorough experiments support our theory on how prompt-tuning attends to the context and how it can potentially outperform full fine-tuning. Our results also suggest many interesting future directions including 1) extension to deeper architectures by characterizing the role of softmax-attention in individual layers, 2) developing a stronger theoretical and empirical understanding of when/if prompt-tuning is superior to fine-tuning, 3) extending our model to include multiple prompt vectors (and perhaps extending (DATA) to include multiple context vectors), and 4) investigating the role of multi-head attention in prompt-tuning.

## Acknowledgement

We thank the reviewers for their feedback and suggesting additional experiments involving fine-tuning only first layer self-attention weights. SO was supported by the NSF grants CCF-2046816 and CCF-2212426, Google Research Scholar award, and Army Research Office grant W911NF2110312. MS is supported by the Packard Fellowship in Science and Engineering, a Sloan Fellowship in Mathematics, an NSF-CAREER under award #1846369, DARPA Learning with Less Labels (LwLL) and FastNICS programs, and NSF-CIF awards #1813877 and #2008443. CT was partially supported by the NSERC Discovery Grant RGPIN-2021-03677 and by the NSF grant CCF-2009030.

## References

- <https://openai.com/blog/chatgpt/>.
- Bahdanau, D., Cho, K., and Bengio, Y. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- Baldi, P. and Vershynin, R. The quarks of attention. *arXiv preprint arXiv:2202.08371*, 2022.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.
- Cao, Y., Chen, Z., Belkin, M., and Gu, Q. Benign overfitting in two-layer convolutional neural networks. *arXiv preprint arXiv:2202.06526*, 2022.
- Chan, W., Jaitly, N., Le, Q. V., and Vinyals, O. Listen, attend and spell. *arXiv preprint arXiv:1508.01211*, 2015.
- Dehghani, M., Gritsenko, A., Arnab, A., Minderer, M., and Tay, Y. Scenic: A jax library for computer vision research and beyond. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 21393–21398, 2022.
- Dong, Y., Cordonnier, J.-B., and Loukas, A. Attention is not all you need: Pure attention loses rank doubly exponentially with depth. In *International Conference on Machine Learning*, pp. 2793–2803. PMLR, 2021.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houshy, N. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- Edelman, B. L., Goel, S., Kakade, S., and Zhang, C. Inductive biases and variable creation in self-attention mechanisms. *arXiv preprint arXiv:2110.10090*, 2021.
- Ergen, T., Neyshabur, B., and Mehta, H. Convexifying transformers: Improving optimization and understanding of transformer networks. *arXiv preprint arXiv:2211.11052*, 2022.
- Frei, S., Chatterji, N. S., and Bartlett, P. L. Random feature amplification: Feature learning and generalization in neural networks. *arXiv preprint arXiv:2202.07626*, 2022.
- Jelassi, S., Sander, M. E., and Li, Y. Vision transformers provably learn spatial structure. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=eMW9AkXaREI>.
- Karp, S., Winston, E., Li, Y., and Singh, A. Local signal adaptivity: Provable feature learning in neural networks beyond kernels. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 24883–24897. Curran Associates, Inc., 2021. URL <https://proceedings.neurips.cc/paper/2021/file/d064bflad039fff366564f352226e7640-Paper.pdf>.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- Ledoux, M. and Talagrand, M. *Probability in Banach Spaces: isoperimetry and processes*, volume 23. Springer Science & Business Media, 1991.
- Lester, B., Al-Rfou, R., and Constant, N. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021.
- Li, H., Wang, M., Liu, S., and Chen, P.-Y. A theoretical understanding of shallow vision transformers: Learning, generalization, and sample complexity. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=jC1Gv3Qjhb>.
- Li, X. L. and Liang, P. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021.
- Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., and Neubig, G. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35, 2023.
- Mohammadi, H., Zare, A., Soltanolkotabi, M., and Jovanović, M. R. Convergence and sample complexity of gradient methods for the model-free linear quadratic regulator problem. *arXiv preprint arXiv:1912.11899*, 2019.
- Narayanan, D., Shoeybi, M., Casper, J., LeGresley, P., Patwary, M., Korthikanti, V., Vainbrand, D., Kashinkunti, P., Bernauer, J., Catanzaro, B., et al. Efficient large-scale language model training on gpu clusters using megatron-lm. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pp. 1–15, 2021.

- Oymak, S. Stochastic gradient descent learns state equations with nonlinear activations. In *conference on Learning Theory*, pp. 2551–2579. PMLR, 2019.
- Pollard, D. Empirical processes: theory and applications. Ims, 1990.
- Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., and Chen, M. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- Reed, S., Zolna, K., Parisotto, E., Colmenarejo, S. G., Novikov, A., Barth-Maron, G., Gimenez, M., Sulsky, Y., Kay, J., Springenberg, J. T., et al. A generalist agent. *arXiv preprint arXiv:2205.06175*, 2022.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., and Fei-Fei, L. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E., Ghasemipour, S. K. S., Ayan, B. K., Mahdavi, S. S., Lopes, R. G., et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022.
- Sahiner, A., Ergen, T., Ozturkler, B., Pauly, J., Mardani, M., and Pilanci, M. Unraveling attention via convex duality: Analysis and interpretations of vision transformers. *International Conference on Machine Learning*, 2022.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., Zemel, R., and Bengio, Y. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pp. 2048–2057. PMLR, 2015.
- Yun, C., Bhojanapalli, S., Rawat, A. S., Reddi, S., and Kumar, S. Are transformers universal approximators of sequence-to-sequence functions? In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=ByxRM0Ntvr>.

## A. Sharp characterization of accuracy under known context

While the discrete dataset model is insightful, incorporating noise is crucial for understanding the fundamental limits of the benefits of context in attention. To this end, let us focus on the core dataset model where we set  $\delta^q = \delta^w = 0$  and explore the role of noise in population accuracy. Also assume that noise is standard normal, i.e. Assumption 1.b.

- *Linear model.* The linear model aggregates tokens to obtain a simple Gaussian mixture distribution. Specifically, aggregated tokens are exactly distributed as  $\frac{1}{T} \mathbf{X}^\top \mathbf{1}_T \sim \mathcal{N}(\zeta \mathbf{w}_*, \frac{1-\zeta}{T} \mathbf{I})$ , leading to the following well-known result.

**Fact A.1.** *For linear models, optimal accuracy obeys  $\min_{\mathbf{w}} \text{ERR}(f^{\text{LIN}}(\mathbf{w})) = Q(\sqrt{\frac{\zeta^2 W^2 T}{1-\zeta}})$  where  $Q(\cdot)$  is the tail function of the standard normal distribution.*

- *Prompt-attention model.* Since prompt-attention strictly generalizes the linear model, its accuracy is at least as good. The theorem below quantifies this and demonstrates how context vector can enable an optimal weighting of relevant and irrelevant tokens to maximize accuracy. A general version of this theorem is proven under a non-asymptotic setting (finite  $T, d$ ) as Theorem 8.

**Theorem 4.** *Consider the prompt-attention model  $f_{\theta}^{\text{ATT}}$ . Suppose  $\mathbf{w}_* \perp \mathbf{q}_*$  and let  $\tau, \bar{\tau} > 0$  be hyperparameters. Consider the following algorithm which uses the hindsight knowledge of  $\mathbf{q}_*$  to estimate  $\mathbf{w}_*$  and make prediction:*

Set  $\hat{\mathbf{w}} = (\mathbf{I} - \bar{\mathbf{q}}_* \bar{\mathbf{q}}_*^\top) \nabla \mathcal{L}_{\mathbf{w}}(0, \tau \bar{\mathbf{q}}_*)$  and  $\theta = (\hat{\mathbf{w}}, \bar{\tau} \bar{\mathbf{q}}_*)$ . Suppose  $\zeta^2 W^2 T, 1 - \zeta, \alpha := n/d, e^{Q^2}, e^\tau$  each lie between two positive absolute constants. Suppose  $T$  is polynomially large in  $n$  and these constants and  $\tilde{\mathcal{O}}(\cdot)$  hides polynomial terms in  $n$ . Define inverse-signal-to-noise-ratio:  $\text{ISNR}(\alpha, \tau) = \frac{(1-\zeta)e^{2\tau(\tau-Q)}}{\alpha \zeta^2 W^2 T}$ . In the limit  $T, d \rightarrow \infty$ , the test error converges in probability to  $\mathcal{Q}\left(\frac{e^{Q\bar{\tau}-\bar{\tau}^2}}{\sqrt{1+\text{ISNR}(\alpha, \tau)}} \cdot \sqrt{\frac{\zeta^2 W^2 T}{1-\zeta}}\right)$ . In this limit, optimal hyperparameters are  $\tau = \bar{\tau} = Q/2$  and leads to optimal  $\text{ISNR}(\alpha) := \frac{(1-\zeta)e^{-Q^2/2}}{\alpha \zeta^2 W^2 T}$  and the error

$$\text{ERR}(\alpha, Q, W) = \mathcal{Q}\left(\frac{e^{Q^2/4}}{\sqrt{1+\text{ISNR}(\alpha)}} \cdot \sqrt{\frac{\zeta^2 W^2 T}{1-\zeta}}\right)$$

Here, a few remarks are in place. Note that  $\text{rate}_{\text{LIN}} := \sqrt{\frac{\zeta^2 W^2 T}{1-\zeta}}$  term is the population error rate of  $f^{\text{LIN}}$  from Fact A.1. In the limit  $\alpha = n/d \rightarrow \infty$ , the rate of  $f^{\text{ATT}}$  is simply given by  $e^{Q/4} \text{rate}_{\text{LIN}}$  demonstrating the strict superiority of prompt-attention. Moreover setting  $Q = 0$  (no prompt info), since feature-output of  $f^{\text{LIN}}$  (i.e.  $\mathbf{X}^\top \phi(\mathbf{X} \mathbf{1})$ ) is (essentially) a binary Gaussian mixture distribution, our error-rate recovers the Bayes-optimal  $f^{\text{LIN}}$  classifier which has a finite-sample rate of  $\text{rate}_{\text{LIN}} / \sqrt{1 + (1-\zeta)/(\alpha \zeta^2 W^2 T^2)}$ . Prompt-tuning also strictly improves this because our  $\text{ISNR}(\alpha)$  introduces an additional  $e^{-Q^2/2}$  multiplier.

## B. Gradient calculations and concentration

In this section, we focus on finite-sample analysis of Algorithm 9. Introduce the following shorthand notation analogous to the population counterparts in Section 4.2:

$$\begin{aligned} \widehat{\mathbf{G}}_q(\mathbf{q}, \mathbf{w}) &:= -\nabla_q \widehat{\mathcal{L}}_{\mathcal{D}}(\theta) = \frac{1}{n} \sum_{i \in [n]} (y_i - f_{\theta}(\mathbf{X}_i)) \mathbf{X}_i^\top \phi'(\mathbf{X}_i \mathbf{q}) \mathbf{X}_i \mathbf{w}, \\ \widehat{\mathbf{G}}_w(\mathbf{q}, \mathbf{w}) &:= -\nabla_w \widehat{\mathcal{L}}_{\mathcal{D}}(\theta) = \frac{1}{n} \sum_{i \in [n]} (y_i - f_{\theta}(\mathbf{X}_i)) \mathbf{X}_i^\top \phi(\mathbf{X}_i \mathbf{q}). \end{aligned} \quad (14)$$

### B.1. Gradient Calculations

We begin with the gradient calculations for the first two steps of the algorithm.

For convenience, we make use of the following shorthands

$$\begin{aligned} R_{\mathbf{q}_*} &:= R_{\mathbf{w}, \mathbf{q}_*} := \mathbf{w}^\top \mathbf{q}_*, & R_{\mathbf{w}_*} &:= R_{\mathbf{w}, \mathbf{w}_*} := \mathbf{w}^\top \mathbf{w}_*, & \alpha_i &:= \alpha(\mathbf{w}, y_i) := R_{\mathbf{q}_*} + y_i R_{\mathbf{w}_*}, \\ \gamma_i &:= \gamma(\mathbf{Z}_i) = \frac{1}{T} \mathbf{Z}_i^\top \mathbf{1}, & \beta_i &:= \beta(\mathbf{Z}_i; \mathbf{w}) = \gamma_i^\top \mathbf{w}, & \hat{\Sigma}_i &:= \frac{1}{T} \mathbf{Z}_i^\top \mathbf{Z}_i, \end{aligned}$$

where  $\mathbf{Z}_i \in \mathbb{R}^{(1-\zeta)T \times d}$  is the matrix of irrelevant tokens  $\mathbf{z}_{i,t}, t \in \mathcal{R}^c$  for sample  $i \in [n]$ .

**Lemma 3.** *Under dataset model (DATA) and Assumption 1.a, we have*

$$\widehat{\mathbf{G}}_{\mathbf{w}}(0, 0) = \zeta \mathbf{w}_* + \zeta \mathbf{q}_* \left( \frac{1}{n} \sum_{i \in [n]} y_i \right) + \frac{1}{n} \sum_i y_i \gamma_i. \quad (16)$$

**Lemma 4.** *Under dataset model (DATA) and Assumption 1.a, we have that*

$$\widehat{\mathbf{G}}_{\mathbf{q}}(0, \mathbf{w}) := \frac{1}{n} \sum_{i=1}^n (y_i - \zeta \alpha_i - \beta_i) \left[ (\zeta - \zeta^2) \alpha_i - \zeta \beta_i \right] (\mathbf{q}_* + y_i \mathbf{w}_*) + \widehat{\Sigma}_i \mathbf{w} - (\zeta \alpha_i + \beta_i) \gamma_i. \quad (17)$$

### B.1.1. PROOF OF LEMMA 3

By direct computation,

$$\begin{aligned} \widehat{\mathbf{G}}_{\mathbf{w}}(0, 0) &:= -\nabla_{\mathbf{w}} \widehat{\mathcal{L}}_{\mathcal{D}}(\boldsymbol{\theta}) = \frac{1}{nT} \sum_{i \in [n]} y_i \mathbf{X}_i^\top \mathbf{1}_T = \frac{1}{nT} \sum_{i \in [n]} \sum_{t \in [T]} y_i \mathbf{x}_{i,t} \\ &= \frac{\zeta}{n} \left( \sum_{i \in [n]} y_i \right) \mathbf{q}_* + \frac{\zeta}{n} \mathbf{w}_* + \frac{1}{n} \sum_{i \in [n]} y_i \gamma_i. \end{aligned}$$

### B.1.2. PROOF OF LEMMA 4

Note that  $\phi'(0) = \frac{1}{T} \mathbb{1} - \frac{1}{T^2} \mathbb{1} \mathbb{1}^\top$ ; hence, for  $\boldsymbol{\theta} = (0, \mathbf{w})$ :

$$\widehat{\mathbf{G}}_{\mathbf{q}}(0, \mathbf{w}) := -\nabla_{\mathbf{q}} \widehat{\mathcal{L}}_{\mathcal{D}}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i \in [n]} \underbrace{\frac{1}{T} (y_i - f_{\boldsymbol{\theta}}(\mathbf{X}_i)) \mathbf{X}_i^\top \mathbf{X}_i \mathbf{w}}_{\text{Term}_{1,i}} - \frac{1}{n} \sum_{i \in [n]} \underbrace{\frac{1}{T^2} (y_i - f_{\boldsymbol{\theta}}(\mathbf{X}_i)) \mathbf{X}_i^\top \mathbb{1} \mathbb{1}^\top \mathbf{X}_i \mathbf{w}}_{\text{Term}_{2,i}}.$$

Moreover, note that,

$$\begin{aligned} f_{\boldsymbol{\theta}}(\mathbf{X}_i) &= \frac{1}{T} \mathbf{w}^\top \mathbf{X}_i^\top \mathbb{1}, \\ \mathbf{X}_i^\top \mathbf{X}_i &= \zeta T (\mathbf{q}_* + y_i \mathbf{w}_*) (\mathbf{q}_* + y_i \mathbf{w}_*)^\top + \mathbf{Z}_i^\top \mathbf{Z}_i, \\ \mathbf{X}_i^\top \mathbb{1} &= \zeta T (\mathbf{q}_* + y_i \mathbf{w}_*) + \mathbf{Z}_i^\top \mathbb{1}, \end{aligned}$$

where recall the notation in Lemma 3 for  $\mathbf{Z}_i$ . Hence, using the lemma's notation (repeated here for convenience)

$$\alpha_i := R_{\mathbf{q}_*} + y_i R_{\mathbf{w}_*}, \quad \beta_i := \beta(\mathbf{Z}_i; \mathbf{w}) := \frac{1}{T} \mathbb{1}^\top \mathbf{Z}_i \mathbf{w}, \quad \gamma_i := \gamma(\mathbf{Z}_i) = \frac{1}{T} \mathbf{Z}_i^\top \mathbb{1}, \quad \widehat{\Sigma}_i := \frac{1}{T} \mathbf{Z}_i^\top \mathbf{Z}_i.$$

we find that

$$\begin{aligned} y_i - f_{\boldsymbol{\theta}}(\mathbf{X}_i) &= y_i - \zeta \alpha_i - \beta_i \\ \frac{1}{T} \mathbf{X}_i \mathbf{X}_i^\top \mathbf{w} &= \zeta \alpha_i \mathbf{q}_* + \zeta y_i \alpha_i \mathbf{w}_* + \widehat{\Sigma}_i \mathbf{w} \\ \frac{1}{T^2} \mathbf{X}_i^\top \mathbb{1} \mathbb{1}^\top \mathbf{X}_i \mathbf{w} &= \zeta (\zeta \alpha_i + \beta_i) \mathbf{q}_* + \zeta (\zeta \alpha_i + \beta_i) y_i \mathbf{w}_* + (\zeta \alpha_i + \beta_i) \gamma_i. \end{aligned}$$

With the above, each one of the two terms becomes:

$$\text{Term}_{1,i} = (y_i - \zeta \alpha_i - \beta_i) \zeta \alpha_i \mathbf{q}_* + (y_i - \zeta \alpha_i - \beta_i) \zeta \alpha_i y_i \mathbf{w}_* + (y_i - \zeta \alpha_i - \beta_i) \widehat{\Sigma}_i \mathbf{w}$$

$$\text{Term}_{2,i} = \zeta (y_i - \zeta \alpha_i - \beta_i) (\zeta \alpha_i + \beta_i) \mathbf{q}_* + \zeta (y_i - \zeta \alpha_i - \beta_i) (\zeta \alpha_i + \beta_i) y_i \mathbf{w}_* + (y_i - \zeta \alpha_i - \beta_i) (\zeta \alpha_i + \beta_i) \gamma_i.$$

Combining the above:

$$\text{Term}_{1,i} - \text{Term}_{2,i} = (y_i - \zeta \alpha_i - \beta_i) \left[ \zeta ((1 - \zeta) \alpha_i - \beta_i) (\mathbf{q}_* + y_i \mathbf{w}_*) + \widehat{\Sigma}_i \mathbf{w} - (\zeta \alpha_i + \beta_i) \gamma_i \right]. \quad (18)$$

## B.2. Concentration of Gradient $\widehat{\mathbf{G}}_q(0, \mathbf{w})$ in the $q$ direction

The main result of this section is the following lemma about concentration of gradient with respect to  $q$ .

**Lemma 5** (Concentration of  $\widehat{\mathbf{G}}_q(0, \mathbf{w})$ ). *Fix any vectors  $\mathbf{v}, \mathbf{w} \in \mathbb{R}^d$ . For convenience define  $R_{v,q_*} := \mathbf{v}^\top \mathbf{q}_*$  and  $R_{v,w_*} := \mathbf{v}^\top \mathbf{w}_*$  and recall  $R_{w_*}, R_{q_*}$  notations from Lemma 4. Then, we can decompose*

$$\mathbf{v}^\top \widehat{\mathbf{G}}_q(0, \mathbf{w}) = \mathbf{v}^\top \mathbf{G}_q(0, \mathbf{w}) + \mathbf{v}^\top \widetilde{\mathbf{G}}_q(0, \mathbf{w}),$$

where the expectation term is given by

$$\begin{aligned} \mathbf{v}^\top \mathbf{G}_q(0, \mathbf{w}) := \mathbb{E}[\mathbf{v}^\top \widehat{\mathbf{G}}_q(0, \mathbf{w})] &= ((\zeta - \zeta^2)(R_{w_*} + \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w}/T) - (\zeta^2 - \zeta^3)(R_{w_*}^2 + R_{q_*}^2)) R_{v,q_*} \\ &\quad + (((\zeta - \zeta^2) - 2(\zeta^2 - \zeta^3)R_{w_*}) R_{q_*}) R_{v,w_*} - ((1 + 2/T)(\zeta - \zeta^2)R_{q_*}) \mathbf{v}^\top \boldsymbol{\Sigma} \mathbf{w}, \end{aligned} \quad (19)$$

and the deviation term obeys

$$\begin{aligned} \mathbf{v}^\top \widetilde{\mathbf{G}}_q(0, \mathbf{w}) &= [((\zeta - \zeta^2) - 2(\zeta^2 - \zeta^3)R_{w_*}) R_{q_*} R_{v,q_*} + ((\zeta - \zeta^2)R_{w_*} - (\zeta^2 - \zeta^3)(R_{w_*}^2 + R_{q_*}^2)) R_{v,w_*}] \left( \frac{1}{n} \sum_{i=1}^n y_i \right) \\ &\quad + [(-\zeta + 2\zeta^2) R_{q_*} R_{v,q_*} + (-\zeta + (-\zeta + 2\zeta^2)R_{w_*}) R_{v,w_*} + (1 - \zeta)\mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{v}] \left( \frac{1}{n} \sum_{i=1}^n (\boldsymbol{\gamma}_i^\top \mathbf{w}) \right) \\ &\quad + \left[ \zeta R_{w_*} - \zeta^2(R_{q_*}^2 + R_{w_*}^2) + \frac{(1 - \zeta)}{T} \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w} \right] \left( \frac{1}{n} \sum_{i=1}^n \mathbf{v}^\top \boldsymbol{\gamma}_i \right) \\ &\quad + [(-\zeta + (-\zeta + 2\zeta^2)R_{w_*}) R_{v,q_*} + (-\zeta + 2\zeta^2) R_{q_*} R_{v,w_*}] \left( \frac{1}{n} \sum_{i=1}^n y_i (\boldsymbol{\gamma}_i^\top \mathbf{w}) \right) \\ &\quad + [\zeta R_{q_*} - 2\zeta^2 R_{q_*} R_{w_*}] \left( \frac{1}{n} \sum_{i=1}^n y_i (\mathbf{v}^\top \boldsymbol{\gamma}_i) \right) \\ &\quad + \zeta R_{v,q_*} \left( \frac{1}{n} \sum_{i=1}^n (\boldsymbol{\gamma}_i^\top \mathbf{w})^2 - \frac{(1 - \zeta)}{T} \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w} \right) \\ &\quad + \zeta R_{v,w_*} \left( \frac{1}{n} \sum_{i=1}^n (\boldsymbol{\gamma}_i^\top \mathbf{w})^2 y_i \right) \\ &\quad + [1 - 2\zeta R_{w_*}] \left( \frac{1}{n} \sum_{i=1}^n y_i (\mathbf{w}^\top \boldsymbol{\gamma}_i) (\mathbf{v}^\top \boldsymbol{\gamma}_i) \right) \\ &\quad + (1 - \zeta R_{w_*}) \left( \frac{1}{n} \sum_{i=1}^n y_i \mathbf{v}^\top \widehat{\boldsymbol{\Sigma}}_i \mathbf{w} \right) \\ &\quad - \zeta R_{q_*} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{v}^\top \widehat{\boldsymbol{\Sigma}}_i \mathbf{w} - (1 - \zeta)\mathbf{v}^\top \boldsymbol{\Sigma} \mathbf{w} \right) \\ &\quad - [2\zeta R_{q_*}] \left( \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \boldsymbol{\gamma}_i) (\mathbf{v}^\top \boldsymbol{\gamma}_i) - \frac{1 - \zeta}{T} \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w} \right) \\ &\quad - \left( \frac{1}{n} \sum_{i=1}^n (\mathbf{v}^\top \boldsymbol{\gamma}_i) \left( (\mathbf{w}^\top \boldsymbol{\gamma}_i)^2 - \frac{(1 - \zeta)}{T} \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w} \right) \right) \\ &\quad - \left( \frac{1}{n} \sum_{i=1}^n (\boldsymbol{\gamma}_i^\top \mathbf{w}) (\mathbf{v}^\top \widehat{\boldsymbol{\Sigma}}_i \mathbf{w} - (1 - \zeta)\mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{v}) \right). \end{aligned} \quad (20)$$

Moreover, all random terms in (20) are zero-mean and concentrate as prescribed by Lemma 6 below.

**Lemma 6** (Main concentration lemma). *Let  $y_i, i \in [n]$  be iid Rademacher random variables. Let  $\mathbf{Z}_i \in \mathbb{R}^{(1-\zeta)T \times d}, i \in [n]$  be iid copies of a random matrix  $\mathbf{Z}$ . Each row  $\mathbf{z}_t, t \in [(1-\zeta)T]$  of  $\mathbf{Z}$  is an iid copy of a random vector  $\mathbf{z}$  satisfying Assumption 1.a. For convenience denote  $\boldsymbol{\gamma}_i := \mathbf{Z}_i^\top \mathbf{1}/T$  and  $\widehat{\boldsymbol{\Sigma}}_i := \mathbf{Z}_i^\top \mathbf{Z}_i/T$ . Then, the following statements are true for all vectors  $\mathbf{w}, \mathbf{v} \in \mathbb{R}^d$ :*

$$\left\| \frac{1}{n} \sum_{i \in [n]} y_i \right\|_{\psi_2} \leq \frac{C}{\sqrt{n}}$$

$$\begin{aligned}
 & \left\| \frac{1}{n} \sum_{i \in [n]} \gamma_i^\top \mathbf{w} \right\|_{\psi_2} \vee \left\| \frac{1}{n} \sum_{i \in [n]} y_i \gamma_i^\top \mathbf{w} \right\|_{\psi_2} \leq \frac{C\sigma\sqrt{1-\zeta}\|\mathbf{w}\|_2}{\sqrt{nT}} \\
 & \left\| \frac{1}{n} \sum_{i \in [n]} (\gamma_i^\top \mathbf{w})(\gamma_i^\top \mathbf{v}) - \frac{1-\zeta}{T} \mathbf{v}^\top \Sigma \mathbf{w} \right\|_{\psi_1} \vee \left\| \frac{1}{n} \sum_{i \in [n]} y_i (\gamma_i^\top \mathbf{w})(\gamma_i^\top \mathbf{v}) \right\|_{\psi_1} \leq \frac{C\sigma^2(1-\zeta)\|\mathbf{w}\|_2\|\mathbf{v}\|_2}{T\sqrt{n}} \\
 & \left\| \frac{1}{n} \sum_{i \in [n]} \mathbf{w}^\top \hat{\Sigma}_i \mathbf{v} - (1-\zeta)\mathbf{w}^\top \Sigma \mathbf{v} \right\|_{\psi_1} \vee \left\| \frac{1}{n} \sum_{i \in [n]} y_i \mathbf{w}^\top \hat{\Sigma}_i \mathbf{v} \right\|_{\psi_1} \leq \frac{C\sigma^2\sqrt{1-\zeta}\|\mathbf{w}\|_2\|\mathbf{v}\|_2}{\sqrt{nT}} \\
 & \left\| \frac{1}{n} \sum_{i \in [n]} \left( (\gamma_i^\top \mathbf{w})^2 - \frac{1-\zeta}{T} \mathbf{w}^\top \Sigma \mathbf{w} \right) \gamma_i^\top \mathbf{v} \right\|_{\psi_{2/3}} \leq \frac{C\sigma^3(1-\zeta)^{3/2}\|\mathbf{w}\|_2^2\|\mathbf{v}\|_2}{T^{3/2}\sqrt{n}} \log n \\
 & \left\| \frac{1}{n} \sum_{i \in [n]} (\gamma_i^\top \mathbf{w})(\mathbf{w}^\top \hat{\Sigma}_i \mathbf{v} - (1-\zeta)\mathbf{w}^\top \Sigma \mathbf{v}) \right\|_{\psi_{2/3}} \leq \frac{C\sigma^3(1-\zeta)\|\mathbf{w}\|_2^2\|\mathbf{v}\|_2}{T\sqrt{n}} \log n.
 \end{aligned}$$

Also, all the random variables that appear above are zero mean.

### B.2.1. PROOF OF LEMMA 5

We split (17) in four terms and handle each of them separately.

- **Term<sub>I</sub>** =  $\frac{1}{n} \sum_{i=1}^n (y_i - \zeta\alpha_i - \beta_i) ((\zeta - \zeta^2)\alpha_i - \zeta\beta_i) \mathbf{q}_*$

We first focus on

$$\text{Term}_I = \frac{1}{n} \sum_{i=1}^n (y_i(1 - \zeta R_{\mathbf{w}_*}) - \beta_i - \zeta R_{\mathbf{q}_*}) (y_i(\zeta - \zeta^2) R_{\mathbf{w}_*} - \zeta\beta_i + (\zeta - \zeta^2) R_{\mathbf{q}_*}) \mathbf{q}_* =: A \mathbf{q}_*.$$

We can express  $A$  above conveniently as follows (recall  $y_i^2 = 1$ ):

$$\begin{aligned}
 A & := -\zeta(\zeta - \zeta^2) R_{\mathbf{q}_*}^2 + (1 - \zeta R_{\mathbf{w}_*})(\zeta - \zeta^2) R_{\mathbf{w}_*} + ((1 - \zeta R_{\mathbf{w}_*})(\zeta - \zeta^2) R_{\mathbf{q}_*} - \zeta(\zeta - \zeta^2) R_{\mathbf{w}_*} R_{\mathbf{q}_*}) \left( \frac{1}{n} \sum_{i=1}^n y_i \right) \\
 & \quad + (-\zeta(\zeta - \zeta^2) R_{\mathbf{q}_*} + \zeta^2 R_{\mathbf{q}_*}) \left( \frac{1}{n} \sum_{i=1}^n \beta_i \right) + (-(1 - \zeta R_{\mathbf{w}_*})\zeta - (\zeta - \zeta^2) R_{\mathbf{w}_*}) \left( \frac{1}{n} \sum_{i=1}^n y_i \beta_i \right) + \zeta \left( \frac{1}{n} \sum_{i=1}^n \beta_i^2 \right) \\
 & = (\zeta - \zeta^2) R_{\mathbf{w}_*} - (\zeta^2 - \zeta^3) (R_{\mathbf{w}_*}^2 + R_{\mathbf{q}_*}^2) + ((\zeta - \zeta^2) - 2(\zeta^2 - \zeta^3) R_{\mathbf{w}_*}) R_{\mathbf{q}_*} \left( \frac{1}{n} \sum_{i=1}^n y_i \right) \\
 & \quad + (-\zeta + 2\zeta^2) R_{\mathbf{q}_*} \left( \frac{1}{n} \sum_{i=1}^n \beta_i \right) + (-\zeta + (-\zeta + 2\zeta^2) R_{\mathbf{w}_*}) \left( \frac{1}{n} \sum_{i=1}^n y_i \beta_i \right) \\
 & \quad + \zeta \left( \frac{1}{n} \sum_{i=1}^n \beta_i^2 - \frac{(1-\zeta)}{T} \mathbf{w}^\top \Sigma \mathbf{w} \right) + \frac{(\zeta - \zeta^2)}{T} \mathbf{w}^\top \Sigma \mathbf{w}.
 \end{aligned}$$

From Lemma 6, all random terms above are zero mean. Hence,

$$\begin{aligned}
 \mathbb{E}[A] & = -(\zeta^2 - \zeta^3) R_{\mathbf{q}_*}^2 + (1 - \zeta R_{\mathbf{w}_*})(\zeta - \zeta^2) R_{\mathbf{w}_*} + (\zeta - \zeta^2) \frac{\mathbf{w}^\top \Sigma \mathbf{w}}{T} \\
 & = -(\zeta^2 - \zeta^3) R_{\mathbf{q}_*}^2 + (\zeta - \zeta^2) (R_{\mathbf{w}_*} + \mathbf{w}^\top \Sigma \mathbf{w} / T) - (\zeta^2 - \zeta^3) R_{\mathbf{w}_*}^2 \\
 & = (\zeta - \zeta^2) (R_{\mathbf{w}_*} + \mathbf{w}^\top \Sigma \mathbf{w} / T) - (\zeta^2 - \zeta^3) (R_{\mathbf{w}_*}^2 + R_{\mathbf{q}_*}^2). \tag{21}
 \end{aligned}$$

- $\mathbf{Term}_{II} = \frac{1}{n} \sum_{i=1}^n (y_i - \zeta \alpha_i - \beta_i) ((\zeta - \zeta^2) \alpha_i - \zeta \beta_i) y_i \mathbf{w}_*$

$$\mathbf{Term}_{II} = \frac{1}{n} \sum_{i=1}^n (y_i (1 - \zeta R_{\mathbf{w}_*}) - \beta_i - \zeta R_{\mathbf{q}_*}) (y_i (\zeta - \zeta^2) R_{\mathbf{w}_*} - \zeta \beta_i + (\zeta - \zeta^2) R_{\mathbf{q}_*}) y_i \mathbf{w}_* = B \mathbf{w}_*.$$

We can express  $B$  above conveniently as:

$$\begin{aligned} B := & ((\zeta - \zeta^2) R_{\mathbf{w}_*} - (\zeta^2 - \zeta^3) (R_{\mathbf{w}_*}^2 + R_{\mathbf{q}_*}^2)) \left( \frac{1}{n} \sum_{i \in [n]} y_i \right) + ((\zeta - \zeta^2) - 2(\zeta^2 - \zeta^3) R_{\mathbf{w}_*}) R_{\mathbf{q}_*} \\ & + (-\zeta + 2\zeta^2) R_{\mathbf{q}_*} \left( \frac{1}{n} \sum_{i=1}^n \beta_i y_i \right) + (-\zeta + (-\zeta + 2\zeta^2) R_{\mathbf{w}_*}) \left( \frac{1}{n} \sum_{i=1}^n \beta_i \right) + \zeta \left( \frac{1}{n} \sum_{i=1}^n \beta_i^2 y_i \right). \end{aligned}$$

All the random terms above are zero-mean. Hence,

$$\mathbb{E}[B] = ((\zeta - \zeta^2) - 2(\zeta^2 - \zeta^3) R_{\mathbf{w}_*}) R_{\mathbf{q}_*}. \quad (22)$$

- $\mathbf{Term}_{III} = \frac{1}{n} \sum_{i=1}^n (y_i - \zeta \alpha_i - \beta_i) \hat{\Sigma}_i \mathbf{w}$

Fix any vector  $\mathbf{v}$ :

$$\begin{aligned} \mathbf{v}^\top \{\mathbf{Term}_{III}\} &= \left( \frac{1}{n} \sum_{i=1}^n y_i \mathbf{v}^\top \hat{\Sigma}_i \mathbf{w} \right) - \zeta (R_{\mathbf{q}_*} + y_i R_{\mathbf{w}_*}) \left( \frac{1}{n} \sum_{i=1}^n \mathbf{v}^\top \hat{\Sigma}_i \mathbf{w} \right) - \left( \frac{1}{n} \sum_{i=1}^n (\gamma_i^\top \mathbf{w}) \mathbf{v}^\top \hat{\Sigma}_i \mathbf{w} \right) \\ &= (1 - \zeta R_{\mathbf{w}_*}) \left( \frac{1}{n} \sum_{i=1}^n y_i \mathbf{v}^\top \hat{\Sigma}_i \mathbf{w} \right) - \zeta R_{\mathbf{q}_*} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{v}^\top \hat{\Sigma}_i \mathbf{w} \right) - \left( \frac{1}{n} \sum_{i=1}^n (\gamma_i^\top \mathbf{w}) \mathbf{v}^\top \hat{\Sigma}_i \mathbf{w} \right) \\ &= (1 - \zeta R_{\mathbf{w}_*}) \left( \frac{1}{n} \sum_{i=1}^n y_i \mathbf{v}^\top \hat{\Sigma}_i \mathbf{w} \right) - \zeta R_{\mathbf{q}_*} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{v}^\top \hat{\Sigma}_i \mathbf{w} - (1 - \zeta) \mathbf{v}^\top \Sigma \mathbf{w} \right) - (\zeta - \zeta^2) R_{\mathbf{q}_*} \mathbf{v}^\top \Sigma \mathbf{w} \\ &\quad - \frac{1}{n} \sum_{i=1}^n (\gamma_i^\top \mathbf{w}) (\mathbf{v}^\top \hat{\Sigma}_i \mathbf{w} - (1 - \zeta) \mathbf{w}^\top \Sigma \mathbf{v}) + (1 - \zeta) (\mathbf{w}^\top \Sigma \mathbf{v}) \frac{1}{n} \sum_{i=1}^n (\gamma_i^\top \mathbf{w}) \end{aligned}$$

From Lemma 6 all random terms above are zero mean. Hence,

$$\mathbb{E}[\mathbf{v}^\top \{\mathbf{Term}_{III}\}] = -(\zeta - \zeta^2) R_{\mathbf{q}_*} \mathbf{w}^\top \Sigma \mathbf{v}. \quad (23)$$

- $\mathbf{Term}_{IV} = \frac{1}{n} \sum_{i=1}^n (y_i - \zeta \alpha_i - \beta_i) (\zeta \alpha_i + \beta_i) \gamma_i$

For fixed vector  $\mathbf{v}$ ,  $\mathbf{v}^\top \{\mathbf{Term}_{IV}\} = \frac{1}{n} \sum_{i=1}^n (y_i - \zeta \alpha_i - \beta_i) (\zeta \alpha_i + \beta_i) \mathbf{v}^\top \gamma_i$ . Reorganizing, note that

$$\begin{aligned} (y_i - \zeta \alpha_i - \beta_i) (\zeta \alpha_i + \beta_i) &= \zeta y_i (R_{\mathbf{q}_*} + y_i R_{\mathbf{w}_*}) - \zeta^2 (R_{\mathbf{q}_*}^2 + R_{\mathbf{w}_*}^2 + 2y_i R_{\mathbf{q}_*} R_{\mathbf{w}_*}) - 2\zeta R_{\mathbf{q}_*} \beta_i - 2\zeta R_{\mathbf{w}_*} y_i \beta_i + y_i \beta_i - \beta_i^2 \\ &= (\zeta R_{\mathbf{q}_*} - 2\zeta^2 R_{\mathbf{q}_*} R_{\mathbf{w}_*}) y_i + (\zeta R_{\mathbf{w}_*} - \zeta^2 (R_{\mathbf{q}_*}^2 + R_{\mathbf{w}_*}^2)) - (2\zeta R_{\mathbf{q}_*}) \beta_i + (1 - 2\zeta R_{\mathbf{w}_*}) y_i \beta_i - \beta_i^2 \end{aligned}$$

Overall,

$$\begin{aligned} \mathbf{v}^\top \{\mathbf{Term}_{IV}\} &= (\zeta R_{\mathbf{w}_*} - \zeta^2 (R_{\mathbf{q}_*}^2 + R_{\mathbf{w}_*}^2)) \left( \frac{1}{n} \sum_{i=1}^n \mathbf{v}^\top \gamma_i \right) + (\zeta R_{\mathbf{q}_*} - 2\zeta^2 R_{\mathbf{q}_*} R_{\mathbf{w}_*}) \left( \frac{1}{n} \sum_{i=1}^n y_i (\mathbf{v}^\top \gamma_i) \right) \\ &\quad + (1 - 2\zeta R_{\mathbf{w}_*}) \left( \frac{1}{n} \sum_{i=1}^n y_i (\mathbf{w}^\top \gamma_i) (\mathbf{v}^\top \gamma_i) \right) \\ &\quad - (2\zeta R_{\mathbf{q}_*}) \left( \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \gamma_i) (\mathbf{v}^\top \gamma_i) - \frac{1 - \zeta}{T} \mathbf{v}^\top \Sigma \mathbf{w} \right) - (2\zeta R_{\mathbf{q}_*}) \frac{1 - \zeta}{T} \mathbf{v}^\top \Sigma \mathbf{w} \\ &\quad - \frac{1}{n} \sum_{i=1}^n (\mathbf{v}^\top \gamma_i) \left( (\mathbf{w}^\top \gamma_i)^2 - \frac{(1 - \zeta)}{T} \mathbf{w}^\top \Sigma \mathbf{w} \right) + \frac{(1 - \zeta)}{T} \mathbf{w}^\top \Sigma \mathbf{w} \left( \frac{1}{n} \sum_{i=1}^n (\mathbf{v}^\top \gamma_i) \right). \end{aligned}$$



According to Lemma 6 all random terms above are zero mean. Thus,

$$\mathbb{E}[\mathbf{v}^\top \{\text{Term}_{\text{IV}}\}] = -2\zeta R_{q_*} \frac{1-\zeta}{T} \mathbf{w}^\top \Sigma \mathbf{v} \quad (24)$$

• **Combined**

The desired identities (19) and (20) follow by combining all the terms above.

B.2.2. PROOF OF LEMMA 6

First bound: Obvious by boundedness (hence, sub-gaussianity) of  $y_i$  and Fact G.1.

Second bound: For convenience set  $\tilde{T} = (1-\zeta)T$  and assume wlog that  $\mathcal{R}^c = [\tilde{T}]$ . Recall that

$$\beta_i = \frac{1}{T} \sum_{t=1}^{\tilde{T}} \mathbf{z}_{i,t}^\top \mathbf{w} = \frac{1-\zeta}{\tilde{T}} \sum_{t=1}^{\tilde{T}} \mathbf{z}_{i,t}^\top \mathbf{w}.$$

Also for all  $t$ :  $\|\mathbf{z}_{i,t}^\top \mathbf{w}\|_{\psi_2} \leq K \|\mathbf{w}\|_2$ . Thus, from Fact G.1:

$$\|\beta_i\|_{\psi_2} \leq \frac{C\sigma(1-\zeta)\|\mathbf{w}\|_2}{\sqrt{\tilde{T}}} = \frac{C\sigma\sqrt{(1-\zeta)}\|\mathbf{w}\|_2}{\sqrt{T}}. \quad (25)$$

The bound then follows by applying Fact G.1 once more.

For the second term in this bound recall that  $y_i \in \{\pm 1\}$  and  $\beta_i = \sum_t \mathbf{z}_{i,t}^\top \mathbf{w}/T$ . Also, for all  $i \in [n]$ :  $y_i, \{\mathbf{z}_{i,t}\}_t$  are zero-mean and independent. Thus (i)  $\mathbb{E}[y_i \beta_i] = 0$ , and (ii)  $\{y_i \mathbf{z}_{i,t} \stackrel{D}{\sim} \mathbf{z}_{i,t} \text{ and } y_i \mathbf{z}_{i,t} \perp y_i \mathbf{z}_{i,t'}\} \implies y_i \beta_i \stackrel{D}{\sim} \beta_i$ . Thus, the same bound as the first term holds.

Third bound: It is easy to compute

$$\mathbb{E}[(\gamma_i^\top \mathbf{w})(\gamma_i^\top \mathbf{v})] = \frac{1}{T^2} \sum_{t=1}^{\tilde{T}} \sum_{t'=1}^{\tilde{T}} \mathbb{E}[\mathbf{w}^\top \mathbf{z}_{i,t} \mathbf{z}_{i,t'}^\top \mathbf{v}] = \frac{1-\zeta}{T} \mathbf{v}^\top \Sigma \mathbf{w}, \quad (26)$$

and, using (25)

$$\|(\gamma_i^\top \mathbf{w})(\gamma_i^\top \mathbf{v}) - \mathbb{E}[(\gamma_i^\top \mathbf{w})(\gamma_i^\top \mathbf{v})]\|_{\psi_1} \leq C \|(\gamma_i^\top \mathbf{w})(\gamma_i^\top \mathbf{v})\|_{\psi_1} \leq C \|\gamma_i^\top \mathbf{w}\|_{\psi_2} \|\gamma_i^\top \mathbf{v}\|_{\psi_2} = \frac{C\sigma^2(1-\zeta)\|\mathbf{w}\|_2 \|\mathbf{v}\|_2}{T}. \quad (27)$$

Since  $(\gamma_i^\top \mathbf{w})(\gamma_i^\top \mathbf{v}), i \in [n]$  are independent, the desired bound on the first term follows from Fact G.1.

Consider now the second term. By independence of  $y_i, \gamma_i$  it holds that  $\mathbb{E}[y_i (\gamma_i^\top \mathbf{w})(\gamma_i^\top \mathbf{v})] = 0$ . Arguing as we did above for the second bound,  $y_i \gamma_i^\top \mathbf{w} \stackrel{D}{\sim} \gamma_i^\top \mathbf{w}$ . Hence, the subexponential bound is the same as for the first term.

Fourth bound: First, it is easy to compute that for all  $i \in [n]$ :

$$\mathbb{E}[\mathbf{w}^\top \Sigma_i \mathbf{v}] = \frac{1}{T} \mathbb{E}[\mathbf{w}^\top \mathbf{Z}_i^\top \mathbf{Z}_i \mathbf{v}] = \frac{1}{T} \sum_{t=1}^{\tilde{T}} \mathbb{E}[\mathbf{w}^\top \mathbf{z}_{i,t}^\top \mathbf{z}_{i,t} \mathbf{v}] = \frac{\tilde{T}}{T} \mathbf{w}^\top \Sigma \mathbf{v} = (1-\zeta) \mathbf{w}^\top \Sigma \mathbf{v}.$$

Thus,

$$\mathbf{w}^\top \hat{\Sigma}_i \mathbf{v} - \mathbb{E}[\mathbf{w}^\top \hat{\Sigma}_i \mathbf{v}] = \mathbf{w}^\top \hat{\Sigma}_i \mathbf{v} - (1-\zeta) \mathbf{w}^\top \Sigma \mathbf{v} = \frac{1}{T} \sum_{t=1}^{\tilde{T}} \left( (\mathbf{z}_{i,t}^\top \mathbf{w})(\mathbf{z}_{i,t}^\top \mathbf{v}) - \mathbf{w}^\top \Sigma \mathbf{v} \right)$$

and so

$$\frac{1}{n} \sum_{i=1}^n \mathbf{w}^\top \hat{\Sigma}_i \mathbf{v} - (1-\zeta) \mathbf{w}^\top \Sigma \mathbf{v} = \frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^{\tilde{T}} \left( (\mathbf{z}_{i,t}^\top \mathbf{w})(\mathbf{z}_{i,t}^\top \mathbf{v}) - \mathbf{w}^\top \Sigma \mathbf{v} \right).$$

Now, each random variable in the double sum above is independent and such that

$$\|(\mathbf{z}_{i,t}^\top \mathbf{w})(\mathbf{z}_{i,t}^\top \mathbf{v}) - \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{v}\|_{\psi_1} \leq 2\|(\mathbf{z}_{i,t}^\top \mathbf{w})(\mathbf{z}_{i,t}^\top \mathbf{v})\|_{\psi_1} \leq 2\|\mathbf{z}_{i,t}^\top \mathbf{w}\|_{\psi_2} \|\mathbf{z}_{i,t}^\top \mathbf{v}\|_{\psi_2} \leq C\sigma^2 \|\mathbf{w}\|_2 \|\mathbf{v}\|_2. \quad (28)$$

Hence, from Fact G.1,

$$\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{w}^\top \hat{\boldsymbol{\Sigma}}_i \mathbf{v} - (1-\zeta) \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{v} \right\|_{\psi_1} \leq \frac{C\sigma^2 \sqrt{1-\zeta} \|\mathbf{w}\|_2 \|\mathbf{v}\|_2}{\sqrt{nT}}.$$

The bound for the second term follows along the same lines. The two key observations are that (i)  $\mathbb{E}[y_i \mathbf{w}^\top \hat{\boldsymbol{\Sigma}}_i \mathbf{v}] = 0$  because  $y_i$  and  $\mathbf{z}_{i,t}$  are independent, and, (ii)

$$y_i \mathbf{w}^\top \hat{\boldsymbol{\Sigma}}_i \mathbf{v} = \frac{1}{T} \mathbf{w}^\top y_i \mathbf{Z}_i^\top \mathbf{Z}_i \mathbf{v} \stackrel{D}{\sim} \frac{1}{T} \mathbf{w}^\top \tilde{\mathbf{Z}}_i^\top \mathbf{Z}_i \mathbf{v} = \frac{1}{T} \sum_{t=1}^{\tilde{T}} (\tilde{\mathbf{z}}_{i,t}^\top \mathbf{w})(\mathbf{z}_{i,t}^\top \mathbf{v})$$

where  $\tilde{\mathbf{Z}}_i$  is an independent copy of  $\mathbf{Z}_i$ .

Fifth bound: From (29) and (26), we have for all  $i \in [n]$  that

$$\left\| (\boldsymbol{\gamma}_i^\top \mathbf{w})(\boldsymbol{\gamma}_i^\top \mathbf{v}) - \frac{1-\zeta}{T} \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w} \right\|_{\psi_1} \leq \frac{C\sigma^2(1-\zeta) \|\mathbf{w}\|_2^2}{T}.$$

Moreover, recall from Eq. (25) that

$$\|\boldsymbol{\gamma}_i^\top \mathbf{v}\|_{\psi_2} \leq \frac{C\sigma \sqrt{1-\zeta} \|\mathbf{v}\|_2}{\sqrt{T}}.$$

Combining the above two displays and applying Fact G.2 we find for all  $i \in [n]$  that

$$\left\| \left( (\boldsymbol{\gamma}_i^\top \mathbf{w})^2 - \frac{1-\zeta}{T} \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w} \right) \boldsymbol{\gamma}_i^\top \mathbf{v} \right\|_{\psi_{2/3}} \leq \frac{C\sigma^3(1-\zeta)^{3/2} \|\mathbf{w}\|_2^2 \|\mathbf{v}\|_2}{T^{3/2}}. \quad (29)$$

The desired bound follows from the above after using Fact G.3.

Sixth bound: From Eq. (28):

$$\|\mathbf{w}^\top \hat{\boldsymbol{\Sigma}}_i \mathbf{v} - (1-\zeta) \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{v}\|_{\psi_1} \leq \frac{C\sigma^2 \sqrt{1-\zeta} \|\mathbf{w}\|_2 \|\mathbf{v}\|_2}{\sqrt{T}}$$

and from Eq. (25)

$$\|\boldsymbol{\gamma}_i^\top \mathbf{w}\|_{\psi_2} \leq \frac{C\sigma \sqrt{1-\zeta} \|\mathbf{w}\|_2}{\sqrt{T}}.$$

Next we use Fact G.2 with  $\alpha = 2$  and  $\beta = 1$  to find that

$$\left\| (\boldsymbol{\gamma}_i^\top \mathbf{w})(\mathbf{w}^\top \hat{\boldsymbol{\Sigma}}_i \mathbf{v} - (1-\zeta) \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{v}) \right\|_{\psi_{2/3}} \leq \frac{CK^3(1-\zeta) \|\mathbf{w}\|_2^2 \|\mathbf{v}\|_2}{T}.$$

Next we use Fact G.3 which allows us to conclude that

$$\left\| \frac{1}{n} \sum_{i=1}^n (\boldsymbol{\gamma}_i^\top \mathbf{w})(\mathbf{w}^\top \hat{\boldsymbol{\Sigma}}_i \mathbf{v} - (1-\zeta) \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{v}) \right\|_{\psi_{2/3}} \leq \frac{CK^3(1-\zeta) \|\mathbf{w}\|_2^2 \|\mathbf{v}\|_2}{T\sqrt{n}} \log n.$$

Finally, the zero-mean property follows since

$$\begin{aligned} \mathbb{E}[(\mathbf{1}^\top \mathbf{Z}_i \mathbf{w})(\mathbf{w}^\top \mathbf{Z}_i \mathbf{Z}_i^\top \mathbf{v})] &= \sum_{t=1}^{\tilde{T}} \sum_{t'=1}^{\tilde{T}} \mathbb{E}[(\mathbf{z}_{i,t}^\top \mathbf{w})(\mathbf{w}^\top \mathbf{z}_{i,t'} \mathbf{z}_{i,t'}^\top \mathbf{v})] = \sum_{t=1}^{\tilde{T}} \sum_{t'=1}^{\tilde{T}} \mathbb{E}[\mathbf{w}^\top \mathbf{z}_{i,t} \operatorname{tr}(\mathbf{z}_{i,t'} \mathbf{z}_{i,t'}^\top \mathbf{v} \mathbf{w}^\top)] \\ &= \tilde{T}^2 \mathbb{E}[\operatorname{tr}(\mathbf{z}^\top \mathbf{w}) \operatorname{tr}(\mathbf{z} \mathbf{z}^\top \mathbf{v} \mathbf{w}^\top)] = \tilde{T}^2 \mathbb{E}[\operatorname{tr}((\mathbf{z}^\top \mathbf{w}) \otimes (\mathbf{z} \mathbf{z}^\top \mathbf{v} \mathbf{w}^\top))] \\ &= \tilde{T}^2 \operatorname{tr}(\mathbb{E}[(\mathbf{z}^\top \otimes \mathbf{z} \mathbf{z}^\top)](\mathbf{w} \otimes \mathbf{v} \mathbf{w}^\top)) = 0, \end{aligned}$$

where the last equality follows by the zero third moment property in Assumption 1.a.

## C. Finite-sample gradient analysis

In the following, we assume without explicit further reference that Ass. 3.a holds (i.e.  $|\rho| < W/Q$ ) and additionally that  $Q > W$ . To simplify the results we further assume  $\sigma \asymp 1$  and  $W \gtrsim 1$ .

### C.1. First gradient step

The lemma below studies the deviation of the first-step of GD  $\widehat{\mathbf{w}}_1$  with respect to its population counterpart  $\mathbf{w}_1$ . Provided that  $n$  and  $n\zeta T/d$  are larger than appropriate functions of other problem parameters, then the deviations are of small multiplicative order.

**Lemma 7** (First gradient step). *Consider the one-step population and finite updates  $\mathbf{w}_1 = \eta \mathbf{G}_w(0, 0)$  and  $\widehat{\mathbf{w}}_1 = \eta \widehat{\mathbf{G}}_w(0, 0)$ , respectively. For convenience denote*

$$R_{\mathbf{w}_*} = \mathbf{w}_1^\top \mathbf{w}_*, \quad R_{\mathbf{q}_*} = \mathbf{w}_1^\top \mathbf{q}_*, \quad \widehat{R}_{\mathbf{w}_*} = \widehat{\mathbf{w}}_1^\top \mathbf{w}_*, \quad \widehat{R}_{\mathbf{q}_*} = \widehat{\mathbf{w}}_1^\top \mathbf{q}_*.$$

For any  $u > 0$  and any small constant  $c_0 > 0$ , there exist absolute constants  $c, c' > 0$  and large enough constant  $C = C(c_0) > 0$  such that if

$$\sqrt{n} \geq Cu \frac{Q}{W} \quad \text{and} \quad \sqrt{n\zeta T} \geq Cu \frac{\sigma}{W} \sqrt{\zeta^{-1} - 1}, \quad (30)$$

then, with probability at least  $1 - c'e^{-cu^2}$

$$|\widehat{R}_{\mathbf{w}_*} - R_{\mathbf{w}_*}| \leq c_0 R_{\mathbf{w}_*} \quad \text{and} \quad |\widehat{R}_{\mathbf{q}_*} - R_{\mathbf{q}_*}| \leq c_0 \eta \zeta QW. \quad (31)$$

Additionally, if

$$\sqrt{n} \geq Cu \frac{Q}{W} \quad \text{and} \quad \sqrt{n\zeta T} \geq C(1+u) \frac{\sigma}{W} \sqrt{\zeta^{-1} - 1} \sqrt{d} \quad (32)$$

then with probability  $1 - e^{-cu^2} - e^{-cdu^2}$

$$\|\widehat{\mathbf{w}}_1\| - \|\mathbf{w}_1\| \leq c_0 \eta \zeta W. \quad (33)$$

*Proof.* Note that the conclusions of the lemma are all homogeneous in  $\eta$ . Hence, without loss of generality, set  $\eta = 1$ . Also recall by Lemma 3 that

$$\widehat{\mathbf{w}}_1 = \widehat{\mathbf{G}}_w(0, 0) = \zeta \mathbf{w}_* + \zeta \mathbf{q}_* \left( \frac{1}{n} \sum_{i \in [n]} y_i \right) + \frac{1}{n} \sum_i y_i \gamma_i, \quad (34)$$

and  $\mathbf{w}_1 = \mathbf{G}_w(0, 0) = \zeta \mathbf{w}_*$ ; thus,  $R_{\mathbf{w}_*} = \zeta W^2$ . From these, and also using Lemma 6, for any  $u > 0$  with probability at least  $1 - 2e^{-cu^2}$

$$\begin{aligned} |\widehat{R}_{\mathbf{w}_*} - R_{\mathbf{w}_*}| &= |\mathbf{w}_*^\top (\widehat{\mathbf{w}}_1 - \mathbf{w}_1)| \leq \zeta |\rho| WQ \left| \frac{1}{n} \sum_{i \in [n]} y_i \right| + \left| \frac{1}{n} \sum_i y_i \gamma_i^\top \mathbf{w}_* \right| \leq \frac{Cu\zeta |\rho| WQ}{\sqrt{n}} + \frac{Cu\sigma \sqrt{1-\zeta} \sqrt{\zeta} W}{\sqrt{n\zeta T}} \\ &\leq c_0 \zeta W^2 = c_0 R_{\mathbf{w}_*}. \end{aligned}$$

where the last inequality follows by assuming  $n, \zeta T$  large enough as in the condition of the lemma and using  $|\rho| \leq 1$ . Similarly,

$$|\widehat{R}_{\mathbf{q}_*} - R_{\mathbf{q}_*}| = |\mathbf{q}_*^\top (\widehat{\mathbf{w}}_1 - \mathbf{w}_1)| \leq \zeta Q^2 \left| \frac{1}{n} \sum_{i \in [n]} y_i \right| + \left| \frac{1}{n} \sum_i y_i \gamma_i^\top \mathbf{q}_* \right| \leq \frac{Cu\zeta Q^2}{\sqrt{n}} + \frac{Cu\sigma \sqrt{1-\zeta} \sqrt{\zeta} Q}{\sqrt{n\zeta T}} \leq c_0 \zeta QW$$

for sufficiently large  $n$  per (30). Finally, with probability at least  $1 - e^{-cu^2} - e^{-cdu^2}$

$$\|\widehat{\mathbf{w}}_1\| - \|\mathbf{w}_1\| \leq \|\widehat{\mathbf{w}}_1 - \mathbf{w}_1\| \leq \frac{Cu\zeta Q}{\sqrt{n}} + \frac{C(1+u)\sigma \sqrt{1-\zeta} \sqrt{\zeta} \sqrt{d}}{\sqrt{n\zeta T}} \leq c_0 \zeta W. \quad (35)$$

where, again, the last inequality follows by assuming  $n, \zeta T$  large enough as stated in the lemma. In the second inequality, we used from Lemma 6 that  $\sum_{i \in [n]} y_i \gamma_i / n$  is  $C\sigma \sqrt{1-\zeta} / \sqrt{nT}$ -subGaussian and applied Fact G.4 to get a high-probability bound on its euclidean norm.  $\square$

## C.2. Second gradient step

Next, we move on to the second gradient update in the direction of  $\widehat{\mathbf{G}}_q(0, \widehat{\mathbf{w}}_1)$ . Recall our goal of controlling the relevance scores of signal and noisy tokens. The first lemma below takes a step in this direction by computing the signal and noise relevance scores assuming access to the population gradient  $\mathbf{G}_q(0, \widehat{\mathbf{w}}_1) := \mathbb{E}[\widehat{\mathbf{G}}_q(0, \widehat{\mathbf{w}}_1)]$ .

**Lemma 8** ( $\widehat{\mathbf{G}}_q(0, \cdot)$  control: Expectation term). *Let  $\mathbf{G}_q(0, \mathbf{w}_1) = \mathbb{E}[\widehat{\mathbf{G}}_q(0, \mathbf{w}_1)]$  be the expectation of a gradient step in the  $\mathbf{q}$ -direction evaluated at  $(0, \widehat{\mathbf{w}}_1)$  and recall that  $\widehat{\mathbf{w}}_1 = \eta \widehat{\mathbf{G}}_w(0, 0)$  for  $\eta > 0$ . Suppose  $\widehat{\mathbf{w}}_1$  satisfies (31) and (33) for sufficiently small enough constant  $c_0 > 0$ . Further assume that the step-size  $\eta$  satisfies the following for sufficiently small absolute constant  $c_\eta > 0$ :*

$$\eta = \frac{c_\eta}{\sigma^2 Q^2}. \quad (36)$$

Then, for  $y \in \{\pm 1\}$  it holds that

$$(\mathbf{q}_* + y\mathbf{w}_*)^\top \mathbf{G}_q(0, \widehat{\mathbf{w}}_1) \geq \eta \zeta (\zeta - \zeta^2) W^2 Q \left( (1/4 - \rho^2/8 - 3c_0) Q - (9/4|\rho| + c_0) W \right), \quad (37)$$

and

$$\|\mathbf{G}_q(0, \widehat{\mathbf{w}}_1)\| \leq \eta \zeta (\zeta - \zeta^2) W^2 Q (13/4 + 2c_0). \quad (38)$$

*Proof.* Fix any  $\mathbf{v}$  and recall the notation of Lemma 7. With these, we have from Lemma 5 that

$$\begin{aligned} \mathbf{v}^\top \mathbf{G}_q(0, \widehat{\mathbf{w}}_1) &= ((\zeta - \zeta^2) (\widehat{R}_{\mathbf{w}_*} + \sigma^2 \|\widehat{\mathbf{w}}_1\|^2/T) - (\zeta^2 - \zeta^3) (\widehat{R}_{\mathbf{w}_*}^2 + \widehat{R}_{\mathbf{q}_*}^2)) \mathbf{v}^\top \mathbf{q}_* \\ &\quad + ((\zeta - \zeta^2) - 2(\zeta^2 - \zeta^3) \widehat{R}_{\mathbf{w}_*}) \widehat{R}_{\mathbf{q}_*} \mathbf{v}^\top \mathbf{w}_* - \sigma^2 (1 + 2/T) (\zeta - \zeta^2) \widehat{R}_{\mathbf{q}_*} \mathbf{v}^\top \widehat{\mathbf{w}}_1 \\ &= (\zeta - \zeta^2) \widehat{R}_{\mathbf{w}_*} (1 + \sigma^2 \|\widehat{\mathbf{w}}_1\|^2/(T \widehat{R}_{\mathbf{w}_*})) - \zeta (\widehat{R}_{\mathbf{w}_*} + \widehat{R}_{\mathbf{q}_*}^2/\widehat{R}_{\mathbf{w}_*}) \mathbf{v}^\top \mathbf{q}_* \\ &\quad + (\zeta - \zeta^2) \widehat{R}_{\mathbf{q}_*} (1 - 2\zeta \widehat{R}_{\mathbf{w}_*}) \mathbf{v}^\top \mathbf{w}_* - \sigma^2 (1 + 2/T) (\zeta - \zeta^2) \widehat{R}_{\mathbf{q}_*} \mathbf{v}^\top \widehat{\mathbf{w}}_1 \\ &= (\zeta - \zeta^2) \widehat{R}_{\mathbf{w}_*} \underbrace{(1 + \sigma^2 \|\widehat{\mathbf{w}}_1\|^2/(T \widehat{R}_{\mathbf{w}_*})) - \zeta (\widehat{R}_{\mathbf{w}_*} + \widehat{R}_{\mathbf{q}_*}^2/\widehat{R}_{\mathbf{w}_*})}_{:= \widehat{C}_1} \mathbf{v}^\top \mathbf{q}_* \\ &\quad + (\zeta - \zeta^2) \widehat{R}_{\mathbf{q}_*} \underbrace{(1 - 2\zeta \widehat{R}_{\mathbf{w}_*} - \eta \zeta \sigma^2 (1 + 2/T))}_{:= \widehat{C}_2} \mathbf{v}^\top \mathbf{w}_* - \underbrace{\eta \sigma^2 (1 + 2/T) (\zeta - \zeta^2) \widehat{R}_{\mathbf{q}_*}}_{:= \widehat{C}_3} \mathbf{v}^\top \delta_1 \end{aligned} \quad (39)$$

where, in the last line we used (34) and set

$$\delta_1 := (\widehat{\mathbf{w}}_1 - \mathbf{w}_1)/\eta = \zeta \mathbf{q}_* \left( \frac{1}{n} \sum_{i \in [n]} y_i \right) + \frac{1}{n} \sum_i y_i \gamma_i.$$

Recall from (31) and (33) that  $R_{\mathbf{w}_*}/2 \leq \widehat{R}_{\mathbf{w}_*} \leq 3R_{\mathbf{w}_*}/2$ ,  $|\widehat{R}_{\mathbf{q}_*} - R_{\mathbf{q}_*}| \leq c_0 \eta \zeta Q W$  and  $\|\widehat{\mathbf{w}}_1\| \leq \|\mathbf{w}_1\| + c_0 \eta \zeta W$ . Also, from (33) we have that  $\|\delta_1\| \leq c_0 \zeta W$ . Here and onwards,  $c_0 > 0$  is a small enough absolute constant (smaller than 1/2) whose value may change from line to line. Further recall  $R_{\mathbf{w}_*} = \eta \zeta W^2$ ,  $R_{\mathbf{q}_*} = \eta \zeta \rho W Q$  and  $\|\mathbf{w}_1\| = \eta \zeta W$ . With these, we can set small enough step size  $\eta \propto \sigma^{-2} Q^{-2}$  such that

$$\widehat{C}_1 \in [1/2, 3/2], \quad \widehat{C}_2 \in [-1/8, 1], \quad \widehat{C}_3 \in [0, (c_3/\zeta)/Q^2],$$

for constant  $c_3 > 0$  to be made small enough later in the proof.

From the above, we can compute

$$\widehat{C}_3 |\widehat{R}_{\mathbf{q}_*}| \|\delta_1\| \leq \frac{c_3}{\zeta Q^2} \eta \zeta (c_0 Q W + \rho W Q) (c_0 \zeta Q) \leq c_3 c_0 \eta \zeta Q.$$

Moreover,

$$\begin{aligned}\hat{C}_1 \widehat{R}_{\mathbf{w}_*} Q^2 &\geq \eta \zeta W^2 Q^2 / 4 \\ \hat{C}_2 \widehat{R}_{\mathbf{q}_*} \rho W Q &= \hat{C}_2 \rho W Q R_{\mathbf{q}_*} + C_2 \rho W Q (\widehat{R}_{\mathbf{q}_*} - R_{\mathbf{q}_*}) \geq -\eta \zeta \rho^2 Q^2 W^2 / 8 - c_0 \eta \zeta Q^2 W^2 \\ &\geq -\eta \zeta \rho^2 Q^2 W^2 / 8 - c_0 \eta \zeta W^2 Q^2.\end{aligned}$$

Putting the above displays together we find that

$$\mathbf{q}_*^\top \mathbf{G}_{\mathbf{q}}(0, \widehat{\mathbf{w}}_1) \geq \eta \zeta (\zeta - \zeta^2) (W^2 Q^2 / 4 - \rho^2 Q^2 W^2 / 8 - c_0 W^2 Q^2 - c_3 c_0 Q^2).$$

Similarly, we can compute that

$$\begin{aligned}|\hat{C}_1 \widehat{R}_{\mathbf{w}_*} \rho Q W| &\leq (9/4) \eta \zeta W^3 Q |\rho| \\ |\hat{C}_2 \widehat{R}_{\mathbf{q}_*} W^2| &\leq \eta \zeta (|\rho| W^3 Q + c_0 Q W^3).\end{aligned}$$

Hence, for  $y \in \{\pm 1\}$

$$|y \mathbf{w}_*^\top \mathbf{G}_{\mathbf{q}}(0, \widehat{\mathbf{w}}_1)| \leq \eta \zeta (\zeta - \zeta^2) 4W^3 Q |\rho| + c_0 Q W^3 + c_3 c_0 Q W.$$

The above two displays put together yield (37). Specifically, we also use the simplifying assumption  $Q \gtrsim W \gtrsim 1$  and pick  $c_3$  small enough so that  $c_3 W^{-2} Q^{-1} \leq 1$ .

The norm bound in (38) follows again starting from (39) and using similar arguments as above to show:

$$\begin{aligned}\|\mathbf{G}_{\mathbf{q}}(0, \widehat{\mathbf{w}}_1)\| &\leq \eta \zeta (\zeta - \zeta^2) (\widehat{C}_1 R_{\mathbf{w}_*} Q + |\widehat{C}_2| |\widehat{R}_{\mathbf{q}_*}| W + \widehat{C}_3 |\widehat{R}_{\mathbf{q}_*}| \|\delta_1\|) \\ &\leq \eta \zeta (\zeta - \zeta^2) ((9/4) W^2 Q + (|\rho| + c_0) W^2 Q + c_3 c_0 Q).\end{aligned}$$

□

The next lemma controls the effect on the relevance scores of the deviation term  $\widetilde{\mathbf{G}}_{\mathbf{q}}(0, \widehat{\mathbf{w}}) = \widehat{\mathbf{G}}_{\mathbf{q}}(0, \widehat{\mathbf{w}}) - \mathbf{G}_{\mathbf{q}}(0, \widehat{\mathbf{w}})$ .

**Lemma 9** ( $\widehat{\mathbf{G}}_{\mathbf{q}}(0, \cdot)$  control: Deviation term). *Let  $\widetilde{\mathbf{G}}_{\mathbf{q}}(0, \widehat{\mathbf{w}}_1) := \widehat{\mathbf{G}}_{\mathbf{q}}(0, \widehat{\mathbf{w}}_1) - \mathbf{G}_{\mathbf{q}}(0, \widehat{\mathbf{w}}_1)$  and suppose  $\widehat{\mathbf{w}}_1$  satisfies (31) and (33). Also assume  $\sigma = 1$ . Fix any  $u > 0$  and any small constant  $c_1 > 0$ . Then, there exists small enough constant  $c_\eta$  (dependent on  $c_1$ ) such that if step-size  $\eta$  is small enough as per (36) the following statements hold.*

*With probability at least  $1 - c' e^{-cu^{2/3}}$  for positive constants  $C, c', c > 0$  it holds for signal tokens that*

$$|(\mathbf{q}_* + y \mathbf{w}_*)^\top \widetilde{\mathbf{G}}_{\mathbf{q}}(0, \widehat{\mathbf{w}}_1)| \leq u C c_1 Q \left( Q \vee \frac{\sigma \vee \sigma^2}{\sqrt{T}} \vee \frac{\sigma^3}{T} \right) \frac{\log(n)}{\sqrt{n}}, \quad y \in \{\pm 1\} \quad (40)$$

*Moreover, with probability at least  $1 - c' d e^{-cu^{2/3}}$  it holds that*

$$\|\widetilde{\mathbf{G}}_{\mathbf{q}}(0, \widehat{\mathbf{w}}_1)\| \leq u C c_1 \left( Q \vee \frac{\sigma \vee \sigma^2}{\sqrt{T}} \vee \frac{\sigma^3}{T} \right) \frac{\log(n) \sqrt{d}}{\sqrt{n}}. \quad (41)$$

*Proof.* We study each one of the terms of  $\widetilde{\mathbf{G}}_{\mathbf{q}}(0, \widehat{\mathbf{w}}_1)$  in (20) separately. We repeat the terms here for convenience, also

noting the substitutions  $R_{w_*} \leftarrow \widehat{R}_{w_*} = \widehat{w}_1^\top w_*$  and  $R_{q_*} \leftarrow \widehat{R}_{q_*} = \widehat{w}_1^\top q_*$ .

$$\begin{aligned}
 v^\top \widetilde{G}_q(0, \widehat{w}_1) &= [((\zeta - \zeta^2) - 2(\zeta^2 - \zeta^3)\widehat{R}_{w_*})\widehat{R}_{q_*}R_{v,q_*} + ((\zeta - \zeta^2)\widehat{R}_{w_*} - (\zeta^2 - \zeta^3)(\widehat{R}_{w_*}^2 + \widehat{R}_{q_*}^2))R_{v,w_*}] \left( \frac{1}{n} \sum_{i=1}^n y_i \right) \\
 &+ [(-\zeta + 2\zeta^2)\widehat{R}_{q_*}R_{v,q_*} + (-\zeta + (-\zeta + 2\zeta^2)\widehat{R}_{w_*})R_{v,w_*} + (1 - \zeta)\widehat{w}_1^\top \Sigma v] \left( \frac{1}{n} \sum_{i=1}^n (\gamma_i^\top \widehat{w}_1) \right) \\
 &+ \left[ \zeta \widehat{R}_{w_*} - \zeta^2(\widehat{R}_{q_*}^2 + \widehat{R}_{w_*}^2) + \frac{(1 - \zeta)}{T} \widehat{w}_1^\top \Sigma \widehat{w}_1 \right] \left( \frac{1}{n} \sum_{i=1}^n v^\top \gamma_i \right) \\
 &+ [(-\zeta + (-\zeta + 2\zeta^2)\widehat{R}_{w_*})R_{v,q_*} + (-\zeta + 2\zeta^2)\widehat{R}_{q_*}R_{v,w_*}] \left( \frac{1}{n} \sum_{i=1}^n y_i (\gamma_i^\top \widehat{w}_1) \right) \\
 &+ [\zeta \widehat{R}_{q_*} - 2\zeta^2 \widehat{R}_{q_*} \widehat{R}_{w_*}] \left( \frac{1}{n} \sum_{i=1}^n y_i (v^\top \gamma_i) \right) \\
 &+ \zeta R_{v,q_*} \left( \frac{1}{n} \sum_{i=1}^n (\gamma_i^\top \widehat{w}_1)^2 - \frac{(1 - \zeta)}{T} \widehat{w}_1^\top \Sigma \widehat{w}_1 \right) \\
 &+ \zeta R_{v,w_*} \left( \frac{1}{n} \sum_{i=1}^n (\gamma_i^\top \widehat{w}_1)^2 y_i \right) \\
 &+ [1 - 2\zeta \widehat{R}_{w_*}] \left( \frac{1}{n} \sum_{i=1}^n y_i (\widehat{w}_1^\top \gamma_i) (v^\top \gamma_i) \right) \\
 &+ (1 - \zeta \widehat{R}_{w_*}) \left( \frac{1}{n} \sum_{i=1}^n y_i v^\top \widehat{\Sigma}_i \widehat{w}_1 \right) \\
 &\quad - \zeta \widehat{R}_{q_*} \left( \frac{1}{n} \sum_{i=1}^n v^\top \widehat{\Sigma}_i \widehat{w}_1 - (1 - \zeta) v^\top \Sigma \widehat{w}_1 \right) \\
 &- [2\zeta \widehat{R}_{q_*}] \left( \frac{1}{n} \sum_{i=1}^n (\widehat{w}_1^\top \gamma_i) (v^\top \gamma_i) - \frac{1 - \zeta}{T} v^\top \Sigma \widehat{w}_1 \right) \\
 &- \left( \frac{1}{n} \sum_{i=1}^n (v^\top \gamma_i) \left( (\widehat{w}_1^\top \gamma_i)^2 - \frac{(1 - \zeta)}{T} \widehat{w}_1^\top \Sigma \widehat{w}_1 \right) \right) \\
 &- \left( \frac{1}{n} \sum_{i=1}^n (\gamma_i^\top \widehat{w}_1) (v^\top \widehat{\Sigma}_i \widehat{w}_1 - (1 - \zeta) \widehat{w}_1^\top \Sigma v) \right). \tag{42}
 \end{aligned}$$

Recall from the lemma assumption that (31) holds and from  $R_{w_*} = \eta \zeta W^2$  and  $R_{q_*} = \eta \zeta \rho W Q$ , that  $1/2\eta \zeta W^2 \leq \widehat{R}_{w_*} \leq 3/2\eta \zeta W^2$  and  $|\widehat{R}_{q_*}| \leq \eta \zeta (|\rho| + c_0) W Q$ . The observation is that we can choose step-size  $\eta$  small enough (as stated in (36)) to bound (in absolute value) all the coefficients in (42) (aka all terms in square brackets) that include  $\widehat{R}_{w_*}, \widehat{R}_{q_*}$ . Therefore, for any small positive constant  $c_1 > 0$  it can be checked that there is sufficiently small constant  $c_\eta$  that determines step-size  $\eta$

in (36) such that

$$\begin{aligned}
 |\mathbf{v}^\top \tilde{\mathbf{G}}_q(0, \hat{\mathbf{w}}_1)| &\leq c_1 [|R_{\mathbf{v}, \mathbf{q}_*}| + |R_{\mathbf{v}, \mathbf{w}_*}|] \left| \frac{1}{n} \sum_{i=1}^n y_i \right| \\
 &+ [c_1 |R_{\mathbf{v}, \mathbf{q}_*}| + (1 + c_1) |R_{\mathbf{v}, \mathbf{w}_*}| + (1 - \zeta) \sigma^2 |\hat{\mathbf{w}}_1^\top \mathbf{v}|] \left| \frac{1}{n} \sum_{i=1}^n (\gamma_i^\top \hat{\mathbf{w}}_1) \right| \\
 &+ \left[ c_1 + \frac{(1 - \zeta) \sigma^2 \|\hat{\mathbf{w}}_1\|^2}{T} \right] \left| \frac{1}{n} \sum_{i=1}^n \mathbf{v}^\top \gamma_i \right| \\
 &+ [(1 + c_1) |R_{\mathbf{v}, \mathbf{q}_*}| + c_1 |R_{\mathbf{v}, \mathbf{w}_*}|] \left| \frac{1}{n} \sum_{i=1}^n y_i (\gamma_i^\top \hat{\mathbf{w}}_1) \right| \\
 &+ [c_1] \left| \frac{1}{n} \sum_{i=1}^n y_i (\mathbf{v}^\top \gamma_i) \right| \\
 &+ \zeta |R_{\mathbf{v}, \mathbf{q}_*}| \left| \frac{1}{n} \sum_{i=1}^n (\gamma_i^\top \hat{\mathbf{w}}_1)^2 - \frac{(1 - \zeta)}{T} \hat{\mathbf{w}}_1^\top \Sigma \hat{\mathbf{w}}_1 \right| \\
 &+ \zeta |R_{\mathbf{v}, \mathbf{w}_*}| \left| \frac{1}{n} \sum_{i=1}^n (\gamma_i^\top \hat{\mathbf{w}}_1)^2 y_i \right| \\
 &+ [1 + c_1] \left| \frac{1}{n} \sum_{i=1}^n y_i (\hat{\mathbf{w}}_1^\top \gamma_i) (\mathbf{v}^\top \gamma_i) \right| \\
 &+ [1 + c_1] \left| \frac{1}{n} \sum_{i=1}^n y_i \mathbf{v}^\top \hat{\Sigma}_i \hat{\mathbf{w}}_1 \right| \\
 &+ c_1 \left| \frac{1}{n} \sum_{i=1}^n \mathbf{v}^\top \hat{\Sigma}_i \hat{\mathbf{w}}_1 - (1 - \zeta) \mathbf{v}^\top \Sigma \hat{\mathbf{w}}_1 \right| \\
 &+ [c_1] \left| \frac{1}{n} \sum_{i=1}^n (\hat{\mathbf{w}}_1^\top \gamma_i) (\mathbf{v}^\top \gamma_i) - \frac{1 - \zeta}{T} \mathbf{v}^\top \Sigma \hat{\mathbf{w}}_1 \right| \\
 &+ \left| \frac{1}{n} \sum_{i=1}^n (\mathbf{v}^\top \gamma_i) \left( (\hat{\mathbf{w}}_1^\top \gamma_i)^2 - \frac{(1 - \zeta)}{T} \hat{\mathbf{w}}_1^\top \Sigma \hat{\mathbf{w}}_1 \right) \right| \\
 &+ \left| \frac{1}{n} \sum_{i=1}^n (\gamma_i^\top \hat{\mathbf{w}}_1) (\mathbf{v}^\top \hat{\Sigma}_i \hat{\mathbf{w}}_1 - (1 - \zeta) \hat{\mathbf{w}}_1^\top \Sigma \mathbf{v}) \right|. \tag{43}
 \end{aligned}$$

Now, we use successively Lemma 6 to bound the random terms. Also note that  $|R_{\mathbf{v}, \mathbf{q}_*}| \leq Q \|\mathbf{v}\|$ ,  $|R_{\mathbf{v}, \mathbf{w}_*}| \leq W \|\mathbf{v}\|$  and  $\|\hat{\mathbf{w}}_1\| \leq \eta \zeta (1 + c_0) W =: \eta \zeta M$ . Here, we denote  $M := (1 + c_0) W$  for convenience. With these, for any  $u > 0$ , with probability at least  $1 - c' e^{-cu^{2/3}}$  we have

$$\begin{aligned}
 |\mathbf{v}^\top \tilde{\mathbf{G}}_q(0, \widehat{\mathbf{w}}_1)| &\leq u \cdot \frac{C}{\sqrt{n}} \|\mathbf{v}\| c_1 (W + Q) \\
 &+ u \cdot \frac{C\sigma\sqrt{1-\zeta}}{\sqrt{nT}} \|\mathbf{v}\| \left( ((1+c_1)(2W+2Q) + c_1\eta\zeta(1-\zeta)\sigma^2 M) \eta\zeta M + (2c_1 + \sigma^2(1-\zeta)\eta^2\zeta^2 M^2/T) \right) \\
 &+ u \cdot \frac{C\sigma^2(1-\zeta)}{T\sqrt{n}} \|\mathbf{v}\| \left( \eta^2\zeta^3 M^2 Q + \eta^2\zeta^3 W M^2 + (1+c_1)\eta\zeta M \right) \\
 &+ u \cdot \frac{C\sigma^2\sqrt{1-\zeta}}{\sqrt{nT}} \|\mathbf{v}\| \left( (1+2c_1)\eta\zeta M \right) \\
 &+ u \frac{C\sigma^2(1-\zeta)}{T\sqrt{n}} \|\mathbf{v}\| \left( c_1\eta\zeta M \right) \\
 &+ u \frac{C\sigma^3(1-\zeta)^{3/2}\log(n)}{T^{3/2}\sqrt{n}} \|\mathbf{v}\| \left( \eta^2\zeta^2 M^2 \right) \\
 &+ u \frac{C\sigma^3(1-\zeta)\log(n)}{T\sqrt{n}} \|\mathbf{v}\| \eta^2\zeta^2 M^2. \tag{44}
 \end{aligned}$$

Now, again using small step size  $\eta$  as per (36) (recall that  $M = (1+c_0)W \lesssim Q$ ), this can be further simplified to the following (here the value of constant  $c_1$  might be different from (44))

$$\begin{aligned}
 |\mathbf{v}^\top \tilde{\mathbf{G}}_q(0, \widehat{\mathbf{w}}_1)| &\leq u \cdot \frac{C}{\sqrt{n}} \|\mathbf{v}\| c_1 (W + Q) + u \cdot \frac{C\sigma\sqrt{1-\zeta}}{\sqrt{nT}} \|\mathbf{v}\| c_1 + u \cdot \frac{C\sigma^2(1-\zeta)}{T\sqrt{n}} \|\mathbf{v}\| c_1 \\
 &+ u \cdot \frac{C\sigma^2\sqrt{1-\zeta}}{\sqrt{nT}} \|\mathbf{v}\| c_1 + u \frac{C\sigma^3(1-\zeta)^{3/2}\log(n)}{T^{3/2}\sqrt{n}} \|\mathbf{v}\| c_1 + u \frac{C\sigma^3(1-\zeta)\log(n)}{T\sqrt{n}} \|\mathbf{v}\| c_1 \\
 &\leq u \cdot \|\mathbf{v}\| \cdot \frac{C}{\sqrt{n}} \cdot c_1 \left( Q \vee \frac{\sigma \vee \sigma^2}{\sqrt{T}} \vee \frac{\sigma^2}{T} \vee \frac{\sigma^3 \log(n)}{T^{3/2}} \vee \frac{\sigma^3 \log(n)}{T} \right) \\
 &\leq u \cdot \|\mathbf{v}\| \cdot \frac{C \log(n)}{\sqrt{n}} \cdot c_1 \left( Q \vee \frac{\sigma \vee \sigma^2}{\sqrt{T}} \vee \frac{\sigma^3}{T} \right). \tag{45}
 \end{aligned}$$

Now, we can compute the deviation of the relevance scores. For signal tokens we have for both  $y \in \{\pm 1\}$ , and all  $u > 0$  with probability at least  $1 - c'e^{-cu^{2/3}}$ , there exist constant  $C > 0$  such that

$$|(\mathbf{q}_* + y\mathbf{w}_*)^\top \tilde{\mathbf{G}}_q(0, \widehat{\mathbf{w}}_1)| \leq u C c_1 Q \left( Q \vee \frac{\sigma \vee \sigma^2}{\sqrt{T}} \vee \frac{\sigma^3}{T} \right) \frac{\log(n)}{\sqrt{n}} \tag{46}$$

Similarly, since (45) holds for all  $\mathbf{v}$ , we can apply it for all standard basis vectors  $\mathbf{v} = \mathbf{e}_j, j \in [n]$  and union bounding yields for all  $u > 0$  with probability at least  $1 - c'de^{-cu^{2/3}}$ ,

$$\|\tilde{\mathbf{G}}_q(0, \widehat{\mathbf{w}}_1)\| \leq u C c_1 \left( Q \vee \frac{\sigma \vee \sigma^2}{\sqrt{T}} \vee \frac{\sigma^3}{T} \right) \frac{\log(n)\sqrt{d}}{\sqrt{n}}.$$

□

### C.3. First and second gradient steps combined: Learning the relevant features

With lemmas 7, 8, and 9 at hand, we are now ready to put things together stating our final bounds for relevance scores. The finding is presented as a stand-alone lemma below.

**Lemma 10** (Put things together). *Consider the finite-sample gradient step  $\widehat{\mathbf{q}}_1 = \widehat{\mathbf{G}}_q(0, \widehat{\mathbf{w}}_1)$ , where recall that  $\widehat{\mathbf{w}}_1 = \eta\widehat{\mathbf{G}}_w(0, 0)$ . Fix any  $u_0, u_1, u_2, u_3 > 0$  and any small constants  $\tilde{c}_0, \tilde{c}_1 > 0$ . Suppose step-size  $\eta$  of first gradient step satisfies (36) for sufficiently small constant  $c_\eta = c_\eta(\tilde{c}_1) > 0$  and further assume*

$$\sqrt{n} \geq u_0 \cdot C_0 \frac{Q}{W} \quad \text{and} \quad \sqrt{\frac{n\zeta T}{d}} \geq (1+u_0) \cdot C_0 \frac{\sigma}{W} \sqrt{1/\zeta - 1}. \tag{47}$$



for some large enough constant  $C_0 = C_0(\tilde{c}_0) > 0$ . Finally, make the following mild assumption (for simplicity),  $\frac{\sigma \vee \sigma^2}{\sqrt{T}} \vee \frac{\sigma^3}{T} \leq Q$ . For a fresh dataset  $(\mathbf{X}_i, y_i)_{i \in [n]}$  consider the signal relevance scores  $\hat{r}_{i,t} := r_{i,0} := (\mathbf{q}_* + y_i \mathbf{w}_*)^\top \widehat{\mathbf{G}}_{\mathbf{q}}(0, \hat{\mathbf{w}}_1)$  for  $t \in \mathcal{R}$ . There exist positive absolute constants  $c_i, c'_i, i = 0, 1, 2, 3$  such that the following statements hold.

- With probability at least  $1 - c'_0 e^{-c_0 u_0^2} - c'_1 e^{c_1 u_1^{2/3}}$ , the signal relevance scores satisfy

$$\min_{i \in [n]} r_{i,0} \geq \eta \zeta (\zeta - \zeta^2) W^2 Q \left( (1/4 - \rho^2/8 - 2\tilde{c}_0) Q - (9/4|\rho| + 2\tilde{c}_0) W \right) - u_1 \tilde{c}_1 Q^2 \frac{\log(n)}{\sqrt{n}} =: B(u_1). \quad (48)$$

- With probability at least  $1 - c'_0 e^{-c_0 u_0^2} - c'_3 d e^{c_3 u_3^{2/3}}$ , the norm of  $\widehat{\mathbf{G}}_{\mathbf{q}}(0, \hat{\mathbf{w}}_1)$  satisfies

$$\|\widehat{\mathbf{G}}_{\mathbf{q}}(0, \hat{\mathbf{w}}_1)\| \leq \eta \zeta (\zeta - \zeta^2) W^2 Q (13/4 + 2\tilde{c}_0) + u_3 \tilde{c}_1 \sigma Q \frac{\log(n) \sqrt{d}}{\sqrt{n}}. \quad (49)$$

*Proof.* The lemma follows by collecting (37), (40), (38), (41), and applying union bound for the noise terms over  $i \in [n], t \in \mathcal{R}^c$ .  $\square$

From the lemma above, we can show that provided  $Q$  is large enough with respect to  $W$  and  $\rho$  is small enough (see (50) below) then the normalized signal relevance score is with high probability over the training set proportional to  $Q$ .

**Corollary 1.** Suppose there exists positive constant  $\alpha \in (0, 3/16)$  such that

$$A_Q := (3/16 - \rho^2/8) Q - (9/4|\rho| + 1/16) W \geq \alpha \cdot Q. \quad (50)$$

Then, for sufficiently small step-size  $\eta \propto Q^{-2}$  and large enough  $n$  there exist constants  $c'_0, c_0, c'_1, c_1, c_3, c'_3$  such that with probability at least

$$\begin{aligned} & 1 - c'_0 \exp\left(-c_0 n \left( (W/Q)^2 \wedge \frac{\zeta^2 W^2 T}{d} \right)\right) - c'_1 \exp\left(-c_1 n^{1/3} (W/Q)^{4/3} \zeta^{4/3} / \log^{2/3}(n)\right) \\ & - c'_3 d \exp\left(-c_3 (n/d)^{1/3} (W/Q)^{4/3} \zeta^{4/3} / \log^{2/3}(n)\right), \end{aligned} \quad (51)$$

it holds that

$$(\mathbf{q}_* + y \mathbf{w}_*)^\top \left( \frac{\widehat{\mathbf{G}}_{\mathbf{q}}(0, \hat{\mathbf{w}}_1)}{\|\widehat{\mathbf{G}}_{\mathbf{q}}(0, \hat{\mathbf{w}}_1)\|} \right) \geq (\alpha/14) \cdot Q.$$

*Proof.* Suppose (50) holds and recall the notation  $A_Q$  defined therein. Further suppose  $n$  is large enough such that (47) holds so that we can invoke Lemma 10. Therein set

$$\begin{aligned} u_1 & \propto \eta W^2 \zeta^2 \frac{\sqrt{n}}{\log(n)} \propto (W/Q)^2 \zeta^2 \frac{\sqrt{n}}{\log(n)}, \\ u_3 & \propto \eta W^2 \zeta^2 \frac{\sqrt{n}}{\log(n) \sqrt{d}} \propto (W/Q)^2 \zeta^2 \frac{\sqrt{n}}{\log(n) \sqrt{d}} \end{aligned}$$

and

$$u_0 \propto \sqrt{n} \left( (W/Q) \wedge \zeta W \sqrt{T/d} \right).$$

Note that the latter condition is consistent with (47). On the other hand, the former two conditions are chosen so that the following two hold.

First,

$$\eta \zeta (\zeta - \zeta^2) W^2 Q \cdot A_Q > 2u_1 \tilde{c}_1 Q^2 \frac{\log n}{\sqrt{n}}.$$

Thus, (48) and setting  $2\tilde{c}_0 = 1/16$  give

$$(\mathbf{q}_* + y \mathbf{w}_*)^\top \widehat{\mathbf{G}}_{\mathbf{q}}(0, \hat{\mathbf{w}}_1) \geq \frac{1}{2} \eta \zeta (\zeta - \zeta^2) W^2 Q A_Q$$

Second,

$$\eta\zeta(\zeta - \zeta^2)W^2Q > u_3\tilde{c}_1Q \frac{\log n\sqrt{d}}{\sqrt{n}}$$

Thus, from (49)

$$\|\widehat{\mathbf{G}}_{\mathbf{q}}(0, \hat{\mathbf{w}}_1)\| \leq 2\eta\zeta(\zeta - \zeta^2)W^2Q(13/4 + 1/16).$$

Combining the above we conclude with the desired: the signal correlation is lower bounded by

$$(\mathbf{q}_* + y\mathbf{w}_*)^\top \widehat{\mathbf{G}}_{\mathbf{q}}(0, \hat{\mathbf{w}}_1) / \|\widehat{\mathbf{G}}_{\mathbf{q}}(0, \hat{\mathbf{w}}_1)\| \geq A_Q/14 \geq (\alpha/14) \cdot Q.$$

□

Having shown that the second gradient step has high signal relevance score, we can use it to show in the lemma below that the attention features are close to the signal features with high-probability. This property will play a key role in finalizing the finite-sample analysis. Indeed, the rate of the event in (55) will turn out to govern the error rate of our algorithm.

**Lemma 11** (From learning the context to learning good features). *Suppose the second gradient step  $\mathbf{q}_1$  has correlation coefficient  $\geq \alpha$  with  $\mathbf{q}_* + y\mathbf{w}_*$  for any  $y \in \{\pm 1\}$ , i.e.,  $(\mathbf{q}_* + y\mathbf{w}_*)^\top \mathbf{q}_1 / \|\mathbf{q}_1\| \geq \alpha Q$ . Then, for any  $\epsilon > 0$  there exists sufficiently large  $\gamma_*(\epsilon)$  and absolute constant  $c > 0$  such that setting  $\mathbf{q}_1^\gamma = \gamma \frac{\mathbf{q}_1}{\|\mathbf{q}_1\|}$  for any  $\gamma \geq \gamma_*(\epsilon)$ , we have*

$$\mathbb{P}_{(\mathbf{X}, y)}(\|\mathbf{X}^\top \phi(\mathbf{X}\mathbf{q}_1^\gamma) - (\mathbf{q}_* + y\mathbf{w}_*)\| \leq \epsilon) \geq 1 - 2Te^{-c\frac{\alpha^2 Q^2}{\sigma^2}}. \quad (52)$$

*Proof.* Recall  $\mathbf{X}$  consists of  $\zeta T$  relevant and  $(1 - \zeta)T$  irrelevant tokens. For relevant tokens, we have by assumption that  $B := (\mathbf{q}_* + y\mathbf{w}_*)^\top \bar{\mathbf{q}}_1 \geq \alpha Q$ , where we denote  $\bar{\mathbf{q}}_1 := \mathbf{q}_1 / \|\mathbf{q}_1\|$ , for convenience. Let  $\mathbf{z}_1, \dots, \mathbf{z}_{(1-\zeta)T}$  denote the irrelevant tokens, which are  $\sigma$ -subgaussian. In order to ensure that softmax-attention perfectly selects the relevant tokens, we require

$$M := \max_{1 \leq t \leq (1-\zeta)T} \mathbf{z}_t^\top \bar{\mathbf{q}}_1 \leq \alpha Q/2 = B/2.$$

Condition on this event, which holds with at least probability  $1 - Te^{-c\frac{\alpha^2 Q^2}{\sigma^2}}$ . Then, the attention coefficients  $a_t = [\phi(\gamma \mathbf{X}\mathbf{q}_1)]_t$  are as follows:

$$t \text{ relevant: } a_t = \frac{1}{\zeta T + (1-\zeta)Te^{\gamma(M-B)}} \geq \frac{1}{\zeta T + (1-\zeta)Te^{-\gamma B/2}} =: \frac{1}{\zeta T} a_R.$$

$$t \text{ irrelevant: } a_t = \frac{1}{\zeta Te^{\gamma(B-M)} + (1-\zeta)T} \leq \frac{1}{\zeta Te^{\gamma B/2} + (1-\zeta)T} =: \frac{1}{(1-\zeta)T} a_I = \frac{1}{(1-\zeta)T} \cdot \frac{1}{1 + \frac{\zeta}{1-\zeta}e^{\gamma B/2}}.$$

Therefore,

$$\|\mathbf{X}^\top \phi(\gamma \mathbf{X}\mathbf{q}_1) - (\mathbf{q}_* + y\mathbf{w}_*)\| \leq (Q + W)(1 - a_R) + a_I \max_{1 \leq t \leq (1-\zeta)T} \|\mathbf{z}_t\|.$$

To continue, further condition on the event

$$\max_{1 \leq t \leq (1-\zeta)T} \|\mathbf{z}_t\| \leq C\sigma\sqrt{d} + \alpha Q \quad (53)$$

which holds with probability at least  $1 - (1 - \zeta)Te^{-c\alpha^2 Q^2/\sigma^2}$ . Also note that  $1 - a_R = a_I$ . Hence, with probability at least  $1 - 2Te^{-c\alpha^2 Q^2/(8\sigma^2)}$

$$\|\mathbf{X}^\top \phi(\gamma \mathbf{X}\mathbf{q}_1) - (\mathbf{q}_* + y\mathbf{w}_*)\| \leq \left((1 + \alpha)Q + W + C\sigma\sqrt{d}\right) a_I.$$

The right hand-side above can be made smaller than  $\epsilon$  by choosing  $\gamma$  large enough (depending on  $\epsilon, \alpha, Q, W, \sigma$  and  $d$ ). This completes the proof. □

Combining Lemma 11 and Corollary 1 we arrive at the following result, which we state as a stand-alone theorem since it summarizes the effect of the first two-gradient steps on learning good (aka relevant) features.

**Theorem 5** (Summary result of first two GD steps). *Suppose  $Q, W$  and  $\rho$  are such that there exists positive constant  $\alpha \in (0, 3/16)$  for which*

$$A_Q := (3/16 - \rho^2/8)Q - (9/4|\rho| + 1/16)W \geq \alpha \cdot Q.$$

*Fix any  $\epsilon > 0$ . For sufficiently small step-size  $\eta \lesssim Q^{-2}$ , large enough  $n$  and sufficiently large  $\gamma_* = \gamma_*(\epsilon)$  such that the following statements are true about the first two gradient steps:*

$$\begin{aligned} \widehat{\mathbf{w}}_1 &:= \eta \widehat{\mathbf{G}}_{\mathbf{w}}(0, 0) \\ \widehat{\mathbf{q}}_1^\gamma &:= \widehat{\mathbf{G}}_{\mathbf{q}}(0, \widehat{\mathbf{w}}_1), \quad \text{for any } \gamma \geq \gamma_*. \end{aligned}$$

*There exist constants  $c'_0, c_0, c'_1, c_1, c_3, c'_3, c$  such that with probability at least*

$$\begin{aligned} 1 - c'_0 \exp\left(-c_0 n \left( (W/Q)^2 \wedge \frac{\zeta^2 W^2 T}{d} \right)\right) - c'_1 \exp\left(-c_1 n^{1/3} (W/Q)^{4/3} \zeta^{4/3} / \log^{2/3}(n)\right) \\ - c'_3 d \exp\left(-c_3 (n/d)^{1/3} (W/Q)^{4/3} \zeta^{4/3} / \log^{2/3}(n)\right), \end{aligned} \quad (54)$$

*it holds for any fresh sample  $(\mathbf{X}, y)$  that*

$$\mathbb{P}(\|\mathbf{X}^\top \phi(\mathbf{X} \widehat{\mathbf{q}}_1^\gamma) - (\mathbf{q}_* + y \mathbf{w}_*)\| \leq \epsilon) \geq 1 - 2T e^{-c \frac{\alpha^2 Q^2}{\sigma^2}}. \quad (55)$$

*Moreover, it holds that*

$$\|\widehat{\mathbf{w}}_1\| \leq \frac{c\zeta W}{Q^2}.$$

#### C.4. Third gradient step

With the characterization of the quality of learnt features in Theorem 5, we are now ready to turn our attention to the third gradient step. For this last step, it turns out all we need is that  $\widehat{\mathbf{w}}_2 := \eta \widehat{\mathbf{G}}_{\mathbf{w}}(\widehat{\mathbf{q}}_1, \widehat{\mathbf{w}}_1)$  has a strictly positive correlation with  $\mathbf{w}_*$ . This is indeed the case and the result is formalized in the lemma below.

**Lemma 12** (Third step). *Suppose that the first and second gradient steps  $\widehat{\mathbf{w}}_1, \widehat{\mathbf{q}}_1$  are such that the following hold. First, for absolute constant  $c_\eta > 0$*

$$\|\widehat{\mathbf{w}}_1\| \leq c_\eta \zeta W / Q^2.$$

*Second, for  $\epsilon > 0$  and any fresh datapoint  $(\mathbf{X}, y)$  there exists  $\delta$  that does not depend on  $\epsilon$  such that*

$$\mathbb{P}_{(\mathbf{X}, y)}(\|\mathbf{X}^\top \phi(\mathbf{X} \widehat{\mathbf{q}}_1) - (\mathbf{q}_* + y \mathbf{w}_*)\| \leq \epsilon) \geq 1 - \delta.$$

*Consider the third gradient step  $\widehat{\mathbf{w}}_2 := \eta \widehat{\mathbf{G}}_{\mathbf{w}}(\widehat{\mathbf{q}}_1, \widehat{\mathbf{w}}_1)$ . There exists absolute constants  $c, C$  such that, for all sufficiently small  $\epsilon$  and  $c_\eta$ , with probability at least  $1 - n\delta - 2e^{c\eta}$ ,*

$$\frac{\widehat{\mathbf{w}}_2^\top \mathbf{w}_*}{\|\widehat{\mathbf{w}}_2\|} \geq C \frac{W^2}{Q}.$$

*Proof.* Note that the lemma's conclusion is insensitive to the choice of step-size  $\eta$  for the third gradient step. Thus, without loss of generality, assume below that  $\eta = 1$ .

Recall the gradient formula

$$\widehat{\mathbf{G}}_{\mathbf{w}}(\mathbf{q}, \mathbf{w}) := -\nabla_{\mathbf{w}} \widehat{\mathcal{L}}_{\mathcal{D}}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i \in [n]} (y_i - f_{\boldsymbol{\theta}}(\mathbf{X}_i)) \mathbf{X}_i^\top \phi((\mathbf{X}_i \mathbf{q}))$$

and denote for convenience

$$\boldsymbol{\epsilon}_i := \mathbf{X}_i^\top \phi(\mathbf{X}_i \mathbf{q}) - (\mathbf{q}_* + y_i \mathbf{w}_*), \quad i \in [n].$$

With this notation, we can conveniently rewrite the third gradient step evaluated at  $\mathbf{q} := \mathbf{q}^\gamma := \gamma \widehat{\mathbf{G}}_{\mathbf{q}}(0, \widehat{\mathbf{w}}_1)$  for any  $\gamma > 0$  as follows:

$$\begin{aligned}
 \widehat{\mathbf{w}}_2 &:= \widehat{\mathbf{G}}_{\mathbf{w}}(\mathbf{q}, \widehat{\mathbf{w}}_1) = \frac{1}{n} \sum_{i \in [n]} y_i \mathbf{X}_i^\top \phi(\mathbf{X}_i \mathbf{q}) - \frac{1}{n} \sum_{i \in [n]} (\widehat{\mathbf{w}}_1^\top \mathbf{X}_i^\top \phi(\mathbf{X}_i \mathbf{q})) \mathbf{X}_i^\top \phi(\mathbf{X}_i \mathbf{q}) \\
 &= \mathbf{w}_* + \mathbf{q}_* \left( \frac{1}{n} \sum_{i \in [n]} y_i \right) + \frac{1}{n} \sum_{i \in [n]} y_i (\mathbf{X}_i^\top \phi(\mathbf{X}_i \mathbf{q}) - (\mathbf{q}_* + y_i \mathbf{w}_*)) \\
 &\quad - \frac{1}{n} \sum_{i \in [n]} (\widehat{\mathbf{w}}_1^\top \mathbf{X}_i^\top \phi(\mathbf{X}_i \mathbf{q})) (\mathbf{q}_* + y_i \mathbf{w}_*) - \frac{1}{n} \sum_{i \in [n]} (\widehat{\mathbf{w}}_1^\top \mathbf{X}_i^\top \phi(\mathbf{X}_i \mathbf{q})) (\mathbf{X}_i^\top \phi(\mathbf{X}_i \mathbf{q}) - (\mathbf{q}_* + y_i \mathbf{w}_*)) \\
 &= \mathbf{w}_* + \mathbf{q}_* \left( \frac{1}{n} \sum_{i \in [n]} y_i \right) + \frac{1}{n} \sum_{i \in [n]} y_i \boldsymbol{\epsilon}_i \\
 &\quad - \frac{1}{n} \sum_{i \in [n]} \widehat{\mathbf{w}}_1^\top \boldsymbol{\epsilon}_i \boldsymbol{\epsilon}_i - \frac{1}{n} \sum_{i \in [n]} \widehat{\mathbf{w}}_1^\top (\mathbf{q}_* + y_i \mathbf{w}_*) \boldsymbol{\epsilon}_i - \frac{1}{n} \sum_{i \in [n]} \widehat{\mathbf{w}}_1^\top \boldsymbol{\epsilon}_i (\mathbf{q}_* + y_i \mathbf{w}_*) - \frac{1}{n} \sum_{i \in [n]} \widehat{\mathbf{w}}_1^\top (\mathbf{q}_* + y_i \mathbf{w}_*) (\mathbf{q}_* + y_i \mathbf{w}_*). \quad (56)
 \end{aligned}$$

In order to control the correlation  $\widehat{\mathbf{w}}_2^\top \mathbf{w}_* / \|\widehat{\mathbf{w}}_2\|$  it suffices to control  $\widehat{\mathbf{w}}_2^\top \mathbf{v}$  for arbitrary  $\mathbf{v} \in \mathbb{R}^d$ . In view of the expression above, it will suffice bounding the two random terms below:

$$\begin{aligned}
 \text{Term}_I &:= \left| \frac{1}{n} \sum_{i \in [n]} y_i \right| \\
 \text{Term}_{II} &:= \|\boldsymbol{\epsilon}_i\| = \|\mathbf{X}_i^\top \phi(\mathbf{X}_i \mathbf{q}) - (\mathbf{q}_* + y_i \mathbf{w}_*)\|
 \end{aligned}$$

For the first term, we have with probability at least  $1 - 2e^{-u_1^2/2}$  that

$$\text{Term}_I \leq u_1 / \sqrt{n}.$$

For the second term, we know by assumption that with probability at least  $1 - n\delta$ ,

$$\text{Term}_{II} \leq \epsilon.$$

Putting the above together, with probability at least  $1 - 2e^{-u^2/2} - n\delta$ , we have that

$$\begin{aligned}
 \mathbf{w}_*^\top \widehat{\mathbf{w}}_2 &\geq W^2 - |\rho| QW \frac{u_1}{\sqrt{n}} - \epsilon W - \frac{3}{2} \eta \zeta \left( \epsilon^2 W^2 + 2\epsilon W^2 (Q + W) + W^2 (Q + W)^2 \right) \\
 &\geq W^2 \left( 1 - |\rho| \frac{Q}{W} \frac{u_1}{\sqrt{n}} - \frac{\epsilon}{W} - \frac{3}{2} \eta \zeta \left( \epsilon^2 + 2\epsilon (Q + W) + (Q + W)^2 \right) \right)
 \end{aligned}$$

where we also used the lemma's assumption on  $\|\widehat{\mathbf{w}}_1\| \leq (3/2)\eta\zeta W$ .

To further lower bound  $\mathbf{w}_*^\top \widehat{\mathbf{w}}_2$ , recall that  $|\rho| \leq W/Q$ ,  $W \gtrsim 1$  and that  $\epsilon$  can be made arbitrarily small constant. Further pick

$$u_1 = c_1 \sqrt{n} \quad \text{and} \quad \eta = c_\eta / Q^2 \quad (57)$$

for sufficiently small constants  $c_1$  and  $c_\eta$ . With these, we guarantee with probability at least  $1 - 2e^{-c_1 n} - n\delta$  that

$$\widehat{\mathbf{w}}_2^\top \mathbf{w}_* \gtrsim W^2. \quad (58)$$

Next, we use similar arguments to bound  $\|\widehat{\mathbf{w}}_2\|$ . Conditioning on the event where the bounds derived above hold for  $\text{Term}_I$  and  $\text{Term}_{II}$ , we have from (56) that

$$\|\widehat{\mathbf{w}}_2\| \leq W + Q \frac{u_1}{\sqrt{n}} + \epsilon + \frac{3}{2} \eta \zeta W \left( \epsilon^2 + 2\epsilon (Q + W) + (Q + W)^2 \right) \lesssim Q,$$

where in the second inequality, we chose  $u_1, \eta$  as in (57) and used again that  $\epsilon$  is arbitrarily small constant, as well as,  $Q > W \gtrsim 1$ .

All the above combined, shows that

$$\frac{\widehat{\mathbf{w}}_2^\top \mathbf{w}_*}{\|\widehat{\mathbf{w}}_2\|} \gtrsim \frac{W^2}{Q}.$$

This completes the proof.  $\square$

### C.5. De-biasing step

**Lemma 13** (Debiasing predictions). *For some  $\epsilon > 0$ , suppose  $\mathbf{q}_1$  is such that a test example  $(y, \mathbf{X})$  satisfies*

$$\mathbb{P}_{(\mathbf{X}, y)} \left( \|\mathbf{X}^\top \phi(\mathbf{X} \hat{\mathbf{q}}_1) - (\mathbf{q}_* + y \mathbf{w}_*)\| \leq \epsilon \right) \geq 1 - \delta$$

and that  $\mathbf{w}_2$  is such that

$$\frac{\mathbf{w}_2^\top \mathbf{w}_*}{\|\mathbf{w}_2\|} > 4\epsilon.$$

Given a fresh dataset  $\mathcal{S} = (y_i, \mathbf{X}_i)_{i=1}^n$ , set  $b = \frac{1}{n} \sum_{i=1}^n f_\theta(\mathbf{X}_i) = \frac{1}{n} \sum_{i=1}^n \mathbf{w}_2^\top \mathbf{v}_i$  where  $\mathbf{v}_i := \mathbf{X}_i^\top \phi(\mathbf{X}_i \mathbf{q}_1)$ . Set the debiased classifier  $f'_\theta(\mathbf{X}) = f_\theta(\mathbf{X}) - b$ . Suppose  $n \geq 8 \log\left(\frac{2}{\delta n}\right)$ . Then, with probability  $1 - 2\delta n$  over  $\mathcal{S}$ , the test error of  $f'_\theta$  obeys

$$\text{ERR}(f'_\theta) \leq \delta.$$

*Proof.* First, let us prove the following intermediate statement: With probability  $1 - 2\delta n$  over  $\mathcal{S}$ , for a new test sample  $(y, \mathbf{X})$ , with probability  $1 - \delta$ ,

$$|y f'_\theta(\mathbf{X}) - \mathbf{w}_2^\top \mathbf{w}_*| \leq \sqrt{\frac{2 \log(2/\delta n)}{n}} \cdot \mathbf{w}_2^\top \mathbf{w}_* + 2\epsilon \|\mathbf{w}_2\|. \quad (59)$$

To see the above, start with observing that, with probability  $1 - n\delta$  over the dataset  $(y_i, \mathbf{X}_i)_{i=1}^n$ , for each  $\mathbf{v}_i$ ,

$$|\mathbf{w}_2^\top \mathbf{v}_i - \mathbf{w}_2^\top (\mathbf{q}_* + y_i \mathbf{w}_*)| \leq \epsilon \|\mathbf{w}_2\|.$$

Set  $\bar{b} = \mathbf{w}_2^\top \mathbf{q}_*$  and  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ . With probability  $1 - \delta n$ ,  $\bar{y} \leq \sqrt{\frac{2 \log(2/\delta n)}{n}}$ . Combining, with overall probability at least  $1 - 2\delta n$ , the classifier bias obeys

$$|b - \bar{b}| \leq |\mathbf{w}_2^\top \mathbf{w}_*| \sqrt{\frac{2 \log(2/\delta n)}{n}} + \epsilon \|\mathbf{w}_2\|.$$

To finalize, for a new sample  $(y, \mathbf{X})$ , with probability  $1 - \delta$ , we have that  $|\mathbf{w}_2^\top \mathbf{v} - \mathbf{w}_2^\top (\mathbf{q}_* + y \mathbf{w}_*)| \leq \epsilon \|\mathbf{w}_2\|$  where  $\mathbf{v} = \mathbf{X}^\top \phi(\mathbf{X} \mathbf{q}_1)$ . Thus, the prediction  $f'(\mathbf{X}) = f(\mathbf{X}) - b$  obeys

$$|y f'_\theta(\mathbf{X}) - y(f_\theta(\mathbf{X}) - \bar{b})| \leq |b - \bar{b}| \leq |\mathbf{w}_2^\top \mathbf{w}_*| \sqrt{\frac{2 \log(2/\delta n)}{n}} + \epsilon \|\mathbf{w}_2\|. \quad (60)$$

To conclude with (59), note that

$$|y(f_\theta(\mathbf{X}) - \bar{b}) - \mathbf{w}_2^\top \mathbf{w}_*| \leq \epsilon \|\mathbf{w}_2\|,$$

and apply triangle inequality with (60).

To prove the statement of the lemma, note that, when  $n \geq 8 \log(2/\delta n)$  and  $\mathbf{w}_2^\top \mathbf{w}_* > 4\epsilon \|\mathbf{w}_2\|$ , a test sample (with  $\geq 1 - \delta$  probability) obeys

$$y f'_\theta(\mathbf{X}) \geq \mathbf{w}_2^\top \mathbf{w}_* - \sqrt{\frac{2 \log(2/\delta n)}{n}} \mathbf{w}_2^\top \mathbf{w}_* - 2\epsilon \|\mathbf{w}_2\| \geq 0.5 \mathbf{w}_2^\top \mathbf{w}_* - 2\epsilon \|\mathbf{w}_2\| > 0. \quad (61)$$

Thus, the classifier makes the correct decision with the same probability.  $\square$

### C.6. Finishing the finite sample analysis

**Theorem 6** (Main theorem: Finite-sample). *Suppose  $Q, W$  and  $\rho$  are such that there exists positive constant  $\alpha \in (0, 3/16)$  for which*

$$(3/16 - \rho^2/8)Q - (9/4|\rho| + 1/16)W \geq \alpha \cdot Q. \quad (62)$$

Fix any  $\epsilon > 0$ . For sufficiently small step-size  $\eta \lesssim Q^{-2}$ , sufficiently large step-size  $\gamma = \gamma(\epsilon)$ , and large enough  $n$ , there exist constants  $c'_j, c_j, j = 0, 1, 2, 3, 4$  such that the following statements hold with probability at least

$$1 - c'_0 \exp\left(-c_0 n \left( (W/Q)^2 \wedge \left( \frac{\zeta^2 W^2 T}{d} \wedge 1 \right) \right)\right) - c'_1 \exp\left(-c_1 n^{1/3} (W/Q)^{4/3} \zeta^{4/3} / \log^{2/3}(n)\right) \\ - c'_3 d \exp\left(-c_3 (n/d)^{1/3} (W/Q)^{4/3} \zeta^{4/3} / \log^{2/3}(n)\right) - 2nT \exp\left(-c_4 \alpha^2 Q^2\right), \quad (63)$$

over the training set:

**1. Prompt attends to relevant tokens:** For any test sample  $(\mathbf{X}, y)$ , with probability at least  $1 - 2T \exp(-c_4 \alpha^2 Q^2)$ , the attention coefficients  $a_t = [\phi(\mathbf{X} \hat{\mathbf{q}}_1)]_t$  after the second gradient step satisfy:

$$a_t \begin{cases} \geq \frac{1-\epsilon}{\zeta T} & t \text{ relevant} \\ \leq \frac{\epsilon}{(1-\zeta)T} & t \text{ irrelevant.} \end{cases} \quad (64)$$

**2. Prompt learns relevant features:** The prompt attention mechanism outputs relevant tokens with the same probability. Concretely,

$$\mathbb{P}_{(\mathbf{X}, y)} \left( \|\mathbf{X}^\top \phi(\mathbf{X} \hat{\mathbf{q}}_1) - (\mathbf{q}_* + y \mathbf{w}_*)\| \leq \epsilon \right) \geq 1 - 2T \exp(-c_4 \alpha^2 Q^2).$$

**3. Test error:** The test error of the model  $f'_\theta$  satisfies

$$\text{ERR}(f'_\theta) \leq 2T \exp(-c_4 \alpha^2 Q^2).$$

*Proof.* The theorem follows by combining Theorem 5, Lemma 12 and Lemma 13.  $\square$

## D. Proofs for Population-gradient Analysis in Section 4.2

This section includes the missing proofs of all the results in Section 4.2 regarding population analysis of Algorithm 9.

### D.1. Proof of Lemma 2

We repeat here the lemma for convenience also stated for general (not necessarily isotropic) noise covariance  $\Sigma$ .

**Lemma 14.** *The second population gradient step  $\mathbf{q}_1 = \gamma \mathbf{G}_w(\mathbf{w}_1, 0)$  satisfies the following for  $\alpha := \eta \zeta$*

$$\mathbf{G}_q(0, \alpha \mathbf{w}_*) = \left( (\zeta - \zeta^2) (\alpha W^2 + \alpha^2 \mathbf{w}_*^\top \Sigma \mathbf{w}_* / T) - \alpha^2 (\zeta^2 - \zeta^3) (W^4 + (\mathbf{w}_*^\top \mathbf{q}_*)^2) \right) \mathbf{q}_* \\ + \left( ((\zeta - \zeta^2) - 2(\zeta^2 - \zeta^3) \alpha W^2) \alpha (\mathbf{w}_*^\top \mathbf{q}_*) \right) \mathbf{w}_* \\ - \left( (1 + 2/T) (\zeta - \zeta^2) \alpha (\mathbf{w}_*^\top \mathbf{q}_*) \right) \alpha \Sigma \mathbf{w}_* \quad (65)$$

*Proof.* The lemma follows immediately from Eqn. (19) of Lemma 5 by recognizing that for  $\mathbf{w} = \alpha \mathbf{w}_*$  it holds  $R_{\mathbf{q}_*} = \alpha \mathbf{q}_*^\top \mathbf{w}_*$  and  $R_{\mathbf{w}_*} = \alpha W^2$ .  $\square$

### D.2. Corollary 2

**Corollary 2.** *Suppose small enough step-size  $\eta$  obeying*

$$\eta (\zeta^2 (W^2 + Q^2) - \zeta \cdot \sigma^2 / T) \leq 1/2. \quad (66a)$$

$$\eta \zeta (2\zeta W^2 + (1 + 2/T) \sigma^2) \leq 5/4. \quad (66b)$$

$$\eta \zeta (\sigma^2 / T) \leq 1/2. \quad (66c)$$

Then, for  $C_1 \in [1/2, 3/2]$  and  $C_2 \in [-1/4, 1]$ , we have that

$$\mathbf{q}_1 = \gamma \eta \zeta (\zeta - \zeta^2) W (C_1 W \mathbf{q}_* + C_2 \rho Q \mathbf{w}_*).$$

In particular,  $\mathbf{q}_*^\top \mathbf{q}_1 = \gamma \eta \zeta (\zeta - \zeta^2) W^2 Q^2 (C_1 + C_2 \rho^2)$  and  $\mathbf{w}_*^\top \mathbf{q}_1 = \gamma \eta \zeta (\zeta - \zeta^2) W^3 Q \rho (C_1 + C_2)$ .

*Proof.* Set  $\alpha = \eta\zeta$  and

$$3/2 \geq C_1 := (1 + \alpha\sigma^2/T) - \alpha\zeta(W^2 + \rho^2Q^2) \geq 1/2. \quad (67a)$$

$$1 \geq C_2 := 1 - 2\alpha\zeta W^2 - (1 + 2/T)\alpha\sigma^2 \geq -1/4. \quad (67b)$$

The gradient formula follows directly from (10). For the lower/upper bounds on  $C_1, C_2$  use (66a), (66b) and (66c).  $\square$

**Remark 2** (Condition on correlation). *To classify correctly the signal tokens, we need*

$$y\mathbf{w}_*^\top(\mathbf{q}_* + y\mathbf{w}_*) > 0 \iff y\rho Q + W > 0 \iff |\rho| < W/Q \quad (68)$$

Note that if (68) holds, then

$$C_1 \geq 1 + \alpha\sigma^2/T - 2\alpha\zeta W^2 = C_2 + (1 + 3/T)\alpha\sigma^2. \quad (69)$$

### D.3. Proof of Theorem 2

We start by showing that for any  $\epsilon > 0$ , and all sufficiently large  $\gamma \geq \gamma_*(\epsilon)$ , it holds

$$\mathbb{P}_{(\mathbf{X}, y) \sim \mathcal{D}}(\|\mathbf{X}^\top \phi(\mathbf{X}\mathbf{q}_1^\gamma) - (\mathbf{q}_* + y\mathbf{w}_*)\| \leq \epsilon) \geq 1 - 2Te^{-c\frac{\alpha^2 Q^2}{\sigma^2}} =: 1 - \delta. \quad (70)$$

We can get this by applying Lemma 11 provided only that we show the correlation of the normalized gradient-step  $\mathbf{q}_1^\gamma / \|\mathbf{q}_1^\gamma\|$  ( $= \bar{\mathbf{q}}_1$  below) with signal-relevant tokens is at least  $\alpha Q$ . Concretely, we have from Corollary 2 that  $\mathbf{q}_1^\gamma = \gamma\mathbf{q}_1 := \gamma\eta\zeta(\zeta - \zeta^2)W(C_1W\mathbf{q}_* + C_2\rho Q\mathbf{w}_*)$  with  $3/2 \geq C_1 \geq 1/2$ ,  $1 \geq C_2 \geq -1/4$ . Define for convenience  $\mathbf{q}_1 := \eta\zeta(\zeta - \zeta^2)W(C_1W\mathbf{q}_* + C_2\rho Q\mathbf{w}_*)$  and consider the normalized gradient step

$$\bar{\mathbf{q}}_1 := \frac{\mathbf{q}_1}{\|\mathbf{q}_1\|_2} = \frac{C_1W\mathbf{q}_* + C_2\rho Q\mathbf{w}_*}{\sqrt{C_1^2W^2Q^2 + C_2^2\rho^2Q^2W^2 + 2C_1C_2\rho^2W^2Q^2}} = \frac{C_1W\mathbf{q}_* + C_2\rho Q\mathbf{w}_*}{QW\sqrt{C_1^2 + C_2\rho^2(C_2 + 2C_1)}}.$$

We can lower-bound its correlation with a signal token  $\mathbf{q}_* + y\mathbf{w}_*$  as follows:

$$\begin{aligned} \bar{\mathbf{q}}_1^\top(\mathbf{q}_* + y\mathbf{w}_*) &= \frac{C_1W(Q^2 + y\rho WQ) + C_2\rho Q(\rho WQ + yW^2)}{QW\sqrt{C_1^2 + C_2\rho^2(C_2 + 2C_1)}} = \frac{C_1(Q + y\rho W) + C_2\rho(\rho Q + yW)}{\sqrt{C_1^2 + C_2\rho^2(C_2 + 2C_1)}} \\ &= \frac{(C_1 + C_2\rho^2)Q + y\rho(C_1 + C_2)W}{\sqrt{C_1^2 + C_2\rho^2(C_2 + 2C_1)}} \\ &\geq \frac{(C_1 + C_2\rho^2)Q + y\rho(C_1 + C_2)W}{C_1\sqrt{1 + 3\rho^2}} \geq \frac{1 + (C_2/C_1)\rho^2}{\sqrt{1 + 3\rho^2}}Q - \frac{|\rho|(1 + C_2/C_1)}{\sqrt{1 + 3\rho^2}} \\ &\geq \frac{(1 + (C_2/C_1)\rho^2)Q - (|\rho|(1 + C_2/C_1))W}{\sqrt{1 + 3\rho^2}} \\ &\geq \frac{(1 - \rho^2/2)Q - 2|\rho|W}{\sqrt{1 + 3\rho^2}} \\ &\geq \alpha Q, \end{aligned}$$

where: (i) the inequality  $\sqrt{C_1^2 + \rho^2C_2(C_2 + 2C_1)} \leq C_1\sqrt{1 + 3\rho^2}$  used in the third line follows because  $C_1 > 0$  and  $C_2 \leq C_1$  from (69); (ii) the penultimate inequality uses  $C_2/C_1 \in [-1/2, 1]$  (for the lower bound recall  $C_2 \geq -1/4, C_1 \geq 1/2$ ); (iii) the last inequality is because of the theorem's assumption in (13).

Next, recall that

$$\mathbf{w}_2^\gamma := \mathbf{G}_w(0, \mathbf{q}_1^\gamma) = \mathbb{E}_{(\mathbf{X}, y) \sim \mathcal{D}}[y\mathbf{X}^\top \phi(\mathbf{X}\mathbf{q}_1^\gamma)] = \mathbb{E}_{(\mathbf{X}, y) \sim \mathcal{D}}[y\mathbf{v}(\mathbf{X})], \quad (71)$$

where we set  $\mathbf{v}(\mathbf{X}) = \mathbf{X}^\top \phi(\mathbf{X}\mathbf{q}_1^\gamma)$  for convenience. Thus, the prediction of the model with parameters  $\theta = (\mathbf{w}_2^\gamma, \mathbf{q}_1^\gamma)$  for test datapoint  $(\tilde{\mathbf{X}}, \tilde{y})$  is

$$\hat{o} := \tilde{y}f_\theta(\tilde{\mathbf{X}}) = \tilde{y}\langle \mathbf{w}_2^\gamma, \mathbf{v}(\tilde{\mathbf{X}}) \rangle = \underbrace{\langle \mathbf{w}_2^\gamma, \tilde{y}\mathbf{q}_* + \mathbf{w}_* \rangle}_{\hat{o}_1} + \underbrace{\tilde{y}\langle \mathbf{w}_2^\gamma, \mathbf{v}(\tilde{\mathbf{X}}) - (\mathbf{q}_* + \tilde{y}\mathbf{w}_*) \rangle}_{\hat{o}_2}. \quad (72)$$

Let  $\mathcal{E}$  denote the event for which  $\|\mathbf{v}(\tilde{\mathbf{X}}) - (\mathbf{q}_* + \tilde{y}\mathbf{w}_*)\| \leq \epsilon$ , which has probability at least  $1 - \delta$  by (70). Note that

$$\text{ERR}(f_\theta) = \mathbb{P}(\hat{\delta} < 0) \leq \mathbb{P}(\hat{\delta} < 0 | \mathcal{E}) + \Pr(\mathcal{E}) = \mathbb{P}(\hat{\delta} < 0 | \mathcal{E}) + \delta. \quad (73)$$

Hence, our goal below is to bound  $\mathbb{P}(\hat{\delta} < 0 | \mathcal{E})$ . In fact, we will show that  $\mathbb{P}(\hat{\delta} < 0 | \mathcal{E}) \leq 0$ , so the error rate is  $\delta$  as stated in the theorem.

To do this, condition on  $\mathcal{E}$  for which  $|\hat{\delta}_2| \leq \epsilon \|\mathbf{w}_2^\gamma\|$ , thus

$$\hat{\delta} \geq \hat{\delta}_1 - \epsilon \|\mathbf{w}_2^\gamma\|.$$

Further note that  $\hat{\delta}_1 = \langle \mathbf{w}_2^\gamma, \mathbf{w}_* \rangle - |\langle \mathbf{w}_2^\gamma, \mathbf{q}_* \rangle|$ . Thus, it suffices to show that

$$\langle \mathbf{w}_2^\gamma, \mathbf{w}_* \rangle \geq |\langle \mathbf{w}_2^\gamma, \mathbf{q}_* \rangle| + \epsilon \|\mathbf{w}_2^\gamma\|. \quad (74)$$

For this, go back to  $\mathbf{w}_2^\gamma$  and write continuing from (71)

$$\begin{aligned} \mathbf{w}_2^\gamma &= \mathbb{E}_{(\mathbf{X}, y) \sim \mathcal{D}} [y\mathbf{q}_* + \mathbf{w}_*] + \mathbb{E}_{(\mathbf{X}, y) \sim \mathcal{D}} [y(\mathbf{v}(\mathbf{X}) - (\mathbf{q}_* + y\mathbf{w}_*))] \\ &= \mathbf{w}_* + \mathbb{E}_{(\mathbf{X}, y) \sim \mathcal{D}} [y(\mathbf{v}(\mathbf{X}) - (\mathbf{q}_* + y\mathbf{w}_*))]. \end{aligned}$$

Denote for convenience  $\mathbf{e} := \mathbf{e}(y, \mathbf{X}) := \mathbf{v}(\mathbf{X}) - (\mathbf{q}_* + y\mathbf{w}_*)$ . Thus,

$$\begin{aligned} \langle \mathbf{w}_2^\gamma, \mathbf{w}_* \rangle &= W^2 + \mathbb{E}[y\langle \mathbf{w}_*, \mathbf{e} \rangle] \\ \langle \mathbf{w}_2^\gamma, \mathbf{q}_* \rangle &= \rho QW + \mathbb{E}[y\langle \mathbf{q}_*, \mathbf{e} \rangle] \\ \|\mathbf{w}_2^\gamma\| &\leq W + \mathbb{E}[\|\mathbf{e}\|]. \end{aligned}$$

where in the last line we used triangle and Jensen's inequalities. We now compute (recall  $\mathcal{E}$  is the event for which  $\|\mathbf{e}\| \leq \epsilon$ ):

$$\begin{aligned} \mathbb{E}[\langle \mathbf{w}_*, \mathbf{e} \rangle] &\leq \mathbb{E}[\langle \mathbf{w}_*, \mathbf{e} \rangle | \mathcal{E}] + \mathbb{E}[\langle \mathbf{w}_*, \mathbf{e} \rangle | \mathcal{E}^c] \mathbb{P}(\mathcal{E}) \\ &\leq \epsilon W + W \delta \mathbb{E}[\|\mathbf{e}\| | \mathcal{E}^c]. \end{aligned} \quad (75)$$

Similarly, we can upper bound  $\mathbb{E}[\langle \mathbf{q}_*, \mathbf{e} \rangle]$  and  $\mathbb{E}[\|\mathbf{e}\|]$ . Combining these with the above displays, the desired Eq. (74) holds provided:

$$W^2 - |\rho|QW \geq \epsilon(W + Q) + \delta(W + Q)B + \epsilon\delta B + \epsilon^2. \quad (76)$$

Above, we have denoted  $B := \mathbb{E}[\|\mathbf{e}\| | \mathcal{E}^c]$ .  $\|\mathbf{w}_2^\gamma\| \leq W + B$ . Note that the LHS of 76 is  $> 0$  because of Assumption 3.a that  $|\rho| \leq W/Q$ . Thus, we can guarantee (76) holds once  $\epsilon$  is small enough (by making  $\gamma$  large enough) and  $\delta$  is also small enough (by making  $\gamma$  large enough). It only remains to bound  $B$ . To do this, note that

$$\|\mathbf{e}\|_2 \leq \|\mathbf{v}(\mathbf{X})\| + Q + W \leq \max_{t \in [T]} \|\mathbf{x}_t\| + Q + W \leq 2(Q + W) + \max_{t \in \mathcal{R}^c} \|\mathbf{z}_t\|.$$

By Lemma 15 we further have that

$$\mathbb{E} \left[ \max_{t \in \mathcal{R}^c} \|\mathbf{z}_t\|_2 | \mathcal{E}^c \right] \mathbb{P}(\mathcal{E}^c) \leq \delta \cdot C\sigma\sqrt{d}\sqrt{\log(2T/\delta)}.$$

Hence,

$$B \leq 2(Q + W) + C\sigma\sqrt{d}\sqrt{\log(2T/\delta)}.$$

### D.3.1. AUXILIARY LEMMA

**Lemma 15** (Subgaussian euclidean-norm tail control). *Let  $\mathbf{z}_i \in \mathbb{R}^d, i \in [N]$  be  $K$ -subgaussian random vectors. Then, for any event  $\mathcal{E}$  with  $\mathbb{P}(\mathcal{E}) = \delta$ , it holds that*

$$\mathbb{E} \left[ \max_{i \in [N]} \|\mathbf{z}_i\| | \mathcal{E}^c \right] \leq 12K\sqrt{d}\sqrt{\log(2N/\delta)}.$$



*Proof.* Set  $Z = \max_{i \in [n]} \|z_i\|$  and define event  $\mathcal{B} = \{Z \geq M\}$  for  $M := 4K\sqrt{d}\sqrt{\log(2N/\delta)}$ . By Fact G.4 for all  $t > 0$ ,  $\mathbb{P}(Z > t) \leq 2Ne^{-t^2/(16dK^2)}$ . Thus, by choice of  $M$ ,  $\mathbb{P}(\mathcal{B}) \leq \mathbb{P}(\mathcal{E}) = \delta$ .

Denote the pdf and cdf complement of  $Z$  by  $f_Z, Q_Z$  respectively. Observe that, we set  $Q_Z(M) \leq \delta$ . Using integration by parts we have,

$$\begin{aligned} \mathbb{E}[Z|\mathcal{B}]\mathbb{P}(\mathcal{B}) &= \int_M^\infty z f_Z(z) dz = - \int_M^\infty z dQ_Z(z) = \int_M^\infty Q_Z(z) dz - [Q_Z(z)z]_M^\infty \\ &= \int_M^\infty Q_Z(z) dz + \delta M \\ &= \delta M + \int_M^\infty \mathbb{P}(Z \geq t) dt \leq \delta M + \int_M^\infty 2Ne^{-t^2/(16dK^2)} dt \\ &\leq \delta M + 2\sqrt{2}K\sqrt{d}(2N) \int_{\sqrt{2\log(2N/\delta)}}^\infty e^{-u^2/2} du \\ &= \delta 4K\sqrt{d}\sqrt{\log(2N/\delta)} + 2\sqrt{\pi}K\sqrt{d}\delta \leq 2\delta M. \end{aligned}$$

We can conclude the proof by noting:

$$\begin{aligned} \mathbb{E}[Z|\mathcal{E}]\mathbb{P}(\mathcal{E}) &= \mathbb{E}[Z|\mathcal{E} \cap \mathcal{B}^c]\mathbb{P}(\mathcal{E} \cap \mathcal{B}^c) + \mathbb{E}[Z|\mathcal{E} \cap \mathcal{B}]\mathbb{P}(\mathcal{E} \cap \mathcal{B}) \\ &\leq M\delta + \mathbb{E}[Z|\mathcal{B}]\mathbb{P}(\mathcal{B}). \end{aligned}$$

□

## E. Proofs of results on discrete datasets

### E.1. Proof of Theorem 1 and Observation 1

• **Proof for Prompt-attention:** Let  $\bar{w}_* = w_*/\|w_*\|$  and  $\bar{q}_* = q_*/\|q_*\|$ .  $q'_*$  be the projection of  $q_*$  to the orthogonal complement of  $w_*$  i.e.  $q'_* = q_* - \bar{w}_*\bar{w}_*^\top q_*$ . Similarly, let  $w'_*$  be the projection of  $w_*$  to the orthogonal complement of  $q_*$  i.e.  $w'_* = w_* - \bar{q}_*\bar{q}_*^\top w_*$ . Denote correlation coefficient between two vectors by  $\rho(a, b) = \frac{a^\top b}{\|a\|\|b\|}$ .

To proceed, observe that,  $q_*^\top q_* = \|q_*\|^2 - (\bar{w}_*^\top q_*)^2 = \|q_*\|^2(1 - \rho(q_*, w_*)^2) > 0$ . The positivity follows from the fact that  $q_*, w_*$  are not parallel, thus, the absolute value of their correlation coefficient is strictly bounded away from 1. Similarly  $w_*^\top w_* = \|w_*\|^2(1 - \rho(q_*, w_*)^2) > 0$ . To proceed, set  $\bar{\rho} := 1 - \rho(q_*, w_*)^2$  and observe that the classifier  $\theta = (w'_*, \Gamma q'_*)$  achieves the attention scores

$$a_i = \phi(\mathbf{X}q'_*)_i = \begin{cases} S^{-1}e^{\|q_*\|^2\Gamma\bar{\rho}} & \text{if } i \text{ relevant,} \\ S^{-1}e^{-\|q_*\|^2\Gamma\delta^q\bar{\rho}} & \text{if } i \text{ irrelevant,} \end{cases}$$

where  $S = T\zeta e^{\|q_*\|^2\Gamma\bar{\rho}} + T(1-\zeta)e^{-\|q_*\|^2\Gamma\delta^q\bar{\rho}}$ . Using orthogonality of  $w'_*$  and  $q_*$ , the final prediction obeys

$$yf_\theta(\mathbf{X}) = \|w'_*\|^2\bar{\rho}S^{-1} \left[ \zeta e^{\|q_*\|^2\Gamma\bar{\rho}} - \delta^w(1-\zeta)e^{-\|q_*\|^2\Gamma\delta^q\bar{\rho}} \right].$$

The classifier achieves perfect accuracy when  $\zeta e^{\|q_*\|^2\Gamma\bar{\rho}} > |\delta^w|(1-\zeta)e^{-\|q_*\|^2\Gamma\delta^q\bar{\rho}}$ . Since we have  $\delta^q \geq 0$  and we have assumed  $\delta^w$  is a  $C$ -bounded variable (i.e.  $|\delta^w| \leq C$ ), thus, the desired inequality can be guaranteed by choosing

$$\Gamma > \frac{1}{\|q_*\|^2\bar{\rho}} \log\left(\frac{C(1-\zeta)}{\zeta}\right).$$

• **Proof for Observation 1:** To prove this, observe that for any  $\delta^q = \Delta^q$ ,  $\delta^w = \Delta^w$  choices, using orthogonality of  $q_*, w'_*$ , for any  $(y, \mathbf{X}) \sim \mathcal{D}$ , we have

$$yf^{\text{LIN}}(w'_*) = \|w'_*\|^2\bar{\rho}(\zeta - (1-\zeta)\delta^w).$$

Thus, as long as  $\delta^w \neq \zeta/(1-\zeta)$ ,  $\text{sign}(yf^{\text{LIN}}(w'_*))$  is always 1 or always -1, resulting in perfect accuracy for  $w'_*$  or  $-w'_*$ .

• **Proof for Self-attention:** The proof is provided under Theorem 7.

• **Proof for Linear Prompt-attention:** Let  $W_1 = \mathbf{w}^\top \mathbf{w}_*$ ,  $W_2 = \mathbf{w}^\top \mathbf{q}_*$ ,  $Q_1 = \mathbf{q}^\top \mathbf{q}_*$ ,  $Q_2 = \mathbf{q}^\top \mathbf{w}_*$ . Since context-irrelevant tokens are of the form  $-\delta^q \mathbf{q}_* - y \delta^w \mathbf{w}_*$ , the model decision is given by  $\frac{1}{T} f(\mathbf{X}) = \frac{1}{T} \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{q} = \zeta \mathbf{w}^\top (y \mathbf{w}_* + \mathbf{q}_*) (y \mathbf{w}_* + \mathbf{q}_*)^\top \mathbf{q} + (1 - \zeta) \mathbf{w}^\top (y \delta^w \mathbf{w}_* + \delta^q \mathbf{q}_*) (y \delta^w \mathbf{w}_* + \delta^q \mathbf{q}_*)^\top \mathbf{q}$  and

$$\begin{aligned} \frac{1}{T} f(\mathbf{X}) &= \zeta (y W_1 + W_2) (Q_1 + y Q_2) + (1 - \zeta) (y \delta^w W_1 + \delta^q W_2) (y \delta^w Q_2 + \delta^q Q_1) \\ &= \zeta y (W_1 Q_1 + W_2 Q_2) + \zeta (W_2 Q_1 + W_1 Q_2) + \\ &\quad (1 - \zeta) y \delta^q \delta^w (W_1 Q_1 + W_2 Q_2) + (1 - \zeta) (\delta^{q^2} W_2 Q_1 + \delta^{w^2} W_1 Q_2). \\ \frac{y f(\mathbf{X})}{T} &= (\zeta + (1 - \zeta) \delta^q \delta^w) (W_1 Q_1 + W_2 Q_2) + y ((\zeta + (1 - \zeta) \delta^{q^2}) W_2 Q_1 + (\zeta + (1 - \zeta) \delta^{w^2}) W_1 Q_2). \end{aligned}$$

To proceed, set  $(\delta^q, \delta^w)$  to be  $(0, 0)$  or  $(\Delta, -\Delta)$  equally-likely for  $\Delta > \sqrt{\zeta/(1-\zeta)}$ . For fixed  $\Delta$ , for any choice of  $W_1, W_2, Q_1, Q_2$  observe that, with  $1/2$  probability the event  $E = \{y((\zeta + (1 - \zeta) \delta^{q^2}) W_2 Q_1 + (\zeta + (1 - \zeta) \delta^{w^2}) W_1 Q_2) \leq 0\}$  happens. On this event (which is over the label  $y$ ), probability that  $(\zeta + (1 - \zeta) \delta^q \delta^w) (W_1 Q_1 + W_2 Q_2) > 0$  is at most  $1/2$  because  $\text{sign}(\zeta + (1 - \zeta) \delta^q \delta^w)$  is Rademacher variable. Combining, we find that  $\mathbb{P}(\frac{y f(\mathbf{X})}{T} \leq 0) \geq 25\%$  as advertised whenever  $\Delta > \sqrt{\zeta/(1-\zeta)}$ .

## E.2. Failure proof for Self-attention

We have the following theorem regarding self-attention.

**Theorem 7.** Fix  $\Delta > 0$  to be sufficiently large. In (DATA), choose  $\delta = (\delta^q, \delta^w)$  to be  $(0, 0)$  or  $(\Delta, \Delta)$  equally-likely, where  $\Delta > 1/(1-\zeta)^2$ .

- For any choice of  $(\mathbf{U} = \mathbf{1} \mathbf{u}^\top, \mathbf{W})$ ,  $f^{\text{SAT}}(\mathbf{1} \mathbf{u}^\top, \mathbf{W})$  achieves 50% accuracy (i.e. random guess).
- For any choice of  $(\mathbf{U}, \mathbf{W})$ , there exists a (DATA) distribution with adversarial relevance set choices such that  $f^{\text{SAT}}(\mathbf{U}, \mathbf{W})$  achieves 50% accuracy.

Here, adversarial relevance set choice means that, the relevance set can be chosen adaptively to the label  $y$ , out-of-context term  $\delta$ , and the self-attention model weights  $(\mathbf{U}, \mathbf{W})$  to cause misclassification.

*Proof.* Let  $\tilde{\mathbf{w}} = \mathbf{W} \mathbf{w}_*$  and  $\tilde{\mathbf{q}} = \mathbf{W} \mathbf{q}_*$ . Also let  $b_w = \mathbf{u}^\top \mathbf{w}_*$  and  $b_q = \mathbf{u}^\top \mathbf{q}_*$ . Since  $\mathbf{W}$  is allowed to be full-rank and arbitrary,  $\tilde{\mathbf{w}}, \tilde{\mathbf{q}}$  are allowed to be arbitrary as well (but fixed given  $\mathbf{W}$ ). In our analysis, the critical terms are the attention weights given by the correlation between the relevant/irrelevant keys/queries.

Setting attention queries as the raw tokens (without losing any generality), relevant queries  $\mathbf{x}_R$  and keys  $\mathbf{k}_R$  become

$$\mathbf{x}_R = y \mathbf{w}_* + \mathbf{q}_*, \quad \mathbf{k}_R = y \tilde{\mathbf{w}} + \tilde{\mathbf{q}}.$$

Thanks to our choice of  $\delta := \delta^w = \delta^q$  to be equally-likely in  $\{0, \Delta\}$ , observe that irrelevant queries and keys are simply

$$\mathbf{x}_I = -\delta \mathbf{x}_R, \quad \mathbf{k}_I = -\delta \mathbf{k}_R.$$

This will greatly help the proof because it will mean that attention weights are highly structured. Specifically, set  $\rho = \mathbf{x}_R^\top \mathbf{k}_R$ . All weights of the attention similarities belong to the set  $(\rho, -\delta \rho, \delta^2 \rho)$ . Consequently, softmax-attention output

$$\mathbf{A} = \phi(\mathbf{X} \mathbf{W} \mathbf{X}^\top) \mathbf{X} = \begin{bmatrix} \mathbf{a}_1^\top \\ \vdots \\ \mathbf{a}_T^\top \end{bmatrix} \text{ is given by}$$

$$\mathbf{a}_i = \begin{cases} \frac{\zeta e^\rho - \delta(1-\zeta)e^{-\delta\rho}}{\zeta e^\rho + (1-\zeta)e^{-\delta\rho}} \cdot \frac{1}{T} \cdot \mathbf{x}_R & \text{if } i \in \mathcal{R} \text{ (relevant)} \\ \frac{\zeta e^{-\delta\rho} - \delta(1-\zeta)e^{\delta^2\rho}}{\zeta e^{-\delta\rho} + (1-\zeta)e^{\delta^2\rho}} \cdot \frac{1}{T} \cdot \mathbf{x}_R & \text{if } i \in \mathcal{R}^c \text{ (irrelevant)} \end{cases} \quad (77)$$

Set  $a_+ = \frac{e^\rho}{\zeta e^\rho + (1-\zeta)e^{-\delta\rho}}$ ,  $a_- = \frac{e^{-\delta\rho}}{\zeta e^\rho + (1-\zeta)e^{-\delta\rho}}$ ,  $b_- = \frac{e^{-\delta\rho}}{\zeta e^{-\delta\rho} + (1-\zeta)e^{\delta^2\rho}}$ ,  $b_+ = \frac{e^{\delta^2\rho}}{\zeta e^{-\delta\rho} + (1-\zeta)e^{\delta^2\rho}}$ . With this, we also set

$$\Delta_R = \frac{\zeta e^\rho - \delta(1-\zeta)e^{-\delta\rho}}{\zeta e^\rho + (1-\zeta)e^{-\delta\rho}} = \zeta a_+ - \delta(1-\zeta)a_-$$

$$\Delta_I = \frac{\zeta e^{-\delta\rho} - \delta(1-\zeta)e^{\delta^2\rho}}{\zeta e^{-\delta\rho} + (1-\zeta)e^{\delta^2\rho}} = \zeta b_- - \delta(1-\zeta)b_+.$$

Also define  $\Delta_i = \Delta_R$  if  $i$  is relevant and  $\Delta_I$  otherwise. With this, we have  $\mathbf{a}_i = \Delta_i \mathbf{x}_R$  based on (77).

The following lemma will be helpful for the downstream analysis. The goal of this lemma is showing that, by choosing  $\delta \in \{0, \Delta\}$ , we can confuse the model output.

**Lemma 16.** Fix a scalar  $\kappa$ . Set  $f_\delta = \kappa \Delta_R + (1-\kappa) \Delta_I$ . Recalling  $\rho = \mathbf{x}_R^\top \mathbf{k}_R$ , the following statements hold:

- Set  $\delta = 0$ . Suppose “ $1 \geq \kappa \geq 0$ ” OR “ $\kappa \geq 1, \rho \geq 0$ ” OR “ $\kappa \leq 0, \rho \leq 0$ ”. Then  $f_\delta > 0$ .
- Fix  $0 \leq \alpha \leq 1$ . Suppose

$$\delta > \Delta_0 := \frac{1}{1-\zeta} \max\left(\frac{\zeta}{\alpha(1-\zeta)}, \frac{1}{1-\alpha}\right).$$

and “ $\kappa \leq \alpha, \rho \geq 0$ ” OR “ $\kappa \geq \alpha, \rho \leq 0$ ”. Then  $f_\delta < 0$ .

*Proof.* Plugging in  $\delta$ , we write

$$f_\delta = \kappa \Delta_R + (1-\kappa) \Delta_I = \kappa \zeta a_+ - \delta \kappa (1-\zeta) a_- + \zeta (1-\kappa) b_- - \delta (1-\zeta) (1-\kappa) b_+ \quad (78)$$

$$= \zeta (\kappa a_+ + (1-\kappa) b_-) - \delta (1-\zeta) (\kappa a_- + (1-\kappa) b_+). \quad (79)$$

- **Suppose  $\delta = 0$ .** In this case, we obtain the first statement of the lemma as follows

$$f_\delta / \zeta = \frac{\kappa e^\rho}{\zeta e^\rho + 1 - \zeta} + 1 - \kappa > 0 \quad \text{whenever} \quad \begin{cases} 1 \geq \kappa \geq 0 & \text{OR} \\ \kappa \geq 1, \rho \geq 0 & \text{OR} \\ \kappa \leq 0, \rho \leq 0 & \end{cases} \quad (80)$$

- **Now suppose  $\delta > \Delta_0$ . First, assume  $\rho \geq 0$  and  $\kappa \leq \alpha$ .** We use the facts

$$1/\zeta \geq a_+ \geq 1, \quad 1 \geq a_- \geq 0, \quad b_+ \geq 1, \quad 1 \geq b_- \geq 0.$$

Observe that, since  $b_+ \geq a_-$  and  $\kappa \leq \alpha$

$$\kappa a_- + (1-\kappa) b_+ \geq \begin{cases} b_+ & \text{if } \kappa \leq 0 \\ (1-\alpha) b_+ & \text{if } \kappa \geq 0 \end{cases} \geq 1 - \alpha.$$

Additionally, if  $\kappa \leq 0$ , we have that

$$\kappa a_- + (1-\kappa) b_+ \geq b_+ \geq b_- \geq \kappa a_+ + (1-\kappa) b_-.$$

If  $\kappa \leq 0$ , we obtain  $f_\delta \leq \zeta b_- - \delta(1-\zeta) b_+$ . Thus,  $f_\delta < 0$  whenever  $\delta > \Delta_0 \geq \zeta / (1-\zeta)$ .

If  $\kappa \geq 0$ , we use  $\kappa a_+ + (1-\kappa) b_- \leq 1/\zeta$  to obtain that whenever  $\delta > \Delta_0 \geq \frac{1}{(1-\zeta)(1-\alpha)}$

$$f_\delta \leq 1 - \delta(1-\zeta)(1-\alpha) < 0.$$

**Now assume  $\rho \leq 0$  and  $\kappa \geq \alpha$ .** We use the facts

$$1 \geq a_+ \geq 0, \quad a_- \geq 1, \quad 1 \geq b_+ \geq 0, \quad \frac{1}{1-\zeta} \geq b_- \geq 1.$$

Observe that, since  $b_+ \leq a_-$  and  $\kappa \geq \alpha$

$$\kappa a_- + (1 - \kappa)b_+ \geq \begin{cases} a_- & \text{if } \kappa \geq 1 \\ \alpha a_- & \text{if } \kappa \leq 1 \end{cases} \geq \alpha.$$

Additionally, if  $\kappa \geq 1$ , we have that

$$\kappa a_- + (1 - \kappa)b_+ \geq a_- \geq a_+ \geq \kappa a_+ + (1 - \kappa)b_-.$$

If  $\kappa \geq 1$ , we obtain  $f_\delta \leq \zeta a_+ - \delta(1 - \zeta)a_-$ . Thus,  $f_\delta < 0$  whenever  $\delta > \Delta_0 \geq \zeta/(1 - \zeta)$ .

If  $\kappa \leq 1$ , we use  $\kappa a_+ + (1 - \kappa)b_- \leq \frac{1}{1 - \zeta}$  to obtain that whenever  $\delta > \Delta_0 \geq \frac{\zeta}{(1 - \zeta)^2 \alpha}$

$$f_\delta \leq \frac{\zeta}{1 - \zeta} - \delta(1 - \zeta)\alpha < 0.$$

□

To proceed, we will conclude with the proof as follows. Set  $\nu_i = y \mathbf{u}_i^\top \mathbf{x}_R$  for  $i \in [T]$  where  $\mathbf{u}_i$  is the  $i$ th row of the output layer weights  $\mathbf{U}$ . Here  $\nu_i$  is obviously  $y$ -dependent. However, we will show that for any choice of  $y$ , the model accuracy is at most 50%. Towards this we fix  $y$  and (mostly) omit it from the notation during the following discussion. Let  $\mathbf{a}_i$  be the  $i$ th token of the attention output. The linear output layer  $\mathbf{U}$  aggregates  $\mathbf{u}_i^\top \mathbf{a}_i$  to obtain

$$yf(\mathbf{U}, \mathbf{W}) = \sum_{i=1}^T \mathbf{u}_i^\top \mathbf{a}_i = \sum_{i=1}^T \nu_i \Delta_i.$$

Aggregating  $\nu_+ = \frac{1}{T} \sum_{i \in \mathcal{R}=\text{relevant}} \nu_i$  and  $\nu_- = \frac{1}{T} \sum_{i \in \mathcal{R}^c=\text{irrelevant}} \nu_i$  and recalling from (77) that over relevant/irrelevant sets attention tokens are given by  $\Delta_R \mathbf{x}_R$  and  $\Delta_I \mathbf{x}_I$ , we find

$$\frac{1}{T} yf(\mathbf{U}, \mathbf{W}) = \nu_R \Delta_R + \nu_I \Delta_I.$$

**Scenario 1: Rows of  $\mathbf{U}$  are identical and we have  $\mathbf{U} = \mathbf{1} \mathbf{u}^\top$ .** In this scenario, we simply have  $\nu_i = \nu$  and  $\nu_R = \zeta \nu$  and  $\nu_I = (1 - \zeta)\nu$ . Thus, we find

$$\frac{1}{T} yf(\mathbf{U}, \mathbf{W}) = T\nu[\zeta \Delta_R + (1 - \zeta)\Delta_I].$$

Set  $f_\delta = \zeta \Delta_R + (1 - \zeta)\Delta_I$ . We claim that  $\text{sign}(f_\delta)$  is Rademacher (given arbitrary  $y$  choice) which will prove that accuracy is at most 50%. Specifically, let us apply Lemma 16 with  $\kappa = \zeta$  and  $\alpha = \zeta$ . When  $\delta = 0$ , we have  $f_\delta > 0$ . When  $\delta = \Delta$ , since the conditions  $\kappa \leq \alpha$  and  $\kappa \geq \alpha$  hold, for any choice of  $\rho$ , for  $\Delta > \Delta_0 := \frac{1}{(1 - \zeta)^2}$  we have that  $f_\delta < 0$ . Thus, we have that  $\mathbb{P}_\delta(f_\delta > 0) = \mathbb{P}_\delta(f_\delta < 0) = 0.5$  as advertised. This follows from the fact that  $f_\delta > 0$  for  $\delta = 0$  and  $f_\delta < 0$  for  $\delta = \Delta$  and  $\delta$  is equally likely over two options.

**Scenario 2:** Suppose rows of  $\mathbf{U}$  are not identical. In this case, we will leverage the fact that relevant set  $\mathcal{R}$  is allowed to be chosen adversarially with respect to the self-attention weights. We will show that by selecting  $\mathcal{R}$  adversarially, on any label  $y$  event, accuracy is a coin flip.

**First consider the scenario  $\nu_{\text{tot}} := \nu_R + \nu_I \leq 0$ :** We will show that model achieves at least 50% error on label  $y$ : Let us denote  $\nu_R$  with  $\nu_R^\mathcal{R}$  which makes the relevance set dependence explicit. Given  $\mathcal{R}$ , fixing  $\delta = 0$ , the model outputs (following (80))

$$\frac{1}{T} yf(\mathbf{U}, \mathbf{W}) = \frac{\nu_R^\mathcal{R} e^\rho}{\zeta e^\rho + 1 - \zeta} + \nu_I^\mathcal{R}.$$

Suppose there is a relevance set  $\mathcal{R}_0$  (that depends on  $y$ ) such that the right hand side is non-positive. Let us select this  $\mathcal{R}_0$  as our relevance set. Then, the model makes 50% error on label  $y$  thanks to the event  $\delta = 0$  (which is exactly what we want). If there is no such  $\mathcal{R}_0$ , then, for all  $\mathcal{R}$ , we have

$$\frac{\nu_R^\mathcal{R} e^\rho}{\zeta e^\rho + 1 - \zeta} + \nu_I^\mathcal{R} > 0$$

By taking average of all relevance sets (“ $T$  choose  $\zeta T$ ” many), all  $v_i$ ’s will be equally-weighted and we obtain  $\nu_{\text{tot}} = \nu_R + \nu_I > 0$ . This contradicts with our initial  $\nu_{\text{tot}} \leq 0$  assumption, thus,  $\mathcal{R}_0$  has to exist.

**Now consider the scenario  $\nu_{\text{tot}} = \nu_R + \nu_I > 0$ :** Let  $\mathcal{D}$  be the uniform distribution over “ $T$  choose  $\zeta T$ ” relevant sets  $\mathcal{R}$ . Clearly  $\mathbb{E}_{\mathcal{D}}[\nu_R^{\mathcal{R}}] = \zeta \nu_{\text{tot}} > 0$ . Thus, there is a relevance set  $\mathcal{R}_+$  such that  $\nu_R^{\mathcal{R}_+} \geq \zeta \nu_{\text{tot}}$  and there is a relevance set  $\mathcal{R}_-$  such that  $\nu_R^{\mathcal{R}_-} \leq \zeta \nu_{\text{tot}}$ . We will make use of these two sets to finalize the proof.

To proceed, set  $\kappa_{\pm} = \nu^{\mathcal{R}_{\pm}} / \nu_{\text{tot}}$  and again set  $\alpha = \zeta$  and  $\Delta_0 = \frac{1}{(1-\zeta)^2}$  in Lemma 16. Here, we are investigating the sign of the prediction

$$\frac{1}{T\nu_{\text{tot}}} yf(\mathbf{U}, \mathbf{W}) = \frac{\nu_R \Delta_R + \nu_I \Delta_I}{\nu_{\text{tot}}} = \kappa_{\pm} \Delta_R + (1 - \kappa_{\pm}) \Delta_I.$$

First, assume that the attention weights are so that  $\rho = \rho_y \geq 0$ . In this case (and for this particular label  $y$ ),

- When  $\delta = 0$ , we choose the relevance set  $\mathcal{R}_+$  which ensures  $\kappa_+ \geq \zeta \geq 0$  and  $f_0 > 0$ .
- When  $\delta = \Delta > \Delta_0$ , we choose the relevance set  $\mathcal{R}_-$  which ensures  $\kappa_- \leq \zeta$  and  $f_{\Delta} < 0$ .

Secondly, assume that the attention weights are so that  $\rho = \rho_y \leq 0$ . In this case,

- When  $\delta = 0$ , we choose the relevance set  $\mathcal{R}_-$  which ensures  $\kappa_- \leq \zeta \leq 1$  and  $f_0 > 0$ .
- When  $\delta = \Delta > \Delta_0$ , we choose the relevance set  $\mathcal{R}_+$  which ensures  $\kappa_+ \geq \zeta$  and  $f_{\Delta} < 0$ .

In either case, by adaptively choosing  $\mathcal{R} \in \{\mathcal{R}_+, \mathcal{R}_-\}$  as a function of  $(\delta, y)$  pair, we ensure accuracy is at most 50% because  $f_{\Delta}$  and  $f_0$  have conflicting signs.  $\square$

### E.3. Success proof for $\mathcal{R}$ -Adaptive Self-Attention

Consider the setting of Theorem 1 and Appendix E.2. We have the following lemma which shows that self-attention can succeed in Theorem 1 if  $\mathbf{U}$  can adapt to the relevance set (rather than  $\mathcal{R}$  being adversarial to  $\mathbf{U}$ ).

**Lemma 17.** *In (DATA), choose  $(\delta^q, \delta^w)$  to be  $(0, 0)$  or  $(\Delta, \Delta)$  equally-likely. Consider the self-attention model  $f^{\text{SAT}}(\mathbf{U}, \mathbf{W})$  where we set*

$$\mathbf{U} = \mathbb{1}_{\mathcal{R}} \mathbf{w}'_{\star} \Gamma \quad \text{and} \quad \mathbf{W} = \Gamma \mathbf{I}.$$

*This model achieves perfect accuracy whenever  $\mathbf{w}'_{\star} = (\mathbf{I} - \bar{\mathbf{q}}_{\star} \bar{\mathbf{q}}_{\star}^{\top}) \mathbf{w}_{\star} \neq 0$  by choosing*

$$\Gamma > \frac{1}{(1 + \Delta)(\|\mathbf{q}_{\star}\| - \|\mathbf{w}_{\star}\|)^2 + \|\mathbf{w}_{\star}\| \|\mathbf{q}_{\star}\| (1 - |\rho(\mathbf{q}_{\star}, \mathbf{w}_{\star})|)} \log\left(\Delta \frac{1 - \zeta}{\zeta}\right).$$

*where  $\rho(\cdot)$  is the correlation coefficient.*<sup>3</sup>

*Proof.* Thanks to the masking  $\mathbb{1}_{\mathcal{R}}$ , we only need to consider the attention scores along relevant tokens. Let  $c = \|\mathbf{y} \mathbf{w}_{\star} + \mathbf{q}_{\star}\|^2$ . For each relevant token, the attention rows are given by

$$\mathbf{a}_i = \begin{cases} e^{\Gamma c} & \text{if } i \in \mathcal{R} \\ e^{-\Delta \Gamma c} & \text{if } i \notin \mathcal{R}. \end{cases}$$

To proceed, attention tokens corresponding to relevant tokens are given by

$$\mathbf{f} = \sum_{i \in \mathcal{R}} \mathbf{a}_i (\mathbf{w}_{\star} + \mathbf{y} \mathbf{q}_{\star}) - \sum_{i \notin \mathcal{R}} \Delta \mathbf{a}_i (\mathbf{y} \mathbf{w}_{\star} + \mathbf{q}_{\star}) \quad (81)$$

$$= (\zeta e^{\Gamma c} - \Delta (1 - \zeta) e^{-\Delta \Gamma c}) (\mathbf{y} \mathbf{w}_{\star} + \mathbf{q}_{\star}). \quad (82)$$

<sup>3</sup>Note that the only instance  $\Gamma$  does not exist is when  $\mathbf{q}_{\star} = c \mathbf{w}_{\star}$  for  $|c| \geq 1$ . In this scenario, classification is impossible using the linear head  $\mathbf{w}'_{\star}$  without a bias term because all tokens are in the  $\text{sign}(c)$  direction regardless of the label  $y$ .

Thus, using  $\mathbf{w}'_* \mathbf{w}_* > 0$ ,

$$\text{sign}(y f^{\text{SATT}}(\mathbf{U}, \mathbf{W})) = \text{sign}(y \mathbf{w}'_* \mathbf{f}) = \text{sign}(\zeta e^{\Gamma c} - \Delta(1 - \zeta) e^{-\Delta \Gamma c}).$$

Thus, we need  $e^{(1+\Delta)\Gamma c} > \Delta \frac{1-\zeta}{\zeta}$  which is implied by  $\Gamma > \frac{1}{(1+\Delta)c} \log(\Delta \frac{1-\zeta}{\zeta})$ . To conclude, note that for both  $y = \pm 1$

$$c \geq \|\mathbf{y} \mathbf{w}_* + \mathbf{q}_*\|^2 \geq \|\mathbf{q}_*\|^2 + \|\mathbf{w}_*\|^2 - 2|\mathbf{q}_*^\top \mathbf{w}_*| \geq (\|\mathbf{q}_*\| - \|\mathbf{w}_*\|)^2 + \|\mathbf{w}_*\| \|\mathbf{q}_*\| (1 - |\rho(\mathbf{q}_*, \mathbf{w}_*)|) > 0.$$

where we used  $|\mathbf{q}_*^\top \mathbf{w}_*| = \|\mathbf{q}_*\| \|\mathbf{w}_*\| |\rho(\mathbf{q}_*, \mathbf{w}_*)|$ .  $\square$

## F. Proofs of sharp population risk formulas (Theorem 4)

Throughout this section, we use slightly different notation from the one stated in the main body for compactness purposes. Specifically, we set  $Q = \|\mathbf{q}_*\|^2$ ,  $W = T\|\mathbf{w}_*\|^2$  rather than  $Q = \|\mathbf{q}_*\|$ ,  $W = \|\mathbf{w}_*\|$ .

**Theorem 8.** Consider the prompt-attention model  $f_{\hat{\theta}}^{\text{ATT}}$ . Set  $Q = \|\mathbf{q}_*\|^2$ ,  $W = T\|\mathbf{w}_*\|^2$ , suppose  $\mathbf{w}_* \perp \mathbf{q}_*$ , and let  $\tau, \bar{\tau} > 0$  be hyperparameters. Consider the following algorithm which uses the hindsight knowledge of  $\mathbf{q}_*$  to estimate  $\mathbf{w}_*$  and make prediction:

1.  $\hat{\mathbf{w}} = (\mathbf{I} - \bar{\mathbf{q}}_* \bar{\mathbf{q}}_*^\top) \nabla \mathcal{L}_{\mathbf{w}}(0, \tau \bar{\mathbf{q}}_*)$ .
2. Set  $\hat{\theta} = (\hat{\mathbf{w}}, \bar{\tau} \bar{\mathbf{q}}_*)$ .

Suppose  $\zeta^2 W, 1 - \zeta, \alpha := n/d, e^Q, e^\tau$  each lie between two positive absolute constants. Suppose  $T$  is polynomially large in  $n$  and these constants and  $\tilde{O}(\cdot)$  hides polynomial terms in  $n$ . Define inverse-signal-to-noise-ratio:  $ISNR(\alpha, \tau) = \frac{(1-\zeta)e^{2\tau(\tau-\sqrt{Q})}}{\zeta^2 W \alpha}$ . With probability  $1 - 2e^{-t^2/2} - \tilde{O}(T^{-1/3})$  over the training data, the test error obeys

$$\text{ERR}(f_{\hat{\theta}}^{\text{ATT}}) = \mathcal{Q} \left( \frac{e^{\sqrt{Q}\bar{\tau}-\bar{\tau}^2}}{\sqrt{1 + (1 \mp \frac{1+t}{\sqrt{d}}) ISNR(\alpha, \tau)}} \cdot \sqrt{\frac{\zeta^2 W}{1 - \zeta}} \right) \pm \tilde{O}(T^{-1/3}).$$

Above,  $\mp, \pm$  highlights the upper/lower range of the test error (see (86) for exact statement). In the limit  $T, d \rightarrow \infty$ , the test error converges in probability to

$$\text{ERR}(\alpha, \zeta, Q, W, \tau, \bar{\tau}) = \mathcal{Q} \left( \frac{e^{\sqrt{Q}\bar{\tau}-\bar{\tau}^2}}{\sqrt{1 + ISNR(\alpha, \tau)}} \cdot \sqrt{\frac{\zeta^2 W}{1 - \zeta}} \right)$$

In this limit, optimal hyperparameters are  $\tau = \bar{\tau} = \sqrt{Q}/2$  and leads to optimal  $ISNR(\alpha) := \frac{(1-\zeta)e^{-Q/2}}{\zeta^2 W \alpha}$  and the error

$$\text{ERR}(\alpha, \zeta, Q, W) = \mathcal{Q} \left( \frac{e^{Q/4}}{\sqrt{1 + ISNR(\alpha)}} \cdot \sqrt{\frac{\zeta^2 W}{1 - \zeta}} \right)$$

*Proof.* Without losing generality, assume first  $\zeta T$  tokens are relevant and remaining tokens are irrelevant. Consider  $\mathbf{X}_I$  of size  $(1 - \zeta)T \times d$  induced by the irrelevant tokens with normal distribution. Observe that  $\mathbf{g} = \mathbf{X}_I \bar{\mathbf{w}}_*$  and  $\mathbf{h} = \mathbf{X}_I \bar{\mathbf{q}}_*$  are two independent i.i.d.  $\mathcal{N}(0, \mathbf{I}_{(1-\zeta)T})$  vectors. Also for standard normal  $g \sim \mathcal{N}(0, 1)$ , recall that moment-generating function is given by  $\mathbb{E}[e^{\tau g}] = e^{\tau^2/2}$ .

**Step 1: Characterizing the distribution of  $\hat{\mathbf{w}}$ .** Note that, the attention weights have the form  $\mathbf{a} = \phi(\tau \begin{bmatrix} \sqrt{Q} \mathbf{1}_{\zeta T} \\ \mathbf{h} \end{bmatrix})$ . Here,

the softmax denominator is  $T \cdot D_T$  where  $D_T := (\zeta e^{\sqrt{Q}\tau} + \frac{1}{T} \sum_{i=1}^{(1-\zeta)T} e^{\tau h_i})$ . Define  $e^{\tau \mathbf{h}}$  to be the numerator corresponding to irrelevant tokens i.e.

$$e^{\tau \mathbf{h}} = [e^{\tau h_1} \dots e^{\tau h_{(1-\zeta)T}}].$$

Define the matrix  $\mathbf{Q}_\perp = \mathbf{I} - \bar{\mathbf{q}}_\star \bar{\mathbf{q}}_\star^\top$ ,  $\mathbf{W}_\perp = \mathbf{I} - \bar{\mathbf{w}}_\star \bar{\mathbf{w}}_\star^\top$ . Set the vector  $\mathbf{v} = \frac{1}{T} \mathbf{Q}_\perp \mathbf{X}_I^\top e^{\tau \mathbf{h}}$  and  $\mathbf{v}_\perp = \frac{1}{T} \mathbf{h}^\top e^{\tau \mathbf{h}} \bar{\mathbf{q}}_\star$ . To proceed, observe that, for a single sample  $(y, \mathbf{X})$ , the gradient has the form

$$\nabla \mathcal{L}_{\mathbf{w}}^{y, \mathbf{X}}(0, \tau \bar{\mathbf{q}}_\star) = y \mathbf{X}^\top \mathbf{a} = \frac{\zeta(\mathbf{w}_\star + y \mathbf{q}_\star) e^{\sqrt{Q}\tau} + \mathbf{v} + \mathbf{v}_\perp}{D_T}. \quad (83)$$

After projection this onto the  $\mathbf{q}_\star$ -complement  $\mathbf{Q}_\perp$ , we get rid of the  $\mathbf{q}_\star$  direction to obtain

$$\hat{\mathbf{w}}_{y, \mathbf{X}} = \mathbf{Q}_\perp \nabla \mathcal{L}_{\mathbf{w}}^{y, \mathbf{X}}(0, \tau \bar{\mathbf{q}}_\star) = \frac{1}{D_T} [\zeta \mathbf{w}_\star e^{\sqrt{Q}\tau} + \mathbf{Q}_\perp \mathbf{X}_I^\top e^{\tau \mathbf{h}} / T].$$

The projected gradient over the full dataset is given by the empirical average

$$\hat{\mathbf{w}} = \mathbf{Q}_\perp \nabla \mathcal{L}_{\mathbf{w}}(0, \tau \bar{\mathbf{q}}_\star) = \frac{1}{n} \sum_{i=1}^n \frac{1}{D_{i,T}} [\zeta \mathbf{w}_\star e^{\sqrt{Q}\tau} + \mathbf{Q}_\perp \mathbf{X}_{i,I}^\top e^{\tau \mathbf{h}_i} / T].$$

Here  $\mathbf{h}_i, \mathbf{X}_{i,I}, D_{i,T}$  denote the random variables induced by the  $i$ th sample. Here, a critical observation is the fact that  $\mathbf{Q}_\perp \mathbf{X}_{i,I}$  is independent of  $\mathbf{h}_i$  (thanks to Gaussian orthogonality), thus,  $\mathbf{Q}_\perp \mathbf{X}_{i,I} e^{\tau \mathbf{h}_i}$  is normal conditioned on  $\mathbf{h}_i$ . To proceed, we apply Chebyshev's inequality over number of tokens  $T$ . Recall that we assumed  $e^\tau \leq C$  for an absolute constant  $C \geq 1$ . This means that  $e^{c\tau} \leq C^{c\tau} \leq C^{c \log C}$  is polynomial in  $C$  and is also upper bounded by a constant. In what follows  $\tilde{\mathcal{O}}(\cdot)$  only reflects the  $T$  dependence and subsumes polynomial dependence on the terms  $n, C$ . For all  $1 \leq i \leq n$ , applying Chebyshev's inequality, for  $T \gtrsim \text{poly}(n, e^{\tau^2})$ , with probability  $1 - T^{-1/3}$ , we have that

- Since  $\|e^{\tau \mathbf{h}_i}\|^2 / T = \frac{1}{T} \sum_{j=1}^T e^{2\tau \mathbf{h}_{ij}}$  thus  $\|e^{\tau \mathbf{h}_i}\|^2 / T - (1 - \zeta) e^{2\tau^2} \leq \tilde{\mathcal{O}}(T^{-1/3})$ ,
- Set  $\mathbb{E}[D_T] = D_\infty := \zeta e^{\sqrt{Q}\tau} + (1 - \zeta) e^{\tau^2/2}$ .  $|D_{i,T} - D_\infty| \leq \tilde{\mathcal{O}}(T^{-1/3})$ .

With these, set  $\mathbf{b}_i = \frac{\sqrt{1 - \zeta} e^{\tau^2}}{\|e^{\tau \mathbf{h}_i}\|} e^{\tau \mathbf{h}_i}$  which is a vector with fixed  $\ell_2$  norm that is perfectly parallel to  $e^{\tau \mathbf{h}_i}$ . Since  $\|\mathbf{b}_i\|^2 = \mathbb{E}[\|e^{\tau \mathbf{h}_i}\|^2 / T] = (1 - \zeta) e^{2\tau^2}$ , from above, observe that,

$$\|\mathbf{b}_i - \frac{1}{\sqrt{T}} e^{\tau \mathbf{h}_i}\| \leq \tilde{\mathcal{O}}(T^{-1/3}).$$

Now, let

$$\bar{\mathbf{v}} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{Q}_\perp \mathbf{X}_{i,I}^\top \mathbf{b}_i.$$

Since  $\mathbf{Q}_\perp \mathbf{X}_{i,I}^\top, \mathbf{b}_i$  are independent and  $\mathbf{b}_i$  has fixed  $\ell_2$  norm, we have that

$$\bar{\mathbf{v}} \sim \mathcal{N}(0, (1 - \zeta) e^{2\tau^2} \mathbf{Q}_\perp).$$

Finally, let  $\mathbf{c} = \zeta e^{\sqrt{Q}\tau} \mathbf{w}_\star$ . Recalling  $\sqrt{T} \|\mathbf{w}_\star\| = W$ , combining the perturbations bounds above, we have that

$$\sqrt{T} \|\mathbf{c} / D_\infty - \frac{1}{n} \sum_{i=1}^n \frac{1}{D_{i,T}} (\zeta \mathbf{w}_\star e^{\sqrt{Q}\tau})\| \leq \tilde{\mathcal{O}}(T^{-1/3}).$$

Combining these observe that

$$\|\sqrt{T} D_\infty \hat{\mathbf{w}} - \sqrt{T} \zeta e^{\sqrt{Q}\tau} \mathbf{w}_\star - \bar{\mathbf{v}} / \sqrt{n}\| \leq \tilde{\mathcal{O}}(T^{-1/3}). \quad (84)$$

Since  $\bar{\mathbf{v}}$  is normally distributed, above also implies that  $\sqrt{T} D_\infty \hat{\mathbf{w}}$  converges to the normal distribution  $\mathcal{N}(\sqrt{T} \zeta e^{\sqrt{Q}\tau} \mathbf{w}_\star, \frac{(1 - \zeta) e^{2\tau^2}}{n} \mathbf{Q}_\perp)$  in the limit  $T \rightarrow \infty$ .

**Lemma 18** (Inverse Signal-to-Noise Ratio (ISNR)). Set  $\mathbf{W}_\perp = \mathbf{I} - \bar{\mathbf{w}}_\star \bar{\mathbf{w}}_\star^\top$ . Define SNR of  $\hat{\mathbf{w}}$  to be

$$ISNR(\hat{\mathbf{w}}) = \frac{\|\mathbf{W}_\perp \hat{\mathbf{w}}\|^2}{\|\bar{\mathbf{w}}_\star^\top \hat{\mathbf{w}}\|^2}.$$

Recall  $\text{ISNR}(\alpha, \tau) = \frac{(1-\zeta)e^{2\tau(\tau-\sqrt{Q})}}{\zeta^2 W \alpha}$ . With probability  $1 - 2e^{-t^2/2} - T^{-1/3}$  over the dataset, we have that

$$\left(1 - \frac{t+1}{\sqrt{d}} - \tilde{O}(T^{-1/3})\right)_+^2 \leq \frac{\text{ISNR}(\hat{\mathbf{w}})}{\text{ISNR}(\alpha, \tau)} \leq \left(1 + \frac{\tau}{\sqrt{d}} + \tilde{O}(T^{-1/3})\right)^2.$$

*Proof.* Let us recall the standard normal concentration: For  $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_{d-1})$ ,  $\sqrt{d-1} \geq \mathbb{E}[\|\mathbf{g}\|] \geq \frac{d-1}{\sqrt{d}}$ . Thus, with probability  $1 - 2e^{-t^2/2}$ , through Lipschitz concentration,

$$\sqrt{d} + t \geq \|\mathbf{g}\| \geq \sqrt{d} - 1 - t.$$

This means that, with the same probability

$$\sqrt{d} + t \geq \frac{\|\bar{\mathbf{v}}\|}{\sqrt{1-\zeta}e^{\tau^2}} \geq (\sqrt{d} - 1 - t)_+.$$

We first upper bound  $\|\mathbf{W}_\perp \hat{\mathbf{w}}\|^2$ . Recalling (84),

$$\|\mathbf{W}_\perp \hat{\mathbf{w}} - \bar{\mathbf{v}}/\sqrt{n}\| \leq \tilde{O}(T^{-1/3}).$$

Thus,

$$(\sqrt{d} + t)^2 + \tilde{O}(T^{-1/3}) \geq \frac{n\|\mathbf{W}_\perp \hat{\mathbf{w}}\|^2}{(1-\zeta)e^{2\tau^2}} \geq (\sqrt{d} - 1 - t)_+^2 - \tilde{O}(T^{-1/3}).$$

Using  $\|\mathbf{w}_*\|^2 T = W$ , We similarly have that

$$\|\bar{\mathbf{w}}_*^\top \hat{\mathbf{w}}\|^2 - W\zeta^2 e^{2\sqrt{Q}\tau} \leq \tilde{O}(T^{-1/3}).$$

To conclude, with probability  $1 - 2e^{-t} - T^{-1/3}$ ,  $\text{ISNR}(\hat{\mathbf{w}})$  obeys

$$\frac{(\sqrt{d} + t)^2 + \tilde{O}(T^{-1/3})}{W\zeta^2 e^{2\sqrt{Q}\tau} - \tilde{O}(T^{-1/3})} \geq \frac{n}{(1-\zeta)e^{2\tau^2}} \text{ISNR}(\hat{\mathbf{w}}) \geq \frac{(\sqrt{d} - 1 - t)_+^2 - \tilde{O}(T^{-1/3})}{W\zeta^2 e^{2\sqrt{Q}\tau} + \tilde{O}(T^{-1/3})}$$

Rewriting this bound, we find

$$\left(1 + \frac{t}{\sqrt{d}} + \tilde{O}(T^{-1/3})\right)^2 \frac{(1-\zeta)e^{2\tau^2}}{\zeta^2 W \alpha e^{2\sqrt{Q}\tau}} \geq \text{ISNR}(\hat{\mathbf{w}}) \geq \left(1 - \frac{1+t}{\sqrt{d}} - \tilde{O}(T^{-1/3})\right)_+^2 \frac{(1-\zeta)e^{2\tau^2}}{\zeta^2 W \alpha e^{2\sqrt{Q}\tau}}.$$

Recalling the definition of  $\text{ISNR}(\alpha, \tau) = \frac{(1-\zeta)e^{2\tau(\tau-\sqrt{Q})}}{\zeta^2 W \alpha}$ , we conclude with the bound.  $\square$

**Step 2: Characterizing the error rate of  $\theta = (\hat{\mathbf{w}}, \bar{\tau}, \mathbf{q}_*)$ .** To achieve this goal, we will leverage Theorem 9. Since conditions of this theorem is satisfied (noticing that their  $\gamma$  is our  $\text{ISNR}(\hat{\mathbf{w}})$  which is upper bounded by a positive constant), for a new test point  $(y, \mathbf{X})$ , we have that

$$\left| \text{ERR}(f_{\theta}^{\text{ATT}}) - \mathcal{Q}\left(\frac{e^{\sqrt{Q}\bar{\tau}-\bar{\tau}^2}}{\sqrt{1 + \text{ISNR}(\hat{\mathbf{w}})}} \cdot \sqrt{\frac{\zeta^2 W}{1-\zeta}}\right) \right| \leq \tilde{O}(T^{-1/3}).$$

Using the Lipschitzness of the Q-function (i.e.  $\mathcal{Q}(x+\epsilon) - \mathcal{Q}(x) = \int_x^{x+\epsilon} e^{-t^2/2} dt \leq \epsilon$ ), as we have done in Theorem 9, we pull out the perturbation term  $\tilde{O}(T^{-1/3})$  within  $\text{ISNR}(\hat{\mathbf{w}})$  to obtain the advertised bound

$$\mathcal{Q}\left(\frac{e^{\sqrt{Q}\bar{\tau}-\bar{\tau}^2}}{\sqrt{1 + (1 + \frac{t}{\sqrt{d}})\text{ISNR}(\alpha, \tau)}} \cdot \sqrt{\frac{\zeta^2 W}{1-\zeta}}\right) - \tilde{O}(T^{-1/3}) \leq \text{ERR}(f_{\theta}^{\text{ATT}}) \leq \tag{85}$$

$$\mathcal{Q}\left(\frac{e^{\sqrt{Q}\bar{\tau}-\bar{\tau}^2}}{\sqrt{1 + (1 - \frac{1+t}{\sqrt{d}})\text{ISNR}(\alpha, \tau)}} \cdot \sqrt{\frac{\zeta^2 W}{1-\zeta}}\right) + \tilde{O}(T^{-1/3}). \tag{86}$$



To emphasize, this bound holds with probability  $1 - 2e^{-t^2/2} - \tilde{\mathcal{O}}(T^{-1/3})$  over a new test datapoint  $(y, \mathbf{X})$ . To see the optimal choices for  $\bar{\tau}, \tau$ , we need to optimize the error bound. This results in

$$\bar{\tau}_* = \arg \min_{\bar{\tau}} \sqrt{Q}\bar{\tau} - \bar{\tau}^2 = \sqrt{Q}/2 \quad (87)$$

$$\tau_* = \arg \min_{\tau} \text{ISNR}(\alpha, \tau) = 2\tau(\tau - \sqrt{Q}) = \sqrt{Q}/2. \quad (88)$$

□

**Theorem 9.** Consider the prompt-attention model  $f_{\theta}^{\text{ATT}}$  where we set  $\theta = (\mathbf{w}_* + \mathbf{p}, \tau \bar{\mathbf{q}}_*)$ . Set  $Q = \|\mathbf{q}_*\|^2, W = T\|\mathbf{w}_*\|^2$ . Here  $\tau$  is a tuning parameter and  $\mathbf{p}$  is a perturbation vector and assume all vectors are perpendicular i.e.  $\mathbf{p} \perp \mathbf{w}_* \perp \mathbf{q}_*$ . Set  $\gamma := \|\mathbf{p}\|^2/\|\mathbf{w}_*\|^2$  and suppose  $1 + \gamma, \zeta^2 W, 1 - \zeta, e^Q, e^\tau$  each lie between two positive absolute constants.  $\tilde{\mathcal{O}}(\cdot)$  subsumes polynomial dependencies in these constants. We have that

$$\left| \text{ERR}(f_{\theta}^{\text{ATT}}) - \mathcal{Q} \left( \frac{e^{\sqrt{Q}\tau - \tau^2}}{\sqrt{1 + \gamma}} \cdot \sqrt{\frac{\zeta^2 W}{1 - \zeta}} \right) \right| \leq \mathcal{O}(T^{-1/3}).$$

Thus, as  $T \rightarrow \infty$ , the optimal tuning obeys  $\tau_* = \sqrt{Q}/2$  and yields an error of  $\mathcal{Q} \left( e^{Q/4} \cdot \sqrt{\frac{\zeta^2 W}{1 - \zeta}} \right)$ .

*Proof.* Let us recap the notation of Theorem 4. Without losing generality, assume first  $\zeta T$  tokens are relevant and remaining tokens are irrelevant. Consider  $\mathbf{X}_I$  of size  $(1 - \zeta)T \times d$  induced by the irrelevant tokens with normal distribution. Using orthogonality of  $\mathbf{q}_*, \mathbf{w}_*, \mathbf{p}$ , observe that  $\mathbf{g} = \mathbf{X}_I(\bar{\mathbf{w}}_* + \frac{\mathbf{p}}{\|\mathbf{w}_*\|}) \sim \mathcal{N}(0, (1 + \gamma)\mathbf{I}_{(1 - \zeta)T})$  and  $\mathbf{h} = \mathbf{X}_I \bar{\mathbf{q}}_* \sim \mathcal{N}(0, \mathbf{I}_{(1 - \zeta)T})$  are independent vectors. Also for standard normal  $g \sim \mathcal{N}(0, 1)$ , recall that moment-generating function is given by  $\mathbb{E}[e^{\tau g}] = e^{\tau^2/2}$ .

Note that, the attention weights have the form  $\mathbf{a} = \phi(\tau \left[ \begin{smallmatrix} \sqrt{Q}\mathbf{1}_{\zeta T} \\ \mathbf{h} \end{smallmatrix} \right])$ . Here, the softmax denominator is  $T \cdot D_T$  where  $D_T := (\zeta e^{\sqrt{Q}\tau} + \frac{1}{T} \sum_{i=1}^{(1 - \zeta)T} e^{\tau h_i})$ . Define  $e^{\tau \mathbf{h}}$  to be the numerator corresponding to irrelevant tokens i.e.

$$e^{\tau \mathbf{h}} = [e^{\tau h_1} \dots e^{\tau h_{(1 - \zeta)T}}].$$

Define the matrix  $\mathbf{Q}_{\perp} = \mathbf{I} - \bar{\mathbf{q}}_* \bar{\mathbf{q}}_*^T, \mathbf{W}_{\perp} = \mathbf{I} - \bar{\mathbf{w}}_* \bar{\mathbf{w}}_*^T$ . To proceed, observe that, the prediction with  $\theta = (\mathbf{w}_* + \mathbf{p}, \tau \bar{\mathbf{q}}_*)$  is given by

$$\frac{D_T}{\sqrt{T}\|\mathbf{w}_*\|} y f_{\theta}^{\text{ATT}}(\mathbf{X}) = \sqrt{T}(\bar{\mathbf{w}}_* + \frac{\mathbf{p}}{\|\mathbf{w}_*\|})^T [\zeta e^{\|\mathbf{q}_*\|^2 \tau} (\mathbf{w}_* + y \mathbf{q}_*) + \frac{\mathbf{X}_I e^{\tau \mathbf{h}}}{T}] \quad (89)$$

$$= \zeta e^{\|\mathbf{q}_*\|^2 \tau} \sqrt{W} + \frac{1}{\sqrt{T}} \mathbf{g}^T e^{\tau \mathbf{h}}. \quad (90)$$

With this, conditioned on  $e^{\tau \mathbf{h}}$  observe that  $\mathbf{g}^T e^{\tau \mathbf{h}} \sim \mathcal{N}(0, \frac{1}{T} \|e^{\tau \mathbf{h}}\|^2)$ , thus,

$$\mathbb{P}_{\mathbf{g}}(y f_{\theta}^{\text{ATT}}(\mathbf{X}) > 0) = 1 - \mathcal{Q} \left( \frac{\zeta e^{\sqrt{Q}\tau} \sqrt{W}}{\sqrt{1 + \gamma} \|e^{\tau \mathbf{h}}\| / \sqrt{T}} \right).$$

To proceed, similar to Theorem 4, we apply Chebyshev's inequality over number of tokens  $T$  to find that with probability  $1 - \tilde{\mathcal{O}}(T^{-1/3})$  over  $\mathbf{h}$ ,

$$\left| \|e^{\tau \mathbf{h}}\|^2 / T - e^{2\tau^2} \right| \leq \tilde{\mathcal{O}}(T^{-1/3}).$$

In aggregate, this implies that, with probability  $1 - \tilde{\mathcal{O}}(T^{-1/3})$  over  $\mathbf{h}$ , we have that

$$1 - \mathcal{Q} \left( (1 + \tilde{\mathcal{O}}(T^{-1/3})) \frac{\zeta e^{\sqrt{Q}\tau} \sqrt{W}}{\sqrt{1 + \gamma} e^{\tau^2}} \right) \geq \mathbb{P}_{\mathbf{g}}(y f_{\theta}^{\text{ATT}}(\mathbf{X}) > 0) \geq 1 - \mathcal{Q} \left( (1 - \tilde{\mathcal{O}}(T^{-1/3})) \frac{\zeta e^{\sqrt{Q}\tau} \sqrt{W}}{\sqrt{1 + \gamma} e^{\tau^2}} \right),$$

Finally, note that since  $\frac{\zeta e^{\sqrt{Q}\tau} \sqrt{W}}{\sqrt{1 + \gamma} e^{\tau^2}}$  is upper/lower bounded by a positive constant, and since  $\mathcal{Q}(x + \epsilon) - \mathcal{Q}(x) = \int_x^{x + \epsilon} e^{-t^2/2} dt \leq \epsilon$ , we can rewrite

$$\left| \mathbb{P}_{\mathbf{g}}(y f_{\theta}^{\text{ATT}}(\mathbf{X}) > 0) - \mathcal{Q} \left( \frac{\zeta e^{\sqrt{Q}\tau} \sqrt{W}}{\sqrt{1 + \gamma} e^{\tau^2}} \right) \right| \leq \tilde{\mathcal{O}}(T^{-1/3}).$$

Union bounding with failure probability over  $\mathbf{h}$ , we conclude with the result. □

## G. Useful facts

For a random variable  $Z$  and  $\alpha > 0$ ,  $\|Z\|_{\psi_\alpha}$  denotes its  $\psi_\alpha$ -norm for Orlicz function  $\psi_\alpha(z) = e^{z^\alpha} - 1$  (Ledoux & Talagrand, 1991).

**Fact G.1.** Let  $X_1, \dots, X_n$  be independent zero-mean sub-gaussian or sub-exponential random variables with  $\|X_i\|_{\psi_m} \leq K$  for all  $i \in [n]$  for either  $m = 2$  or  $m = 1$ . Then,

$$\left\| \frac{1}{n} \sum_{i \in [n]} X_i \right\|_{\psi_m} \leq \frac{CK}{\sqrt{n}}.$$

**Fact G.2.** (Pollard, 1990) The following identity holds for Orlicz norms

$$\|XY\|_{\psi_{\frac{\alpha\beta}{\alpha+\beta}}} \leq c \|X\|_{\psi_\alpha} \cdot \|Y\|_{\psi_\beta} \quad (91)$$

for a fixed numerical constant  $c$ .

Next we state a Lemma from Talagrand quoted directly from Lemma 22 of (Mohammadi et al., 2019).

**Fact G.3.** (Ledoux & Talagrand, 1991) For any scalar  $\alpha \in (0, 1]$ , there exists a constant  $C_\alpha$  such that for any sequence of independent random variables  $\xi_1, \xi_2, \dots, \xi_N$  we have

$$\left\| \sum_i \xi_i - \mathbb{E} \left[ \sum_i \xi_i \right] \right\|_{\psi_\alpha} \leq C_\alpha \left( \max_i \|\xi_i\|_{\psi_\alpha} \right) \sqrt{N} \log N.$$

**Fact G.4.** Let  $\mathbf{z} \in \mathbb{R}^d$  be a  $K$ -subgaussian vector i.e.  $\mathbf{z}^\top \mathbf{v}$  is  $K$ -subgaussian for fixed  $\|\mathbf{v}\| = 1$ . Then, the following are true for a constant  $c > 0$

$$\mathbb{P} \left( \|\mathbf{z}\| \geq cK(\sqrt{d} + t) \right) \leq e^{-t^2}.$$

*Proof.* For completeness, we provide a proof. Repeating Lemma 31 of (Oymak, 2019), we can pick a  $1/2$  cover  $\mathcal{C}$  of the unit Euclidean ball in  $\mathbb{R}^d$  with size  $\log |\mathcal{C}| \leq 2d$ . For any  $\mathbf{v} \in \mathcal{C}$  subgaussianity implies  $\mathbb{P}(\mathbf{v}^\top \mathbf{z} \geq t) \leq \exp(-ct^2/K^2)$ . Setting  $K' = K/\sqrt{c}$ ,  $t = K'(\sqrt{d} + \tau)$  and union bounding over all  $\mathbf{v} \in \mathcal{C}$ , we find

$$\mathbb{P}(\sup_{\mathbf{v} \in \mathcal{C}} \mathbf{v}^\top \mathbf{z} \geq K'(\sqrt{d} + \tau)) \leq \exp(-\tau^2).$$

To proceed, set  $\mathbf{v}(\mathbf{z}) \in \mathcal{C}$  to be the nearest point to  $\bar{\mathbf{z}} = \mathbf{z}/\|\mathbf{z}\|$  in  $\mathcal{C}$ . Since  $\|\mathbf{v}(\mathbf{z}) - \bar{\mathbf{z}}\| \leq 0.5$ , note that

$$\|\mathbf{z}\| = \bar{\mathbf{z}}^\top \mathbf{z} = \mathbf{v}(\mathbf{z})^\top \mathbf{z} + (\bar{\mathbf{z}} - \mathbf{v}(\mathbf{z}))^\top \mathbf{z} \leq \mathbf{v}(\mathbf{z})^\top \mathbf{z} + 0.5\|\mathbf{z}\|.$$

Thus,  $\|\mathbf{z}\| \leq 2K'(\sqrt{d} + \tau)$  with probability at least  $1 - \exp(-\tau^2)$ . □

## H. Additional experimental results and details

### H.1. Additional details for image classification experiments

**Dataset.** As mentioned in Section 5.2, we construct three datasets by modifying the original images in CIFAR-10:

- **FULL-TILED.** Each examples consists of a 64x64 images obtained by arranging a 32x32 image from CIFAR-10 in a tiling pattern with four tiles (cf. Fig. 3a).
- **PARTIAL-TILED.** This is dataset is similar to FULL-TILED with the exception that each image has at-least  $T$  out of 4 tiles replaced by patches of i.i.d. random Gaussian noise with mean zero and variance 0.2. Note that, for each example in the dataset,  $T \in \{1, 2, 3\}$  is a random number as well as the location of the noisy tiles (cf. Fig. 3b).
- **EMBED-IN-IMAGENET** (Karp et al., 2021). We construct an example by simply embedding a 32x32 image from CIFAR-10 at a random location in a 64x64 background corresponding to a randomly selected image from ImageNet (Russakovsky et al., 2015). We also add i.i.d. random Gaussian noise with mean zero and variance 0.2 to the background (cf. Fig. 3c).

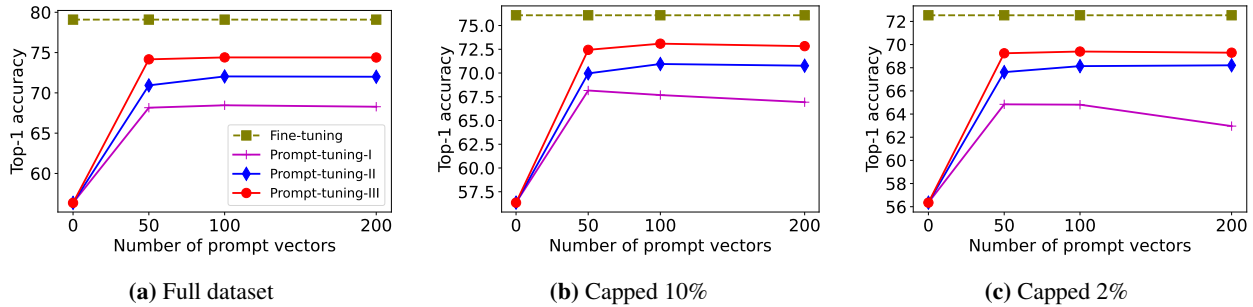


Figure 5. Performance of fine-tuning vs. prompt-tuning on 10-way classification tasks defined by PARTIAL-TILED dataset. Full dataset has 50K training examples. Capped 10% and 2% correspond to sub-sampled *train* sets where we select exactly 500 and 100 examples per class from the full dataset. Note that number of prompt vectors equal to 0 corresponds to *zero-shot* performance.

Table 1. Comparison between prompt-tuning and fine-tuning only first layer self-attention weights (Full dataset).

Method (trainable parameters)	EMBED-IN-IMAGENET	PARTIAL-TILED
PROMPT-TUNING-I w/ 100 prompt vectors (19.2K)	31.89	68.46
PROMPT-TUNING-II w/ 100 prompt vectors (19.2K)	38.48	72.04
PROMPT-TUNING-III w/ 100 prompt vectors (230.8K)	47.81	74.40
Fine-tuning only first layer attention weights (148.2K)	30.35	73.95

By construction, each dataset has 50,000 train and 10,000 test examples corresponding to train and test set of CIFAR-10. We also consider data-limited settings where we keep the test set intact but subsample the train set by selecting a fixed number of images for each class. Note that all three datasets define 10-way multiclass classification tasks with CIFAR-10 classes as potential labels.

**Model architecture.** We utilize a tiny variant of the Vision transformer model (Dosovitskiy et al., 2021) for our experiments. This variant has 12 transformer layers with its hidden dimension, MLP intermediate dimension, and number of heads per attention layer being equal to 192, 768, and 3, respectively. The patch size in our study is set to be 4x4. The model itself (without counting the trainable parameters/weights during prompt-tuning) has approximately 5.44M parameters. We rely on the CLS token to obtain the classification logits.

**Training.** We rely on Scenic library (Dehghani et al., 2022)<sup>4</sup> to conduct our experiments on image classification. Following the default settings in the library along with a coarse grid search, we employ Adam optimizer (Kingma & Ba, 2014) with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ , weight decay = 0.1, and batch size = 128 while training a *randomly initialized* model. Furthermore, we employ a linear warm-up of learning rate followed cosine learning rate schedule with base learning rate  $3e-3$ . As for the fine-tuning and prompt-tuning experiments that (partially) initialize from a *pre-trained* model, we rely on SGD with momentum parameter 0.9 and batch size = 128 to update trainable parameters. Again, we utilize a linear warm-up of learning rate followed by cosine learning rate schedule. Throughout our experiments, the base learning rates for fine-tuning and prompt-tuning are  $1e-3$  and 0.1, respectively.

## H.2. Additional results on image classification

Figure 5 showcases the performance of fine-tuning and various prompt-tuning strategies on PARTIAL-TILED.

**Comparison with fine-tuning the first self-attention layer.** In Tables 1, 2, and 3, we explored fine-tuning only first layer self-attention parameters for the underlying ViT model. This setting aligns well with the single-layer nature of our theoretical results. Similar to Fig. 4 (corresponding to EMBED-IN-IMAGENET dataset) and Fig. 5 (corresponding to PARTIAL-TILED dataset), we considered three settings: 1) Full dataset; 2) Capped 10%; and 3) Capped 2%, which progressively corresponds to smaller amount of (training) data during fine-tuning and prompt-tuning.

The key takeaways are:

<sup>4</sup><https://github.com/google-research/scenic>

Table 2. Comparison between prompt-tuning and fine-tuning only first layer self-attention weights (Capped 10%).

Method (trainable parameters)	EMBED-IN-IMAGENET	PARTIAL-TILED
PROMPT-TUNING-I w/ 100 prompt vectors (19.2K)	28.06	67.68
PROMPT-TUNING-II w/ 100 prompt vectors (19.2K)	36.76	70.95
PROMPT-TUNING-III w/ 100 prompt vectors (230.8K)	42.96	73.09
Fine-tuning only first layer attention weights (148.2K)	18.53	70.44

Table 3. Comparison between prompt-tuning and fine-tuning only first layer self-attention weights (Capped 2%).

Method (trainable parameters)	EMBED-IN-IMAGENET	PARTIAL-TILED
PROMPT-TUNING-I w/ 100 prompt vectors (19.2K)	20.52	64.81
PROMPT-TUNING-II w/ 100 prompt vectors (19.2K)	33.62	68.14
PROMPT-TUNING-III w/ 100 prompt vectors (230.8K)	36.13	69.40
Fine-tuning only first layer attention weights (148.2K)	15.46	65.04

1. When there is a significant distribution-shift between from pre-training data (in case of EMBED-IN-IMAGENET), even the simplest prompt-tuning, namely PROMPT-TUNING-I, significantly outperforms the fine-tuning first layer self-attention parameters.
2. When the distribution-shift is small, prompt-tuning variants realize a better *accuracy vs. training cost* trade-off, e.g. PROMPT-TUNING-II outperforms fine-tuning first layer self-attention parameters in the Capped 10% and Capped 2% setting (while training only 19.2K rather than 148.2K parameters).

### H.3. Illustration of attention weights for prompt vectors

Fig. 6 presents a representative example where we show evolution of average attention weights from prompt vectors to image tokens/patches across transformer layers, when we employ PROMPT-TUNING-III. It is evident from the figure that prompt-attention helps distinguish the relevant tokens/patches from the irrelevant patches.

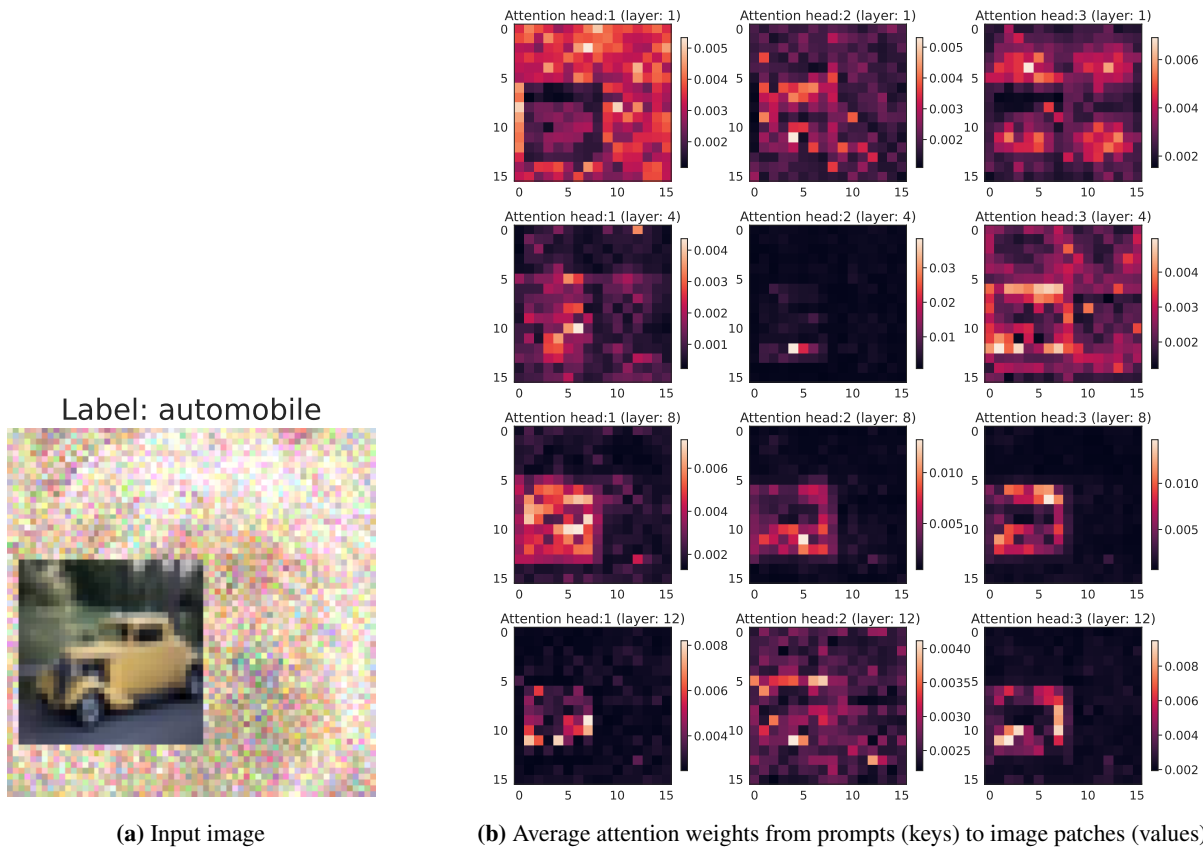


Figure 6. Illustration of how attention weights progressive change from the first layer (Figure 6b-top) to the last layer (Figure 6b-bottom) in the transformer model for a given input image (Figure 6a) when we employ PROMPT-TUNING-III. We plot average attention weights from 50 prompt vectors (keys) to 256 image patches (values). The attention weights for each attention head are naturally arranged in a 16 x 16 grid corresponding to the original locations of the patches in the image. Note that the attention weights in the early layer have a tiling pattern similar to that in FULL-TILED—the dataset utilized by the pre-trained model. However, as we progress deeper into the transformer, attention weights begin to capture the relevant patch locations in the dataset of interest, i.e., EMBED-IN-IMAGENET.