
TEST TIME TRAINING ENHANCES IN-CONTEXT LEARNING OF NONLINEAR FUNCTIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Test-time training (TTT) enhances model performance by explicitly updating designated parameters prior to each prediction to adapt to the test data. While TTT has demonstrated considerable empirical success, its theoretical underpinnings remain limited, particularly for nonlinear models. In this paper, we investigate the combination of TTT with in-context learning (ICL), where the model is given a few examples from the target distribution at inference time. We analyze this framework in the setting of single-index models $y = \sigma_*(\langle \beta, x \rangle)$, where the feature vector β is drawn from a hidden low-dimensional subspace. For single-layer transformers trained with gradient-based algorithms and adopting TTT, we establish an upper bound on the prediction risk. Our theory reveals that TTT enables the single-layer transformers to adapt to both the feature vector β and the link function σ_* , which vary across tasks. This creates a sharp contrast with ICL alone, which is theoretically difficult to adapt to shifts in the link function. Moreover, we provide the convergence rate with respect to the data length, showing the predictive error can be driven arbitrarily close to the noise level as the context size and the network width grow.

1 INTRODUCTION

In-context learning (ICL) is a powerful capability of pretrained transformers to solve tasks using a few labeled examples provided as input, without updating their weights. This ability has gained increasing attention with the advent of models with massive context windows, as more examples lead to significantly improved performance (Agarwal et al., 2024). This approach has also proven effective for multimodal tasks (Jiang et al., 2024). From a theoretical perspective, transformers are known to implement algorithms such as linear regression. Recent studies have extended this understanding to nonlinear settings, showing that transformers can learn nonlinear single-index models (Oko et al., 2024) and that softmax attention facilitates data-efficient feature learning (Nishikawa et al., 2025). Nevertheless, ICL faces its inherent limitations: the performance of ICL is fundamentally constrained by factors like the pretraining data (Bigoulaeva et al., 2025) and model architecture (Naim & Asher, 2025).

Test-time training (TTT) has emerged as a promising strategy to overcome these barriers. TTT adapts the model by updating its parameters on the test data before each prediction. This adaptive mechanism has led to strong empirical success across various fields, including large language models (Hu et al., 2025) and video object segmentation (Bertrand et al., 2023). In the context of ICL, TTT can be seamlessly integrated by using the in-context examples as data for task-specific adaptation. For instance, Akyürek et al. (2025) demonstrated that this combination achieves notable improvements on few-shot reasoning benchmarks.

Despite its empirical success, the theoretical foundations of TTT remain underdeveloped. A notable work by Gozeten et al. (2025) established the statistical efficiency of TTT over standard ICL, but their analysis was restricted to linear regression with a linear transformer. This simplified model fails to capture the true potential of TTT’s power on nonlinear, complex tasks. This gap motivates our primary research question:

Does test-time training improve in-context learning in nonlinear settings?

Furthermore, existing theoretical works commonly analyze performance in high-dimensional regimes, showing that the prediction loss is $o_d(1)$ and thus vanishes as the data dimension d grows. In practice, however, this dimension is fixed, making it crucial to understand how the loss behaves as the number of data points n increases. This raises our second key question:

How does the test loss behave when we fix the dimension and increase the data size?

1.1 OUR CONTRIBUTION

To address the questions above, we analyze the performance of transformers on learning single-index models, a simple type of nonlinear function. In our setting, each task is defined as

$$\mathbf{x}_1^t, \dots, \mathbf{x}_N^t, \mathbf{x}^t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_d), y_i^t \approx \sigma_*^t(\langle \beta^t, \mathbf{x}_i^t \rangle) \quad (i = 1, \dots, n)$$

where σ_*^t is an unknown polynomial that varies across tasks, and β^t is drawn from a fixed r -dimensional subspace. Using the in-context data $\mathbf{x}_1, \dots, \mathbf{x}_N$, our model constructs the predictor for the new query \mathbf{x} . We establish a rigorous upper-bound on the predictive risk for a transformer that utilizes TTT. Our main result is as follows:

Theorem 1 (Informal). *Consider learning single-index polynomial $\sigma_*^t(\langle \beta^t, \mathbf{x} \rangle)$ with the transformer trained via Algorithm 1. Then, with probability at least 0.99, we can construct the model f_{TF} for each prompt that satisfies $\mathbb{E}[|f_{\text{TF}}(\mathbf{x}) - y|] = \tau + \tilde{O}(m^{-1/2}) + \tilde{O}(\sqrt{\frac{r\sqrt{r}}{N_{\text{test}}}})$, where N_{pt} and T_{pt} are the context length and the number of tasks in pretraining, respectively, N_{test} is the context length in test-time, m is the network width and $\tau = O(1)$ is the noise level, if $N_{\text{pt}} = T_{\text{pt}} = (d^{\Omega(\text{ie}(\sigma_*^t))} r^2)$ and $N_{\text{test}} = \tilde{\Omega}(r^{\Omega(\text{ge}(\sigma_*^{\text{test}}))})$, where $\text{ie}(\sigma_*)$ and $\text{ge}(\sigma_*)$ are the information exponent and the general exponent of the polynomial σ_* , respectively.*

This theorem ensures the effectiveness of TTT in learning nonlinear single-index models, extending its known applicability to linear models. Our analysis reveals several strengths of our approach:

- **Efficient sample complexity:** Theorem 1 implies that $N_{\text{test}} = \tilde{\Omega}(r^{2\nu\Theta(\text{ge}(\sigma_*^{\text{test}}))})$ to ensure low predictive loss $o(1)$, which does not depend on the entire dimension d . This shows that transformers can adapt to the low-dimensionality of β . In addition, this statistical complexity does not depend on either the degree of the polynomial $\text{deg}(\sigma_*^{\text{test}})$ or $\text{ie}(\sigma_*^{\text{test}})$, which outperforms CSQ learners and shows a comparable performance with that of SQ learners.
- **Flexibility for varying nonlinearity:** Our framework allows the link function σ_* to vary across tasks. This adaptability is enabled by TTT, which fine-tunes the MLP layer each time using task-specific test data.
- **Statistical guarantee for practical settings:** We provide an explicit convergence rate with respect to the context length N_{test} . This result offers a statistical guarantee in practical scenarios where the dimensions d, r are large but fixed.

1.2 RELATED WORKS

In-context learning and its theoretical analysis In-context learning (Brown et al., 2020) is transformer’s ability to adapt to the specific task using few labeled examples, without updating any parameters. Agarwal et al. (2024) demonstrated that many in-context examples lead to considerably improved performance, and Jiang et al. (2024) confirmed that many-shot ICL is also beneficial for multimodal tasks. The theoretical background of ICL is extensively studied. For example, a wide array of works (Garg et al., 2023; von Oswald et al., 2023; Zhang et al., 2024; Gatmiry et al., 2024) have shown that linear transformers can be trained to perform linear regression in-context. As for nonlinear transformers, Cheng et al. (2024) demonstrated that nonlinear transformers learn to perform gradient descent and thus learn nonlinear functions. Also, Nichani et al. (2024) analyzed learning of causal structure by softmax transformer. Recently, Dherin et al. (2025) showed that ICL in a single-transformer block (a self-attention layer and subsequent MLP layer) corresponds to low-rank update in MLP layer. Regarding limitations of ICL, Bigoulaeva et al. (2025) argued that pretraining datasets impose a fundamental limit on the model’s capability with ICL. Furthermore, (Naim & Asher, 2025) demonstrated that the transformer’s ICL ability to generalize functions is limited to certain input values, and found that this limitation comes from layer normalization and softmax attention.

Test-time training Test-time training (Sun et al. (2020), Liu et al. (2021)) updates the model using test data before making predictions, thereby addressing distribution shifts. TTT achieved success in many fields. For example, Bertrand et al. (2023) applied TTT for the video object segmentation task and achieved a significant improvement in the performance. Moreover, Hu et al. (2025) analyzed test-time learning of large language models and achieved at least 20% higher performance on domain knowledge adaptation. Furthermore, Zhang et al. (2025) proposed adopting a large chunk update, and validated the effectiveness of their approach to long-context data through tasks like image sets and language model. As for TTT combined with few-shot prediction, Akyürek et al. (2025) reported that introducing TTT with in-context examples resulted in 6 times higher accuracy in the Abstraction and Reasoning Corpus and a 7.3 percent higher score on BIG-Bench Hard. Finally, regarding the theory behind TTT for in-context learning, Gozeten et al. (2025) analyzed the linear transformer with a single gradient step and characterized the prediction risk of the model with TTT, showing that TTT can mitigate distribution shift.

2 PRELIMINARIES AND PROBLEM SETTINGS

Notations Let $\text{He}_i(z) = (-1)^i e^{\frac{z^2}{2}} \frac{d^i}{dz^i} e^{-\frac{z^2}{2}}$ be the degree- i (probabilist's) Hermite polynomial. \mathbb{S}^{d-1} denotes the unit sphere in \mathbb{R}^d . For matrix \mathbf{A} , we denote its ℓ_2 operator norm and Frobenius norm as $\|\mathbf{A}\|_2$ and $\|\mathbf{A}\|_F$, respectively. For a set S , $\text{Unif}(S)$ denotes the uniform distribution over S . $\tilde{O}, \tilde{\Omega}, \tilde{\Theta}$ means O, Ω, Θ where polylogarithmic terms of d and $1/\varepsilon$ are hidden. O_d, Ω_d, Θ_d means the order with respect to the dimension d, r , while Θ_ε denotes the order in terms of ε .

2.1 IN-CONTEXT LEARNING AND TEST-TIME TRAINING

We consider the basic setting in ICL, which is introduced by Garg et al. (2023) (see Oko et al. (2024) and Lee et al. (2024)). In this setting, the model is given a sequence $(\mathbf{x}_1, y_1, \dots, \mathbf{x}_N, y_N, \mathbf{x})$ called prompt. The labeled pairs $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$ are called contexts, and $\mathbf{x} \in \mathbb{R}^d$ is referred to as query. The model is asked to predict the output that corresponds to \mathbf{x} based on the context. The context is sometimes abbreviated as $(\mathbf{X}_n, \mathbf{y}_n)$ where $\mathbf{X}_n = (\mathbf{x}_1, \dots, \mathbf{x}_n)$, $\mathbf{y}_n = (y_1, \dots, y_n)$. In this work, we assume that the \mathbf{x} and y are generated as follows:

$$\mathbf{x}_1^t, \dots, \mathbf{x}_N^t, \mathbf{x}^t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_d), y_i^t = f_*^t(\mathbf{x}_i^t) + \zeta_i, \zeta_i \sim \text{Unif}(-\tau, \tau).$$

ICL aims to predict $y = f_*(\mathbf{x}) + \zeta$ by mere observation of the context, without updating parameters for each prompt. However, we introduce test-time training to further enhance the model's accuracy. In pretraining, we train the model with T_{pt} distinct datasets $(\mathbf{x}_1^t, y_1^t, \dots, \mathbf{x}_{N_{pt}}^t, y_{N_{pt}}^t, \mathbf{x}^t, y^t)_{t=1}^{T_{pt}}$, with each prompts consisting of N_{pt} queries. In test-time, we divide the context into four groups with i -th group's length N_i . This means the test-time prompt is $(\mathbf{X}_{N_1}, \mathbf{y}_{N_1}, \mathbf{X}_{N_2}, \mathbf{y}_{N_2}, \mathbf{X}_{N_3}, \mathbf{y}_{N_3}, \mathbf{X}_{N_4}, \mathbf{y}_{N_4}, \mathbf{x})$. Each group of data plays a different role in test-time training: See Subsection 2.4 for the detail. Let $N_{test} = \sum_{i=1}^4 N_i$ be the total number of contexts test-time.

For the evaluation of the model $f(\mathbf{x}, \theta)$ with parameter θ , we define prediction risk as

$$\mathcal{R}_f(\theta) = \mathbb{E}[|f(\mathbf{x}, \theta) - y|],$$

where $y = f_*(\mathbf{x}) + \zeta$ and the expectation is taken over the prompt $\mathbf{x}_i, \mathbf{x} \sim \mathcal{D}_{\mathbf{x}}, f_* \sim \mathcal{D}_{f_*}, \zeta_i, \zeta \sim \mathcal{D}_{\zeta}$. Note that $\mathcal{R}_f(\theta)$ not only depends on the dimensions d, r but also on the context length N .

2.2 SINGLE INDEX MODEL

We consider single-index models for the input-output relationship, where the output depends solely on the direction of the feature vector β . To predict this relationship accurately, the models are expected to learn the target direction β from the high-dimensional data in \mathbb{R}^d . Consequently, single-index models have been extensively studied in machine learning theory (Bai & Lee, 2020; Ba et al., 2022; Bietti et al., 2022; Mousavi-Hosseini et al., 2023; Berthier et al., 2024), particularly to examine adaptability to the low-dimensional subspace. The target function $f_*^t(\mathbf{x})$ is determined as:

$$f_*^t(\mathbf{x}) = \sigma_*^t(\langle \beta^t, \mathbf{x} \rangle).$$

- 162 1. **Feature vector** The feature vector β is chosen uniformly from the unit sphere in an r -dimensional
 163 subspace. Let S_r be an r -dimensional linear subspace in \mathbb{R}^d . For each prompt, β^t is uniformly
 164 drawn from its support $\text{Supp}(\beta) = \{\beta \mid \beta \in S_r, \|\beta\| = 1\}$.
 165
 166 2. **Link function** We take $\sigma_*^t(z) = \sum_{i=Q}^P \frac{c_i^t}{i!} \text{He}_i(z)$ where $1 \leq Q \leq P$. We assume P and Q are
 167 constants of $O(1)$. The coefficients c_i^t are drawn from any distribution \mathcal{D}_{σ_*} that satisfies

$$168 \mathbb{E}[c_Q^t] = \Theta(1) \neq 0, \sum_{i=Q}^P (c_i^t)^2 \leq R_c = \Theta(1) \text{ (a.s.) and } (c_Q^t, \dots, c_P^t) \neq (0, \dots, 0) \text{ (a.s.)}. \\ 169 \\ 170 \\ 171$$

172 **Remark 2.** We draw a new feature vector β^t and a new link function $\sigma_*^t(z)$ for each prompt. When
 173 $r \ll d$, β^t has low-dimensional support. We will later demonstrate that our model can leverage this
 174 low-dimensionality. Note that our analysis is also valid when $r = d$. As for the link function, we do
 175 not assume a specific distribution of the coefficients c_i : our framework allows different distributions
 176 as long as they satisfy the assumptions above.

177 For simplicity, we assume that $N_{pt}, T_{pt}, N_1, N_2, N_3, N_4 = \text{poly}(d)$, and there exists $\alpha_r > 0$ that
 178 $r = \Omega(d^{\alpha_r})$, which means that r grows faster than $\text{polylog}(d)$.
 179

180 The complexity of learning single-index model is governed by three key quantities: the degree of
 181 the polynomial $\text{deg}(\sigma_*)$, the information exponent $\text{ie}(\sigma_*)$, and the general exponent $\text{ge}(\sigma_*)$.

- 182 • The information exponent (Dudeja & Hsu, 2018; Arous et al., 2021) $\text{ie}(\sigma_*)$ is the smallest non-
 183 zero degree of the hermite expansion of σ_* .
 184
- 185 • The general exponent (Damian et al., 2024) $\text{ge}(\sigma_*)$ is the minimum of $\text{ie}(f \circ \sigma_*)$ with respect to
 186 all the L_2 -measurable transformation f .
 187

188 By definition, $\text{deg}(\sigma_*) \geq \text{ie}(\sigma_*) \geq \text{ge}(\sigma_*)$ holds. This relationship characterizes the statistical
 189 complexity required by different algorithms.

- 190 • The kernel method, which relies on pre-determined feature maps, requires the sample complexity
 191 $n \gtrsim d^{\Theta(\text{deg}(\sigma_*))}$ to ensure low predictive error (Ghorbani et al., 2020; Donhauser et al., 2021).
 192
- 193 • The models with access to correlational statistical query (CSQ), of the form $\mathbb{E}[\phi(\mathbf{x})y]$, require a
 194 sample size of $n \gtrsim d^{\Theta(\text{ie}(\sigma_*))}$, known as CSQ lower bound (Damian et al., 2022; Abbe et al.,
 195 2023). This improved sample complexity, independent of $\text{deg}(\sigma_*)$, has been achieved by models
 196 like two-layer neural network with online SGD (Arous et al., 2021; Damian et al., 2023) or one-
 197 step gradient descent (Damian et al., 2022; Dandi et al., 2025).
- 198 • For the broader class of statistical queries (SQ), which takes the form $\mathbb{E}[\phi(\mathbf{x}, y)]$, the sample com-
 199 plexity is further improved to $n \gtrsim d^{\Theta(\text{ge}(\sigma_*))}$. This advantage comes from applying nonlinear
 200 transformations to the label y , thereby reducing $\text{ie}(\sigma_*)$. Recent works have achieved this com-
 201 plexity by reusing the data (Lee et al., 2024; Arnaboldi et al., 2025) or adjusting the loss function
 202 (Joshi et al., 2024). Furthermore, when σ_* is a polynomial, the following result holds:

203 **Lemma 3** (Lee et al. (2024), Proposition 6). *It holds that $\text{ge}(\sigma_*) = \begin{cases} 1 & (\text{if } \sigma_* \text{ is not even}) \\ 2 & (\text{if } \sigma_* \text{ is even}) \end{cases}$.*

204 *Moreover, $\text{ge}(\sigma_*) = \min_{j \geq 1} \text{ie}(\sigma_*^j)$ holds.*
 205

206 This implies that, for any polynomial σ_* , $\text{ge}(\sigma_*) \leq 2$ is a small constant, regardless of $\text{deg}(\sigma_*)$.
 207

208 Our goal is to achieve the test-time sample complexity of $N_{test} = r^{\Theta(\text{ge}(\sigma_*))}$, which is independent
 209 of the entire dimension d and surpasses the CSQ limit.
 210

211 2.3 STUDENT MODEL

212 To facilitate our theoretical analysis, we need to establish a concrete architecture for the student
 213 model. We employ a single-layer transformer model, formulated as follows. In pretraining, we use
 214 the same model as Nishikawa et al. (2025). We first construct the embedding E as
 215

216 $\mathbf{E} = \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_N & \mathbf{x} \\ y_1 & \cdots & y_N & 1 \end{bmatrix} \in \mathbb{R}^{(d+1) \times (N+1)}$. Then, we apply the softmax attention layer as

217
218
219
$$\text{Attn}(\mathbf{E}) = \mathbf{W}^V \mathbf{E} \cdot \text{softmax}(\text{Mask}(\rho^{-1} \cdot (\mathbf{W}^K \mathbf{E})^\top \mathbf{W}^Q \mathbf{E})),$$

220
221 where ρ is temperature and $\mathbf{W}^V, \mathbf{W}^K, \mathbf{W}^Q \in \mathbb{R}^{(d+1) \times (d+1)}$ are the parameters for attention layer.
222 The softmax is applied to each column, while the Mask function converts all the elements in the
223 final row into $-\infty$ to prevent the model from focusing on the uninformative final row. We further
224 apply a multi-layer perceptron (MLP) layer. For the activation function, we use ReLU function
225 $\sigma(z) = \max\{0, z\}$ throughout the paper. With the parameters $\mathbf{W}^F \in \mathbb{R}^{m \times (d+1)}$, $\mathbf{b} \in \mathbb{R}^m$, and
226 $\mathbf{a} \in \mathbb{R}^m$ where m is the network width, the model’s output is

227
$$f_{\text{IC}}(\Gamma, \mathbf{X}_N, \mathbf{y}_N, \mathbf{x}) = \text{MLP} \circ \text{Attn}(\mathbf{E})_{:,N+1} = \mathbf{a}^\top \sigma(\mathbf{W}^F \text{Attn}(\mathbf{E})_{:,N+1} + \mathbf{b}),$$

228 where σ is applied entry-wise. Finally, we adopt some simplification. Let $\mathbf{W}^{KQ} = (\mathbf{W}^K)^\top \mathbf{W}^Q \in$
229 $\mathbb{R}^{(d+1) \times (d+1)}$ and $\mathbf{W}^{FV} = \mathbf{W}^F \mathbf{W}^V \in \mathbb{R}^{(m+1) \times (d+1)}$. We use the following parametrization:

230
231
$$\mathbf{W}^{KQ} = \begin{bmatrix} \Gamma & \mathbf{0}_{d \times 1} \\ \mathbf{0}_{1 \times d} & 1 \end{bmatrix}, \mathbf{W}^{FV} = [\mathbf{O}_{(m+1) \times d} \quad \mathbf{v}],$$

232 for $\Gamma \in \mathbb{R}^{d \times d}$ and $\mathbf{v} \in \mathbb{R}^{(d+1) \times 1}$. This kind of simplification, specifying some of the parameters as
233 zero, is often adopted in many theoretical works on transformers (Zhang et al., 2023; Huang et al.,
234 2023; Kim & Suzuki, 2025). Overall, the model’s output is written as

235
236
237
$$f_{\text{IC}}(\Gamma, \mathbf{X}_N, \mathbf{y}_N, \mathbf{x}) = \sum_{j=1}^m a_j \sigma \left(v_j \frac{\sum_{i=1}^N y_i e^{y_i / \rho} e^{\mathbf{x}_i^\top \Gamma \mathbf{x} / \rho}}{\sum_{i=1}^N e^{y_i / \rho} e^{\mathbf{x}_i^\top \Gamma \mathbf{x} / \rho}} + b_j \right).$$

240 See Appendix A in Nishikawa et al. (2025) for how to derive this equation. At test-time, we adopt
241 low-rank adaptation(LoRA). Specifically, we change the attention matrix Γ^* as $\Gamma_u = \Gamma^* + \mathbf{u}^\top \mathbf{u}$,
242 where Γ^* is fixed during test-time and $\mathbf{u} \in \mathbb{R}^d$ is a trainable parameter vector. Our goal is to find
243 $\hat{\mathbf{u}} \approx \beta$ using test-time context data. See section 2.4 for how test data is used to find $\hat{\mathbf{u}}$. The final
244 prediction of the model is as follows:

245
246
$$f_{\text{TF}}(\mathbf{x}, \hat{\mathbf{u}}, \mathbf{v}, \mathbf{a}, \mathbf{b}) = \sum_{j=1}^m a_j \sigma(v_j \langle \hat{\mathbf{u}}, \mathbf{x} \rangle + b_j).$$

247
248 **Remark 4.** *The use of in-context data in attention output leads to data-efficient prediction,*
249 *but it prevents the model from achieving near-zero error. This is because of the high or-*
250 *der term of $\langle \beta, \mathbf{x} \rangle$. When N is sufficiently large, the denominator of the attention output*
251 *$N^{-1} \sum_{i=1}^N e^{y_i / \rho} e^{\mathbf{x}_i^\top \Gamma \mathbf{x} / \rho}$ approximates $\mathbb{E}[\exp \sigma_*(\langle \beta, \mathbf{x}_1 \rangle / \rho) \exp(\langle \Gamma_* \mathbf{x}, \mathbf{x}_1 \rangle / \rho)]$, as explained in*
252 *Nishikawa et al. (2025). The key of their work is, since $\exp(\langle \Gamma_* \mathbf{x}, \mathbf{x}_1 \rangle / \rho)$ contains all the Hermite*
253 *coefficients and $\exp \sigma_*(\langle \beta, \mathbf{x}_1 \rangle / \rho)$ has nonzero coefficient for $\text{He}_{\text{ge}(\sigma_*)}$, this attention layer can*
254 *compute $\langle \beta, \mathbf{x} \rangle^{\text{ge}(\sigma_*)}$. However, this means that the output of the attention module contains $\langle \beta, \mathbf{x} \rangle^i$*
255 *for $i \geq \text{ge}(\sigma_*) + 1$ as well. This causes inevitable prediction error, even when $N \rightarrow \infty$. Mean-*
256 *while, when we find $\hat{\mathbf{u}}$ that is sufficiently near the feature vector β , we can achieve near-zero error.*
257 *Therefore we do not use in-context data in the final prediction f_{TF} .*

258
259
260

2.4 TRAINING ALGORITHM

261 We employ a gradient-based training algorithm, as specified in Algorithm 1. The training procedure
262 consists of pretraining and test-time training, with the latter divided into three stages.

- 263
264 • **Pretraining:** We optimize Γ via one-step gradient descent over T_{pt} prompts. This scheme is origi-
265 nally taken from Nishikawa et al. (2025). The effectiveness of one-step gradient descent is con-
266 firmed by works such as Ba et al. (2022) and Damian et al. (2022), which demonstrate that one-step
267 gradient already captures the key feature.
- 268 • **TTT stage I:** We apply a single gradient descent step to $\mathbf{u}^{(0)}$ with L_2 -regularization. This stage
269 prevents catastrophic forgetting, a phenomenon where the parameter change caused by TTT de-
prives the model of the fundamental ability to adapt to the original task. The goal of this stage is

Algorithm 1: Pretraining and test-time training of transformer

Input : Learning rate $\eta_{pt}, \eta_1, \eta_2$, regularization rate λ_{pt}, λ_1 , initialization scale $\alpha_{pt}, \alpha_1, \alpha_2$, dimensions d, r , temperature ρ .

1 **Initialize** $\Gamma(0) \sim \mathbf{I}_d/\sqrt{d}$, $\mathbf{v}(0) \sim \text{Unif}(\{\pm 1\}^m)$, $\mathbf{b}(0) = \mathbf{0}_m$, $\mathbf{a}(0) = \alpha_{pt}\mathbf{1}_m$.

2 **Pretraining: One-step Gradient descent on Attention Matrix**

3 $\Gamma^* \leftarrow \Gamma(0) - \eta_{pt} \frac{1}{2T_{pt}} \sum_{t=1}^{T_{pt}} \nabla_{\Gamma} ((f_{\text{IC}}(\Gamma, \mathbf{X}_{N_{pt}}, \mathbf{y}_{N_{pt}}, \mathbf{x}^t) - y^t)^2 + \lambda_{pt} \|\Gamma\|_F^2)$.

4 **Test-Time Training**

5 $\mathbf{a}(0) = \alpha_1 \mathbf{1}_m$.

6 **for** $i = 1$ **to** N_1 **do**

7 $\mathbf{x}_i \leftarrow \sqrt{r} \Gamma^* \mathbf{x}_i$.

8 $\mathbf{u}^{(0)} \sim \mathcal{N}(0, \mathbf{I}_d)$, $\mathbf{u}^{(0)} \leftarrow \sqrt{r} \Gamma^* \mathbf{u}^{(0)}$, $\mathbf{u}^{(0)} \leftarrow \frac{\mathbf{u}^{(0)}}{\sqrt{r} \|\mathbf{u}^{(0)}\|}$.

9 **Stage I: Initialization of \mathbf{u} with signals from Original Model**

10 Draw $\mathbf{w}_1, \dots, \mathbf{w}_{N_{new}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_d)$.

11 $\mathbf{b} = \frac{1}{N_{new}} \sum_{i=1}^{N_{new}} g(\Gamma^*, \mathbf{X}_{N_1}, \mathbf{y}_{N_1}, \mathbf{w}_i)$,

12 $\mathbf{u}^{(1)} = \mathbf{u}^{(0)} - \eta_1 \{ \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} \frac{1}{2N_{new}} \sum_{i=1}^{N_{new}} (f_{\text{IC}}(\Gamma_{\mathbf{u}}, \mathbf{X}_{N_1}, \mathbf{y}_{N_1}, \mathbf{w}_i) - (g(\Gamma^*, \mathbf{X}_{N_1}, \mathbf{y}_{N_1}, \mathbf{w}_i) - \mathbf{b}))^2 + \frac{\lambda_1}{2} \nabla_{\mathbf{u}} \|\mathbf{u}\|^2 \}$,

13 $\mathbf{u}^{(1)} \leftarrow \frac{\mathbf{u}^{(1)}}{\|\mathbf{u}^{(1)}\|}$.

14 $\mathbf{a}(0) = \alpha_2 \mathbf{1}_m$.

15 **for** $i = N_1 + 1$ **to** $N_1 + N_2$ **do**

16 $\mathbf{x}_i \leftarrow \sqrt{r} \Gamma^* \mathbf{x}_i$.

17 **Stage II: Strong Recovery**

18 **for** $t = 1$ **to** $t = N_3$ **do**

19 $\mathbf{u}^{(t+1)} = \mathbf{u}^{(t)} - \eta_2 \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} (\frac{1}{2} (f_{\text{IC}}(\Gamma_{\mathbf{u}}, \mathbf{X}_{N_2}, \mathbf{y}_{N_2}, \mathbf{x}_{N_1+N_2+t}) - y_{N_1+N_2+t})^2)$.

20 $\mathbf{u}^{(t+1)} \leftarrow \frac{\mathbf{u}^{(t+1)}}{\|\mathbf{u}^{(t+1)}\|}$.

21 **Initialize** $\mathbf{b}_j^* \sim \text{Unif}([-\log^2 d, \log^2 d])$, $\mathbf{v}^* = \mathbf{v}(0)$.

22 **Stage III: Training of MLP Layer**

$\mathbf{a}^* \leftarrow \arg \min_{\mathbf{a}} \frac{1}{2N_4} \sum_{t=N_1+N_2+N_3+1}^{N_1+N_2+N_3+N_4} (f_{\text{TF}}(\mathbf{x}_t, \mathbf{u}^{(N_3+1)}, \mathbf{v}^*, \mathbf{a}, \mathbf{b}^*) - y_t)^2 + \frac{\lambda_2}{2} \|\mathbf{a}\|^2$.

Output: Prediction $f_{\text{TF}}(\mathbf{x}, \mathbf{u}^{(N_3+1)}, \mathbf{v}^*, \mathbf{a}^*, \mathbf{b}^*)$.

to find a good initial value of $\mathbf{u}^{(1)}$ that satisfies $\langle \beta, \mathbf{u}^{(1)} \rangle \geq 1/\text{polylog}(d)$. For this stage, we use the output from the attention layer of the original model as a teacher signal, which is defined as $g(\Gamma^*, \mathbf{X}_{N_1}, \mathbf{y}_{N_1}, \mathbf{x}) = \frac{N_1^{-1} \sum_{i=1}^{N_1} y_i e^{y_i/\rho} e^{\mathbf{x}_i^T \Gamma^* \mathbf{x}}}{N_1^{-1} \sum_{i=1}^{N_1} e^{y_i/\rho} e^{\mathbf{x}_i^T \Gamma^* \mathbf{x}}}$. In this stage, the query \mathbf{x} do not require ground-truth label y . Therefore, we use N_{new} newly generated vectors $w_1, \dots, w_{N_{new}} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_d)$ for the query. The necessity of this stage is further discussed in Subsection 3.1.2.

- TTT stage II: We continue to optimize \mathbf{u} by applying multi-step online SGD scheme that originates from Lee et al. (2024), using $(\mathbf{X}_{N_2}, \mathbf{y}_{N_2})$ as contexts. This stage aims to align \mathbf{u} more closely to the target β . As we will see, increasing the number of the SGD steps leads to the convergence of \mathbf{u} to β .
- TTT stage III: Finally, we train the MLP layer to fit to the nonlinear link function σ_*^{test} . Specifically, we randomize \mathbf{v} and \mathbf{b} and optimize \mathbf{a} with ridge regression, following the recipe in Nishikawa et al. (2025). This problem is convex with respect to \mathbf{a} , ensuring that the global optimum can be readily found.

3 MAIN RESULT

Based on the problem settings above, we are now ready to present our main result.

Theorem 1 (Formal). We denote the link function drawn in inference-time as σ_*^{test} . Suppose that $T_{pt}, N_{pt} = \tilde{\Omega}(r^2 d^{Q+2})$, $N_1 = \tilde{\Omega}(r^{\text{ge}(\sigma_*^{\text{test})+2})$, $N_{new} = \tilde{\Omega}(r^{\text{ge}(\sigma_*^{\text{test})+2})$, $N_2 = \tilde{\Theta}(r^2)$. Moreover, we assume that $m, N_{pt}, T_{pt}, N_1, N_2, N_3, N_4 = O(\text{poly}(d))$ and there exists $\alpha_r > 0$ that $r = \Omega(d^{\alpha_r})$. When we fix d large enough, then there exists $\lambda_{pt}, \lambda_1, \eta_{pt}, \eta_1, \eta_2$ such that the prediction risk is low with probability at least 0.99 over the training data and random initialization. Concretely, for the model trained via Algorithm 1, we have that

$$|\mathcal{R}_{f_{\text{TF}}}(\mathbf{u}, \mathbf{v}^*, \mathbf{a}^*, \mathbf{b}^*) - \tau| = \tilde{O}(N_4^{-1/2}) + \tilde{O}(m^{-1/2}) + \tilde{O}\left(\sqrt{\frac{r\sqrt{r}}{N_3}}\right),$$

with probability at least 0.99.

The proof is given in Appendix G. Our result has the following advantages compared to previous works.

(i) Nonlinearity While Gozeten et al. (2025) also investigates test-time training combined with ICL, their analysis is restricted to linear datasets and a transformer with linear attention. By contrast, we consider a more complex and general problem of prediction for a single-index model with a nonlinear polynomial. In addition, the student model in our work utilizes nonlinear softmax attention, thereby extending the analysis beyond linear transformers.

(ii) The adaptability to link function Our framework has the flexibility of using a task-specific link function σ_*^t . In the previous work by Nishikawa et al. (2025), the link function is fixed throughout all the tasks (only β varies across different tasks) because the training of the relevant layer is done only in pretraining. In contrast, the use of test-time training allows the model to adapt to the characteristics of each task. Therefore, our algorithm is effective even when the underlying link function varies from one task to another.

(iii) The convergence rate of predictive loss with respect to n The analysis by Nishikawa et al. (2025) guarantees that the ICL risk is $o_d(1)$, but it provides no guarantee that the risk diminishes as $n \rightarrow \infty$. In fact, this is an inherent consequence of adopting softmax attention, as we discussed in Remark 4. In contrast, our result overcomes that limitation. When d is a sufficiently large constant, our theory ensures that the prediction risk can be made arbitrarily close to the inevitable noise τ by increasing the context size N_{test} and the network width m . This is because the increase in the test-time context length N_3 allows the vector \mathbf{u} to converge to the ground truth β .

In addition, it is worth noting that we achieve efficient sample complexity. Theorem 1 implies that $N_{test} = \tilde{\Omega}(r^{2+\text{ge}(\sigma_*^{\text{test})})$ to ensure low predictive loss $o(1)$, yielding two key benefits. First, our sample complexity does not depend on the entire dimension d , which means this model also adapts to the low-dimensionality of β . Second, our sample complexity is independent of either $\text{deg}(\sigma_*^{\text{test}})$ or $\text{ie}(\sigma_*^{\text{test}})$, which breaks the CSQ upper bound. Moreover, for polynomials where $\text{ge}(\sigma_*) \leq 2$ holds, the required number of samples is small. In conclusion, our sample complexity is nearly optimal. It is instructive to compare this with Nishikawa et al. (2025). Although they also achieve sample complexity with these two features ($N_{test} = \tilde{\Omega}(r^{3\text{ge}(\sigma_*)/2})$), their result guarantees convergence with respect to d , not N_{test} . This implies that their result is valid only in the asymptotic limit where $d \rightarrow \infty$. Moreover, they do not provide the convergence rate with respect to d . In contrast, our result holds for any sufficiently large fixed d .

3.1 PROOF SKETCH

The outline of the proof of our main theorem is as follows: using the output of the original model, we can achieve weak recovery i.e. nontrivial overlap between β and \mathbf{u} . Once weak recovery is completed, more training using the signal y leads to strong recovery i.e. $\langle \beta, \mathbf{u} \rangle \geq 1 - \varepsilon$. Finally, the training of MLP layer allows the model to fit to the nonlinearity of the link function.

3.1.1 EXPLOITING THE PRETRAINED ATTENTION MATRIX

As shown in Nishikawa et al. (2025), the attention matrix Γ^* after pretraining captures the r -dimensional Span of β :

Lemma 5 (Informal, Correspond to Proposition 22 in Nishikawa et al. (2025)). *After running the pretraining in Algorithm 1 with $T_{pt}, N_{pt} = \tilde{\Omega}(r^2 d^{Q+2})$, it holds that*

$$\Gamma^* \approx c_r \mathbb{E}_\beta[\beta\beta^\top]$$

with high probability, where $c_r = \tilde{\Theta}(\sqrt{r})$.

This means that Γ^* can project vectors into the r -dimensional subspace $\text{Supp}(\beta)$. Therefore, by multiplying $\sqrt{r}\Gamma^*$ to \mathbf{x}_i and the gradient (the coefficient \sqrt{r} is just for adjusting the scale), we can make the problem virtually r -dimensional, even though the entire dimension is d . This leads to the sample complexity only scaling up with r , not the whole dimension d .

3.1.2 WEAK RECOVERY

First, we initialize \mathbf{u} by one step gradient descent using the output from the original model $g(\Gamma_*, \mathbf{X}_{N_1}, \mathbf{y}_{N_1}, w_i)$. We do not use the signal y in this process. Intuitively, this self-distillation can be seen as a prevention of catastrophic forgetting. Taking the information from the original model makes the LoRA model similar to the original model, thereby preserving the desired features of the original model.

To clarify why we need to use the original model as teacher, consider training the LoRA model with the true signal y . Then, by calculating the gradient, we can get the following:

Lemma 6 (Informal). *The following holds with high probability:*

$$\frac{1}{2} \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} (f_{\text{IC}}(\Gamma, \mathbf{X}_N, \mathbf{y}_N, \mathbf{x}) - y)^2 \approx \tilde{\Theta}(\alpha m) \langle \beta, \mathbf{u} \rangle \{ (\sqrt{r})^{-(\text{ie}(\sigma_*)-1)} + \langle \beta, \mathbf{u} \rangle^{2\text{ie}(\sigma_*)-2} \} \beta.$$

See Appendix C for details. At initialization $\langle \beta, \mathbf{u} \rangle = \tilde{O}(1/r)$ holds, so the signal strength is $r^{\Theta(\text{ie}(\sigma_*^{\text{test}}))}$. Therefore, when we train \mathbf{u} from the signal y , we need at least $r^{\Theta(\text{ie}(\sigma_*^{\text{test}}))}$ data.

However, using the original model as teacher signal reduces this data length to $r^{\Theta(\text{ge}(\sigma_*^{\text{test}}))}$. As Nishikawa et al. (2025) showed, the attention layer after pretraining can compute $\langle \beta, \mathbf{x} \rangle^{\text{ge}(\sigma_*^{\text{test}})}$ in-context. Then, noting that $\text{ie}(\text{He}_{\text{ge}(\sigma_*)}) = \text{ge}(\sigma_*)$, when we learn from the original model, the signal strength becomes $r^{\Theta(\text{ge}(\sigma_*^{\text{test}}))}$. This improved signal strength results in the required data length $N_{\text{test}} = r^{\Theta(\text{ge}(\sigma_*^{\text{test}}))}$, which surpasses CSQ limit.

3.1.3 STRONG RECOVERY

Weak recovery is insufficient for reliably predicting $\langle \beta, \mathbf{x} \rangle$ using only \mathbf{u} . Therefore, we further optimize \mathbf{u} until we achieve strong recovery, defined as $\langle \beta, \mathbf{u} \rangle \geq 1 - \varepsilon$ for a small error tolerance $\varepsilon > 0$. After achieving $\langle \beta, \mathbf{u}^{(1)} \rangle > 1/\text{polylog}(d)$, the signal strength is $O(1/\text{polylog}(d))$, independent of $\text{ge}(\sigma_*^{\text{test}})$. Consequently, the sample size required to achieve strong recovery does not depend on $\text{ge}(\sigma_*^{\text{test}})$. Moreover, we achieve a smaller sample complexity for achieving strong recovery compared to the prior work by Lee et al. (2024). While their work suggested a linear increase in $\langle \beta, \mathbf{u}^{(n)} \rangle$, which resulted in the required data length $N = \Theta_\varepsilon(\varepsilon^{-2})$, we find that beyond a certain point, the error $1 - \langle \beta, \mathbf{u}^{(n)} \rangle$ converges to 0 geometrically. This accelerated geometric convergence yields an improved sample complexity of $N = \Theta_\varepsilon(\frac{1}{\varepsilon} \log \frac{1}{\varepsilon})$. See Appendix E for further details.

3.1.4 ESTIMATION OF THE LINK FUNCTION

Finally, we train MLP layer to predict the link function σ_*^{test} . Following prior work Nishikawa et al. (2025), we show that training MLP layer leads to small empirical loss. Moreover, we derive the upper bound of the predictive risk using Rademacher complexity. See Appendix F for details.

4 SYNTHETIC EXPERIMENT

We conducted a numerical experiment to examine the effectiveness of TTT compared to ICL. To evaluate whether our theory holds with real-world settings, we trained a GPT-2 model (Radford

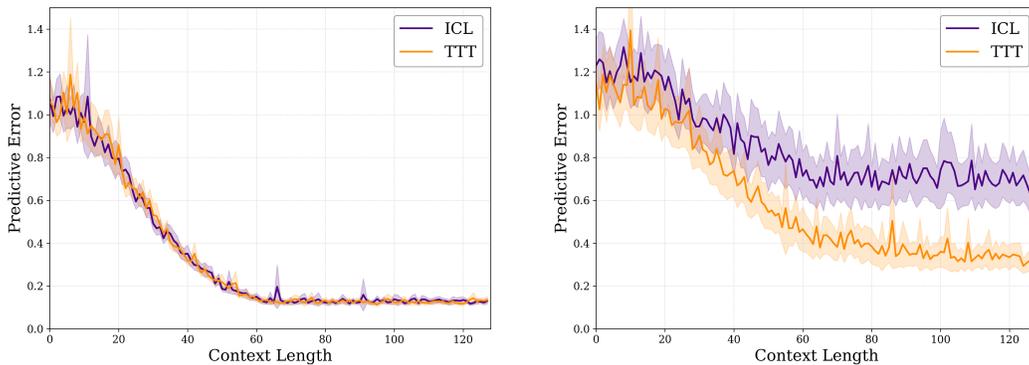
et al., 2019) to learn the Gaussian single-index functions. See Appendix I for further experimental details.

The dimensions were set to $d = r = 16$. For each task t , the data was generated as follows: $\mathbf{x}_1^t, \dots, \mathbf{x}_N^t, \mathbf{x}^t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \mathbf{I}_d)$, $\beta^t \sim \text{Unif}(\mathbb{S}^{r-1})$ and $y_i^t = \sigma_*^t(\langle \beta, \mathbf{x}_i \rangle) = \frac{1}{\sqrt{3!}} \text{He}_3(\langle \beta, \mathbf{x}_i \rangle) + \frac{c^t}{\sqrt{4!}} \text{He}_4(\langle \beta, \mathbf{x}_i \rangle)$, where $c^t \sim \text{Unif}(-0.5, 0.5)$. We compared the following two settings:

1. In-context learning, configured as Garg et al. (2023), Oko et al. (2024) and Nishikawa et al. (2025). We first pretrained the model with the data $(\mathbf{X}^t, \mathbf{y}^t, \mathbf{x}^t, y^t)_{t=1}^{T_{pt}}$, then we calculated the predictive loss.
2. Test-time training: We adapted the same pre-trained model as we used in the evaluation of ICL by fine-tuning the LoRA parameters applied to the attention and MLP layers.

Figure 1a highlights the result. Contrary to our initial expectation, the pretrained model’s ICL ability is powerful enough to adapt to the varying link function, while TTT does not lead to significant improvement of the prediction. This discrepancy is likely due to the difference between our theoretical setting and the practical experimental setup. While our theoretical analysis assumed a single-layer transformer consisting of the attention layer followed by the MLP layer, the experiment utilized a 6-layer GPT-2 model. As noted by Oko et al. (2024), MLP layers followed by linear attention (a structure present in GPT-2) are capable of adapting to nonlinear link functions. It is plausible that the GPT-2 model learned to fit the task-specific link function during pretraining, leaving little room for improvement via TTT. Nevertheless, we observe that TTT at least does not degrade performance even when ICL functions effectively.

This result spurred us to analyze the limitation of ICL and the true potential of TTT in more challenging tasks: we changed the distribution of link function as $y_i^t = \frac{c_3^t}{\sqrt{3!}} \text{He}_3(\langle \beta, \mathbf{x}_i \rangle) + \frac{c_4^t}{\sqrt{4!}} \text{He}_4(\langle \beta, \mathbf{x}_i \rangle)$, where $c_3^t \sim \text{Unif}(0.5, 1.5)$, $c_4^t \sim \text{Unif}(-0.5, 0.5)$, only in test-time. The result, summarized in Figure 1b, clearly shows the advantage of TTT over ICL. As the data length grows, the accuracy of TTT steadily improves, substantially surpassing the performance of ICL. Taken together, these results highlight TTT as a robust strategy that matches ICL in standard settings while significantly outperforming it under distribution shifts.



(a) In-distribution setting: The distribution of the link function is consistent between pretraining and test-time.

(b) Out-of-distribution setting: The distribution of the link function at test-time differs from that of pretraining

Figure 1: The predictive error of In-context learning (ICL) and Test-time training (TTT) for a pre-trained GPT-2 model on single-index polynomials (see Section 4 for details). While (a) shows that ICL effectively adapts to varying link functions for in-distribution tasks, (b) demonstrates that TTT outperforms ICL under distribution shift, showing its superior adaptability.

To further assess the capabilities of TTT, we examined its ability to leverage low-dimensional task structures. We fixed the ambient dimension at $d = 16$ and compared the predictive error in the in-distribution setting for a low intrinsic dimension ($r = 4$) versus a full-rank intrinsic dimension ($r = 16$). As illustrated in Figure 2, the predictive error is substantially lower for $r = 4$ than for

$r = 16$, even though the ambient dimension d remains the same. This finding highlights that the TTT-equipped model effectively adapts to the task’s intrinsic low-dimensionality.

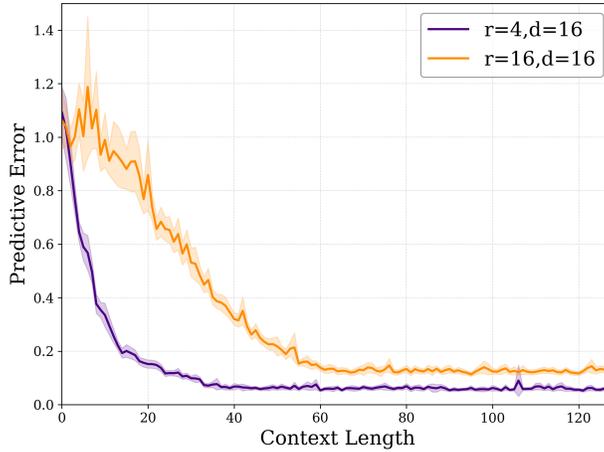


Figure 2: The predictive error of Test-time training for a pretrained GPT-2 model on single-index polynomials, comparing different intrinsic dimensions r (with fixed ambient dimension $d = 16$). Notably, when the intrinsic dimension is low ($r = 4$), the error remains minimal even with very short context lengths.

5 CONCLUSION

We have investigated test-time training combined with in-context learning. We provided an upper bound of the predictive loss in terms of m and N . Our result shows that for a large test-time sample size N_{test} , the predictive loss approaches the inevitable noise τ , even when d remains finite.

Future work and limitation We outline some limitations and future research directions.

- The result of the numerical experiment demonstrates that the ICL model’s loss decreases as we increase n in in-distribution cases. This finding is contrary to our theoretical expectation that ICL cannot adapt to the change in link function. Future work could clarify whether ICL without TTT can manage varying link functions.
- We only considered the single-index model where the link function σ_* is a polynomial. Investigating a more general class of input-output relationships is a possible extension of this work.
- In this work, we assumed that the test-time distribution of the feature vector β is the same as in the pretraining. Considering a distribution shift, such as investigating the situation where $\text{Supp}(\beta)_{test}$ is slightly different from $\text{Supp}(\beta)_{pt}$ is another interesting direction.
- Algorithm 1 divides the test-time training into 3 stages, training different layers sequentially. This differs from the typical situation where the entire model is trained at once, as was done in the experiment. Whether a similar upper bound of the predictive risk can be established in such settings remains to be examined.

LLM USAGE STATEMENT

Large language models are used for three purposes: to proofread and polish English writing, to help us find related works, and to write the code for the synthetic experiment. We did not use any LLM assistant for designing the problem settings and constructing the proofs.

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593

REFERENCES

- Emmanuel Abbe, Enric Boix Adserà, and Theodor Misiakiewicz. Sgd learning on neural networks: leap complexity and saddle-to-saddle dynamics. In Gergely Neu and Lorenzo Rosasco (eds.), *Proceedings of Thirty Sixth Conference on Learning Theory*, volume 195 of *Proceedings of Machine Learning Research*, pp. 2552–2623. PMLR, 12–15 Jul 2023. URL <https://proceedings.mlr.press/v195/abbe23a.html>.
- Rishabh Agarwal, Avi Singh, Lei M. Zhang, Bernd Bohnet, Luis Rosias, Stephanie Chan, Biao Zhang, Ankesh Anand, Zaheer Abbas, Azade Nova, John D. Co-Reyes, Eric Chu, Feryal Behbahani, Aleksandra Faust, and Hugo Larochelle. Many-shot in-context learning, 2024. URL <https://arxiv.org/abs/2404.11018>.
- Ekin Akyürek, Mehul Damani, Adam Zweiger, Linlu Qiu, Han Guo, Jyothish Pari, Yoon Kim, and Jacob Andreas. The surprising effectiveness of test-time training for few-shot learning, 2025. URL <https://arxiv.org/abs/2411.07279>.
- Luca Arnaboldi, Yatin Dandi, Florent Krzakala, Luca Pesce, and Ludovic Stephan. Repetita iuvant: Data repetition allows sgd to learn high-dimensional multi-index functions, 2025. URL <https://arxiv.org/abs/2405.15459>.
- Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. Online stochastic gradient descent on non-convex losses from high-dimensional inference. *Journal of Machine Learning Research*, 22(106):1–51, 2021. URL <http://jmlr.org/papers/v22/20-1288.html>.
- Jimmy Ba, Murat A. Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-dimensional asymptotics of feature learning: How one gradient step improves the representation, 2022. URL <https://arxiv.org/abs/2205.01445>.
- Yu Bai and Jason D. Lee. Beyond linearization: On quadratic and higher-order approximation of wide neural networks, 2020. URL <https://arxiv.org/abs/1910.01619>.
- Raphaël Berthier, Andrea Montanari, and Kangjie Zhou. Learning time-scales in two-layers neural networks. *Foundations of Computational Mathematics*, August 2024. ISSN 1615-3383. doi: 10.1007/s10208-024-09664-9. URL <http://dx.doi.org/10.1007/s10208-024-09664-9>.
- Juliette Bertrand, Giorgos Kordopatis Zilos, Yannis Kalantidis, and Giorgos Toliás. Test-time training for matching-based video object segmentation. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 20918–20941. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/4267d84ca2f6fbb4aa5172b76b433aca-Paper-Conference.pdf.
- Alberto Bietti, Joan Bruna, Clayton Sanford, and Min Jae Song. Learning single-index models with shallow neural networks, 2022. URL <https://arxiv.org/abs/2210.15651>.
- Irina Bigoulaeva, Harish Tayyar Madabushi, and Iryna Gurevych. The inherent limits of pretrained llms: The unexpected convergence of instruction tuning and in-context learning capabilities, 2025. URL <https://arxiv.org/abs/2501.08716>.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL <https://arxiv.org/abs/2005.14165>.
- Xiang Cheng, Yuxin Chen, and Suvrit Sra. Transformers implement functional gradient descent to learn non-linear functions in context, 2024. URL <https://arxiv.org/abs/2312.06528>.

-
- 594 Alex Damian, Jason D. Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations
595 with gradient descent, 2022. URL <https://arxiv.org/abs/2206.15144>.
596
- 597 Alex Damian, Eshaan Nichani, Rong Ge, and Jason D. Lee. Smoothing the landscape boosts the
598 signal for sgd: Optimal sample complexity for learning single index models, 2023. URL <https://arxiv.org/abs/2305.10633>.
599
- 600 Alex Damian, Loucas Pillaud-Vivien, Jason D. Lee, and Joan Bruna. Computational-statistical gaps
601 in gaussian single-index models, 2024. URL <https://arxiv.org/abs/2403.05529>.
602
- 603 Yatin Dandi, Florent Krzakala, Bruno Loureiro, Luca Pesce, and Ludovic Stephan. How two-layer
604 neural networks learn, one (giant) step at a time, 2025. URL <https://arxiv.org/abs/2305.18270>.
605
- 606 Benoit Dherin, Michael Munn, Hanna Mazzawi, Michael Wunder, and Javier Gonzalvo. Learning
607 without training: The implicit dynamics of in-context learning, 2025. URL <https://arxiv.org/abs/2507.16003>.
608
609
- 610 Konstantin Donhauser, Mingqi Wu, and Fanny Yang. How rotational invariance of common kernels
611 prevents generalization in high dimensions. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 2804–2814. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/donhauser21a.html>.
612
613
614
- 615 Rishabh Dudeja and Daniel Hsu. Learning single-index models in gaussian space. In Sébastien
616 Bubeck, Vianney Perchet, and Philippe Rigollet (eds.), *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learning Research*, pp. 1887–1930. PMLR, 06–09 Jul 2018. URL <https://proceedings.mlr.press/v75/dudeja18a.html>.
617
618
619
- 620 Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What can transformers learn
621 in-context? a case study of simple function classes, 2023. URL <https://arxiv.org/abs/2208.01066>.
622
623
- 624 Khashayar Gatmiry, Nikunj Saunshi, Sashank J. Reddi, Stefanie Jegelka, and Sanjiv Kumar. Can
625 looped transformers learn to implement multi-step gradient descent for in-context learning?, 2024. URL <https://arxiv.org/abs/2410.08292>.
626
627
- 628 Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Linearized two-layers
629 neural networks in high dimension, 2020. URL <https://arxiv.org/abs/1904.12191>.
- 630 Halil Alperen Gozeten, M. Emrullah Ildiz, Xuechen Zhang, Mahdi Soltanolkotabi, Marco Mondelli,
631 and Samet Oymak. Test-time training provably improves transformers as in-context learners,
632 2025. URL <https://arxiv.org/abs/2503.11842>.
633
- 634 Jinwu Hu, Zhitian Zhang, Guohao Chen, Xutao Wen, Chao Shuai, Wei Luo, Bin Xiao, Yuanqing Li,
635 and Mingkui Tan. Test-time learning for large language models, 2025. URL <https://arxiv.org/abs/2505.20633>.
636
- 637 Yu Huang, Yuan Cheng, and Yingbin Liang. In-context convergence of transformers, 2023. URL
638 <https://arxiv.org/abs/2310.05249>.
639
- 640 Yixing Jiang, Jeremy Irvin, Ji Hun Wang, Muhammad Ahmed Chaudhry, Jonathan H. Chen, and
641 Andrew Y. Ng. Many-shot in-context learning in multimodal foundation models, 2024. URL
642 <https://arxiv.org/abs/2405.09798>.
- 643 Nirmal Joshi, Theodor Misiakiewicz, and Nathan Srebro. On the complexity of learning sparse func-
644 tions with statistical and gradient queries, 2024. URL <https://arxiv.org/abs/2407.05622>.
645
646
- 647 Juno Kim and Taiji Suzuki. Transformers provably solve parity efficiently with chain of thought,
2025. URL <https://arxiv.org/abs/2410.08633>.

648 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL
649 <https://arxiv.org/abs/1412.6980>.
650

651 Jason D. Lee, Kazusato Oko, Taiji Suzuki, and Denny Wu. Neural network learns low-dimensional
652 polynomials with sgd near the information-theoretic limit, 2024. URL <https://arxiv.org/abs/2406.01581>.
653

654 Yuejiang Liu, Parth Kothari, Bastien van Delft, Baptiste Bellot-Gurlet, Taylor Mordan, and Alexan-
655 dre Alahi. Ttt++: When does self-supervised test-time training fail or thrive? In *Neural*
656 *Information Processing Systems*, 2021. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:247410579)
657 [CorpusID:247410579](https://api.semanticscholar.org/CorpusID:247410579).

658 Alireza Mousavi-Hosseini, Sejun Park, Manuela Girotti, Ioannis Mitliagkas, and Murat A. Erdogdu.
659 Neural networks efficiently learn low-dimensional representations with sgd, 2023. URL <https://arxiv.org/abs/2209.14863>.
660
661

662 Omar Naim and Nicholas Asher. Analyzing limits for in-context learning, 2025. URL <https://arxiv.org/abs/2502.03503>.
663

664 Eshaan Nichani, Alex Damian, and Jason D. Lee. How transformers learn causal structure with
665 gradient descent, 2024. URL <https://arxiv.org/abs/2402.14735>.
666

667 Naoki Nishikawa, Yujin Song, Kazusato Oko, Denny Wu, and Taiji Suzuki. Nonlinear transform-
668 ers can perform inference-time feature learning. In *Forty-second International Conference on*
669 *Machine Learning*, 2025. URL <https://openreview.net/forum?id=xQTSvP57C3>.

670 Kazusato Oko, Yujin Song, Taiji Suzuki, and Denny Wu. Pretrained transformer efficiently learns
671 low-dimensional target functions in-context, 2024. URL <https://arxiv.org/abs/2411.02544>.
672
673

674 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever.
675 Language models are unsupervised multitask learners. 2019. URL <https://api.semanticscholar.org/CorpusID:160025533>.
676

677 Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei A. Efros, and Moritz Hardt. Test-time
678 training with self-supervision for generalization under distribution shifts, 2020. URL <https://arxiv.org/abs/1909.13231>.
679
680

681 Johannes von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordv-
682 intsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient
683 descent, 2023. URL <https://arxiv.org/abs/2212.07677>.

684 Ruiqi Zhang, Spencer Frei, and Peter L. Bartlett. Trained transformers learn linear models in-
685 context, 2023. URL <https://arxiv.org/abs/2306.09927>.
686

687 Ruiqi Zhang, Jingfeng Wu, and Peter L. Bartlett. In-context learning of a linear transformer block:
688 Benefits of the mlp component and one-step gd initialization, 2024. URL <https://arxiv.org/abs/2402.14951>.
689

690 Tianyuan Zhang, Sai Bi, Yicong Hong, Kai Zhang, Fujun Luan, Songlin Yang, Kalyan Sunkavalli,
691 William T. Freeman, and Hao Tan. Test-time training done right, 2025. URL <https://arxiv.org/abs/2505.23884>.
692
693
694
695
696
697
698
699
700
701

| | | |
|-----|--|-----------|
| 702 | CONTENTS | |
| 703 | | |
| 704 | | |
| 705 | 1 Introduction | 1 |
| 706 | 1.1 Our contribution | 2 |
| 707 | | |
| 708 | 1.2 Related works | 2 |
| 709 | | |
| 710 | 2 Preliminaries and problem settings | 3 |
| 711 | | |
| 712 | 2.1 In-context learning and test-time training | 3 |
| 713 | | |
| 714 | 2.2 Single index model | 3 |
| 715 | | |
| 716 | 2.3 Student model | 4 |
| 717 | | |
| 718 | | |
| 719 | 3 Main result | 6 |
| 720 | | |
| 721 | 3.1 Proof sketch | 7 |
| 722 | 3.1.1 Exploiting the pretrained attention matrix | 7 |
| 723 | 3.1.2 Weak recovery | 8 |
| 724 | | |
| 725 | 3.1.3 Strong recovery | 8 |
| 726 | | |
| 727 | 3.1.4 Estimation of the link function | 8 |
| 728 | | |
| 729 | 4 Synthetic experiment | 8 |
| 730 | | |
| 731 | 5 Conclusion | 10 |
| 732 | | |
| 733 | | |
| 734 | A Definition of high probability events | 15 |
| 735 | | |
| 736 | B Pretraining | 15 |
| 737 | | |
| 738 | | |
| 739 | C Gradient calculation | 16 |
| 740 | | |
| 741 | D One step gradient descent for weak recovery | 22 |
| 742 | | |
| 743 | | |
| 744 | E Strong recovery | 24 |
| 745 | | |
| 746 | F Training MLP layer | 30 |
| 747 | | |
| 748 | | |
| 749 | G Proof of the Theorem 1 | 32 |
| 750 | | |
| 751 | H Additional Lemmas | 32 |
| 752 | | |
| 753 | H.1 Pretrained matrix | 33 |
| 754 | | |
| 755 | I Experimental details | 35 |

A DEFINITION OF HIGH PROBABILITY EVENTS

Definition 7. We say that an event A occurs with high probability when the following holds:

$$1 - P(A) \leq O(d^{-C_*}),$$

where C_* is a sufficiently large constant that is independent of d and r .

Note that we can redefine C_* to be sufficiently large if needed. A basic example is the Gaussian tail bound. When $x \sim \mathcal{N}(0, 1)$ and $t > 0$, we have

$$P(|x| > t) \leq 2 \exp(-t^2/2).$$

Thus, by letting $t = \sqrt{2C_* \log d}$, we see that

$$P(|x| > \sqrt{2C_* \log d}) \leq 2d \exp(-C_*) = O(d^{-C_*}).$$

In such a situation, we say that $|x| = O(\sqrt{\log d})$ with high probability.

When A_1, \dots, A_M occurs with high probability where $M = \text{poly}(d)$, then the intersection $A_1 \cap A_2 \cap \dots \cap A_M$ occurs with high probability. In this paper, we assume that $N_{pt}, T_{pt}, N_1, N_2, N_3, N_4, m = O(\text{poly}(d))$. Because of this assumption, when we take the intersection of high probability events A_1, \dots, A_M , $M = \text{poly}(d)$ is satisfied.

B PRETRAINING

We follow the pretraining algorithm that was considered in Nishikawa et al. (2025). In the original paper, the link function σ_* is fixed across all tasks, whereas in this paper, the link function varies depending on the tasks. We should take this difference into account and ensure that the pretraining algorithm works properly even in our setting.

Lemma 5 (Formal). After running the pretraining in Algorithm 1 with $T_{pt}, N_{pt} = \tilde{\Omega}(r^2 d^{Q+2})$, it holds that

$$\Gamma^* = \frac{1}{r^{1/2} \log^{C_\kappa} d} (r \mathbb{E}_\beta[\beta \beta^\top] + \mathbf{N})$$

with high probability over the data distribution, where $\|\mathbf{N}\|_F = O_d(1/\sqrt{d})$ holds, where C_κ can be taken to be sufficiently large.

Proof. We only consider the difference between our settings and Nishikawa et al. (2025). See Section C of Nishikawa et al. (2025) for the original argument.

If we fix t , we may apply Lemma 19 and Lemma 20 in Nishikawa et al. (2025). In addition to that, since the norm upper bound remains unchanged, we can bound the difference between the expectation and the actual value just like Eq.(C.3) and Eq.(C.4) in Lemma 21 of Nishikawa et al. (2025). What remains is to calculate $\mathbb{E}_{\mathbf{x}, \beta, t}[yz^k \beta \mathbf{x}^\top] = (\rho\sqrt{d})^{-k} \mathbb{E}_\beta[\mathbb{E}_{\mathbf{x}, t}[\sigma_*^t(\langle \beta, \mathbf{x} \rangle) (\langle \beta, \mathbf{x} \rangle)^k] \beta \beta^\top] + (\rho\sqrt{d})^{-k} k \mathbb{E}_\beta[\mathbb{E}_{\mathbf{x}, t}[\sigma_*^t(\langle \beta, \mathbf{x} \rangle) (\langle \beta, \mathbf{x} \rangle)^{k-1}] \beta \beta^\top]$. Because of the definition of $\text{ie}(\sigma_*)$, $\mathbb{E}_{\mathbf{x}, \beta, t}[yz^k \beta \mathbf{x}^\top] = 0$ when $k < Q - 1$. As we assume that $\mathbb{E}_t[c_Q] = \Theta(1)$, when $k = Q - 1$, we see that

$$\begin{aligned} \mathbb{E}_\beta[\mathbb{E}_{\mathbf{x}, t}[\sigma_*^t(\langle \beta, \mathbf{x} \rangle) (\langle \beta, \mathbf{x} \rangle)^k] \beta \beta^\top] &\asymp \mathbb{E}_\beta[\mathbb{E}_{\mathbf{x}}[\mathbb{E}_t[c_Q] \text{He}_{Q-1}(\langle \beta, \mathbf{x} \rangle) \langle \beta, \mathbf{x} \rangle^{Q-1}] \beta \beta^\top] \\ &\asymp (\rho\sqrt{d})^{-(Q-1)} \mathbb{E}_\beta[\beta \beta^\top], \end{aligned}$$

and this is the main term that is proportional to $\mathbb{E}_\beta[\beta \beta^\top]$. This means that we can use the same argument as Nishikawa et al. (2025) with $\text{ie}(\sigma_*) = Q$. \square

Let $\kappa = \log^{-C_\kappa} d$. In other words, κ satisfies $\kappa = \Theta(\log^{-C_\kappa} d)$ where C_κ can be taken sufficiently large.

Lemma 8. Suppose that the context length satisfies $N_1 = \tilde{\Omega}(r^{\text{ge}(\sigma_*)+1})$. Then, it holds with high probability that

$$g(\Gamma^*, \mathbf{X}_{N_1}, \mathbf{y}_{N_1}, \mathbf{x}) = P'_1 + P'_2 \left(\frac{\langle \mathbf{x}, \beta \rangle}{\sqrt{r}} \right)^{\text{ge}(\sigma_*)} + n_3,$$

where $P'_1 = o_d(1)$, $P'_2 = \Theta_d((\log d)^{-C_{P_2}})$ and $n_3 = o_d(P'_2 r^{-\text{ge}(\sigma_*)/2-1/2} \log^{-2 \deg(\sigma_*)+2} d)$.

810 *Proof.* When $N_1 = \tilde{\Omega}(r^{\text{ge}(\sigma_*)+1})$, the (h.o.t.) in the proof of Proposition 11 in Nishikawa et al.
811 (2025) can be evaluated as $o(\rho^{-1}P_2r^{-\text{ge}(\sigma_*)/2-1/2} \log^{-2\text{deg}(\sigma_*)+2} d)$ by carefully examining the
812 term. Using this fact, the same argument as the proof of Proposition 11 in Nishikawa et al. (2025)
813 yields the result. \square

815 C GRADIENT CALCULATION

816 For the sake of simplicity, we write σ_*^{test} as σ_* . Let $\Gamma_u = \Gamma^* + uu^\top$, $\mathbf{x}_i^* = \sqrt{r}\Gamma^*\mathbf{x}_i$, $\beta^* =$
817 $\sqrt{r}\Gamma^*\beta$ and $f_{\text{IC}}(\mathbf{x}) = \sum_{j=1}^m a_j \sigma \left(v_j \frac{\sum_{i=1}^{N_1} y_i e^{y_i/\rho} e^{\mathbf{x}_i^{*\top} \Gamma_u \mathbf{x} / \rho}}{\sum_{i=1}^{N_1} e^{y_i/\rho} e^{\mathbf{x}_i^{*\top} \Gamma_u \mathbf{x} / \rho}} + b_j \right)$. Moreover, let us define A_i, B_i
818 as $\bar{\sigma}_*(z) \exp((\bar{\sigma}_*(z)) = \sum_{i \geq 0} \frac{A_i}{i!} z^i$ and $\exp((\bar{\sigma}_*(z)) = \sum_{i \geq 0} \frac{B_i}{i!} z^i$, where

$$819 \bar{\sigma}_*(z) := \begin{cases} \frac{\sigma_*(z)}{\rho} & \left(\text{if } \left| \frac{\sigma_*(z)}{\rho} \right| \leq \frac{1}{\log d} \right) \\ 0 & (\text{otherwise}) \end{cases} .$$

820 **Lemma 9.** Take $\rho = \Theta((\log d)^{C_\rho})$ where C_ρ is a constant sufficiently large. Suppose that \mathbf{u} satisfies
821 $\|\mathbf{u}\| = C_u \leq 1$ and $\langle \mathbf{u}, \mathbf{x}_i^* \rangle = C_u \tilde{O}_d(1)$. Then

$$822 \sqrt{r}\Gamma^* \frac{\partial f_{\text{IC}}}{\partial \mathbf{u}} = \alpha L_m \left\{ \sqrt{r}\Gamma^* \mathbf{x} \cdot (\gamma(\mathbf{x}, y) \langle \beta^*, \mathbf{u} \rangle + \sqrt{r}\Gamma^* \gamma(\mathbf{x}, y) \langle \mathbf{x}, \mathbf{u} \rangle \beta^* + \sqrt{r}\Gamma^* \mathbf{x} \cdot \mathbf{n}_1 + \langle \mathbf{x}, \mathbf{u} \rangle \mathbf{n}_2) \right\},$$

823 where

$$824 \gamma(\mathbf{x}, y) = P_0 + P_1 z + \dots + P_{\text{ie}(\sigma)-1} z^{\text{ie}(\sigma)-1} + \dots$$

825 is satisfied with $z = \langle \beta, \sqrt{r}\Gamma^{*\top} \Gamma_u \mathbf{x} \rangle / \rho$, and $\mathbf{n}_1 = \tilde{O} \left(C_u \sqrt{\frac{1}{N_1}} \right)$ and $\mathbf{n}_2 = \tilde{O} \left(\sqrt{\frac{r}{N_1}} \right)$ hold.

826 Note that from Lemma 29, the condition $\langle \mathbf{u}, \mathbf{x}_i^* \rangle = C_u \tilde{O}_d(1)$ is satisfied at the initialization.

827 *Proof.* Let

$$828 \pi_1(\mathbf{x}, y) = \frac{1}{N_1} \sum_{i=1}^{N_1} \frac{y_i}{\rho} e^{y_i/\rho} e^{\mathbf{x}_i^{*\top} \Gamma_u \mathbf{x} / \rho},$$

$$829 \pi_2(\mathbf{x}, y) = \frac{1}{N_1} \sum_{i=1}^{N_1} e^{y_i/\rho} e^{\mathbf{x}_i^{*\top} \Gamma_u \mathbf{x} / \rho},$$

$$830 \xi_1(\mathbf{x}, y) = \frac{1}{N_1} \sum_{i=1}^{N_1} \frac{y_i}{\rho} e^{y_i/\rho} e^{\mathbf{x}_i^{*\top} \Gamma_u \mathbf{x} / \rho} \{ \langle \mathbf{x}_i^*, \mathbf{u} \rangle \mathbf{x} + \langle \mathbf{u}, \mathbf{x} \rangle \mathbf{x}_i^* \},$$

$$831 \xi_2(\mathbf{x}, y) = \frac{1}{N_1} \sum_{i=1}^{N_1} e^{y_i/\rho} e^{\mathbf{x}_i^{*\top} \Gamma_u \mathbf{x} / \rho} \{ \langle \mathbf{x}_i^*, \mathbf{u} \rangle \mathbf{x} + \langle \mathbf{u}, \mathbf{x} \rangle \mathbf{x}_i^* \}.$$

832 Then, we can write the gradient of f_{IC} with respect to \mathbf{u} as

$$833 \frac{\partial f_{\text{IC}}}{\partial \mathbf{u}} = \sum_{j=1}^m a_j \sigma' \left(v_j \frac{\sum_{i=1}^{N_1} y_i e^{y_i/\rho} e^{\mathbf{x}_i^{*\top} \Gamma_u \mathbf{x} / \rho}}{\sum_{i=1}^{N_1} e^{y_i/\rho} e^{\mathbf{x}_i^{*\top} \Gamma_u \mathbf{x} / \rho}} + b_j \right) \cdot \frac{\xi_1 \pi_2 - \xi_2 \pi_1}{\pi_2^2}.$$

834 First, from the assumption $\langle \mathbf{u}, \mathbf{x}_i^* \rangle = C_u \tilde{O}_d(1)$ with $C_u \leq 1$ and Lemma 25, we can see that

$$835 \mathbf{x}_i^* \mathbf{u} \mathbf{u}^\top \mathbf{x} = \langle \mathbf{x}_i^*, \mathbf{u} \rangle \langle \mathbf{u}, \mathbf{x} \rangle = \tilde{O}_d(1).$$

836 Therefore, by taking C_ρ sufficiently large, we can say $\mathbf{x}_i^{*\top} \Gamma_u \mathbf{x} / \rho = o_d(1)$ with high probability.
837 Hence, following the same argument as Proposition 11 in Nishikawa et al. (2025), we can say that
838 the content of $\rho'(\cdot)$ is $P_1(1 + o(1))$ using the positive constant P_1 . Therefore, $\rho'(\cdot) = 1$ if and only
839 $v_j = 1$, which means that

$$840 \frac{\partial f_{\text{IC}}}{\partial \mathbf{u}} = \alpha L_m \frac{\xi_1 \pi_2 - \xi_2 \pi_1}{\pi_2^2}.$$

Next, we evaluate the term $\frac{\xi_1 \pi_2 - \xi_2 \pi_1}{\pi_2^2}$. For that purpose, we first derive the expectation of each component of this term, and then bound the deviation of its actual value from the expectation. Let $k_\rho = \frac{1}{2}(\exp(\tau/\rho) + \exp(-\tau/\rho))$, then the following holds:

$$\begin{aligned}
\mathbb{E}_{\mathbf{x}_1, \zeta_1} [e^{\bar{\sigma}_*(\langle \beta, \mathbf{x}_1 \rangle)/\rho + \zeta_1/\rho} e^{\mathbf{x}_1^* \top \Gamma_u \mathbf{x}/\rho} \mathbf{x}_1^*] &= k_\rho \mathbb{E}_{\mathbf{x}_1} [e^{\sigma_*(\langle \beta, \mathbf{x}_1 \rangle)} e^{\mathbf{x}_1^* \top \Gamma_u \mathbf{x}/\rho} \mathbf{x}_1^*] \\
&= k_\rho \mathbb{E}_{\mathbf{x}_1} \left[\sum_{i \geq 0} \frac{B_i}{i!} \text{He}_i(\langle \beta, \mathbf{x}_1 \rangle) e^{\mathbf{x}_1^* \top \Gamma_u \mathbf{x}/\rho} \mathbf{x}_1^* \right] \\
&= \sqrt{r} \Gamma^* k_\rho \beta \mathbb{E}_{\mathbf{x}_1} \left[\sum_{i \geq 0} \frac{B_{i+1}}{i!} \text{He}_i(\langle \beta, \mathbf{x}_1 \rangle) e^{\mathbf{x}_1^* \top \Gamma_u \mathbf{x}/\rho} \right] \\
&\quad + \sqrt{r} \Gamma^* k_\rho \frac{\Gamma_u \mathbf{x}}{\rho} \mathbb{E}_{\mathbf{x}_1} \left[\sum_{i \geq 0} \frac{B_i}{i!} \text{He}_i(\langle \beta, \mathbf{x}_1 \rangle) e^{\mathbf{x}_1^* \top \Gamma_u \mathbf{x}/\rho} \right], \tag{1}
\end{aligned}$$

(recall that $z = \langle \beta, \sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x} \rangle / \rho$). Here, we used Stein's lemma to derive the final equality. The expectation in the second term of the RHS can be calculated as

$$\begin{aligned}
&\mathbb{E}_{\mathbf{x}_1} \left[\sum_{i \geq 0} \frac{B_i}{i!} \text{He}_i(\langle \beta, \mathbf{x}_1 \rangle) e^{\mathbf{x}_1^* \top \Gamma_u \mathbf{x}/\rho} \right] \\
&= \mathbb{E}_{\mathbf{x}_1} \left[\sum_{i \geq 0} \frac{B_i}{i!} \text{He}_i(\langle \beta, \mathbf{x}_1 \rangle) \sum_{j \geq 0} \frac{1}{j!} \left(\frac{\|\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}\|}{\rho} \right)^j \exp\left(\frac{\|\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}\|^2}{2\rho^2}\right) \text{He}_j\left(\left\langle \mathbf{x}_1, \frac{\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}}{\|\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}\|} \right\rangle\right) \right] \\
&= \sum_{i \geq 0} \frac{B_i}{i!} \left(\frac{\|\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}\|}{\rho} \right)^i \exp\left(\frac{\|\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}\|^2}{2\rho^2}\right) \left\langle \beta, \frac{\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}}{\|\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}\|} \right\rangle^i \\
&= \exp\left(\frac{\|\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}\|^2}{2\rho^2}\right) \left(\sum_{i \geq 0} \frac{B_i}{i!} \langle \beta, \sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x} / \rho \rangle^i \right).
\end{aligned}$$

The expectation in the first term of Eq. (1) can be evaluated in a similar way, and we obtain:

$$\begin{aligned}
&\mathbb{E}_{\mathbf{x}_1, \zeta_1} [e^{\bar{\sigma}_*(\langle \mathbf{x}_1, \beta \rangle)/\rho + \zeta_1/\rho} e^{\mathbf{x}_1^* \top \Gamma_u \mathbf{x}/\rho} \mathbf{x}_1^*] \\
&= \sqrt{r} \Gamma^* k_\rho \beta \exp\left(\frac{\|\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}\|^2}{2\rho^2}\right) \left(\sum_{i \geq 0} \frac{B_{i+1}}{i!} z^i \right) \\
&\quad + \sqrt{r} \Gamma^* \frac{k_\rho}{\rho} \Gamma_u \mathbf{x} \exp\left(\frac{\|\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}\|^2}{2\rho^2}\right) \left(\sum_{i \geq 0} \frac{B_i}{i!} z^i \right). \tag{2}
\end{aligned}$$

Similarly, using $k'_\rho := \frac{\tau}{2\rho}(\exp(\tau/\rho) - \exp(-\tau/\rho))$, we have that

$$\begin{aligned}
&\mathbb{E}_{\mathbf{x}_1, \zeta_1} \left[\frac{\bar{\sigma}_*(\langle \mathbf{x}_1, \beta \rangle) + \zeta_1}{\rho} e^{\bar{\sigma}_*(\langle \mathbf{x}_1, \beta \rangle)/\rho + \zeta_1/\rho} e^{\mathbf{x}_1^* \top \Gamma_u \mathbf{x}/\rho} \mathbf{x}_1^* \right] \\
&= \sqrt{r} \Gamma^* k_\rho \exp\left(\frac{\|\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}\|^2}{2\rho^2}\right) \left(\sum_{i \geq 0} \frac{A_{i+1}}{i!} z^i \right) \beta + \sqrt{r} \Gamma^* \frac{k_\rho}{\rho} \exp\left(\frac{\|\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}\|^2}{2\rho^2}\right) \left(\sum_{i \geq 0} \frac{A_i}{i!} z^i \right) \Gamma_u \mathbf{x} \\
&\quad + \sqrt{r} \Gamma^* k'_\rho \exp\left(\frac{\|\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}\|^2}{2\rho^2}\right) \left(\sum_{i \geq 0} \frac{B_{i+1}}{i!} z^i \right) \beta + \sqrt{r} \Gamma^* \frac{k'_\rho}{\rho} \exp\left(\frac{\|\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}\|^2}{2\rho^2}\right) \left(\sum_{i \geq 0} \frac{B_i}{i!} z^i \right) \Gamma_u \mathbf{x}.
\end{aligned}$$

Let $F_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} \left[\frac{y_i}{\rho} e^{y_i/\rho} e^{\mathbf{x}_i^* \top \Gamma_u \mathbf{x}/\rho} \mathbf{x}_i^* \right]$ and $F_2 = \frac{1}{N_1} \sum_{i=1}^{N_1} [e^{y_i/\rho} e^{\mathbf{x}_i^* \top \Gamma_u \mathbf{x}/\rho} \mathbf{x}_i^*]$. Then, the expectation of ξ_1 and ξ_2 can be written as $\mathbb{E}_{\mathbf{x}_1, \zeta_1} [\xi_1] = \langle \mathbb{E}[F_1], \mathbf{u} \rangle \mathbf{x} + \langle \mathbf{u}, \mathbf{x} \rangle \mathbb{E}[F_1]$ and $\mathbb{E}_{\mathbf{x}_1, \zeta_1} [\xi_2] = \langle \mathbb{E}[F_2], \mathbf{u} \rangle \mathbf{x} + \langle \mathbf{u}, \mathbf{x} \rangle \mathbb{E}[F_2]$ respectively.

Based on the argument above, the term $\frac{\xi_1\pi_2 - \xi_2\pi_1}{\pi_2^2}$ is concentrated on

$$(\mathbb{E}[\pi_2])^{-2} \{ \langle \mathbb{E}[F_1] \mathbb{E}[\pi_2], \mathbf{u} \rangle \mathbf{x} + \langle \mathbf{u}, \mathbf{x} \rangle \mathbb{E}[F_1] \mathbb{E}[\pi_2] - \langle \mathbb{E}[F_2] \mathbb{E}[\pi_1], \mathbf{u} \rangle \mathbf{x} - \langle \mathbf{u}, \mathbf{x} \rangle \mathbb{E}[F_2] \mathbb{E}[\pi_1] \}.$$

Also, following the same argument to obtain Eq. (2), we have

$$\mathbb{E}_{\mathbf{x}_1, \zeta_1} \left[e^{\bar{\sigma}_*(\langle \mathbf{x}_1, \beta \rangle) / \rho + \zeta_1 / \rho} e^{\mathbf{x}_1^* \top \Gamma_u \mathbf{x} / \rho} \right] = k_\rho \exp \left(\frac{\|\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}\|^2}{2\rho^2} \right) \left(\sum_{i \geq 0} \frac{B_i}{i!} z^i \right),$$

and

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}_1, \zeta_1} \left[\frac{\bar{\sigma}_*(\langle \mathbf{x}_1, \beta \rangle) + \zeta_1}{\rho} e^{\bar{\sigma}_*(\langle \mathbf{x}_1, \beta \rangle) / \rho + \zeta_1 / \rho} e^{\mathbf{x}_1^* \top \Gamma_u \mathbf{x} / \rho} \right] \\ &= k_\rho \exp \left(\frac{\|\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}\|^2}{2\rho^2} \right) \left(\sum_{i \geq 0} \frac{A_i}{i!} z^i \right) + k'_\rho \exp \left(\frac{\|\sqrt{r} \Gamma^* \top \Gamma_u \mathbf{x}\|^2}{2\rho^2} \right) \left(\sum_{i \geq 0} \frac{B_i}{i!} z^i \right). \end{aligned}$$

Using these result, we obtain that

$$\begin{aligned} & (\mathbb{E}_{\mathbf{x}_1, \zeta_1}[\pi_2])^{-2} (\mathbb{E}_{\mathbf{x}_1, \zeta_1}[F_1] \mathbb{E}_{\mathbf{x}_1, \zeta_1}[\pi_2] - \mathbb{E}_{\mathbf{x}_1, \zeta_1}[F_2] \mathbb{E}_{\mathbf{x}_1, \zeta_1}[\pi_1]) \\ &= \sqrt{r} \Gamma^* \left(k_\rho \sum_{i \geq 0} \frac{B_i}{i!} z^i \right)^{-2} \times \\ & \left\{ \left[k_\rho \left(\sum_{i \geq 0} \frac{A_{i+1}}{i!} z^i \right) \beta + \frac{k_\rho}{\rho} \left(\sum_{i \geq 0} \frac{A_i}{i!} z^i \right) \Gamma_u \mathbf{x} + k'_\rho \left(\sum_{i \geq 0} \frac{B_{i+1}}{i!} z^i \right) \beta + \frac{k'_\rho}{\rho} \left(\sum_{i \geq 0} \frac{B_i}{i!} z^i \right) \Gamma_u \mathbf{x} \right] \cdot k_\rho \left(\sum_{i \geq 0} \frac{B_i}{i!} z^i \right) \right. \\ & \left. - \left[k_\rho \beta \left(\sum_{i \geq 0} \frac{B_{i+1}}{i!} z^i \right) + \frac{k_\rho}{\rho} \Gamma_u \mathbf{x} \left(\sum_{i \geq 0} \frac{B_i}{i!} z^i \right) \right] \cdot \left[k_\rho \left(\sum_{i \geq 0} \frac{A_i}{i!} z^i \right) + k'_\rho \left(\sum_{i \geq 0} \frac{B_i}{i!} z^i \right) \right] \right\} \\ &= \sqrt{r} \Gamma^* \left(k_\rho \sum_{i \geq 0} \frac{B_i}{i!} z^i \right)^{-2} \left\{ \left[k_\rho \left(\sum_{i \geq 0} \frac{A_{i+1}}{i!} z^i \right) \beta + k'_\rho \left(\sum_{i \geq 0} \frac{B_{i+1}}{i!} z^i \right) \beta \right] \cdot k_\rho \left(\sum_{i \geq 0} \frac{B_i}{i!} z^i \right) \right. \\ & \left. - k_\rho \beta \left(\sum_{i \geq 0} \frac{B_{i+1}}{i!} z^i \right) \cdot \left[k_\rho \left(\sum_{i \geq 0} \frac{A_i}{i!} z^i \right) + k'_\rho \left(\sum_{i \geq 0} \frac{B_i}{i!} z^i \right) \right] \right\}. \end{aligned}$$

Therefore, we can expand this value using z . Specifically,

$$(\mathbb{E}_{\mathbf{x}_1, \zeta_1}[\pi_2])^{-2} (\mathbb{E}_{\mathbf{x}_1, \zeta_1}[F_1] \mathbb{E}_{\mathbf{x}_1, \zeta_1}[\pi_2] - \mathbb{E}_{\mathbf{x}_1, \zeta_1}[F_2] \mathbb{E}_{\mathbf{x}_1, \zeta_1}[\pi_1]) = \sqrt{r} \Gamma^* (P_0 + P_1 z + \dots) \beta$$

holds. By letting $\gamma(\mathbf{x}, y) = P_0 + P_1 z + \dots$, we obtain

$$\frac{\mathbb{E}[\xi_1] \mathbb{E}[\pi_2] - \mathbb{E}[\xi_2] \mathbb{E}[\pi_1]}{\mathbb{E}[\pi_2]^2} = \mathbf{x} \cdot (\gamma(\mathbf{x}, y) \langle \sqrt{r} \Gamma^* \beta, \mathbf{u} \rangle) + \gamma(\mathbf{x}, y) \langle \mathbf{x}, \mathbf{u} \rangle \sqrt{r} \Gamma^* \beta.$$

Following the same argument as Lemma 20 in Nishikawa et al. (2025) yields the order of P_i .

Next we deal with the deviation from the expectation. Using the same technique as Lemma 18 in Nishikawa et al. (2025), we have

$$\frac{1}{N_1} \sum_{i=1}^{N_1} \frac{y_i}{\rho} e^{y_i / \rho} e^{\mathbf{x}_i^* \top \Gamma_u \mathbf{x} / \rho} = \mathbb{E}_{\mathbf{x}_1, \zeta_1} \left[\frac{\bar{\sigma}_*(\langle \mathbf{x}_1, \beta \rangle) + \zeta_1}{\rho} e^{\bar{\sigma}_*(\langle \mathbf{x}_1, \beta \rangle) / \rho + \zeta_1 / \rho} e^{\mathbf{x}_1^* \top \Gamma_u \mathbf{x} / \rho} \right] + \tilde{O}(N_1^{-1/2}),$$

and

$$\frac{1}{N_1} \sum_{i=1}^{N_1} e^{y_i / \rho} e^{\mathbf{x}_i^* \top \Gamma_u \mathbf{x} / \rho} = \mathbb{E}_{\mathbf{x}_1, \zeta_1} \left[e^{\bar{\sigma}_*(\langle \mathbf{x}_1, \beta \rangle) / \rho + \zeta_1 / \rho} e^{\mathbf{x}_1^* \top \Gamma_u \mathbf{x} / \rho} \right] + \tilde{O}(N_1^{-1/2}).$$

Moreover, noting that $\sqrt{r}\Gamma^* \mathbf{x}_i^* = \tilde{O}(\sqrt{r})$ (see Lemma 28), from Lemma 30, we have

$$\begin{aligned} & \frac{1}{N_1} \sum_{i=1}^{N_1} \sqrt{r}\Gamma^* \frac{y_i}{\rho} e^{y_i/\rho} e^{\mathbf{x}_i^{*\top} \Gamma_u \mathbf{x}/\rho} \mathbf{x}_i^* \\ &= \sqrt{r}\Gamma^* \mathbb{E}_{\mathbf{x}_1, \zeta_1} \left[\frac{\bar{\sigma}_*(\langle \mathbf{x}_1, \beta \rangle) + \zeta_1}{\rho} e^{\bar{\sigma}_*(\langle \mathbf{x}_1, \beta \rangle)/\rho + \zeta_1/\rho} e^{\mathbf{x}_1^{*\top} \Gamma_u \mathbf{x}/\rho} \mathbf{x}_1^* \right] + \tilde{O}(r^{1/2} N_1^{-1/2}), \end{aligned}$$

and

$$\frac{1}{N_1} \sum_{i=1}^{N_1} \sqrt{r}\Gamma^* e^{y_i/\rho} e^{\mathbf{x}_i^{*\top} \Gamma_u \mathbf{x}/\rho} \mathbf{x}_i^* = \sqrt{r}\Gamma^* \mathbb{E}_{\mathbf{x}_1, \zeta_1} \left[e^{\bar{\sigma}_*(\langle \mathbf{x}_1, \beta \rangle)/\rho + \zeta_1/\rho} e^{\mathbf{x}_1^{*\top} \Gamma_u \mathbf{x}/\rho} \mathbf{x}_1^* \right] + \tilde{O}(r^{1/2} N_1^{-1/2}).$$

Hence, we obtain that

$$\sqrt{r}\Gamma^* \pi_2^{-2} (F_1 \pi_2 - F_2 \pi_1) = \frac{\text{poly}(z)\beta + \delta_2}{\text{poly}(z) + \delta_1},$$

for some δ_1 and δ_2 which satisfy $\delta_1 = \tilde{O}(1/\sqrt{N_1})$ and $\delta_2 = \tilde{O}\left(\sqrt{\frac{r}{N_1}}\right)$ with high probability. Let $\mathbf{n}_2 = \sqrt{r}\Gamma^* \{\pi_2^{-2} (F_1 \pi_2 - F_2 \pi_1) - \gamma(\mathbf{x}, y)\beta\}$, then $\mathbf{n}_2 = \tilde{O}\left(\sqrt{\frac{r}{N_1}}\right)$ holds.

Also, noting that $\langle \mathbf{x}_i^*, \mathbf{u} \rangle = C_u \tilde{O}_d(1)$ with high probability, it holds that

$$\begin{aligned} & \frac{1}{N_1} \sum_{i=1}^{N_1} \frac{y_i}{\rho} e^{y_i/\rho} e^{\mathbf{x}_i^{*\top} \Gamma_u \mathbf{x}/\rho} \langle \mathbf{x}_i^*, \mathbf{u} \rangle \\ &= \mathbb{E}_{\mathbf{x}_1, \zeta_1} \left[\frac{\bar{\sigma}_*(\langle \mathbf{x}_1, \beta \rangle) + \zeta_1}{\rho} e^{\bar{\sigma}_*(\langle \mathbf{x}_1, \beta \rangle)/\rho + \zeta_1/\rho} e^{\mathbf{x}_1^{*\top} \Gamma_u \mathbf{x}/\rho} \langle \mathbf{x}_1^*, \mathbf{u} \rangle \right] + \tilde{O}(C_u N_1^{-1/2}), \end{aligned}$$

and

$$\frac{1}{N_1} \sum_{i=1}^{N_1} e^{y_i/\rho} e^{\mathbf{x}_i^{*\top} \Gamma_u \mathbf{x}/\rho} \langle \mathbf{x}_i^*, \mathbf{u} \rangle = \mathbb{E}_{\mathbf{x}_1, \zeta_1} \left[e^{\bar{\sigma}_*(\langle \mathbf{x}_1, \beta \rangle)/\rho + \zeta_1/\rho} e^{\mathbf{x}_1^{*\top} \Gamma_u \mathbf{x}/\rho} \langle \mathbf{x}_1^*, \mathbf{u} \rangle \right] + \tilde{O}(C_u N_1^{-1/2}).$$

Therefore, letting $n_1 = \pi_2^{-2} \{\langle F_1, \mathbf{u} \rangle \pi_2 - \langle F_2, \mathbf{u} \rangle \pi_1\} - \gamma(\mathbf{x}, y) \langle \beta^*, \mathbf{u} \rangle$, we obtain that $n_1 = \tilde{O}(C_u \sqrt{\frac{1}{N_1}})$.

In summary, we obtain that

$$\sqrt{r}\Gamma^* \frac{\xi_1 \pi_2 - \xi_2 \pi_1}{\pi_2^2} = \sqrt{r}\Gamma^* [\mathbf{x} \cdot (\gamma(\mathbf{x}, y) \langle \beta^*, \mathbf{u} \rangle) + \gamma(\mathbf{x}, y) \langle \mathbf{x}, \mathbf{u} \rangle \beta^*] + \sqrt{r}\Gamma^* \mathbf{x} \cdot n_1 + \langle \mathbf{x}, \mathbf{u} \rangle n_2,$$

where $n_1 = \tilde{O}\left(C_u \sqrt{\frac{1}{N_1}}\right)$ and $n_2 = \tilde{O}\left(\sqrt{\frac{r}{N_1}}\right)$. \square

Lemma 10. *The norm of the vector $\|\sqrt{r}\Gamma^* \frac{\partial f_{\text{IC}}}{\partial \mathbf{u}}\|$ has sub-Weibull tail with tail index $1/(P+2)$.*

Proof. It suffices to show that $\left\| \sqrt{r}\Gamma^* \frac{\xi_1 \pi_2 - \xi_2 \pi_1}{\pi_2^2} \right\|$ has sub-Weibull tail. Note that $\frac{\xi_1 \pi_2 - \xi_2 \pi_1}{\pi_2^2} = \frac{\xi_1}{\pi_2} - \frac{\xi_2 \pi_1}{\pi_2^2}$. Therefore we examine $\frac{\xi_1}{\pi_2}$, $\frac{\xi_2}{\pi_2}$ and $\frac{\pi_1}{\pi_2}$. Let us define

$$p_i = \frac{e^{\mathbf{x}_i^{*\top} \Gamma_u \mathbf{x}/\rho}}{\sum_{i=1}^{N_1} e^{y_i/\rho} e^{\mathbf{x}_i^{*\top} \Gamma_u \mathbf{x}/\rho}},$$

then the following holds:

$$\frac{\xi_1}{\pi_2} = \sum_{i=1}^{N_1} p_i \frac{y_i}{\rho} \{\langle \mathbf{x}_i^*, \mathbf{u} \rangle \mathbf{x} + \langle \mathbf{u}, \mathbf{x} \rangle \mathbf{x}_i^*\}$$

$$\frac{\xi_2}{\pi_2} = \sum_{i=1}^{N_1} p_i \{\langle \mathbf{x}_i^*, \mathbf{u} \rangle \mathbf{x} + \langle \mathbf{u}, \mathbf{x} \rangle \mathbf{x}_i^*\}$$

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

$$\frac{\pi_1}{\pi_2} = \sum_{i=1}^{N_1} p_i \frac{y_i}{\rho}$$

By the triangle inequality, we have

$$\left\| \frac{\xi_1}{\pi_2} \right\| \leq \sum_{i=1}^{N_1} p_i \frac{y_i}{\rho} \|\langle \mathbf{x}_i^*, \mathbf{u} \rangle \mathbf{x} + \langle \mathbf{u}, \mathbf{x} \rangle \mathbf{x}_i^*\|.$$

By further applying the Cauchy-Schwarz inequality, we have

$$\left\| \frac{\xi_1}{\pi_2} \right\| \leq \left(\sum_{i=1}^{N_1} p_i^2 \right) \left(\sum_{i=1}^{N_1} \frac{y_i^2}{\rho^2} \|\langle \mathbf{x}_i^*, \mathbf{u} \rangle \mathbf{x} + \langle \mathbf{u}, \mathbf{x} \rangle \mathbf{x}_i^*\|^2 \right).$$

Note that $\sum_{i=1}^{N_1} p_i^2 \leq 1$. Next, we can see that

$$y_i^2 \|\langle \mathbf{x}_i^*, \mathbf{u} \rangle \mathbf{x} + \langle \mathbf{u}, \mathbf{x} \rangle \mathbf{x}_i^*\|^2 \leq 2y_i^2 \|\langle \mathbf{x}_i^*, \mathbf{u} \rangle \mathbf{x}\|^2 + 2y_i^2 \|\langle \mathbf{u}, \mathbf{x} \rangle \mathbf{x}_i^*\|^2$$

Since $y_i = \sigma_*(\langle \beta, \mathbf{x} \rangle) + \zeta_i$, the tail probability of each term is almost equal to that of $\text{poly}(t)$ where $t \sim \mathcal{N}(0, 1)$, which means $y_i^2 \|\langle \mathbf{x}_i^*, \mathbf{u} \rangle \mathbf{x}\|^2$ and $y_i^2 \|\langle \mathbf{u}, \mathbf{x} \rangle \mathbf{x}_i^*\|^2$ have sub-Weibull tail with tail index $2/(2P + 4)$. This means that $\left\| \frac{\xi_1}{\pi_2} \right\|$ has sub-Weibull tail with tail index $1/(P + 2)$. By applying the same argument, you can see that $\frac{\xi_2}{\pi_2} \frac{\pi_1}{\pi_2}$ has sub-Weibull tail. Therefore, considering that $\|\sqrt{r}\Gamma^*\|_2$ is a finite constant, $\|\sqrt{r}\Gamma^* \frac{\partial f_{\text{IC}}}{\partial \mathbf{u}}\|$ has sub-Weibull tail with tail index $1/(P + 2)$. \square

Lemma 6 (Formal). *It holds that*

$$\begin{aligned} & \frac{1}{2} \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} (f_{\text{IC}}(\mathbf{x}) - y)^2 \\ &= \Theta(\alpha m) \langle \beta, \mathbf{u} \rangle \frac{\beta}{\kappa^{\text{ie}(\sigma_*)+1} \rho^{\text{ie}(\sigma_*)}} \{ (\kappa \sqrt{r})^{-(\text{ie}(\sigma_*)-1)} (1 + O(1/\sqrt{d})) + \langle \beta, \mathbf{u} \rangle^{2\text{ie}(\sigma_*)-2} (1 + O(1/\sqrt{d})) \} \\ & \quad + \tilde{O}(\alpha^2 m^2 C_u \sqrt{r}) + \tilde{O} \left(\alpha m C_u \sqrt{\frac{r}{N_1}} \right) + \nu, \end{aligned}$$

with high probability, where ν satisfies $\nu = \tilde{O}(\alpha m C_u \sqrt{r})$ with high probability and $\mathbb{E}_{\mathbf{x}}[\nu] = 0$. Moreover, $\|\nu\|$ has sub-Weibull tail.

Proof. Note that

$$\frac{1}{2} \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} (f_{\text{IC}}(\mathbf{x}) - y)^2 = f_{\text{IC}}(\mathbf{x}) \cdot \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} f_{\text{IC}}(\mathbf{x}) - y \cdot \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} f_{\text{IC}}(\mathbf{x}). \quad (3)$$

First, we analyze the first term of the RHS. Since $y_i = \tilde{O}_p(1)$, we have that

$$|f_{\text{IC}}(\mathbf{x})| \leq \alpha m \left| \frac{\sum_{i=1}^{N_1} y_i e^{y_i/\rho} e^{\mathbf{x}_i^* \Gamma_{\mathbf{u}} \mathbf{x} / \rho}}{\sum_{i=1}^{N_1} e^{y_i/\rho} e^{\mathbf{x}_i^* \Gamma_{\mathbf{u}} \mathbf{x} / \rho}} \right| = \tilde{O}_p(\alpha m).$$

Moreover, from Lemma 9, it holds that

$$\sqrt{r} \Gamma^* \frac{\partial f_{\text{IC}}}{\partial \mathbf{u}} = \alpha L_m \{ \sqrt{r} \Gamma^* \mathbf{x} \cdot \langle \gamma(\mathbf{x}, y) \rangle \langle \beta^*, \mathbf{u} \rangle + \sqrt{r} \Gamma^* \gamma(\mathbf{x}, y) \langle \mathbf{x}, \mathbf{u} \rangle \beta^* + \sqrt{r} \Gamma^* \mathbf{x} \cdot \mathbf{n}_1 + \langle \mathbf{x}, \mathbf{u} \rangle \mathbf{n}_2 \},$$

where $\sqrt{r} \Gamma^* \mathbf{x} = \tilde{O}(\sqrt{r})$ (See Lemma 28) and $\gamma(\mathbf{x}, y) = \tilde{O}(1)$ with high probability, which means that $\sqrt{r} \Gamma^* \frac{\partial f_{\text{IC}}}{\partial \mathbf{u}} = \tilde{O}_p(\alpha m C_u \sqrt{r})$. Therefore, we obtain

$$f_{\text{IC}}(\mathbf{x}) \cdot \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} f_{\text{IC}}(\mathbf{x}) = \tilde{O}_p(\alpha^2 m^2 C_u \sqrt{r}).$$

Next, we evaluate the second term of Eq. (3). For that purpose, let

$$\nu = y \cdot \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} f_{\text{IC}}(\mathbf{x}) - \mathbb{E}_{\mathbf{x}}[y \cdot \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} f_{\text{IC}}(\mathbf{x})].$$

By noticing that $y = \tilde{O}_p(1)$ and $\sqrt{r} \Gamma^* \nabla_{\mathbf{u}} f_{\text{IC}}(\mathbf{x}) = \tilde{O}(\alpha m C_u \sqrt{r})$ with high probability, we have

$$\nu = \tilde{O}_p(\alpha m C_u \sqrt{r}).$$

Moreover, $\mathbb{E}_{\mathbf{x}}[\nu] = 0$ by definition, and $\|\nu\|$ has sub-Weibull tail from Lemma 10. Thus, what remains is to calculate $\mathbb{E}_{\mathbf{x}}[y \cdot \sqrt{r}\Gamma^* \nabla_{\mathbf{u}} f_{\text{IC}}(\mathbf{x})]$:

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}}[y\alpha L_m \{ \sqrt{r}\Gamma^* \mathbf{x} \cdot \gamma(\mathbf{x}, y) \langle \beta^*, \mathbf{u} \rangle + \sqrt{r}\Gamma^* \gamma(\mathbf{x}, y) \langle \mathbf{x}, \mathbf{u} \rangle \beta^* \}] \\ & = \alpha L_m \langle \beta^*, \mathbf{u} \rangle \mathbb{E}_{\mathbf{x}}[\sqrt{r}\Gamma^* \mathbf{x} \sigma_*(\langle \beta, \mathbf{x} \rangle)(P_0 + P_1 z + \dots)] \\ & \quad + \alpha L_m \sqrt{r}\Gamma^* \beta^* \mathbb{E}_{\mathbf{x}}[\langle \mathbf{u}, \mathbf{x} \cdot \sigma_*(\langle \beta, \mathbf{x} \rangle)(P_0 + P_1 z + \dots) \rangle]. \end{aligned} \quad (4)$$

Now, we calculate $\mathbb{E}_{\mathbf{x}}[\sqrt{r}\Gamma^* \mathbf{x} \sigma_*(\langle \beta, \mathbf{x} \rangle) z^k]$ in order to investigate $\mathbb{E}_{\mathbf{x}}[\sqrt{r}\Gamma^* \mathbf{x} \sigma_*(\langle \beta, \mathbf{x} \rangle)(P_0 + P_1 z + \dots)]$:

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}}[\sqrt{r}\Gamma^* \mathbf{x} \sigma_*(\langle \beta, \mathbf{x} \rangle) z^k] \\ & = \sqrt{r}\Gamma^* \{ \mathbb{E}_{\mathbf{x}}[\sigma'_*(\langle \beta, \mathbf{x} \rangle) z^k] \beta + \mathbb{E}_{\mathbf{x}}[\sigma_*(\langle \beta, \mathbf{x} \rangle) k z^{k-1}] \Gamma_u^\top \beta \} \end{aligned} \quad (5)$$

holds from Stein's lemma, and when $k < \text{ie}(\sigma_*) - 1$, $\mathbb{E}_{\mathbf{x}}[\sigma'_*(\langle \beta, \mathbf{x} \rangle) z^k] = 0$ from the definition of the information exponent. When $k = \text{ie}(\sigma_*) - 1$, we have that

$$\sqrt{r}\Gamma^* \mathbb{E}_{\mathbf{x}}[\sigma'_*(\langle \beta, \mathbf{x} \rangle) z^k] \simeq \sqrt{r}\Gamma^* \beta \left(\frac{\langle \beta, \sqrt{r}\Gamma_u^\top \Gamma^* \beta \rangle}{\rho} \right)^{\text{ie}(\sigma_*) - 1},$$

and from Lemma 31, it holds that

$$\beta^* = \sqrt{r}\Gamma^* \beta = \kappa^{-1}(\beta + O(1/\sqrt{d})),$$

and

$$\begin{aligned} \langle \beta, \sqrt{r}\Gamma_u^\top \Gamma^* \beta \rangle & = \langle \Gamma_u \beta, \sqrt{r}\Gamma^* \beta \rangle \\ & = \langle \Gamma^* \beta, \sqrt{r}\Gamma^* \beta \rangle + \langle \beta, \mathbf{u} \rangle \langle \mathbf{u}, \sqrt{r}\Gamma^* \beta \rangle \\ & = (\sqrt{r}\kappa^2)^{-1}(1 + O(1/\sqrt{d})) + \kappa^{-1} \langle \beta, \mathbf{u} \rangle \{ \langle \beta, \mathbf{u} \rangle + O(1/\sqrt{d}) \} \\ & = (\sqrt{r}\kappa^2)^{-1}(1 + O(1/\sqrt{d})) + \kappa^{-1} \langle \beta, \mathbf{u} \rangle^2 \{ 1 + O(1/\sqrt{d}) \} \end{aligned}$$

with high probability. By ignoring small terms, this yields that

$$\begin{aligned} & \sqrt{r}\Gamma^* \beta \left(\frac{\langle \beta, \sqrt{r}\Gamma_u^\top \Gamma^* \beta \rangle}{\rho} \right)^{\text{ie}(\sigma_*) - 1} \\ & = \frac{P_{\text{ie}(\sigma_*) - 1} \beta}{\kappa^{\text{ie}(\sigma_*)} \rho^{\text{ie}(\sigma_*) - 1}} \{ (\kappa\sqrt{r})^{-(\text{ie}(\sigma_*) - 1)} (1 + O(1/\sqrt{d})) + \langle \beta, \mathbf{u} \rangle^{2\text{ie}(\sigma_*) - 2} (1 + O(1/\sqrt{d})) \}. \end{aligned}$$

We can see that the other term in Eq. (5) is negligible. Using $P_{\text{ie}(\sigma_*) - 1} = \Theta(1/\rho)$, plugging this result into Eq. (4) yields

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}}[y\alpha L_m \{ \sqrt{r}\Gamma^* \mathbf{x} \cdot \gamma(\mathbf{x}, y) \langle \beta^*, \mathbf{u} \rangle + \sqrt{r}\Gamma^* \gamma(\mathbf{x}, y) \langle \mathbf{x}, \mathbf{u} \rangle \beta^* \}] \\ & = \alpha L_m \langle \beta, \mathbf{u} \rangle \frac{\beta}{\kappa^{\text{ie}(\sigma_*) + 1} \rho^{\text{ie}(\sigma_*)}} \{ (\kappa\sqrt{r})^{-(\text{ie}(\sigma_*) - 1)} (1 + O(1/\sqrt{d})) + \langle \beta, \mathbf{u} \rangle^{2\text{ie}(\sigma_*) - 2} (1 + O(1/\sqrt{d})) \} \end{aligned}$$

with high probability. Finally, we evaluate the term $\mathbb{E}_{\mathbf{x}}[\alpha L_m \sigma_*(\langle \beta, \mathbf{x} \rangle)(\sqrt{r}\Gamma^* \mathbf{x} \cdot \mathbf{n}_1 + \langle \mathbf{x}, \mathbf{u} \rangle \mathbf{n}_2)]$. From Corollary 17 of Oko et al. (2024), we have $\sigma_*(\langle \beta, \mathbf{x} \rangle) = O_p(\log^{\text{deg}(\sigma_*)/2} d)$. Moreover, from Lemma 9, it holds that $\mathbf{n}_1 = \tilde{O}(C_u \sqrt{\frac{1}{N_1}})$ and $\mathbf{n}_2 = \tilde{O}(\sqrt{\frac{r}{N_1}})$ with high probability. Therefore, we arrive at

$$\mathbb{E}_{\mathbf{x}}[\alpha L_m \sigma_*(\langle \beta, \mathbf{x} \rangle)(\sqrt{r}\Gamma^* \mathbf{x} \cdot \mathbf{n}_1 + \langle \mathbf{x}, \mathbf{u} \rangle \mathbf{n}_2)] = \tilde{O}_p \left(\alpha m C_u \sqrt{\frac{r}{N_1}} \right).$$

Combining all results above completes the proof. \square

The following lemma will be useful in the analysis of weak recovery (the next section).

Lemma 11. *When $\text{ge}(\sigma_*) = 2$, $\gamma(\mathbf{x}, y) = \tilde{O}(1/\sqrt{r})$ holds with high probability.*

1134 *Proof.* Recall that

$$\begin{aligned}
1135 \quad \gamma(\mathbf{x}, y) &= P_0 + P_1 z + \dots \\
1136 \quad &= \left(k_\rho \sum_{i \geq 0} \frac{B_i}{i!} z^i \right)^{-2} \left\{ \left[k_\rho \left(\sum_{i \geq 0} \frac{A_{i+1}}{i!} z^i \right) + k'_\rho \left(\sum_{i \geq 0} \frac{B_{i+1}}{i!} z^i \right) \right] \cdot k_\rho \left(\sum_{i \geq 0} \frac{B_i}{i!} z^i \right) \right. \\
1137 \quad & \left. - k_\rho \left(\sum_{i \geq 0} \frac{B_{i+1}}{i!} z^i \right) \cdot \left[k_\rho \left(\sum_{i \geq 0} \frac{A_i}{i!} z^i \right) + k'_\rho \left(\sum_{i \geq 0} \frac{B_i}{i!} z^i \right) \right] \right\}.
\end{aligned}$$

1140 By expanding the RHS in powers of z and comparing their coefficients, we obtain

$$\begin{aligned}
1141 \quad P_0 &= \frac{(k_\rho A_0 + k'_\rho B_0) \cdot k_\rho B_1 - k_\rho B_0 (k_\rho A_1 + k'_\rho B_1)}{(k_\rho B_0)^2} \\
1142 \quad &= \frac{A_0 B_1 - B_0 A_1}{B_0^2}.
\end{aligned}$$

1143 When $\text{ge}(\sigma_*)$ is 2, σ_* is even, and thus $\exp(\bar{\sigma}_*)$ and $\bar{\sigma}_* \exp(\bar{\sigma}_*)$ are also even. This means that $\exp(\bar{\sigma}_*)$ and $\bar{\sigma}_* \exp(\bar{\sigma}_*)$ can be expanded in polynomial of z^2 . This means that the coefficients of z in the hermite expansion of $\exp(\bar{\sigma}_*)$ and $\bar{\sigma}_* \exp(\bar{\sigma}_*)$ are 0, namely $A_1 = B_1 = 0$, which yields $P_0 = 0$. Hence, we have

$$1144 \quad \gamma(\mathbf{x}, y) = P_1 z + P_2 z^2 + \dots,$$

1145 where $P_1 = \tilde{O}(1)$ and $z = \tilde{O}_p(1/\sqrt{r})$, which yields the assertion. \square

1156 D ONE STEP GRADIENT DESCENT FOR WEAK RECOVERY

1157 **Lemma 12.** Set $C_u = 1/\sqrt{r}$, $N_1 = \tilde{\Omega}(r^{\text{ge}(\sigma_*)+1})$ and $N_{new} = \tilde{\Omega}(r^{\text{ge}(\sigma_*)+2})$. Let $w_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d)$ for $i = 1, 2, \dots, N_{new}$ and $\mathbf{h} = \frac{1}{2} \sqrt{r} \Gamma^* \frac{1}{N_{new}} \sum_{i=1}^{N_{new}} \nabla_{\mathbf{u}} (f_{\text{IC}}(\mathbf{w}_i) - (g(\Gamma^*, \mathbf{w}_i) - b))^2$, then

$$\begin{aligned}
1161 \quad \mathbf{h} &= P'_2 \Theta(\alpha_1 m) \langle \beta, \mathbf{u} \rangle \frac{\beta}{\kappa^{2\text{ge}(\sigma_*)} \rho^{\text{ge}(\sigma_*)}} (\sqrt{r})^{-(2\text{ge}(\sigma_*)-1)} (1 + O(1/\sqrt{d})) \\
1162 \quad &+ \tilde{O}(\alpha_1^2 m^2 r^{-1/2} r^{-\frac{\text{ge}(\sigma_*)-1}{2}}) + \tilde{O} \left(\alpha_1 m \sqrt{\frac{r^{-\text{ge}(\sigma_*)}}{N_1}} \right) + \tilde{O} \left(\alpha_1 m \sqrt{\frac{r^{-\text{ge}(\sigma_*)}}{N_{new}}} \right) + \tilde{O}(\alpha_1 m r^{-\text{ge}(\sigma_*)-1/2})
\end{aligned}$$

1163 holds with high probability.

1164 *Proof.* From Lemma 8, we have that

$$1165 \quad g(\Gamma^*, \mathbf{x}) = P'_1 + P'_2 \left(\frac{\langle \mathbf{x}, \beta \rangle}{\sqrt{r}} \right)^{\text{ge}(\sigma_*)} + n_3,$$

1166 where $P'_1 = o_d(1)$, $P'_2 = \Theta_d((\log d)^{-C_{P_2}})$ and $n_3 = o_d(P'_2 r^{-\text{ge}(\sigma_*)/2-1/2} \log^{-2 \text{deg}(\sigma_*)+2} d)$.

1167 Noticing that $g(\Gamma^*, \mathbf{x}) = O(1)$ with high probability, then, letting $b = \frac{1}{N_{new}} \sum_{i=1}^{N_{new}} g(\Gamma^*, \mathbf{w}_i)$ and applying Lemma 26, it holds that

$$1168 \quad g(\Gamma^*, \mathbf{w}_i) - b = P'_2 r^{-\frac{\text{ge}(\sigma_*)}{2}} \text{He}_{\text{ge}(\sigma_*)}(\langle \beta, \mathbf{w}_i \rangle) + n_3 + \tilde{O}(1/\sqrt{N_{new}}).$$

1169 When $N_{new} = O(r^{\text{ge}(\sigma_*)+2})$, the term $\tilde{O}(1/\sqrt{N_{new}})$ can be dominated by n_3 , which means that

$$1170 \quad g(\Gamma^*, \mathbf{w}_i) - b = P'_2 r^{-\frac{\text{ge}(\sigma_*)}{2}} \text{He}_{\text{ge}(\sigma_*)}(\langle \beta, \mathbf{w}_i \rangle) + n_3,$$

1171 and $n_3 = o_d(P'_2 r^{-\text{ge}(\sigma_*)/2-1/2} \log^{-2 \text{deg}(\sigma_*)+2} d)$ with high probability.

1172 In this stage, we use $g(\Gamma^*, \mathbf{w}_i) - b$ as a teacher signal. Thus, as $C_u = 1/\sqrt{r}$ and $g(\Gamma^*, \mathbf{w}_i) - b = \tilde{O}(r^{-\frac{\text{ge}(\sigma_*)}{2}})$, with high probability, we have

$$\begin{aligned}
1173 \quad &\frac{1}{2} \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} (f_{\text{IC}}(\mathbf{w}_i) - (g(\Gamma^*, \mathbf{w}_i) - b))^2 \\
1174 \quad &= \tilde{O}(\alpha_1^2 m^2 r^{-1/2} r^{-\frac{\text{ge}(\sigma_*)-1}{2}}) - \left\{ P'_2 r^{-\frac{\text{ge}(\sigma_*)}{2}} \text{He}_{\text{ge}(\sigma_*)}(\langle \beta, \mathbf{w}_i \rangle) + n_3 \right\} \cdot \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} f_{\text{IC}}(\mathbf{w}_i).
\end{aligned}$$

Following the same argument as Lemma 6 yields

$$\begin{aligned}
& P_2' r^{-\frac{\text{ge}(\sigma_*)}{2}} \text{He}_{\text{ge}(\sigma_*)}(\langle \beta, \mathbf{w}_i \rangle) \cdot \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} f_{\text{IC}}(\mathbf{w}_i) \\
&= P_2' \Theta(\alpha_1 m) \langle \beta, \mathbf{u} \rangle \frac{r^{-\frac{\text{ge}(\sigma_*)}{2}} \beta}{\kappa^{\text{ge}(\sigma_*)+1} \rho^{\text{ge}(\sigma_*)}} \{(\kappa \sqrt{r})^{-(\text{ge}(\sigma_*)-1)} (1 + O(1/\sqrt{d})) + \langle \beta, \mathbf{u} \rangle^{2\text{ge}(\sigma_*)-2} (1 + O(1/\sqrt{d}))\} \\
&+ \nu + \tilde{O} \left(\alpha_1 m \sqrt{\frac{r^{-\text{ge}(\sigma_*)}}{N_1}} \right)
\end{aligned}$$

with high probability. Since $\langle \beta, \mathbf{u} \rangle = \tilde{O}_p(1/r)$ at the initialization, the first term (regarding $(\kappa \sqrt{r})^{-(\text{ge}(\sigma_*)-1)}$) dominates the second term (regarding $\langle \beta, \mathbf{u} \rangle^{2\text{ge}(\sigma_*)-2}$). Taking the average over $i = 1, \dots, N_{\text{new}}$ yields that

$$\begin{aligned}
& P_2' r^{-\frac{\text{ge}(\sigma_*)}{2}} \frac{1}{N_{\text{new}}} \sum_{i=1}^{N_{\text{new}}} \text{He}_{\text{ge}(\sigma_*)}(\langle \beta, \mathbf{w}_i \rangle) \cdot \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} f_{\text{IC}}(\mathbf{w}_i) \\
&= P_2' \Theta(\alpha_1 m) \langle \beta, \mathbf{u} \rangle \frac{\beta}{\kappa^{2\text{ge}(\sigma_*)} \rho^{\text{ge}(\sigma_*)}} (\sqrt{r})^{-(2\text{ge}(\sigma_*)-1)} (1 + O(1/\sqrt{d})) \\
&+ \frac{1}{N_{\text{new}}} \sum_{i=1}^{N_{\text{new}}} \nu_i + \tilde{O} \left(\alpha_1 m \sqrt{\frac{r^{-\text{ge}(\sigma_*)}}{N_1}} \right),
\end{aligned}$$

where ν_i is a series of i.i.d. mean-zero random variable vectors which satisfy $\nu_i = \tilde{O}_p(\alpha_1 m \cdot r^{-\text{ge}(\sigma_*)/2})$. Then, from Hoeffding's inequality, we have

$$\frac{1}{N_{\text{new}}} \sum_{i=1}^{N_{\text{new}}} \nu_i = \tilde{O} \left(\alpha_1 m \sqrt{\frac{r^{-\text{ge}(\sigma_*)}}{N_{\text{new}}}} \right),$$

with high probability. Next we investigate the effect of n_3 . From Lemma 9, $\sqrt{r} \Gamma^* \nabla_{\mathbf{u}} f_{\text{IC}}(\mathbf{w}_i) = \tilde{O}(\alpha_1 m C_u) = \tilde{O}(\alpha_1 m r^{-1/2})$ holds. This leads to $n_3 \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} f_{\text{IC}}(\mathbf{w}_i) = \tilde{O}(\alpha_1 m r^{-\text{ge}(\sigma_*)/2-1})$. Moreover, from Lemma 11, when $\text{ge}(\sigma_*) = 2$ we have $\gamma(\mathbf{x}, \mathbf{y}) = O(1/\sqrt{r})$. Combining this fact with $N_1 = \tilde{\Omega}(r^{\text{ge}(\sigma_*)+1})$ yields $\sqrt{r} \Gamma^* \nabla_{\mathbf{u}} f_{\text{IC}}(\mathbf{w}_i) = \tilde{O}(\alpha_1 m/r)$. Therefore, whether $\text{ge}(\sigma_*) = 1$ or $\text{ge}(\sigma_*) = 2$, we obtain

$$n_3 \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} f_{\text{IC}}(\mathbf{w}_i) = \tilde{O}(\alpha_1 m r^{-\text{ge}(\sigma_*)-1/2}).$$

Taking the average over $i = 1, \dots, N_{\text{new}}$ completes the proof. \square

Lemma 13. *Under the condition of Lemma 12, it holds that*

$$\begin{aligned}
\langle \beta, \mathbf{h} \rangle &= P_2' \Theta(\alpha_1 m) \langle \beta, \mathbf{u} \rangle \frac{1}{\kappa^{2\text{ge}(\sigma_*)} \rho^{\text{ge}(\sigma_*)}} (\sqrt{r})^{-(2\text{ge}(\sigma_*)-1)} (1 + O(1/\sqrt{d})) \\
&+ \tilde{O} \left(\alpha_1^2 m^2 r^{-\frac{\text{ge}(\sigma_*)}{2}} \right) + \tilde{O} \left(\alpha_1 m \sqrt{\frac{r^{-\text{ge}(\sigma_*)}}{N_1}} \right) + \tilde{O} \left(\alpha_1 m \sqrt{\frac{r^{-\text{ge}(\sigma_*)}}{N_{\text{new}}}} \right) + o(\alpha_1 m r^{-\text{ge}(\sigma_*)-1/2}),
\end{aligned}$$

with high probability.

Proof. As $\|\beta\| = 1$, using the result of Lemma 12 and taking the inner product of \mathbf{h} and β yields

$$\begin{aligned}
\langle \beta, \mathbf{h} \rangle &= P_2' \Theta(\alpha_1 m) \langle \beta, \mathbf{u} \rangle \frac{1}{\kappa^{2\text{ge}(\sigma_*)} \rho^{\text{ge}(\sigma_*)}} (\sqrt{r})^{-(2\text{ge}(\sigma_*)-1)} (1 + O(1/\sqrt{d})) \\
&+ \tilde{O}(\alpha_1^2 m^2 r^{-\frac{\text{ge}(\sigma_*)}{2}}) + \tilde{O} \left(\alpha_1 m \sqrt{\frac{r^{-\text{ge}(\sigma_*)}}{N_1}} \right) + \tilde{O} \left(\alpha_1 m \sqrt{\frac{r^{-\text{ge}(\sigma_*)}}{N_{\text{new}}}} \right) + \tilde{O}(\alpha_1 m r^{-\text{ge}(\sigma_*)-1/2}).
\end{aligned}$$

However, by carefully investigating the last term ($\tilde{O}(\alpha_1 m r^{-\text{ge}(\sigma_*)-1/2})$) in the RHS, we can obtain a tighter bound of that term. Indeed, this term arises from

$$\begin{aligned} & n_3 \left\langle \beta, \frac{1}{N_{\text{new}}} \sum_{i=1}^{N_{\text{new}}} \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} f_{\text{IC}}(\mathbf{w}_i) \right\rangle \\ &= n_3 \alpha_1 L_m \frac{1}{N_{\text{new}}} \sum_{i=1}^{N_{\text{new}}} \left\{ \langle \beta, \sqrt{r} \Gamma^* \mathbf{x} \rangle \cdot \gamma(\mathbf{x}, y) \langle \beta^*, \mathbf{u} \rangle + \gamma(\mathbf{x}, y) \langle \mathbf{x}, \mathbf{u} \rangle \langle \beta^*, \sqrt{r} \Gamma^* \beta \rangle \right. \\ & \quad \left. + \langle \beta, \sqrt{r} \Gamma^* \mathbf{x} \rangle \cdot n_1 + \langle \mathbf{x}, \mathbf{u} \rangle \langle \beta, \mathbf{n}_2 \rangle \right\}. \end{aligned}$$

We notice that the leading term of the right hand side is $n_3 \alpha_1 L_m P_0 \langle \mathbf{u}, \mathbf{x} \rangle$ when $\text{ge}(\sigma_*) = 1$, and $n_3 \alpha_1 m P_1 \langle \beta, \sqrt{r} \Gamma_u^\top \Gamma_u \mathbf{x} \rangle / \rho \cdot \langle \mathbf{x}, \mathbf{u} \rangle$ when $\text{ge}(\sigma_*) = 2$. Note that $\mathbb{E}_{\mathbf{x}}[P_0 \langle \mathbf{u}, \mathbf{x} \rangle] = 0$ and $\mathbb{E}_{\mathbf{x}}[P_1 \langle \sqrt{r} \Gamma_u^\top \Gamma_u \mathbf{x} \rangle / \rho \cdot \langle \mathbf{x}, \mathbf{u} \rangle] = P_1 \langle \sqrt{r} \Gamma_u^\top \Gamma_u \beta, \mathbf{u} \rangle / \rho = \tilde{O}(r^{-3/2})$ because of $\sqrt{r} \Gamma_u^\top \Gamma_u \beta = \tilde{O}(1/\sqrt{r})\beta + \tilde{O}(1/r) + \langle \beta, \mathbf{u} \rangle \mathbf{u}$ (see Lemma 32). Moreover, we have $\mathbb{E}_{\mathbf{x}}[P_0^2 \langle \mathbf{u}, \mathbf{x} \rangle^2] = \tilde{O}(1)$ and $\mathbb{E}_{\mathbf{x}}[P_1^2 \langle \sqrt{r} \Gamma_u^\top \Gamma_u \mathbf{x} \rangle^2 / \rho^2 \cdot \langle \mathbf{x}, \mathbf{u} \rangle^2] \leq P_1^2 / \rho^2 \cdot \mathbb{E}[\langle \sqrt{r} \Gamma_u^\top \Gamma_u \beta, \mathbf{x} \rangle^4]^{1/2} \mathbb{E}[\langle \mathbf{x}, \mathbf{u} \rangle^4]^{1/2} = \tilde{O}(1)$. Therefore, Lemma 26 yields that

$$\begin{aligned} & \frac{1}{N_{\text{new}}} \sum_{i=1}^{N_{\text{new}}} P_0 \langle \mathbf{u}, \mathbf{w}_i \rangle = \tilde{O} \left(\sqrt{\frac{1}{r N_{\text{new}}}} \right), \\ & \frac{1}{N_{\text{new}}} \sum_{i=1}^{N_{\text{new}}} P_1 \langle \beta, \sqrt{r} \Gamma_u \Gamma_u^* \mathbf{w}_i \rangle / \rho \langle \mathbf{w}_i, \mathbf{u} \rangle = \tilde{O}(r^{-3/2}) + \tilde{O}(r^{-1} N_{\text{new}}^{-1/2}). \end{aligned}$$

Since $n_3 = o_d(P'_2 r^{-\text{ge}(\sigma_*)/2-1/2} \log^{-2 \deg(\sigma_*)+2} d)$ and $N_{\text{new}} = \tilde{\Omega}(r^{\text{ge}(\sigma_*)+2})$, this implies that

$$n_3 \left\langle \beta, \frac{1}{N_{\text{new}}} \sum_{i=1}^{N_{\text{new}}} \sqrt{r} \Gamma^* \nabla_{\mathbf{u}} f_{\text{IC}}(\mathbf{w}_i) \right\rangle = o(\alpha_1 m r^{-\text{ge}(\sigma_*)-1/2}).$$

This yields the assertion. \square

Setting the parameters directly yields the following lemma:

Lemma 14. *Set $N_1 = \tilde{\Omega}(r^{\text{ge}(\sigma_*)+2})$, $N_{\text{new}} = \tilde{\Omega}(r^{\text{ge}(\sigma_*)+2})$, $\alpha_1 m = \tilde{\Theta}(r^{-\text{ge}(\sigma_*)/2-1})$, $\eta_1 = \tilde{\Theta}(r^{3\text{ge}(\sigma_*)/2+3/2})$ and $\lambda_1 = \eta_1^{-1}$, then we have*

$$\mathbf{u}^1 = \tilde{\Theta}(1)\beta + \tilde{O}(1),$$

and

$$\langle \beta, \mathbf{u}^1 \rangle = \tilde{\Theta}(1) + o(1),$$

with high probability. In particular, we have that

$$\left\langle \beta, \frac{\mathbf{u}^1}{\|\mathbf{u}^1\|} \right\rangle \geq 1/\text{polylog}(d).$$

E STRONG RECOVERY

After achieving weak recovery, we train the vector \mathbf{u} , regarding the in-context examples (\mathbf{x}_i, y_i) as training data. The idea of proof in this section is taken from Lee et al. (2024), but we achieve better sample complexity.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

Lemma 15. Let $\mathbf{h} = \frac{1}{2}\sqrt{r}\Gamma^*\nabla_{\mathbf{u}}(f_{\text{IC}}(\mathbf{x}) - y)^2$. Then, we have that

$$\begin{aligned} \mathbf{h} &= \Theta(\alpha_2 m) \langle \beta, \mathbf{u} \rangle^{2ie(\sigma_*)-1} \frac{\beta}{\kappa^{ie(\sigma_*)+1} \rho^{ie(\sigma_*)}} (1 + O(1/\sqrt{d})) \\ &\quad + \tilde{O}(\alpha_2^2 m^2 \sqrt{r}) + \tilde{O}\left(\alpha_2 m \sqrt{\frac{r}{N_2}}\right) + \nu, \\ \langle \beta, \mathbf{h} \rangle &= \Theta(\alpha_2 m) \langle \beta, \mathbf{u} \rangle^{2ie(\sigma_*)-1} \frac{1}{\kappa^{ie(\sigma_*)+1} \rho^{ie(\sigma_*)}} (1 + O(1/\sqrt{d})) \\ &\quad + \tilde{O}(\alpha_2^2 m^2 \sqrt{r}) + \tilde{O}\left(\alpha_2 m \sqrt{\frac{r}{N_2}}\right) + \nu_6, \\ \langle \mathbf{u}, \mathbf{h} \rangle &= \Theta(\alpha_2 m) \langle \beta, \mathbf{u} \rangle^{2ie(\sigma_*)} \frac{1}{\kappa^{ie(\sigma_*)+1} \rho^{ie(\sigma_*)}} (1 + O(1/\sqrt{d})) \\ &\quad + \tilde{O}(\alpha_2^2 m^2 \sqrt{r}) + \tilde{O}\left(\alpha_2 m \sqrt{\frac{r}{N_2}}\right) + \nu_7, \end{aligned}$$

with high probability, where ν_6 and ν_7 are mean-zero sub-Weibull random variables. Moreover, $\nu_6 = \tilde{O}(\alpha_2 m)$ and $\nu_7 = \tilde{O}(\alpha_2 m)$ hold with high probability.

Proof. Remember that, in the stage of strong recovery, the initialization scale is α_2 , the context length is N_2 and $C_u = 1$. By using the same argument to derive Lemma 6, we have

$$\begin{aligned} \mathbf{h} &= \Theta(\alpha_2 m) \langle \beta, \mathbf{u} \rangle \frac{\beta}{\kappa^{ie(\sigma_*)+1} \rho^{ie(\sigma_*)}} \{(\kappa\sqrt{r})^{-(ie(\sigma_*)-1)} (1 + O(1/\sqrt{d})) + \langle \beta, \mathbf{u} \rangle^{2ie(\sigma_*)-2} (1 + O(1/\sqrt{d}))\} \\ &\quad + \tilde{O}(\alpha_2^2 m^2 \sqrt{r}) + \tilde{O}\left(\alpha_2 m \sqrt{\frac{r}{N_2}}\right) + \nu. \end{aligned}$$

Now, since $\langle \beta, \mathbf{u} \rangle \geq 1/\text{polylog}(d)$, the term $\langle \beta, \mathbf{u} \rangle^{2ie(\sigma_*)-2}$ dominates over the term $(\kappa\sqrt{r})^{-(ie(\sigma_*)-1)}$. Therefore, with high probability, we have

$$\begin{aligned} \mathbf{h} &= \Theta(\alpha_2 m) \langle \beta, \mathbf{u} \rangle^{2ie(\sigma_*)-1} \frac{\beta}{\kappa^{ie(\sigma_*)+1} \rho^{ie(\sigma_*)}} (1 + O(1/\sqrt{d})) \\ &\quad + \tilde{O}(\alpha_2^2 m^2 \sqrt{r}) + \tilde{O}\left(\alpha_2 m \sqrt{\frac{r}{N_2}}\right) + \nu. \end{aligned}$$

Next, since $\|\beta\| = 1$, we have

$$\begin{aligned} \langle \beta, \mathbf{h} \rangle &= \Theta(\alpha_2 m) \langle \beta, \mathbf{u} \rangle^{2ie(\sigma_*)-1} \frac{1}{\kappa^{ie(\sigma_*)+1} \rho^{ie(\sigma_*)}} (1 + O(1/\sqrt{d})) \\ &\quad + \tilde{O}(\alpha_2^2 m^2 \sqrt{r}) + \tilde{O}\left(\alpha_2 m \sqrt{\frac{r}{N_2}}\right) + \langle \nu, \beta \rangle, \end{aligned}$$

with high probability, where $\nu = \tilde{O}(\alpha_2 m \sqrt{r})$. Then, since we see that

$$\begin{aligned} \nu &= y \cdot \alpha_2 L_m \{ \sqrt{r}\Gamma^* \mathbf{x} \cdot \gamma(\mathbf{x}, y) \langle \beta^*, \mathbf{u} \rangle + \sqrt{r}\Gamma^* \gamma(\mathbf{x}, y) \langle \mathbf{x}, \mathbf{u} \rangle \beta^* + \sqrt{r}\Gamma^* \mathbf{x} \cdot \mathbf{n}_1 + \langle \mathbf{x}, \mathbf{u} \rangle \mathbf{n}_2 \} \\ &\quad - \mathbb{E}_{\mathbf{x}}[y \cdot \alpha_2 L_m \{ \sqrt{r}\Gamma^* \mathbf{x} \cdot \gamma(\mathbf{x}, y) \langle \beta^*, \mathbf{u} \rangle + \sqrt{r}\Gamma^* \gamma(\mathbf{x}, y) \langle \mathbf{x}, \mathbf{u} \rangle \beta^* + \sqrt{r}\Gamma^* \mathbf{x} \cdot \mathbf{n}_1 + \langle \mathbf{x}, \mathbf{u} \rangle \mathbf{n}_2 \}], \end{aligned}$$

it holds that

$$\begin{aligned} &\langle \beta, \nu \rangle \\ &= y \cdot \alpha_2 L_m \{ \langle \beta, \sqrt{r}\Gamma^* \mathbf{x} \rangle \cdot \gamma(\mathbf{x}, y) \langle \beta^*, \mathbf{u} \rangle + \gamma(\mathbf{x}, y) \langle \mathbf{x}, \mathbf{u} \rangle \langle \beta, \sqrt{r}\Gamma^* \beta^* \rangle + \langle \beta, \sqrt{r}\Gamma^* \mathbf{x} \rangle \cdot \mathbf{n}_1 + \langle \mathbf{x}, \mathbf{u} \rangle \langle \beta, \mathbf{n}_2 \rangle \} \\ &\quad - \mathbb{E}_{\mathbf{x}}[y \cdot \alpha_2 L_m \{ \langle \beta, \sqrt{r}\Gamma^* \mathbf{x} \rangle \cdot \gamma(\mathbf{x}, y) \langle \beta^*, \mathbf{u} \rangle + \gamma(\mathbf{x}, y) \langle \mathbf{x}, \mathbf{u} \rangle \langle \beta, \sqrt{r}\Gamma^* \beta^* \rangle + \langle \beta, \sqrt{r}\Gamma^* \mathbf{x} \rangle \cdot \mathbf{n}_1 + \langle \mathbf{x}, \mathbf{u} \rangle \langle \beta, \mathbf{n}_2 \rangle \}]. \end{aligned}$$

Here, the first term of the right hand side satisfies

$$\begin{aligned} &y \cdot \alpha_2 L_m \{ \langle \beta, \sqrt{r}\Gamma^* \mathbf{x} \rangle \cdot \gamma(\mathbf{x}, y) \langle \beta^*, \mathbf{u} \rangle + \gamma(\mathbf{x}, y) \langle \mathbf{x}, \mathbf{u} \rangle \langle \beta, \sqrt{r}\Gamma^* \beta^* \rangle \\ &\quad + \langle \beta, \sqrt{r}\Gamma^* \mathbf{x} \rangle \cdot \mathbf{n}_1 + \langle \mathbf{x}, \mathbf{u} \rangle \langle \beta, \mathbf{n}_2 \rangle \} \\ &= \tilde{O}(\alpha_2 m), \end{aligned}$$

with high probability, and $\langle \beta, \nu \rangle$ is the difference between this random variable and its expectation. Thus, by defining $\nu_6 = \langle \beta, \nu \rangle$, Hoeffding's inequality yields that $\nu_6 = \tilde{O}(\alpha_2 m)$ with high probability. Moreover, since $|\langle \beta, \nu \rangle| \leq \|\nu\|$, $\nu_6 = \langle \beta, \nu \rangle$ has sub-Weibull tail.

Likewise, we have that

$$\begin{aligned} \langle \mathbf{u}, \mathbf{h} \rangle &= \Theta(\alpha_2 m) \langle \beta, \mathbf{u} \rangle^{2ie(\sigma_*)} \frac{1}{\kappa^{ie(\sigma_*)+1} \rho^{ie(\sigma_*)}} (1 + O(1/\sqrt{d})) \\ &\quad + \tilde{O}(\alpha_2^2 m^2 \sqrt{r}) + \tilde{O}\left(\alpha_2 m \sqrt{\frac{r}{N_2}}\right) + \nu_7, \end{aligned}$$

where ν_7 is a mean-zero sub-Weibull random variable satisfying $\nu_7 = \tilde{O}(\alpha_2 m)$. \square

Lemma 16. *Let $a^{(n)} = \langle \beta, \mathbf{u}^{(n)} \rangle$. Suppose that $c_1 \leq a^{(n)} \leq 1 - \varepsilon$ where $c_1 = \tilde{\Theta}(1)$. Set $\alpha_2 m = \tilde{\Theta}(\varepsilon/r)$, $N_2 = \tilde{\Theta}(r^2)$ and $\eta_2 = \tilde{\Theta}(1/\sqrt{r})$. Then, there exists $c_3 = \tilde{\Theta}(1)$ which satisfies*

$$a^{(n+1)} \geq a^{(n)} + \frac{c_3 \varepsilon}{r\sqrt{r}} a^{(n)2ie(\sigma_*)-1} (1 - a^{(n)2}) (1 - O(1/\sqrt{d})) + \nu_8 - \tilde{O}(\varepsilon/r^2),$$

with high probability, where ν_8 is a mean-zero sub-Weibull random variable which satisfies $\nu_8 = \tilde{O}(\varepsilon/r\sqrt{r})$ with high probability.

Proof. Using the projection matrix $\mathbf{P}_u = \mathbf{I} - uu^\top$, online SGD update can be written as

$$\mathbf{u}^{(n+1)} = \mathbf{u}^{(n)} + \frac{\mathbf{u}^{(n)} + \eta_2 \mathbf{P}_{\mathbf{u}^{(n)}} \mathbf{h}}{\|\mathbf{u}^{(n)} + \eta_2 \mathbf{P}_{\mathbf{u}^{(n)}} \mathbf{h}\|},$$

which gives

$$\begin{aligned} &\langle \beta, \mathbf{u}^{(n+1)} \rangle \\ &= \langle \beta, \mathbf{u}^{(n)} \rangle + \left\langle \beta, \frac{\mathbf{u}^{(n)} + \eta_2 \mathbf{P}_{\mathbf{u}^{(n)}} \mathbf{h}}{\|\mathbf{u}^{(n)} + \eta_2 \mathbf{P}_{\mathbf{u}^{(n)}} \mathbf{h}\|} \right\rangle \\ &= \langle \beta, \mathbf{u}^{(n)} \rangle + \eta_2 \langle \beta, \mathbf{P}_{\mathbf{u}^{(n)}} \mathbf{h} \rangle - \frac{1}{2} \eta_2^2 \|\mathbf{P}_{\mathbf{u}^{(n)}} \mathbf{h}\|^2 + O(\eta_2^3) \\ &= \langle \beta, \mathbf{u}^{(n)} \rangle + \eta_2 \langle \beta, \mathbf{h} \rangle - \eta_2 \langle \beta, \mathbf{u}^{(n)} \rangle \langle \mathbf{u}^{(n)}, \mathbf{h} \rangle - \frac{1}{2} \eta_2^2 \|\mathbf{P}_{\mathbf{u}^{(n)}} \mathbf{h}\|^2 + O((\eta_2 \|\mathbf{h}\|)^3). \end{aligned}$$

By the definition of $a^{(n)}$, i.e., $a^{(n)} = \langle \beta, \mathbf{u}^{(n)} \rangle$, we have

$$a^{(n+1)} = a^{(n)} + \eta_2 \langle \beta, \mathbf{h} \rangle - \eta_2 a^{(n)} \langle \mathbf{u}^{(n)}, \mathbf{h} \rangle - \frac{1}{2} \eta_2^2 \|\mathbf{P}_{\mathbf{u}^{(n)}} \mathbf{h}\|^2 + O((\eta_2 \|\mathbf{h}\|)^3).$$

Now, from Lemma 15, we also have

$$\langle \beta, \mathbf{h} \rangle = \frac{\varepsilon}{r} a^{(n)2ie(\sigma_*)-1} (\tilde{\Theta}(1) + O(1/\sqrt{d})) + \tilde{O}(\varepsilon/r\sqrt{r}) + \tilde{O}(\varepsilon/r\sqrt{r}) + \nu_6$$

and

$$\langle \mathbf{u}, \mathbf{h} \rangle = \frac{\varepsilon}{r} a^{(n)2ie(\sigma_*)} (\tilde{\Theta}(1) + O(1/\sqrt{d})) + \tilde{O}(\varepsilon/r\sqrt{r}) + \tilde{O}(\varepsilon/r\sqrt{r}) + \nu_7.$$

Moreover, since $\|\mathbf{P}_{\mathbf{u}^{(n)}} \mathbf{h}\| = \tilde{O}(\alpha_2 m \sqrt{r}) = \tilde{O}(\varepsilon/\sqrt{r})$, it holds that

$$\frac{1}{2} \eta_2^2 \|\mathbf{P}_{\mathbf{u}^{(n)}} \mathbf{h}\|^2 = \tilde{O}(\varepsilon^2/r^2).$$

Therefore, ignoring the term $O((\eta_2 \|\mathbf{h}\|)^3)$, we arrive at

$$\begin{aligned} a^{(n+1)} &= a^{(n)} + \eta_2 \langle \beta, \mathbf{h} \rangle - \eta_2 a^{(n)} \langle \mathbf{u}, \mathbf{h} \rangle - \frac{1}{2} \eta_2^2 \|\mathbf{P}_{\mathbf{u}^{(n)}} \mathbf{h}\|^2 \\ &= a^{(n)} + \frac{\eta_2 \varepsilon}{r} a^{(n)2ie(\sigma_*)-1} (1 - a^{(n)2}) (\tilde{\Theta}(1) + O(1/\sqrt{d})) + \eta_2 (\nu_6 - a^{(n)} \nu_7) + \eta_2 \tilde{O}(\varepsilon/r\sqrt{r}) + \tilde{O}(\varepsilon/r^2) \\ &\geq a^{(n)} + \frac{c_3 \varepsilon}{r\sqrt{r}} a^{(n)2ie(\sigma_*)-1} (1 - a^{(n)2}) (1 - O(1/\sqrt{d})) + \nu_8 - \tilde{O}(\varepsilon/r^2), \end{aligned}$$

where ν_8 is a mean-zero sub-Weibull random variable satisfying $\nu_8 = \tilde{O}(\varepsilon/r\sqrt{r})$. \square

Lemma 17. Suppose that $\mathbf{u}^{(1)}$ satisfies $\langle \beta, \mathbf{u}^{(1)} \rangle \geq c_1$, where $c_1 \geq 1/\text{polylog}(d)$. Then, there exists $N_3 = \tilde{\Theta}(\frac{r\sqrt{r}}{\varepsilon} \log \frac{1}{\varepsilon})$ such that

$$\langle \beta, \mathbf{u}^{(N_3+1)} \rangle \geq 1 - \varepsilon$$

with high probability.

Proof. Before going into the proof, we first explain the main idea of the proof. Following the same argument as Lemma 19 in Lee et al. (2024), we can see that after $\tilde{\Theta}(r\sqrt{r}/\varepsilon)$ steps, $a^{(n)}$ becomes larger than a constant $\sqrt{\frac{k+1}{k+2}}$, where $k = 2ie(\sigma_*) - 1$. Then, by applying the Mean Value theorem, we can observe that $1 - a^{(i+1)} \lesssim (1 - \frac{c_3\varepsilon}{r\sqrt{r}}C_k)(1 - a^{(i)})$ where $C_k = \Theta(1)$. This means $1 - a^{(i)}$ converges to 0 geometrically, which yields the required data length $N_3 = \tilde{\Theta}(\frac{r\sqrt{r}}{\varepsilon} \log \frac{1}{\varepsilon})$.

Let $\nu_9^{(i)} = r\sqrt{r}\nu_8^{(i)}/\varepsilon$ and $k = 2ie(\sigma_*) - 1$. Because $\nu_9^{(i)}$ is a sequence of mean-zero sub-Weibull random variables with $\nu_9^{(i)} = \tilde{O}_p(1)$, we have

$$\left| \sum_{i=j}^{j+l} \nu_9^{(i)} \right| = C\sqrt{l}, \quad (6)$$

for any $1 \leq j, l \leq N_3$ with high probability, where $C = \tilde{O}(1)$. Let $c_k = 1 - \sqrt{\frac{k+1}{k+2}}$. Note that c_k is a constant that only depends on k , that is $c_k = O(1)$. Suppose that $c_1 \leq a^{(i)} \leq 1 - \frac{1}{3}c_k$ for all $i = 1, 2, \dots, N_3$. Then, from Lemma 16, we have

$$a^{(n+1)} \geq a^{(n)} + \frac{c_3\varepsilon}{r\sqrt{r}}a^{(n)2ie(\sigma_*)-1}(1 - a^{(n)2})(1 - O(1/\sqrt{d})) + \nu_8 - \tilde{O}(\varepsilon/r^2),$$

for $i \leq N_3$. The term $O(1/\sqrt{d})$ is smaller than 1, thus we may ignore this term. Moreover, the term $\tilde{O}(\varepsilon/r^2)$ is dominated by $\frac{c_3\varepsilon}{r\sqrt{r}}a^{(n)2ie(\sigma_*)-1}(1 - a^{(n)2}) = \tilde{O}(\frac{\varepsilon}{r\sqrt{r}})$. Let $c_2 = c_1^{2ie(\sigma_*)-1}$. By ignoring these terms and using $1 - a^{(i)2} \geq \frac{1}{3}c_k$, we have

$$\begin{aligned} a^{(i+1)} &\geq a^{(i)} + \frac{c_2c_3}{r\sqrt{r}}\varepsilon \cdot c_k/3 + \frac{\varepsilon}{r\sqrt{r}} \cdot \nu_9^{(i)} \\ &\geq a^{(1)} + \frac{c_2c_3c_k}{3r\sqrt{r}}\varepsilon i - \frac{\varepsilon}{r\sqrt{r}} \left| \sum_{j=1}^i \nu_9^{(j)} \right| \\ &\geq a^{(1)} + \left(\frac{c_2c_3c_k}{3r\sqrt{r}}\varepsilon \right) i - \frac{\varepsilon}{r\sqrt{r}}C\sqrt{i}. \end{aligned}$$

If $i \leq \frac{r^3c_1^2}{4\varepsilon^2C^2}$, then $\frac{\varepsilon}{r\sqrt{r}}C\sqrt{i} \leq c_1/2$, and if $i \geq (\frac{6C}{c_2c_3c_k})^2$, then $\frac{\varepsilon}{r\sqrt{r}}C\sqrt{i} \leq \frac{c_2c_3c_k}{6r\sqrt{r}}\varepsilon i$. By observing the order in terms of r , we have $\frac{r^3c_1^2}{4\varepsilon^2C^2} \geq (\frac{6C}{c_2c_3c_k})^2$ when r is sufficiently large. Therefore, it holds that

$$a^{(i+1)} \geq \frac{c_1}{2} + \frac{c_2c_3c_k}{6r\sqrt{r}}\varepsilon i. \quad (7)$$

When $i = \frac{6r\sqrt{r}}{c_2c_3c_k\varepsilon} = \tilde{\Theta}(r\sqrt{r}/\varepsilon)$, then the RHS of Eq. (7) exceeds 1. Therefore, there exists $i^* \leq N_3 = \tilde{\Theta}(\frac{r\sqrt{r}}{\varepsilon} \log \frac{1}{\varepsilon})$ such that $a^{(i^*)} \geq 1 - c_k/3$. Next we prove that $a^{(i)} \geq 1 - c_k = \sqrt{\frac{k+1}{k+2}}$ for all $i = i^*, i^* + 1, \dots, N_3$. In this setting, $a^{(i+1)} - a^{(i)} = \tilde{O}(\frac{\varepsilon}{r\sqrt{r}})$ holds. Therefore, if there exists $i \geq i^*$ such that $a^{(i-1)} \geq 1 - c_k/3$ and $a^{(i)} \leq 1 - c_k/3$, we have $a^{(i)} \geq 1 - 2c_k/3$ with high probability. Also, if $a^{(i+l)} \leq 1 - c_k/3$ holds for all $l = 0, 1, \dots, j - 1$, we have

$$a^{(i+j)} \geq 1 - \frac{2c_k}{3} + \frac{c_2c_3c_k}{3r\sqrt{r}}\varepsilon j - \frac{\varepsilon}{r\sqrt{r}}C\sqrt{j}.$$

1458 If $j \leq \frac{r^3 c_k^2}{9\varepsilon^2 C^2} = \tilde{O}(r^3/\varepsilon^2)$, then $\frac{\varepsilon}{r\sqrt{r}} C\sqrt{j} \leq c_k/3$, and if $j \geq (\frac{3C}{c_2 c_3 c_k})^2 = \tilde{O}(1)$, then $\frac{\varepsilon}{r\sqrt{r}} C\sqrt{j} \leq$
1459 $\frac{c_2 c_3 c_k}{3r\sqrt{r}} \varepsilon j$. Since $\frac{r^3 c_k^2}{9\varepsilon^2 C^2} \geq (\frac{3C}{c_2 c_3 c_k})^2$ when r is sufficiently large, we have

$$1460 \quad a^{(i+j)} \geq 1 - c_k$$

1461
1462 until $i + j = N_3 + 1$ or $a^{(i+j)} \geq 1 - c_k/3$ holds again. By repeating this argument if necessary, we
1463 get $a^{(i+1)} \geq 1 - c_k = \sqrt{\frac{k+1}{k+2}}$ for all $i^* \leq i \leq N_3$.

1464 We continue by showing that, after we achieve $a^{(i^*)} \leq 1 - c_k$, the number of remaining steps needed
1465 to ensure $a^{(i)} \geq 1 - \varepsilon$ is $\tilde{O}(\frac{r\sqrt{r}}{\varepsilon} \log \frac{1}{\varepsilon})$. Let $F(x) = x + \frac{c_3 \varepsilon}{r\sqrt{r}} x^k (1 - x^2)$. Then,

$$1466 \quad a^{(i^*+i+1)} = F(a^{(i^*+i)}) + \frac{\varepsilon}{r\sqrt{r}} \nu_9^{(i^*+i)} - \tilde{O}(\varepsilon/r^2).$$

1467 By the Mean Value theorem, there exists γ such that

$$1468 \quad \frac{1 - F(a^{(i^*+i)})}{1 - a^{(i^*+i)}} = F'(\gamma),$$

1469 and $a^{(i^*+i)} < \gamma < 1$. Now $F'(x) = 1 + \frac{c_2 c_3 \varepsilon}{r\sqrt{r}} x^{k-1} (k - (k+2)x^2)$, and since $\gamma > a^{(i^*+i)} \geq \sqrt{\frac{k+1}{k+2}}$,
1470 we have $k - (k+2)\gamma^2 < -1$, which leads to

$$1471 \quad F'(\gamma) \leq 1 - \frac{c_3 \varepsilon}{r\sqrt{r}} \gamma^{k-1} < 1 - \frac{c_3 \varepsilon}{r\sqrt{r}} \left(\sqrt{\frac{k+1}{k+2}} \right)^{k-1}.$$

1472 Let $C_k = \left(\sqrt{\frac{k+1}{k+2}} \right)^{k-1}$, then $1 - F(a^{(i^*+i)}) < (1 - \frac{c_3 \varepsilon}{r\sqrt{r}} C_k)(1 - a^{(i^*+i)})$, which yields that

$$1473 \quad 1 - a^{i^*+i+1} < \left(1 - \frac{c_3 \varepsilon}{r\sqrt{r}} C_k \right) (1 - a^{(i^*+i)}) + \frac{\varepsilon}{r\sqrt{r}} \nu_9^{(i^*+i)} + \tilde{O}(\varepsilon/r^2).$$

1474 Noting that $(1 - \frac{c_3 \varepsilon}{r\sqrt{r}} C_k) < 1$, taking the sum leads to

$$1475 \quad 1 - a^{(N_3+1)} < \left(1 - \frac{c_3 \varepsilon}{r\sqrt{r}} C_k \right)^{N_3+1-i^*} (1 - a^{(i^*)}) + \sum_{i=i^*}^{N_3} \frac{\varepsilon}{r\sqrt{r}} \nu_9^{(i)} + \tilde{O}(\varepsilon N_3/r^2).$$

1476 Since $\lim_{t \rightarrow \infty} (1 - \frac{1}{t})^t = 1/e$, if we set $N_3 = i^* - 1 + \frac{r\sqrt{r}}{c_3 C_k \varepsilon} \log(1/\varepsilon)$, then we have

$$1477 \quad \left(1 - \frac{c_3 \varepsilon}{r\sqrt{r}} C_k \right)^{N_3+1-i^*} (1 - a^{(i^*)}) < \left(1 - \frac{c_3 \varepsilon}{r\sqrt{r}} C_k \right)^{\frac{r\sqrt{r}}{c_3 C_k \varepsilon} \log(1/\varepsilon)}$$

$$1478 \quad \approx (1/e)^{\log(1/\varepsilon)} = \varepsilon.$$

1479 Finally, by noticing that

$$1480 \quad \left| \sum_{i=i^*}^{N_3} \frac{\varepsilon}{r\sqrt{r}} \nu_9^{(i)} \right| \leq \frac{\varepsilon}{r\sqrt{r}} C \sqrt{N_3 - i^* + 1} = \tilde{O}(r^{-3/4}),$$

1481 we can see that this term is negligible. Moreover, the third term $\tilde{O}(\varepsilon N_3/r^2) = \tilde{O}(\frac{\log 1/\varepsilon}{\sqrt{r}})$ is also
1482 negligible. By summarizing the argument above, we conclude that

$$1483 \quad 1 - a^{(N_3+1)} < \varepsilon.$$

1484 □

1485 In Lemma 9, we assumed that $\langle \mathbf{u}, \mathbf{x}_i \rangle = \tilde{O}_d(1)$. Now $\mathbf{u}^{(n)}$ and \mathbf{x}_i are not independent of each other,
1486 so we need to ensure that $\langle \mathbf{u}^{(j)}, \mathbf{x}_i \rangle = \tilde{O}_d(1)$ for all $j = 1, 2, \dots, N_3 + 1$.

1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

Lemma 18. *It holds that*

$$\begin{aligned} \langle \mathbf{x}_i^*, \mathbf{h} \rangle &= \Theta(\alpha_2 m) \langle \beta, \mathbf{u} \rangle^{2ie(\sigma_*)-1} \frac{\langle \mathbf{x}_i^*, \beta \rangle}{(\kappa\rho)^{ie(\sigma_*)}} (1 + O(1/\sqrt{d})) \\ &\quad + \tilde{O}(\alpha_2^2 m^2 r) + \tilde{O}(\alpha_2 m r / \sqrt{N_2}) + \nu_*, \end{aligned}$$

where ν_* is a mean-zero random variable satisfying $\nu_* = \tilde{O}(\alpha_2 m \sqrt{r})$.

Proof. From Lemma 15, we have

$$\begin{aligned} \mathbf{h} &= \Theta(\alpha_2 m) \langle \beta, \mathbf{u} \rangle^{2ie(\sigma_*)-1} \frac{\beta}{(\kappa\rho)^{ie(\sigma_*)}} (1 + O(1/\sqrt{d})) \\ &\quad + \tilde{O}(\alpha_2^2 m^2 \sqrt{r}) + \tilde{O}(\alpha_2 m \sqrt{\frac{r}{N_2}}) + \nu. \end{aligned}$$

Considering $\mathbf{x}_i^* = \tilde{O}(\sqrt{r})$, taking the inner product of \mathbf{x}_i^* and \mathbf{h} leads to

$$\begin{aligned} \langle \mathbf{x}_i^*, \mathbf{h} \rangle &= \Theta(\alpha_2 m) \langle \beta, \mathbf{u} \rangle^{2ie(\sigma_*)-1} \frac{\langle \mathbf{x}_i^*, \beta \rangle}{(\kappa\rho)^{ie(\sigma_*)}} (1 + O(1/\sqrt{d})) \\ &\quad + \tilde{O}(\alpha_2^2 m^2 r) + \tilde{O}(\alpha_2 m r / \sqrt{N_2}) + \langle \mathbf{x}_i^*, \nu \rangle. \end{aligned}$$

Following the same argument as Lemma 15 yields $\langle \mathbf{x}_i^*, \nu \rangle = \tilde{O}(\alpha_2 m \sqrt{r})$. Thus, by letting $\nu_* = \langle \mathbf{x}_i^*, \nu \rangle$, we obtain the assertion. \square

Lemma 19. *Set $\alpha_2 m = \tilde{\Theta}(\varepsilon/r)$, $N_2 = \tilde{\Theta}(r^2)$, $N_3 = \tilde{\Theta}(\frac{r\sqrt{r}}{\varepsilon} \log \frac{1}{\varepsilon})$ and $\eta_2 = \tilde{\Theta}(1/\sqrt{r})$ (This is exactly the same situation as Lemma 16 and Lemma 17). Then we have*

$$\langle \mathbf{x}_i^*, \mathbf{u}^{(j)} \rangle = \tilde{O}_d(1),$$

for all $j = 1, 2, \dots, N_3 + 1$.

Proof. When $j = 1$, \mathbf{x}_i and $\mathbf{u}^{(1)}$ are independent of each other, thus Lemma 29 yields the desired result. Substituting the parameters into the result of Lemma 18 yields

$$\eta_2 \langle \mathbf{x}_i^*, \mathbf{h} \rangle = \tilde{O} \left(\frac{\varepsilon}{r\sqrt{r}} (1 + O(1/\sqrt{d})) \right) + \tilde{O} \left(\frac{\varepsilon}{r\sqrt{r}} \right) + \eta_2 \nu_*,$$

where $\eta_2 \nu_* = \tilde{O}(\varepsilon/r)$. First, following the same argument to derive Eq. (6) yields

$$\left| \sum_{k=1}^j \eta_2 \nu_* \right| = \tilde{O}(\varepsilon \sqrt{j}/r),$$

for $j = 1, 2, \dots, N_3$. Since $\varepsilon \sqrt{N_3}/r = \tilde{O}_d(r^{-1/4}) = o_d(1)$, we can ignore the effect of $\eta_2 \nu_*$. Then, we can see that

$$\begin{aligned} \langle \mathbf{x}_i^*, \mathbf{u}^{(j+1)} \rangle &= \langle \mathbf{x}_i^*, \mathbf{u}^{(j)} \rangle + \eta_2 \langle \mathbf{x}_i^*, \mathbf{h} \rangle - \eta_2 \langle \mathbf{h}, \mathbf{u}^{(j)} \rangle \langle \mathbf{u}^{(j)}, \mathbf{x}_i^* \rangle - \frac{1}{2} \eta_2^2 \|\mathbf{P}_{\mathbf{u}^{(n)}} \mathbf{h}\|^2 + O((\eta_2 \|\mathbf{h}\|)^3) \\ &= \langle \mathbf{x}_i^*, \mathbf{u}^{(j)} \rangle + \langle \beta, \mathbf{u}^{(j)} \rangle^{2ie(\sigma_*)-1} \tilde{O} \left(\frac{\varepsilon}{r\sqrt{r}} \right) \left(\langle \beta, \mathbf{x}_i^* \rangle - \langle \mathbf{u}^{(j)}, \mathbf{x}_i^* \rangle \langle \beta, \mathbf{u}^{(j)} \rangle \right) \\ &\quad - \tilde{O} \left(\frac{\varepsilon}{r\sqrt{r}} \right) - \tilde{O} \left(\frac{\varepsilon^2}{r^2} \right) \end{aligned}$$

holds. Now $\langle \beta, \mathbf{x}_i^* \rangle = \tilde{O}_p(1)$ and $\langle \beta, \mathbf{u}^{(j)} \rangle \leq 1$. Therefore, if $\langle \mathbf{u}^{(j)}, \mathbf{x}_i^* \rangle = \tilde{O}_p(1)$ holds, then $\langle \mathbf{x}_i^*, \mathbf{u}^{(j+1)} \rangle = \langle \mathbf{x}_i^*, \mathbf{u}^{(j)} \rangle + \tilde{O} \left(\frac{\varepsilon}{r\sqrt{r}} \right)$. Therefore, by induction, we have

$$\langle \mathbf{x}_i^*, \mathbf{u}^{(N_3+1)} \rangle = \tilde{O} \left(\log \frac{1}{\varepsilon} \right),$$

which does not depend on d . \square

1566 This lemma allows us to use Lemma 9. Note that $N_3 = \tilde{\Theta}(\frac{r\sqrt{r}}{\varepsilon} \log \frac{1}{\varepsilon})$ means, by ignoring the term
 1567 $\log \frac{1}{\varepsilon}$, $\varepsilon = \tilde{\Theta}(\frac{r\sqrt{r}}{N_3})$. Now we can estimate $\langle \beta, \mathbf{x} \rangle$ using $\langle \mathbf{u}^{(N_3+1)}, \mathbf{x} \rangle$. For notational simplicity, we
 1568 write $\mathbf{u}^{(N_3+1)}$ as \mathbf{u} from now on.

1570 **Lemma 20.** *Let $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$, and \mathbf{u} be a vector independent of \mathbf{x} satisfying $\|\mathbf{u}\| = 1$ and*
 1571 *$\langle \beta, \mathbf{u} \rangle \geq 1 - \varepsilon$. Then, we have*

$$1572 \quad \langle \beta, \mathbf{x} \rangle = \langle \mathbf{u}, \mathbf{x} \rangle + O_p(\sqrt{2\varepsilon \log d}).$$

1575 *Proof.* First, note that

$$1576 \quad \langle \beta, \mathbf{x} \rangle = \langle \mathbf{u}, \mathbf{x} \rangle + \langle \beta - \mathbf{u}, \mathbf{x} \rangle.$$

1577 Since we have $\langle \beta, \mathbf{u} \rangle \geq 1 - \varepsilon$, we have that

$$1578 \quad \|\beta - \mathbf{u}\|^2 = \|\beta\|^2 + \|\mathbf{u}\|^2 - 2\langle \beta, \mathbf{u} \rangle \leq 2\varepsilon,$$

1580 which yields

$$1581 \quad \|\beta - \mathbf{u}\| \leq \sqrt{2\varepsilon}.$$

1582 Then, combining with Lemma 25 yields $\langle \beta - \mathbf{u}, \mathbf{x} \rangle = O_p(\sqrt{2\varepsilon \log d})$. □

1585 F TRAINING MLP LAYER

1587 In this section, we show that the MLP layer can fit the polynomial σ_*^{test} . Most of the argument in
 1588 this section is taken from Nishikawa et al. (2025), but we do not omit the proof for readability.

1589 **Lemma 21.** *Suppose that there exists $g(\mathbf{x})$ such that*

$$1591 \quad |g(\mathbf{x}) - \langle \beta, \mathbf{x} \rangle| \leq \delta.$$

1592 *Then, there exists $\pi(v, b)$ such that*

$$1593 \quad |\mathbb{E}_{v \sim \text{Unif}\{\pm 1\}, b \sim [-\log^2 d, \log^2 d]}[\pi(v, b)\sigma(v \cdot g(\mathbf{x}) + b)] - \sigma_*(\langle \beta, \mathbf{x} \rangle)| = O(\delta(\log d)^{2\deg(\sigma_*)-2})$$

1596 *with high probability over $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$. Moreover, $\sup_{v, b} |\pi(v, b)| = \tilde{O}(1)$ holds.*

1598 *Proof.* Let $\sigma_*(z) = \sum_{k=0}^{\deg(\sigma_*)} s_k z^k$. Now from Lemma 9 in Damian et al. (2022), there exists
 1599 $\pi'_k(v, b)$ such that $\sup_{v, b} |\pi'_k(v, b)| = O(1)$ and

$$1600 \quad \mathbb{E}_{v \sim \text{Unif}\{\pm 1\}, b \sim [-1, 1]}[\pi'_k(v, b)\sigma(vz + b)] = z^k$$

1603 for any $|z| \leq 1$. If we define

$$1604 \quad \pi(v, b) = \sum_{k=0}^{\deg(\sigma_*)} s_k \frac{\pi'_k(v, b \log^{-2} d)}{\log^2 d} \log^{2k} d,$$

1608 then, we have $\sup_{v, b} |\pi(v, b)| = O(\log^{2\deg(\sigma_*)-2} d)$. Moreover, when $|z \log^{-2} d| \leq 1$, we have

$$1610 \quad \mathbb{E}_{v \sim \text{Unif}\{\pm 1\}, b \sim [-\log^2 d, \log^2 d]}[\pi(v, b)\sigma(vz + b)]$$

$$1611 \quad = \sum_{k=0}^{\deg(\sigma_*)} s_k \mathbb{E}_{v \sim \text{Unif}\{\pm 1\}, b \sim [-\log^2 d, \log^2 d]} \left[\frac{\pi'_k(v, b \log^{-2} d)}{\log^2 d} \log^{2k} d \sigma(vz + b) \right]$$

$$1612 \quad = \sum_{k=0}^{\deg(\sigma_*)} s_k \log^{2k-2} d \mathbb{E}_{v \sim \text{Unif}\{\pm 1\}, b \sim [-1, 1]} [\pi'_k(v, b)\sigma(\log^2 d(vz \log^{-2} d + b))]$$

$$1613 \quad = \sum_{k=0}^{\deg(\sigma_*)} s_k z^k = \sigma_*(z).$$

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

Therefore, it holds that

$$\begin{aligned}
& \mathbb{E}_{v \sim \text{Unif}\{\pm 1\}, b \sim [-\log^2 d, \log^2 d]} [\pi(v, b) \sigma(v \cdot g(\mathbf{x}) + b)] \\
&= \mathbb{E}_{v \sim \text{Unif}\{\pm 1\}, b \sim [-\log^2 d, \log^2 d]} [\pi(v, b) \sigma(v \langle \beta, \mathbf{x} \rangle + v\delta + b)] \\
&= \mathbb{E}_{v \sim \text{Unif}\{\pm 1\}, b \sim [-\log^2 d, \log^2 d]} [\pi(v, b) \sigma(v \langle \beta, \mathbf{x} \rangle + b)] + O(\delta \log^{2 \deg(\sigma_*) - 2} d) \\
&= \sigma_*(\langle \beta, \mathbf{x} \rangle) + O(\delta \log^{2 \deg(\sigma_*) - 2} d).
\end{aligned}$$

Here, we used the fact that $|\langle \beta, \mathbf{x} \rangle| \leq \log^2 d$ with high probability. □

Lemma 22. *Under the condition of Lemma 21, there exists $\mathbf{a}' \in \mathbb{R}^m$ such that*

$$\left| \sum_{j=1}^m a'_j \sigma(v_j \cdot g(\mathbf{x}) + b_j) - \sigma_*(\langle \beta, \mathbf{x} \rangle) \right| = \tilde{O}(m^{-\frac{1}{2}}) + O(\delta \log^{2 \deg(\sigma_*) - 2} d)$$

holds with high probability over $\mathbf{x} \sim (0, \mathbf{I}_d)$. Moreover, $\|\mathbf{a}'\|^2 = \tilde{O}(m^{-1/2})$ holds with high probability.

Proof. Using $\pi(v, b)$ in Lemma 21, define $a'_j = m^{-1} \pi(v_j, b_j)$. Since $\sup_{v, b} |\pi(v, b)| = \tilde{O}(1)$, from Hoeffding's inequality, it holds that

$$\left| m^{-1} \sum_{j=1}^m \pi(v_j, b_j) \sigma(v_j \cdot g(\mathbf{x}) + b_j) - \mathbb{E}_{v, b} [\pi(v, b) \sigma(v \cdot g(\mathbf{x}) + b)] \right| = \tilde{O}(m^{-1/2}),$$

with high probability. Hence we have

$$\left| \sum_{j=1}^m a'_j \sigma(v_j \cdot g(\mathbf{x}) + b_j) - \sigma_*(\langle \beta, \mathbf{x} \rangle) \right| = \tilde{O}(m^{-\frac{1}{2}}) + O(\delta \log^{2 \deg(\sigma_*) - 2} d).$$

As for the upper bound of the norm $\|\mathbf{a}'\|^2 = m^{-1} \cdot m^{-1} \sum_{j=1}^m \pi(v_j, b_j)^2$, Hoeffding's inequality yields

$$\left| m^{-1} \sum_{j=1}^m \pi(v_j, b_j)^2 - \mathbb{E}_{v, b} [\pi(v, b)^2] \right| = \tilde{O}(m^{-1/2}).$$

Since $\sup_{v, b} |\pi(v, b)| = \tilde{O}(1)$, we can say that $\mathbb{E}_{v, b} [\pi(v, b)^2] = \tilde{O}(1)$, which completes the proof. □

Lemma 23. *Let \mathbf{a}^* be the parameter trained via Algorithm 1. Then there exists λ_3 such that*

$$\frac{1}{N_4} \sum_{t=N_1+N_2+N_3+1}^{N_1+N_2+N_3+N_4} |y_t - f_{\text{TF}}(\mathbf{x}_t, \mathbf{u}, \mathbf{v}^*, \mathbf{a}^*, \mathbf{b}^*)| = \tau + \tilde{O}(m^{-1/2}) + O(\delta \log^{2 \deg(\sigma_*) - 2} d)$$

with high probability. Moreover, $\|\mathbf{a}^\|^2 \leq \tilde{O}_p(m^{-1/2})$ is satisfied.*

Proof. Let $M = N_1 + N_2 + N_3$ and \mathbf{a}' be the output parameter constructed in Lemma 22: from the equivalence between ℓ_2 -regularized and norm-constrained optimization algorithms, if we carefully choose λ_3 , then we have

$$\begin{aligned}
\left(\frac{1}{N_4} \sum_{t=M+1}^{M+N_4} |y_t - f_{\text{TF}}(\mathbf{x}_t, \mathbf{u}, \mathbf{v}^*, \mathbf{a}^*, \mathbf{b}^*)| \right)^2 &\leq \frac{1}{N_4} \sum_{t=M+1}^{M+N_4} (y_t - f_{\text{TF}}(\mathbf{x}_t, \mathbf{u}, \mathbf{v}^*, \mathbf{a}^*, \mathbf{b}^*))^2 \\
&\leq \frac{1}{N_4} \sum_{t=M+1}^{M+N_4} (y_t - f_{\text{TF}}(\mathbf{x}_t, \mathbf{u}, \mathbf{v}^*, \mathbf{a}', \mathbf{b}^*))^2 \\
&\leq (\tau + \tilde{O}(m^{-1/2}) + O(\delta \log^{2 \deg(\sigma_*) - 2} d))^2
\end{aligned}$$

(recall that we assumed $\tau = \Theta(1)$) with high probability, which yields the first assertion. Moreover, from Lemma 22, we can see that

$$\|\mathbf{a}^*\|^2 \leq \|\mathbf{a}'\|^2 \leq \tilde{O}_p(m^{-1/2}),$$

which completes the proof. \square

Lemma 24. *We fix \mathbf{b} and \mathbf{v} . Let $\mathcal{F}_A = \{\mathbf{a} \mapsto \sum_{j=1}^m a_j \sigma(v_j \langle \mathbf{x}, \mathbf{u} \rangle + b_j) \mid \|\mathbf{a}\| \leq A\}$ the set of transformers where the norm of the MLP parameter \mathbf{a} is upper bounded. When $\|\mathbf{b}\| \leq B$ and $v_j = \pm 1$, then it holds that*

$$\text{Rad}_N(\mathcal{F}_A) = \tilde{O}\left(\frac{A(B + \sqrt{m})}{\sqrt{N_4}}\right).$$

Proof. Note that $\mathbb{E}[g(\mathbf{x})^2] = \mathbb{E}[\langle \mathbf{u}, \mathbf{x} \rangle^2] = O(\log d)$, considering $\|\mathbf{u}\| = 1$. Then, following the same argument as Lemma 25 in Nishikawa et al. (2025) yields the assertion. \square

G PROOF OF THE THEOREM 1

Now, we are ready to give the proof of Theorem 1.

Proof. First note that

$$\begin{aligned} \mathcal{R}_{f_{\text{TF}}}(\mathbf{u}, \mathbf{v}^*, \mathbf{a}^*, \mathbf{b}^*) - \tau &= \frac{1}{N_4} \sum_{i=N_1+N_2+N_3+1}^{N_1+N_2+N_3+N_4} |y_i - f_{\text{TF}}(\mathbf{x}_i, \mathbf{u}, \mathbf{v}^*, \mathbf{a}^*, \mathbf{b}^*)| - \tau \\ &+ \mathcal{R}_{f_{\text{TF}}}(\mathbf{u}, \mathbf{v}^*, \mathbf{a}^*, \mathbf{b}^*) - \frac{1}{N_4} \sum_{i=N_1+N_2+N_3+1}^{N_1+N_2+N_3+N_4} |y_i - f_{\text{TF}}(\mathbf{x}_i, \mathbf{u}, \mathbf{v}^*, \mathbf{a}^*, \mathbf{b}^*)|. \end{aligned}$$

From Lemma 23, the first two terms $\frac{1}{N_4} \sum_{i=N_1+N_2+N_3+1}^{N_1+N_2+N_3+N_4} |y_i - f_{\text{TF}}(\mathbf{x}_i, \mathbf{u}, \mathbf{v}^*, \mathbf{a}^*, \mathbf{b}^*)| - \tau$ are bounded by $\tilde{O}(m^{-1/2}) + O(\delta \log^2 \deg(\sigma_*)^{-2} d)$. Moreover, from Lemma 22 and the definition of \mathbf{b}^* , we have $\|\mathbf{a}^*\| = \tilde{O}(m^{-1/2})$ and $\|\mathbf{b}^*\| = \tilde{O}(\sqrt{m})$. Therefore, from Lemma 24, we have

$$\text{Rad}_N(\mathcal{F}_A) = \tilde{O}(N_4^{-1/2}).$$

Using the standard technique yields (see Appendix D.3 in Oko et al. (2024))

$$|\mathcal{R}_{f_{\text{TF}}}(\mathbf{u}, \mathbf{v}^*, \mathbf{a}^*, \mathbf{b}^*) - \tau| = \tilde{O}(N_4^{-1/2}) + \tilde{O}(m^{-1/2}) + O(\delta \log^2 \deg(\sigma_*)^{-2} d),$$

with probability at least 0.995. Note that all the desired events occur with high probability, except for this event. Therefore, when d is large enough, all the events occur with probability at least 0.99.

Finally, we rewrite the term $O(\delta \log^2 \deg(\sigma_*)^{-2} d)$. From Lemma 20, $\delta = O(\sqrt{2\varepsilon \log d})$ holds. Also, $N_3 = \tilde{\Theta}(\frac{r\sqrt{r}}{\varepsilon} \log \frac{1}{\varepsilon})$, which means, when we ignore the term $\log \frac{1}{\varepsilon}$, $\varepsilon = \tilde{\Theta}(\frac{r\sqrt{r}}{N_3})$. Therefore

$$O(\delta \log^2 \deg(\sigma_*)^{-2} d) = \tilde{O}\left(\sqrt{\frac{r\sqrt{r}}{N_3}}\right)$$

is satisfied. This completes the proof. \square

H ADDITIONAL LEMMAS

Lemma 25. *Suppose $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$ and \mathbf{y} is a vector satisfying $\|\mathbf{y}\| = C$ independent of \mathbf{x} . Then*

$$\langle \mathbf{x}, \mathbf{y} \rangle = O(C\sqrt{\log d})$$

holds with high probability.

Lemma 26. Let z_1, \dots, z_n be i.i.d. random variables which satisfy $\|z_i\| \leq C$ with high probability and $\mathbb{E}[z_1^2] = O(d^\alpha)$, where α is a constant independent of d . When $n = \text{poly}(d)$,

$$\frac{1}{n} \sum_{i=1}^n z_i - \mathbb{E}[z_i] = O\left(C\sqrt{\frac{1}{n}}\right)$$

holds with high probability.

Proof. Let $z'_i = z_i \mathbf{1}_{|z_i| \leq C}$. Applying the uniform bound argument, we may consider that $z_i = z'_i$ for all $i = 1, \dots, n$ with high probability because our assumption $n = \text{poly}(d)$ yields

$$P(\max_i |z_i| \leq C) \geq 1 - O(nd^{-C_*}) \geq 1 - O(d^{-C_*}),$$

where C'_* is a constant determined appropriately. Then, using Hoeffding's inequality yields

$$\frac{1}{n} \sum_{i=1}^n z_i - \mathbb{E}[z'_i] = O\left(C\sqrt{\frac{1}{n}}\right).$$

We complete the proof by showing that $|\mathbb{E}[z_i] - \mathbb{E}[z'_i]|$ is sufficiently small. This can be shown by

$$\begin{aligned} \mathbb{E}[z_i] - \mathbb{E}[z'_i] &= \mathbb{E}[z_i \mathbf{1}_{|z_i| > C_z}] \\ &\leq \mathbb{E}[\mathbf{1}_{|z_i| > C_z}]^{1/2} \mathbb{E}[z_i^2]^{1/2} \\ &\leq O(d^{\alpha - C_*}), \end{aligned}$$

where C_* can be taken sufficiently large from the definition of high probability event. Since we assumed that $n = \text{poly}(d)$, this term is smaller than the main term, which completes the proof. \square

Lemma 27 (Damian et al. (2023), Property 1). Let $\alpha, \beta \in \mathbb{S}^{d-1}$, then

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)}[\text{He}_i(\langle \alpha, \mathbf{x} \rangle) \text{He}_i(\langle \beta, \mathbf{x} \rangle)] = \mathbf{1}_{i=j} i! \langle \alpha, \beta \rangle^i.$$

H.1 PRETRAINED MATRIX

Let Γ^* be a matrix obtained after pretraining. Then, from Nishikawa et al. (2025), Γ^* can be written as

$$\Gamma^* = \frac{r\mathbb{E}_\beta[\beta\beta^\top] + \mathbf{N}}{\kappa\sqrt{r}}$$

for some matrix \mathbf{N} satisfying $\|\mathbf{N}\|_F = O(1/\sqrt{d})$ and a number $\kappa = \Theta(\log^{C_\kappa} d)$. Also, from Nishikawa et al. (2025) and the assumption on the support of β ,

$$U\beta \sim \text{Unif}\{(\alpha_1, \alpha_2, \dots, \alpha_r, 0, \dots, 0) \mid \alpha_1^2 + \dots + \alpha_r^2 = 1\}.$$

holds for some orthogonal matrix U . Then, using $\mathbf{D} = \text{diag}(\underbrace{1, \dots, 1}_r, \underbrace{0, \dots, 0}_{d-r})$, we have that

$$r\mathbb{E}_\beta[\beta\beta^\top] = U^\top \mathbf{D} U.$$

Lemma 28. Let $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$, then

$$\|\sqrt{r}\kappa\Gamma^*\mathbf{x}\| = \tilde{O}(\sqrt{r})$$

holds with high probability.

Proof. By the argument above, we can write that

$$\sqrt{r}\kappa\Gamma^*\mathbf{x} = (U^\top \mathbf{D} U + \mathbf{N})\mathbf{x}.$$

From rotational invariance, when we define \mathbf{y} as $\mathbf{y} = U\mathbf{x}$, $\mathbf{y} \sim \mathcal{N}(0, \mathbf{I}_d)$ is satisfied. From the definition of \mathbf{D} , the first r components of $\mathbf{D}\mathbf{y}$ follow standard normal distribution i.i.d., and the other $(n - r)$ components are equal to zero. Since applying U^\top to a vector does not change the norm of the vector, we obtain

$$U^\top \mathbf{D} U \mathbf{x} = U^\top \mathbf{D} \mathbf{y} = \tilde{O}(\sqrt{r}),$$

with high probability. Finally, since $\|\mathbf{N}\|_F = O(1/\sqrt{d})$, $\mathbf{N}\mathbf{x} = \tilde{O}(1)$ holds. Hence we can ignore the term $\mathbf{N}\mathbf{x}$. \square

1782 **Lemma 29.** Let $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$, and \mathbf{u} be a vector independent of \mathbf{x} satisfying $\|\mathbf{u}\| = 1$. then

$$1783 \langle \sqrt{r}\kappa\Gamma^* \mathbf{x}, \mathbf{u} \rangle = O(\sqrt{\log d})$$

1784 holds with high probability.

1785 *Proof.* As in the previous lemma, we know that

$$1786 \sqrt{r}\kappa\Gamma^* \mathbf{x} = (\mathbf{U}^\top \mathbf{D}\mathbf{U} + \mathbf{N})\mathbf{x}$$

1787 holds. Again, from rotational invariance, when we define \mathbf{y} as $\mathbf{y} = \mathbf{U}\mathbf{x}$, $\mathbf{y} \sim \mathcal{N}(0, \mathbf{I}_d)$ is satisfied.
 1788 From the definition of \mathbf{D} , the first r components of $\mathbf{D}\mathbf{y}$ follow standard normal distribution i.i.d.,
 1789 and the other $(n - r)$ components are equal to zero. Also, when we define $\mathbf{u}' = \mathbf{U}\mathbf{u}$, we have
 1790 $\|\mathbf{u}'\| = 1$, which means $\sqrt{u_1'^2 + \dots + u_r'^2} \leq 1$. Therefore, from Corollary 30 in Nishikawa et al.
 1791 (2025), we have

$$1792 \langle \mathbf{U}^\top \mathbf{D}\mathbf{U}\mathbf{x}, \mathbf{u} \rangle = \langle \mathbf{D}\mathbf{y}, \mathbf{u}' \rangle = O_p(\sqrt{\log d}).$$

1793 As for $\mathbf{N}\mathbf{x}$, again from Corollary 30 in Nishikawa et al. (2025), we have

$$1794 \langle \mathbf{N}\mathbf{x}, \mathbf{u} \rangle = O_p\left(\sqrt{\frac{\log d}{d}}\right).$$

1795 This is smaller than the main term. □

1800 **Lemma 30.** Let $\mathbf{x}_1 \dots \mathbf{x}_n \sim \mathcal{N}(0, \mathbf{I}_d)$, and z_1, \dots, z_n are i.i.d. random variables which satisfy
 1801 $|z_i| \leq C_z$ with high probability. When $n = \text{poly}(d)$, it holds that

$$1802 \left\| \frac{1}{n} \sum_{i=1}^n (z_i \sqrt{r}\kappa\Gamma^* \mathbf{x}_i) - \mathbb{E}[z_1 \sqrt{r}\kappa\Gamma^* \mathbf{x}_1] \right\| = \tilde{O}\left(C_z \sqrt{\frac{r}{n}}\right),$$

1803 with high probability.

1804 *Proof.* We use the same proof strategy as Lemma 31 in Nishikawa et al. (2025). Let $z'_i = z_i \mathbf{1}_{|z_i| \leq C_z}$.
 1805 First, we can confirm that

$$1806 \begin{aligned} \mathbb{E}[z_1 \sqrt{r}\kappa\Gamma^* \mathbf{x}_1] - \mathbb{E}[z'_1 \sqrt{r}\kappa\Gamma^* \mathbf{x}_1] &= \mathbb{E}[\mathbf{1}_{|z_1| > C_z} \sqrt{r}\kappa\Gamma^* \mathbf{x}_1] \\ 1807 &\leq \mathbb{E}[\mathbf{1}_{|z_1| > C_z}^2]^{1/2} \mathbb{E}[\|\sqrt{r}\kappa\Gamma^* \mathbf{x}_1\|^2]^{1/2} \\ 1808 &\leq O(d^{-C}). \end{aligned}$$

1809 Where C can be taken sufficiently large from the definition of high probability event. Because
 1810 $n = \text{poly}(d)$, by redefining C if necessary, we can make this term smaller than $\tilde{O}(C_z \sqrt{\frac{r}{n}})$.
 1811 As $\frac{1}{n} \sum_{i=1}^n (z_i \sqrt{r}\kappa\Gamma^* \mathbf{x}_i) = \frac{1}{n} \sum_{i=1}^n (z'_i \sqrt{r}\kappa\Gamma^* \mathbf{x}_i)$ with high probability, we continue to evalu-
 1812 ate $\frac{1}{n} \sum_{i=1}^n (z'_i \sqrt{r}\kappa\Gamma^* \mathbf{x}_i) - \mathbb{E}[z'_1 \sqrt{r}\kappa\Gamma^* \mathbf{x}_1]$. Let $\mathbf{y}_i = \mathbf{U}\mathbf{x}_i$. Then $\mathbf{y}_i \sim \mathcal{N}(0, \mathbf{I}_d)$ holds. $z'_i \mathbf{D}\mathbf{y}_i$
 1813 is a sub-Gaussian vector, and $(\mathbf{D}\mathbf{y}_i)_k = 0$ when $k > r$ from the definition of \mathbf{D} . Then, applying a
 1814 standard concentration bound for a sub-Gaussian vector to the r -dimensional vector $(\mathbf{D}\mathbf{y}_i)_{1:r}$ yields

$$1815 \frac{1}{n} \sum_{i=1}^n z'_i \mathbf{D}\mathbf{y}_i - \mathbb{E}[z'_1 \mathbf{D}\mathbf{y}_1] = \tilde{O}\left(C_z \sqrt{\frac{r}{n}}\right).$$

1816 As multiplying the orthogonal matrix \mathbf{U}^\top does not change the norm of the vector, we have

$$1817 \frac{1}{n} \sum_{i=1}^n z'_i \mathbf{U}^\top \mathbf{D}\mathbf{U}\mathbf{x}_i - \mathbb{E}[z'_1 \mathbf{U}^\top \mathbf{D}\mathbf{U}\mathbf{x}_1] = \tilde{O}\left(C_z \sqrt{\frac{r}{n}}\right).$$

1818 Again from the standard concentration bound for a sub-Gaussian vector, we have

$$1819 \frac{1}{n} \sum_{i=1}^n z'_i \mathbf{x}_i - \mathbb{E}[z'_1 \mathbf{x}_1] = \tilde{O}\left(C_z \sqrt{\frac{d}{n}}\right),$$

with high probability. Since $\|\mathbf{N}\|_F = O(1/\sqrt{d})$, we have

$$\frac{1}{n} \sum_{i=1}^n z'_i \mathbf{N} \mathbf{z}_i - \mathbb{E}[z'_1 \mathbf{N} \mathbf{x}_1] = \tilde{O} \left(C_z \sqrt{\frac{1}{n}} \right).$$

In summary, we arrive at

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n (z_i \sqrt{r\kappa} \Gamma^* \mathbf{x}_i) - \mathbb{E}[z_1 \sqrt{r\kappa} \Gamma^* \mathbf{x}_1] \\ &= \frac{1}{n} \sum_{i=1}^n (z_i (\mathbf{U}^\top \mathbf{D} \mathbf{U} + \mathbf{N}) \mathbf{x}_i) - \mathbb{E}[z_1 (\mathbf{U}^\top \mathbf{D} \mathbf{U} + \mathbf{N}) \mathbf{x}_1] = \tilde{O} \left(C_z \sqrt{\frac{r}{n}} \right), \end{aligned}$$

with high probability, which completes the proof. \square

Lemma 31. *We have that*

$$\sqrt{r\kappa} \Gamma^* \beta = \beta + O(1/\sqrt{d}).$$

Proof. First note that

$$\sqrt{r\kappa} \Gamma^* \beta = (\mathbf{U}^\top \mathbf{D} \mathbf{U} + \mathbf{N}) \beta.$$

From the definition of \mathbf{U} , we have $\mathbf{U} \beta = (\alpha_1, \alpha_2, \dots, \alpha_r, 0, \dots, 0)^\top$ for some $\alpha_1, \alpha_2, \dots, \alpha_r \in \mathbb{R}$. Therefore, we obtain that

$$\mathbf{U}^\top \mathbf{D} \mathbf{U} \beta = \mathbf{U}^\top \mathbf{U} \beta = \beta.$$

Finally, by noticing $\|\mathbf{N}\|_F = O(1/\sqrt{d})$ and $\|\beta\| = 1$, we can see that $\mathbf{N} \beta = O(1/\sqrt{d})$ holds. \square

Lemma 32. *We have that*

$$\sqrt{r\kappa} \Gamma^* (\sqrt{r\kappa} \Gamma^* \beta) = \beta + O(1/\sqrt{d}).$$

Proof. From lemma 31, we have $\sqrt{r\kappa} \Gamma^* \beta = \beta + \mathbf{d}$ where $\mathbf{d} = O(1/\sqrt{d})$. Therefore, it holds that

$$\sqrt{r\kappa} \Gamma^* (\sqrt{r\kappa} \Gamma^* \beta) = \sqrt{r\kappa} \Gamma^* (\beta + \mathbf{d}).$$

Again from Lemma 31, we see that $\sqrt{r\kappa} \Gamma^* \beta = \beta + O(1/\sqrt{d})$. Also we have $\sqrt{r\kappa} \Gamma^* \mathbf{d} = \mathbf{U}^\top \mathbf{D} \mathbf{U} \mathbf{d} + \mathbf{N} \mathbf{d}$, and since $\|\mathbf{U}^\top \mathbf{D} \mathbf{U}\|_2 = 1$ from the definition of \mathbf{U} and \mathbf{D} , we have $\mathbf{U}^\top \mathbf{D} \mathbf{U} \mathbf{d} = O(1/\sqrt{d})$. Finally, since $\|\mathbf{N}\|_F = O(1/\sqrt{d})$, we have $\mathbf{N} \mathbf{d} = O(1/d)$. This indicates that $\sqrt{r\kappa} \Gamma^* \mathbf{d} = O(1/\sqrt{d})$. \square

I EXPERIMENTAL DETAILS

The GPT-2 model architecture we used in Section 4 originates from Garg et al. (2023). Given the $(N + 1)$ -length prompt $\{(\mathbf{x}_i, y_i)\}_{i=1}^{N+1}$, we first construct the embedding as

$$\mathbf{E} = [\mathbf{x}_1, \tilde{\mathbf{y}}_1, \dots, \mathbf{x}_{N+1}, \mathbf{y}_{N+1}] \in \mathbb{R}^{d \times (2N+2)},$$

where $\tilde{\mathbf{y}}_i = [y_i, 0, \dots, 0]^\top$. Next, the read-in layer transforms this embedding into $\tilde{\mathbf{E}} \in \mathbb{R}^{D \times (2N+2)}$, where $D = 128$. This mapped embedding $\tilde{\mathbf{E}}$ goes through a 6-layer GPT-2 backbone with 4 attention heads, following the configuration by Garg et al. (2023). Finally, the output of GPT-2 backbone is transformed by the read-out layer into the vector $[z_1, z_2, \dots, z_{2N+1}, z_{2N+2}]$. Here, z_{2i-1} is the prediction of y_i given the context $(\mathbf{x}_1, y_1, \dots, \mathbf{x}_{i-1}, y_{i-1}, \mathbf{x}_i)$. We used Adam optimizer (Kingma & Ba, 2017) with a learning rate of 0.0001. To reduce the pretraining cost, we adopted the curriculum learning strategy, which is also used in Garg et al. (2023). The training started with the dimension $d = 4$, and the dimension was increased by two until it reached the target dimension.

For the model with TTT, we used the same base model as ICL evaluation. We introduced low-rank adaptation (LoRA) to the attention projection layers (c.attn and c.proj) and the feedforward layer (c.fc). The rank of LoRA was set to 4, and the parameters LoRA_alpha and LoRA_dropout were set to 8 and 0.1, respectively. The LoRA parameters were updated 10 times for each query.

1890 The inference-time learning rate was 0.01 and 0.05 in the experiment in Figure 1a and Figure 1b,
1891 respectively. To prevent information leakage when evaluating the prediction of y_i , the model used
1892 only the preceding data $(\mathbf{x}_1, y_1, \dots, \mathbf{x}_{k-1}, y_{k-1}, \mathbf{x}_k)$ to update the weight. The LoRA weights were
1893 reset before proceeding to the prediction of the next target, y_{i+1} .

1894 For the additional experiment (shown in Figure 2), we reused the TTT results from the in-distribution
1895 setting (Figure 1a) for the $r = d = 16$ case. For the $r = 4, d = 16$ case, the inference-time learning
1896 rate was also set to 0.01, consistent with the $r = 16$ setting.

1897 The test loss was averaged over 256 runs, with each run containing 256 independent queries.
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943