
Test-Time Training Provably Improves Transformers as In-context Learners

Halil Alperen Gozeten^{* 1} Emrullah Ildiz^{* 1} Xuechen Zhang¹ Mahdi Soltanolkotabi² Marco Mondelli³
Samet Oymak¹

Abstract

Test-time training (TTT) methods explicitly update the weights of a model to adapt to the specific test instance, and they have found success in a variety of settings, including most recently language modeling and reasoning. To demystify this success, we investigate a gradient-based TTT algorithm for in-context learning, where we train a transformer model on the in-context demonstrations provided in the test prompt. Specifically, we provide a comprehensive theoretical characterization of linear transformers when the update rule is a single gradient step. Our theory (i) delineates the role of alignment between pretraining distribution and target task, (ii) demystifies how TTT can alleviate distribution shift, and (iii) quantifies the sample complexity of TTT including how it can significantly reduce the eventual sample size required for in-context learning. As our empirical contribution, we study the benefits of TTT for TabPFN, a tabular foundation model. In line with our theory, we demonstrate that TTT significantly reduces the required sample size for tabular classification (3 to 5 times fewer) unlocking substantial inference efficiency with a negligible training cost.

1. Introduction

Modern language models can be viewed as general-purpose, addressing a diverse range of user queries in a zero-shot fashion (Kojima et al., 2022; Wei et al., 2022). However, as the complexity or novelty of the query increases, such as in complex multi-step reasoning scenarios, the pretrained model

may falter. This has motivated two popular approaches: in-context learning and test-time computation. In-context learning (ICL) incorporates demonstrations related to the query as part of the prompt facilitating better inference by the model (Brown et al., 2020; Min et al., 2022). Test-time computation methods explicitly increase the inference-time compute to elicit higher quality responses (Snell et al., 2024; Jaech et al., 2024). An important instance of test-time computation is test-time training (TTT) where the weights of a model are explicitly updated to adapt to specific test instances (Sun et al., 2020; Liu et al., 2021). A gradient-based TTT approach can be described as follows: Given a test prompt and a pretrained sequence model, we update the pretrained weights by performing a few gradient iterations on a suitable self-supervised objective (e.g. next-token prediction objective). We can then use the resulting model to perform the inference on the test prompt. This procedure is applicable beyond language models, and it is conceptually similar to meta-learning approaches such as model-agnostic meta-learning (Finn et al., 2017).

Recently, TTT has found significant success in the context of language modeling by boosting the accuracy and reasoning capability of state-of-the-art models (Akyürek et al., 2024; Sun et al., 2024). This success is in part due to the fact that TTT can be naturally integrated within in-context learning: We can fine-tune the transformer model to fit to the labels or chain-of-thought rationales provided as part of the in-context demonstrations, and then re-apply ICL on the adapted model. In fact, recent work (Akyürek et al., 2024) utilizes a variation of this procedure and additional data augmentation to obtain remarkable improvements in the ARC reasoning benchmark. This motivates the central question of our work:

What are the provable benefits of test-time training of transformers, specifically, for enhancing in-context learning?

Contributions. As our central contribution, we address this question by providing a comprehensive theoretical study of test-time training for one-layer linear transformers. Focusing on prompts following a linear dataset model, we provide

^{*}Equal contribution ¹University of Michigan, Ann Arbor
²University of Southern California ³Institute of Science and Technology Austria. Correspondence to: Halil Alperen Gozeten <alperen@umich.edu>.

a precise risk characterization of TTT with one gradient step update rule. This characterization is established in terms of three ingredients: (i) *context length (number of in-context examples during inference)*, (ii) *target sample size available for TTT* in Section 4, and (iii) *alignment between the pre-trained model and target task* in Section 5. We show that, as sample size increases, TTT can alleviate the distribution shift bottlenecks that arise in standard ICL. This also reveals regimes where TTT with zero or small initialization is preferable to TTT with the pre-trained model (i.e. cold vs. warm start). Interestingly, our experiments with the GPT2 architecture (Radford et al., 2019) show that multi-layer transformers exhibit behavior in line with our theory. Our theory and experiments demonstrate that one step of TTT yields significant performance gains with less computation, consistent with recent empirical observations (Akyürek et al., 2024) that only a few gradient steps offer substantial test-time improvements. Additionally, while standard ICL requires $\Omega(d)$ context length under an isotropic task prior, with d denoting feature dimension, we prove that TTT can succeed with $o(d)$ context length by effectively memorizing the target task. Our technical novelty stems from accurately capturing the statistical benefit of TTT during in-context learning. Specifically, while the transformer is trained on a single prompt, we characterize how the sample complexity benefit of TTT is proportional to the number of target examples within the prompt.

As empirical corroboration of our theory, we explore tabular learning and TabPFN (Hollmann et al., 2023; 2025) – a state-of-the-art tabular foundation model pretrained with structural causal model priors. TabPFN is well-aligned with our theoretical setting with similar token encodings but different prior distributions. An important drawback of TabPFN lies in its inference cost, as it uses a full tabular dataset as context during inference. In line with our theory, we demonstrate that TTT can convert TabPFN into a task-specific tabular model that works equally well with significantly less data (up to 5 times fewer). This in turn implies substantial inference gains given that the complexity of softmax-attention is quadratic in the sequence length.

2. Related Work

We organize our discussion of related work into two main areas: in-context learning and test-time training.

In-context learning. In-context learning has received significant interest during the past few years (Brown et al., 2020; Liu et al., 2023; Agarwal et al., 2024). This has also motivated works toward a stronger theoretical understanding of ICL (Zhang et al., 2024; Mahankali et al., 2024; Ahn et al., 2023b; Xie et al., 2022; Garg et al., 2022; Li et al., 2023). Closer to us, Mahankali et al. (2024); Ahn et al. (2023b); Zhang et al. (2024) study the optimization

landscape of a one-layer linear attention model and show that the optimized model implements a single step of projected gradient descent over the in-context demonstrations. More recent works (Lu et al., 2024; Wu et al., 2023) extend these to characterize the pretraining task sample complexity of ICL. While these works focus on pretraining capabilities, we focus on the adaptation of the pretrained model to the target task. Additionally, rather than empirical risk minimization, we use gradient descent for adaptation and characterize its risk when the sample size is determined by the target prompt length. In our theoretical model, each token represents a data point containing input features and the target label. Remarkably, TabPFN (Hollmann et al., 2023; 2025) shows that, using this simple encoding and pretraining the model with sufficiently rich data priors, transformers can accomplish state-of-the-art classification on tabular datasets via in-context learning. There have been also efforts (Thomas et al., 2024) to fine-tune in-context tabular models like TabPFN at the dataset level using retrieval-based local context selection.

Test-time training. Test-time training (Sun et al., 2020; Liu et al., 2021; Gandelsman et al., 2022) and related test-time adaptation (Wang et al., 2021; Niu et al., 2022; Yuan et al., 2023) methods aim to overcome distribution shift bottlenecks. This is typically accomplished by adapting the model with the test example using self-supervised or unsupervised objectives. These methods can work with just a single text example or admit a minibatch of examples. For sequence/language modeling tasks, one can utilize TTT on the test sequence via the next-token prediction objective (Sun et al., 2024; Hardt & Sun, 2024; Hübötter et al., 2025). Specifically, during in-context learning, the query is unlabeled (e.g. math problem to solve), but we are provided with related examples and associated labels (e.g. through retrieval). Thus, we can fit to these labels as the source of supervision (essentially fine-tuning the model on the dataset of in-context examples). For instance, the training objective in Akyürek et al. (2024) utilizes this approach boosted by additional data augmentation. Finally, the computational efficiency of TTT is an important consideration which motivates our investigation of TTT with a single gradient update.

3. Problem Setup

Notation. Let $[n]$ denote the set $\{1, \dots, n\}$ for an integer $n \geq 1$. We denote vectors and matrices using bold lower-case and upper-case letters, respectively, such as \mathbf{x} for vectors and \mathbf{X} for matrices. The element x_i refers to the i -th entry of the vector \mathbf{x} . We represent the zero vector of size n by $\mathbf{0}_n$ and the zero matrix of size $m \times n$ by $\mathbf{0}_{m \times n}$. The operator $\text{tr}(\mathbf{X})$ represents the trace of \mathbf{X} , \mathbf{X}^\dagger is its Moore–Penrose pseudoinverse, and $\|\mathbf{X}\|_2$ denotes the spectral norm of \mathbf{X} . Given $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, the notation $[\mathbf{x} \ \mathbf{y}] \in \mathbb{R}^{d \times 2}$ represents the

row-wise concatenation, while $[\mathbf{x}; \mathbf{y}] \in \mathbb{R}^{2d}$ represents the column-wise concatenation.

In-context learning and test-time training. In-context learning is the ability of the model to learn from demonstrations in the prompt, i.e. *in-context*. Specifically, given a sequence of demonstrations of desired input/output pairs $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n) \in \mathbb{R}^d \times \mathbb{R}$ followed by a query input \mathbf{x} in the prompt, the model can guess the corresponding output query y . Concretely, we can define the context tokens $\mathbf{z}_i = [\mathbf{x}_i; y_i] \in \mathbb{R}^{d+1}$ for $i \in [n]$ and the query token $\mathbf{z} = [\mathbf{x}; 0] \in \mathbb{R}^{d+1}$. Then the input prompt can be written in the form

$$\mathbf{Z} = [\mathbf{z}_1 \ \dots \ \mathbf{z}_n \ \mathbf{z}]^\top = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_n & \mathbf{x} \\ y_1 & y_2 & \dots & y_n & 0 \end{bmatrix}^\top \in \mathbb{R}^{(n+1) \times (d+1)}.$$

To estimate the output query y , we focus on a sequence model $\text{SM}(\mathbf{Z}, \mathbf{W})$ with \mathbf{Z} the input prompt and \mathbf{W} the model parameters. We assume that the sequence model is *pre-trained* with a token distribution $\mathcal{D}_z^{\text{PT}}$ where $(\mathbf{z}_i)_{i=1}^n, [\mathbf{x}; y] \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_z^{\text{PT}}$ with the corresponding optimal parameters given by

$$\mathbf{W}^* = \arg \min_{\mathbf{W}} \mathbb{E}_{(\mathbf{z}_i)_{i=1}^n, [\mathbf{x}; y] \sim \mathcal{D}_z^{\text{PT}}} [(y - \text{SM}(\mathbf{Z}, \mathbf{W}))^2]. \quad (1)$$

During inference, we test the sequence model on another distribution $\mathcal{D}_z^{\text{TT}}$ and observe k samples of this test distribution $\mathcal{S}_{\text{TT}} = \{(\mathbf{Z}_j, y_j)\}_{j=1}^k$ where y_j is the label of the query token inside \mathbf{Z}_j . The main idea behind *Test-Time Training* (TTT) is to refine the model's parameters using the test data \mathcal{S}_{TT} before performing inference. Concretely, the empirical loss of the sequence model on the test set \mathcal{S}_{TT} with an arbitrary model parameter \mathbf{W} is given by

$$\hat{\mathcal{L}}_{\mathcal{S}_{\text{TT}}}(\mathbf{W}) = \sum_{j=1}^k (y_j - \text{SM}(\mathbf{Z}_j, \mathbf{W}))^2. \quad (2)$$

One can thus refine \mathbf{W}^* by optimizing the above test-time empirical loss. In this paper, we focus on a TTT strategy involving a single gradient descent step over $\hat{\mathcal{L}}_{\mathcal{S}_{\text{TT}}}$, i.e. $\mathbf{W}_{\text{TT}} := \mathbf{W}^* - \eta \nabla \hat{\mathcal{L}}_{\mathcal{S}_{\text{TT}}}(\mathbf{W}^*)$. The error incurred by the model can then be calculated using the population loss via

$$\mathcal{L}(\mathbf{W}) = \mathbb{E}_{(\mathbf{z}_i)_{i=1}^n, [\mathbf{x}; y] \sim \mathcal{D}_z^{\text{TT}}} [(y - \text{SM}(\mathbf{Z}, \mathbf{W}))^2]. \quad (3)$$

Now, we will define the expected loss of the weights \mathbf{W}_{TT} we achieve after test-time training over all test-time training sets \mathcal{S}_{TT} :

$$\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}}) := \mathbb{E}_{\mathcal{S}_{\text{TT}}} [\mathcal{L}(\mathbf{W}_{\text{TT}})]. \quad (4)$$

Our goal is to characterize the loss $\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}})$ obtained after test-time training as a function of the number of context samples n , the number of data points in test-time training

k , the distributions $(\mathcal{D}_z^{\text{PT}}, \mathcal{D}_z^{\text{TT}})$, and the pretrained starting point \mathbf{W}^* .

Architecture. We study the one-layer linear attention model as the sequence model:

$$\text{SM}(\mathbf{Z}, \mathbf{W}) = [\mathbf{Z} \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{Z}^\top \mathbf{Z} \mathbf{W}_V]_{n+1, d+1} = \mathbf{x}^\top \mathbf{W} \mathbf{X}^\top \mathbf{y}, \quad (5)$$

where $[\cdot]_{n+1, d+1}$ denotes the entry $(n+1, d+1)$ of the corresponding matrix. Here, the query, key, and value matrices $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{(d+1) \times (d+1)}$ are defined as

$$\mathbf{W}_Q \mathbf{W}_K^\top = \begin{bmatrix} \mathbf{W} & \mathbf{0}_{d \times 1} \\ \mathbf{0}_{1 \times d} & 0 \end{bmatrix} \quad \mathbf{W}_V = \begin{bmatrix} \mathbf{0}_{d \times d} & \mathbf{0}_{d \times 1} \\ \mathbf{0}_{1 \times d} & 1 \end{bmatrix},$$

and we set the data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ and the corresponding labels $\mathbf{y} \in \mathbb{R}^n$ to be:

$$\mathbf{X} = [\mathbf{x}_1 \ \dots \ \mathbf{x}_n]^\top \quad \mathbf{y} = [y_1 \ \dots \ y_n]^\top.$$

Here, we collapse the query and key matrices by identifying the top-left $d \times d$ block of $\mathbf{W}_Q \mathbf{W}_K^\top \in \mathbb{R}^{(d+1) \times (d+1)}$ with a single matrix $\mathbf{W} \in \mathbb{R}^{d \times d}$, and choose the value matrix to retrieve the output of the query token with one-layer linear attention. We note that similar models have been used in prior work (Zhang et al., 2023; Mahankali et al., 2023; Ahn et al., 2023a; Li et al., 2024; 2025) for the theoretical analysis of a variety of phenomena, albeit not for characterizing the benefits of TTT.

Data model. We consider a linear model with Gaussian data. Specifically, during pre-training, we assume the context tokens $\mathbf{z}_i = [\mathbf{x}_i; y_i]$ and the query input/output vector $[\mathbf{x}; y]$ to be sampled i.i.d. from the distribution $\mathcal{D}_z^{\text{PT}}(\boldsymbol{\Sigma}_\beta)$. Concretely, the outputs are generated according to the linear model $y_i = \mathbf{x}_i^\top \boldsymbol{\beta}_{\text{PT}} + \xi_i$, with task parameter $\boldsymbol{\beta}_{\text{PT}} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_\beta)$, input features $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_x)$, and noise terms $\xi_i \sim \mathcal{N}(0, \sigma^2)$ for $i \in [n]$.

During inference, we test the sequence model on a new task parameter $\boldsymbol{\beta}_{\text{TT}}$ with the input prompts generated following the test-time distribution $\mathcal{D}_z^{\text{TT}}(\boldsymbol{\beta}_{\text{TT}})$. This test-time distribution is governed by a similar linear setting. Concretely, the outputs are generated according to the linear model $y_i = \mathbf{x}_i^\top \boldsymbol{\beta}_{\text{TT}} + \xi_i$, where the input features \mathbf{x}_i are sampled i.i.d. from the same feature distribution $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_x)$ and the noise terms are sampled i.i.d. from the same distribution $\mathcal{N}(0, \sigma^2)$.

We construct the test-time set $\mathcal{S}_{\text{TT}} = \{(\mathbf{Z}_j, y_j)\}_{j=1}^k$ by following the test-time distribution $\mathcal{D}_z^{\text{TT}}(\boldsymbol{\beta}_{\text{TT}})$. We first sample $(n+k)$ query input/output pairs $\{(\bar{\mathbf{x}}_i, \bar{y}_i)\}_{i=1}^{(n+k)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_z^{\text{TT}}(\boldsymbol{\beta}_{\text{TT}})$. Then, we designate the first n samples as fixed context tokens across the set, whereas we use the latter k samples as queries in the set. In formulas,

$$\mathbf{Z}_j = \begin{bmatrix} \bar{\mathbf{x}}_1 & \bar{\mathbf{x}}_2 & \dots & \bar{\mathbf{x}}_n & \bar{\mathbf{x}}_{j+n} \\ \bar{y}_1 & \bar{y}_2 & \dots & \bar{y}_n & 0 \end{bmatrix}^\top, \quad y_j = \bar{y}_{j+n} \quad \forall j \in [k]. \quad (6)$$

Remark: Our procedure can be implemented using a single forward-backward pass by putting all examples within the prompt and using a suitable attention mask to ensure that only the first n examples (context examples with labels) attend the last k examples (query examples without labels) and vice versa. This allows for efficient parallel training. This also means that we use a fixed context-query split where the same n examples are used as context during TTT. This can be generalized to a K-fold procedure where all examples are used as both context and query during TTT (as in Akyürek et al. (2024)); however, we opted for the 1-fold option, which is more amenable to a precise statistical analysis.

Single-step GD for TTT. We now turn our attention to deriving the single-step gradient update over the test set \mathcal{S}_{TT} . To this aim, we denote by $\mathbf{X}_{\text{context}}$ and $\mathbf{y}_{\text{context}}$ the context inputs and their outputs, whereas we denote by $\mathbf{X}_{\text{train}}$ and $\mathbf{y}_{\text{train}}$ the training query inputs and their outputs:

$$\begin{aligned}\mathbf{X}_{\text{context}} &= [\bar{\mathbf{x}}_1 \dots \bar{\mathbf{x}}_n]^\top, & \mathbf{y}_{\text{context}} &= [\bar{\mathbf{y}}_1 \dots \bar{\mathbf{y}}_n]^\top, \\ \mathbf{X}_{\text{train}} &= [\bar{\mathbf{x}}_{n+1} \dots \bar{\mathbf{x}}_{n+k}]^\top, & \mathbf{y}_{\text{train}} &= [\bar{\mathbf{y}}_{n+1} \dots \bar{\mathbf{y}}_{n+k}]^\top.\end{aligned}$$

The following proposition (proved in Appendix A) characterizes the single-step GD TTT update.

Proposition 3.1. *Consider the linear attention model with parameters $\mathbf{W} \in \mathbb{R}^{d \times d}$. Suppose the test-time training loss function is defined as in (2) and define $\mathbf{u}_{\text{context}} := \mathbf{X}_{\text{context}}^\top \mathbf{y}_{\text{context}} \in \mathbb{R}^d$. Then, for any step size $\eta > 0$, the new parameter \mathbf{W}_{TT} after one gradient-descent step from \mathbf{W} is given by the rank-1 update*

$$\mathbf{W}_{TT} = \mathbf{W} + 2\eta \mathbf{X}_{\text{train}}^\top (\mathbf{y}_{\text{train}} - \mathbf{X}_{\text{train}} \mathbf{W} \mathbf{u}_{\text{context}}) \mathbf{u}_{\text{context}}^\top.$$

In the coming sections, we will start our analysis when the covariance matrices (Σ_β, Σ_x) are isotropic (Section 4), and then we provide an extension to more general covariance matrices (Section 5). Finally, we corroborate our theoretical results with empirical evidence (Section 6).

4. Analysis for Isotropic Covariances

In this section, we study the loss $\mathcal{L}_{TT}(\mathbf{W}_{TT})$ induced by a single-step gradient update in the isotropic scenario $(\Sigma_\beta, \Sigma_x) = (\mathbf{I}, \mathbf{I})$, for an arbitrary test-time parameter β_{TT} . The assumption of isotropic covariance matrices enables us to explicitly characterize the behavior of the loss as a function of n, d , and k . We first analyze how $\mathcal{L}_{TT}(\mathbf{W}_{TT})$ varies with the number of in-context examples (n) and the embedding dimension (d). We then compare two distinct choices of \mathbf{W}^* , as a function of the test-time training set size k : (i) \mathbf{W}^* is obtained from (1), and (ii) $\mathbf{W}^* = \mathbf{0}_{d \times d}$. The former represents a test-time training approach, whereas the latter corresponds to training from scratch with the single-step gradient update.

We start with the characterization of \mathbf{W}^* and its loss $\mathcal{L}(\mathbf{W}^*)$.

Proposition 4.1. *Let $\Sigma_\beta, \Sigma_x = \mathbf{I}, \sigma^2 = 0$ and let β_{TT} be an arbitrary new task parameter. Then, the optimal pre-trained model’s parameter \mathbf{W}^* and its population loss are given by*

$$\mathbf{W}^* = \frac{1}{n+d+1} \mathbf{I}, \quad \mathcal{L}(\mathbf{W}^*) = \|\beta_{TT}\|^2 \frac{d+1}{n+d+1}.$$

The characterization of \mathbf{W}^* is obtained from Li et al. (2024, Corollary 1) and its loss with respect to the new task parameter β_{TT} is provided in Appendix B for the more general setting with noise term σ^2 . Due to the technical difficulty of analyzing the single-step gradient in the noisy setting and to have manageable expressions, we will focus on the noiseless setting.

Given k, n , and d , the loss of the pre-trained parameters, $\hat{\mathcal{L}}_{\mathcal{S}_{TT}}(\mathbf{W}^*)$, is a random variable with respect to the test-time set \mathcal{S}_{TT} . Consequently, the loss $\mathcal{L}(\mathbf{W}_{TT})$ is also a random variable dependent on \mathcal{S}_{TT} . When applying a single-step update to the pre-trained model’s parameters, we define the expectation of the loss $\mathcal{L}(\mathbf{W}_{TT})$ with respect to \mathcal{S}_{TT} in (4), treating it as a function of the step size η . The optimal step size is then selected to minimize the expected test-time training loss. In the following theorem (proved in Appendix B), we present the optimal learning rate and the corresponding improvement in the loss induced by test-time training.

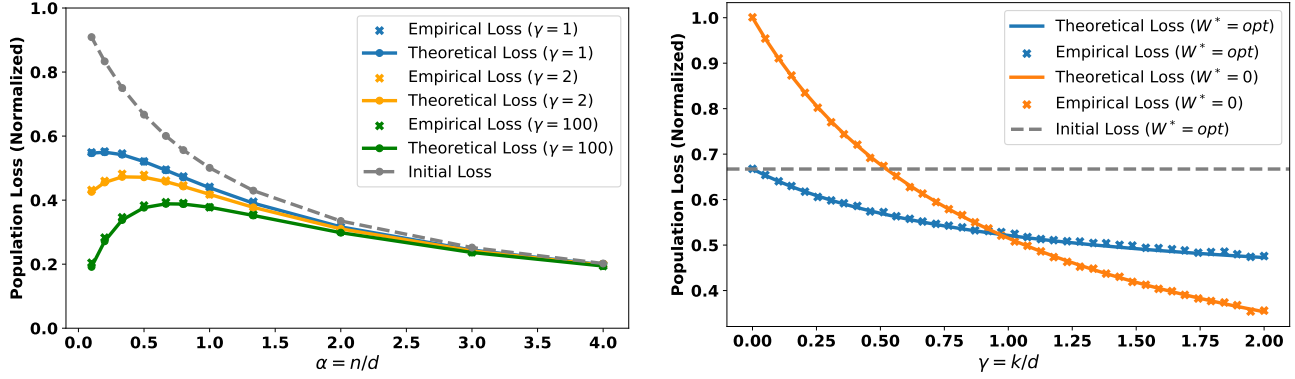
Theorem 4.2. *Let $n/d = \Theta(1)$. Recall the definition of the expected loss $\mathcal{L}_{TT}(\mathbf{W}_{TT})$ with respect to \mathcal{S}_{TT} in (4). In the isotropic covariance and noiseless setting ($\sigma^2 = 0$), the optimal step-size that minimizes $\mathcal{L}_{TT}(\mathbf{W}_{TT})$ is*

$$\eta^* \approx \frac{d}{2(k+d)n^2(n+d)\|\beta_{TT}\|_2^2}.$$

With this optimal step-size η^* , the improvement in the loss due to the test-time training is

$$\mathcal{L}(\mathbf{W}^*) - \mathcal{L}_{TT}(\mathbf{W}_{TT}) \approx \frac{k}{k+d} \frac{d^3}{(n+d)^3} \|\beta_{TT}\|_2^2.$$

Remark 4.3. We note that in the proportional n, d regime, the Gaussian approximation in Lemma B.1, which is utilized in the proof of Theorem 4.2, introduces an error of $O(n^{-1})$. Additionally, during the derivation, we omit lower-order terms (such as $2n$ in $n^2 + 2n$), which introduces errors of one degree lower relative to the leading terms. Consequently, after carrying these approximations through the entire derivation, the overall approximation error remains at most $O(n^{-1})$, since the final loss expression is zeroth-order in n, d . Hence, our result is robust to these errors as $n, d \rightarrow \infty$ while maintaining the ratio $\alpha = n/d$. The same reasoning also applies to the non-isotropic covariance case, which will be discussed in Section 5.



(a) Non-monotonic behavior of the loss $\mathcal{L}(\mathbf{W}_{\text{TT}})$ after the test-time-training update with respect to $\alpha = n/d$ for various $\gamma = k/d$ values.

(b) Losses $\mathcal{L}(\mathbf{W}_{\text{TT}})$ after the test-time-training update as a function of $\gamma = k/d$ when using optimal pre-trained weights vs. $\mathbf{0}_{d \times d}$ (null).

Figure 1. Plot of the normalized population losses after test-time training when $\Sigma_x = \Sigma_\beta = \mathbf{I}$, $\sigma^2 = 0$, and the new task is $\beta_{\text{TT}} = \mathbf{1}_d$. Solid lines denote theoretical predictions which match the empirical (markers) results. (a): The figure illustrates a non-monotonic trend as the ratio $\frac{n}{d}$ shifts, see Corollary 4.4. **Setting:** $n = 300$; the ratio $\frac{n}{d}$ changes between 0.1 and 4; the three lines depict $\gamma = 1$, $\gamma = 2$, and $\gamma = 100$ where $\gamma = k/d$; the step size η is selected optimally as per Theorem 4.2. (b): The figure reveals a threshold in k at which the preferable initialization switches from the optimal pre-trained W^* to the null model $\mathbf{0}_{d \times d}$ as k increases. The transition threshold k aligns well with Corollary 4.6. **Setting:** $n = 200$; $d = 400$; k changing between 0 and $4n$ with equal increments.

Theorem 4.2 establishes the improvement produced by the test-time training as a function of k, n , and d . When n and d are fixed, such improvement is proportional to $\frac{k}{k+d}$. Additionally, if $k/d = \gamma$ is fixed, the loss $\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}})$ exhibits an intriguing behavior as a function of $n/d = \alpha$, as described by the next corollary (proved in Appendix B).

Corollary 4.4. Recall the definitions of $\gamma = k/d$, $\alpha = n/d$, and consider the loss $\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}})$ as a function of α in the isotropic covariance and noiseless setting. If $\gamma > \frac{1}{2}$, then the loss $\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}})$ is non-monotonic in α . Specifically, the loss $\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}})$ is increasing for $\alpha < \sqrt{\frac{3\gamma}{\gamma+1}} - 1$ and decreasing for $\alpha > \sqrt{\frac{3\gamma}{\gamma+1}} - 1$. Conversely, if $\gamma < \frac{1}{2}$, then $\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}})$ is monotonic in α .

The phase transition points described in Corollary 4.4 are approximate but become exact as k, n , and d approach infinity while maintaining the ratios $\gamma = k/d$ and $\alpha = n/d$. The non-monotonic behavior identified in Corollary 4.4 is observed as a result of two opposing effects, which are the initial loss and the improvement by TTT. As n grows against d (i.e. α increases), the pre-trained model does better initially and already has a lower loss before TTT, which makes it harder to be further reduced by TTT as the rank-1 update is unable to correct all directions. On the other hand, when d grows against n (α decreases), the initial loss is high, providing more room (error) to be corrected by rank-1 update, and thus, the improvement by TTT is larger. This intuition aligns with Theorem 4.2, which establishes that the TTT improvement scales as $(\frac{d}{n+d})^3$. Together, these two trends result in the non-monotonic behavior.

We plot the behavior of the loss $\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}})$ as a function of $\alpha = n/d$ for various values of $\gamma = k/d$ in Figure 1a, making the following three observations: (i) As stated in Theorem 4.2, the improvement achieved through test-time training is approximately proportional to $1/(\alpha + 1)^3$ for large n and d . This behavior is evident in Figure 1a, which shows that for every value of γ , the improvement diminishes as $\alpha = n/d$ increases. (ii) The non-monotonic behavior described Corollary 4.4 is clearly observed in Figure 1a. The increasing/decreasing regions in the loss as a function of α for different values of γ are consistent with the theoretical results. (iii) As γ increases, the improvement from the gray curve by test-time training is proportional to the ratio $\gamma/(\gamma + 1)$ for a fixed $\alpha = n/d$. This observation aligns with Theorem 4.2.

In addition to the scenario where we begin from an optimally pre-trained model W^* , a natural question is how much improvement can be achieved when the initial weight matrix $W^* = \mathbf{0}_{d \times d}$ (zero initialization). This initialization corresponds to training the model from scratch using a single-step gradient descent. In the following theorem, we quantify the improvement gained by applying a single-step gradient descent from this initialization.

Theorem 4.5. Consider the isotropic covariance and noiseless setting ($\sigma^2 = 0$). Suppose the initial weight matrix is $W^* = \mathbf{0}_{d \times d}$. Then, the optimal step-size that minimizes $\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}})$ is

$$\eta^* = \frac{1}{2(k+d+1)(n^2+4n+3+d)\|\beta_{\text{TT}}\|_2^2}.$$

With this optimal step-size η^* , the improvement in the loss

due to test-time training is

$$\mathcal{L}(\mathbf{W}^*) - \mathcal{L}_{TT}(\mathbf{W}_{TT}) = \frac{k}{k+d+1} \frac{n^2}{n^2+4n+3+d} \|\boldsymbol{\beta}_{TT}\|_2^2,$$

where $\mathcal{L}(\mathbf{W}^*) = \|\boldsymbol{\beta}_{TT}\|_2^2$.

We provide the proof of Theorem 4.5 in Appendix B, where we also solve the more general noisy setting with initial weight $\mathbf{W}^* = \mathbf{0}_{d \times d}$. If $\sigma^2 = O(\|\boldsymbol{\beta}_{TT}\|_2^2)$ and k grows sufficiently faster than d (e.g., $k/d \rightarrow \infty$), then the test-time training update can reduce the loss to near 0 as $k, n, d \rightarrow \infty$ with $\alpha = n/d = \Theta(1)$.

Armed with Theorem 4.2 and 4.5, we can now establish when it is better to initialize with the pre-trained \mathbf{W}^* , as opposed to the zero (null) initialization, as the size k of the test-time training set varies.

Corollary 4.6. *Recall the definitions of $\alpha = n/d$ and $\gamma = k/d$. Consider the setting in Theorem 4.2 and test-time training described in Proposition 3.1 with both pre-trained and zero initializations. Then, under the isotropic covariance and noiseless setting, there exists a threshold γ^* given by $\gamma^* \approx \frac{(\alpha+1)^2}{(\alpha+2)}$ such that $\gamma < \gamma^*$ if and only if it is better to utilize the pre-trained initialization over the zero initialization $\mathbf{0}_{d \times d}$.*

Corollary 4.6 identifies a phase transition point as a function of γ and α that distinguishes the region where test-time training outperforms training from scratch. In Figure 1b, we illustrate the loss $\mathcal{L}_{TT}(\mathbf{W}_{TT})$ as a function of $\gamma = k/d$ for pre-trained and zero initializations, with a fixed $\alpha = n/d = 1/2$. Based on Corollary 4.6, the phase transition point for γ is $\frac{(3/2)^2}{5/2} = 9/10$, which aligns with the empirical results.

The pre-trained initialization provides a significant advantage when the test-time training set is small, as it introduces a strong prior that facilitates better performance. However, as the test-time training set grows, this initialization becomes a limitation because the rank-one update to the pre-trained matrix is insufficient to achieve a near-zero loss. In contrast, with zero initialization, the rank-one update becomes highly effective when the test-time training set is sufficiently large, enabling it to achieve near-zero error. This demonstrates that while pre-trained initialization is beneficial in data-scarce scenarios, zero initialization is better suited for scenarios with sufficient test-time training data.

5. Analysis for General Covariance

In this section, we extend the previous analysis to the general case where $\boldsymbol{\Sigma}_x$ and $\boldsymbol{\Sigma}_\beta$ are any jointly diagonalizable covariances. Specifically, we assume the task covariance $\boldsymbol{\Sigma}_\beta$ to be a non-isotropic diagonal matrix while keeping the

feature covariance matrix $\boldsymbol{\Sigma}_x$ isotropic.¹ This assumption enables us to characterize the loss function $\mathcal{L}_{TT}(\mathbf{W}_{TT})$ with respect to the alignment between the pre-trained initialization matrix \mathbf{W}^* and the new task parameter $\boldsymbol{\beta}_{TT}$.

We start with the definition of two quantities, which will be used to express the optimal step size and the loss improvement via the test-time training.

Definition 5.1. For an arbitrary $\mathbf{W} \in \mathbb{R}^{d \times d}$, define $A = \boldsymbol{\beta}_{TT}^\top (\mathbf{I} - n\mathbf{W})^2 \boldsymbol{\beta}_{TT}$ and $B = n \|\boldsymbol{\beta}_{TT}\|_2^2 \|\mathbf{W}\|_F^2$.

Here, the term A represents the *misalignment* between the test-time task parameter $\boldsymbol{\beta}_{TT}$ and the initial parameter \mathbf{W} , whereas the term B represents the total signal power of the one-layer linear attention system with parameter \mathbf{W} . Note that the terms A and B are a function of the new task parameter $\boldsymbol{\beta}_{TT}$ and the initial weight parameter \mathbf{W} .

Using the definitions of A and B , we are ready to generalize Proposition 4.1 to the non-isotropic and rank-deficient covariance matrices. Note that for the rank-deficient covariance matrix, we take the optimal \mathbf{W}^* as the minimum Frobenius-norm matrix that minimizes the loss in (1).

Proposition 5.2. *Let $\boldsymbol{\Sigma}_x = \mathbf{I}$, $\sigma^2 = 0$ and suppose $\boldsymbol{\Sigma}_\beta$ is an arbitrary diagonal covariance matrix with the first r diagonal entries being non-zero. Then, the minimizer \mathbf{W}^* of the pre-training loss (1) that has minimal Frobenius norm and its population loss are given by*

$$\mathbf{W}^* = ((n+1)\mathbf{I}' + \text{tr}(\boldsymbol{\Sigma}_\beta)\boldsymbol{\Sigma}_\beta^\dagger)^\dagger, \quad \mathcal{L}(\mathbf{W}^*) \approx A + B,$$

where $\mathbf{I}' = \text{diag}(\mathbf{1}_r, \mathbf{0}_{d-r})$.

We generalize the characterization of \mathbf{W}^* obtained from Li et al. (2024, Corollary 1) to rank-deficient covariance matrices, and its loss with respect to the new task parameter $\boldsymbol{\beta}_{TT}$ is provided in Appendix C.

Note that the eigenvalues of \mathbf{W}^* are smaller than $1/(n+1)$ as $\boldsymbol{\Sigma}_\beta$ is a positive semi-definite matrix. This condition ensures the stability of the linear attention model described in (5) as the magnitude of the SM output scales with the context length n when \mathbf{W} is fixed. Assuming a similar criterion for the set of \mathbf{W} , we provide the optimal step size parameter and the corresponding test-time training loss.

Theorem 5.3. *Let $n/d = \Theta(1)$, $\boldsymbol{\Sigma}_x = \mathbf{I}$, and $\sigma^2 = 0$. Suppose that the initial weight matrix \mathbf{W} is a diagonal matrix whose eigenvalues are in the interval $[0, \frac{1}{n+1}]$.² Then, the*

¹This assumption is equivalent to a more general case where $\boldsymbol{\Sigma}_x$ and $\boldsymbol{\Sigma}_\beta$ are jointly diagonalizable, as proved in Appendix C. Joint diagonalizability refers to the case where there exists a unitary matrix $\mathbf{Q} \in \mathbb{R}^{d \times d}$ such that the covariance matrices $\boldsymbol{\Sigma}_\beta = \mathbf{Q}\boldsymbol{\Lambda}_\beta\mathbf{Q}^\top$ and $\boldsymbol{\Sigma}_x = \mathbf{Q}\boldsymbol{\Lambda}_x\mathbf{Q}^\top$ can be expressed with diagonal $\boldsymbol{\Lambda}_\beta$ and $\boldsymbol{\Lambda}_x$ matrices.

²In the joint diagonalizability case, the restriction on \mathbf{W} covers the space of \mathbf{W} that is jointly diagonal with $\boldsymbol{\Sigma}_\beta$ and $\boldsymbol{\Sigma}_x$.

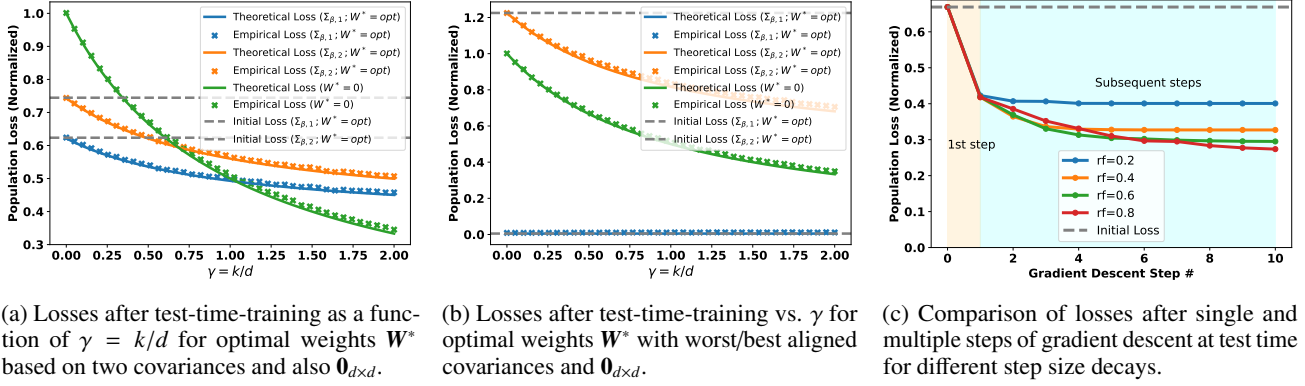


Figure 2. Population losses in single and multi-step update scenarios. $\Sigma_x = I$, $\sigma^2 = 0$ and step sizes are optimal based on the formula in Theorem 5.3. $W^* = \text{opt}$ and $W^* = 0$ represent the pre-trained matrix in (1) and zero (null) initialization, respectively. **(a): Setting:** $n = 300$; $d = 600$; k changing between 0 and $4n$ with equal increments; $\Sigma_{\beta,1} = \text{diag}(I_{250}, 0.5 \cdot I_{250})$; $\Sigma_{\beta,2} = \text{diag}(0.5 \cdot I_{250}, I_{250})$; $\beta_{TT} = [I_{250}; 0.5 \cdot I_{250}]$. Solid lines denote theoretical predictions which match the empirical (markers) results. **(b): Setting:** $n = 250$; $d = 500$; k changing between 0 and $4n$ with equal increments; $\Sigma_{\beta,1} = \text{diag}(1, 0_{499})$; $\Sigma_{\beta,2} = \text{diag}(0, I_{499})$; $\beta_{TT} = [1; 0_{499}]$. **(c):** The figure illustrates the empirical results, which imply that with the initially optimal step size, a significant improvement can be gained by a single gradient step. Each line depicts a different reduce factor on the step size η after each gradient step. **Setting:** $n = 50$; $d = 100$; $k = 50 \times d$; $\Sigma_\beta = I$; $\beta_{TT} = I_{100}$; W^* is optimal.

optimal step size that minimizes the population loss given in (3) after the test-time training update is

$$\eta^* \approx \frac{A}{2(k+d)n^2 \|\beta_{TT}\|_2^2 (A+B)}.$$

With this optimal step-size η^* , the improvement in the loss due to test-time training and the initial loss are approximately

$$\mathcal{L}(W) - \mathcal{L}_{TT}(W_{TT}) \approx \frac{k}{k+d} \frac{A^2}{A+B}, \quad \mathcal{L}(W) \approx A+B.$$

A more general and equivalent version of this result, which applies to any pair of jointly diagonalizable (not necessarily diagonal) covariance matrices Σ_x and Σ_β and to any W that is also jointly diagonalizable with them, can be found in Appendix C. The above approximations are robust and become exact as $n, d \rightarrow \infty$ while having a fixed ratio n/d ; see Remark 4.3. Theorem 5.3 characterizes the improvement due to test-time training for a rather general class of W . It shows that, for fixed $\|W\|_F^2$, the loss improvement increases monotonically with A . In other words, when the misalignment measure A is larger, test-time training achieves greater performance gains. Consequently, considering the optimal pre-trained W^* , if the eigenvalue spectrum of Σ_β aligns with the entries of the task β_{TT} (i.e., larger eigenvalues match larger entries), then A is larger and thus yields a bigger improvement.

Furthermore, we note that the optimal $\mathcal{L}(W)$ over all $W \in \mathbb{R}^{d \times d}$ is $O(n^{-1})$ by Lemma C.1. This optimal loss behavior can be achieved through test-time training with the initial

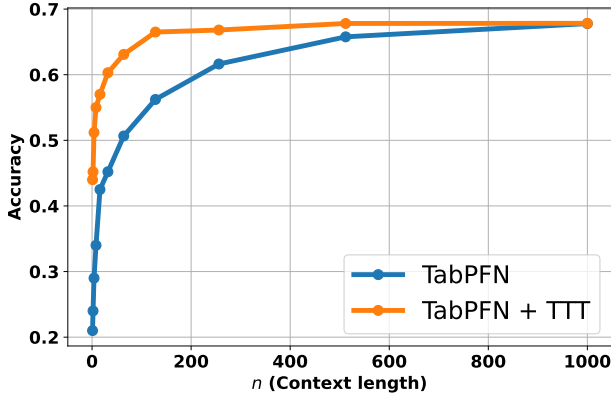
parameter satisfying $n\|W\|_2 < 1 - \delta(n, d)$ for some fixed $\delta(n, d) > 0$ as k/d approaches infinity by Theorem 5.3. This means that test-time training achieves the optimal solution when we have zero or small initialization of W as $k/d \rightarrow \infty$.

For the remainder of this section, we focus on the optimal pre-trained W^* constructed in Proposition 5.2, i.e. the minimum-Frobenius-norm solution of the pre-training loss (1). When we initialize the pre-trained model weight as W^* , the loss $\mathcal{L}_{TT}(W_{TT})$ can be computed by combining Proposition 5.2 and Theorem 5.3. When we initialize the weights as $W = 0_{d \times d}$, then the loss $\mathcal{L}_{TT}(W_{TT})$ is obtained by using the fact that $\mathcal{L}(0_{d \times d}) = \|\beta_{TT}\|_2^2$. By combining these results, we identify the phase transition point in the non-isotropic task covariance setting, which determines when test-time training outperforms training from scratch.

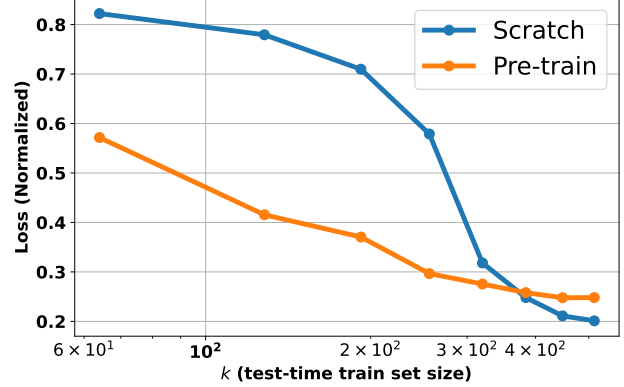
Corollary 5.4. Consider the setting in Theorem 5.3. Recall the definitions of $\alpha = n/d$ and $\gamma = k/d$ and define $c_1 = \frac{A}{\|\beta_{TT}\|_2^2}$, $c_2 = \frac{B}{\|\beta_{TT}\|_2^2}$. Then, there exists a phase transition point $\gamma^* \approx \frac{(c_1 + c_2) - (c_1 + c_2)^2}{c_2(2c_1 + c_2)}$ such that $\gamma < \gamma^*$ if and only if it is better to utilize the pre-trained W^* over the null initialization $0_{d \times d}$.

In a similar way as Corollary 4.4, the phase transition point mentioned above is approximate for finite k, n , and d , and it becomes exact when k, n , and d approach infinity while keeping the ratios α and γ fixed.

We illustrate the findings of Corollary 5.4 in Figures 2a and 2b and make the following three observations: (i) The



(a) Accuracy of TabPFN v2 model with and without test-time-training as a function of number of in-context samples n .



(b) Losses after test-time-training for GPT2 comparing pre-trained and scratch models as a function of test-time training set size k .

Figure 3. (a): Accuracy of TabPFN v2 with and without test-time training. **Setting:** $k = 1000$; n changing between 1 and 1000. (b): Normalized loss after test-time-training with pre-training and zero initialization. **Setting:** $d = 60$; $n = 40$; k changing between 64 and 512 in increments of 64; $\sigma^2 = 0.01$; $\Sigma_\beta = \text{diag}(0.1 \cdot I_{25}, 0.5 \cdot I_{10}, I_{25})$; $\Sigma_x = I$. We sample β_{TT} from the distribution $\mathcal{N}(0, \Sigma_{\beta_{\text{TT}}})$ where $\Sigma_{\beta_{\text{TT}}} = \text{diag}(I_{25}, 0.5 \cdot I_{10}, 0.1 \cdot I_{25})$.

phase transition point is monotonic as a function of c_1 , while keeping c_2 fixed, which implies that it is monotonic as a function of the misalignment term A . This means that, as we increase the alignment between the task covariance Σ_β and the new task parameter β_{TT} , the crossing point of the green and blue curves occurs later than that of the green and orange curves, as observed in Figure 2a. (ii) The worst-aligned case over the set of all possible $(\Sigma_\beta, \beta_{\text{TT}})$ pairs is the one that satisfies $\beta_{\text{TT}}^\top \Sigma_\beta \beta_{\text{TT}} = 0$. In this case, the corresponding c_1 is equal to 1, which implies that the phase transition point in Corollary 5.4 is less than 0 as $c_2 > 0$. This means that training from scratch is always better than utilizing the test-time training for the worst-aligned case, which is depicted in Figure 2b as the orange curve. (iii) The best-aligned case over the set of all possible $(\Sigma_\beta, \beta_{\text{TT}})$ pairs is the one that satisfies $\beta_{\text{TT}}^\top \Sigma_\beta \beta_{\text{TT}} / \text{tr}(\Sigma_\beta) = \|\beta_{\text{TT}}\|_2^2$. In this case, the corresponding c_1 scales on the order of n^{-2} and c_2 on the order of n^{-1} , which in turn makes the phase transition point from Corollary 5.4 approximately n . When k, n , and d all grow while maintaining the ratios $\alpha = n/d$ and $\gamma = k/d$ fixed, this phase transition point approaches infinity. As a result, test-time training is always better than training from scratch in the best-aligned case, which is depicted in Figure 2b as the blue curve.

6. Experiments

In this section, we provide empirical results in light of our theoretical insights from previous sections. The first discussion of this section focuses on whether further computation via multi-step gradient descent at test time yields significant improvements beyond a single-step update. Then, we focus

on tabular learning by demonstrating how test-time training reduces context length for in-context learning on TabPFN. Finally, we demonstrate how the pre-trained model behaves against the scratch model under the test-time-training update for the GPT-2 model in the presence of a distribution shift.

6.1. Multi-Step Gradient Updates

In Figure 2c, we investigate the benefits of multiple gradient steps by choosing the optimal single-step size η based on Theorem 4.2, and then applying different decay factors on step size after each subsequent step. We note that reducing the step size each time is necessary to ensure continued improvement of the loss. The results in Figure 2c confirm that a single-step test-time-training update can capture significant improvement, whereas additional steps yield diminishing returns compared to the initial step. Consequently, the single-step gradient update offers a favorable trade-off with less computation yet with the performance near multi-step finetuning.

6.2. Experiments on TabPFN

We demonstrate that test-time training (TTT) can significantly reduce the context length required for a given performance on TabPFN v2 (Hollmann et al., 2025). Specifically, we evaluate the TabPFN v2 model on the The Tremendous TabLib Trawl (T4) dataset (Gardner et al., 2024) for a more comprehensive evaluation, which is a large-scale high-quality collection of tabular benchmarks. Following the official TabPFN v2 implementation and our theoretical setup, we select the datasets containing at least 1,250 samples (with 1,000 for training, using an 80–20 split), limit

the number of classes to 10 by choosing the most frequent ones, and restrict datasets to 100 randomly selected features. We also convert regression tasks into 10-class classification tasks based on quantiles to maintain consistency across training and evaluation. For each selected dataset from T4, we evaluate TabPFN v2 with context length n varying from 1 to 1000. For the TabPFN (blue curve in Figure 3a), we directly load the pre-trained model and vary the context window length during evaluation. In contrast, for TabPFN+TTT (orange curve in Figure 3a), we finetune the model using different context lengths with $k = 1000$ samples. As the context length n decreases, the samples are divided into $1000/n$ groups, where each group undergoes 50 training iterations. We then report the average performance across all datasets with enough training samples.

In the previous sections, we have theoretically shown that TTT reduces the required in-context sample size for queries from a new task by performing only a single adaptation step, thus enhancing the efficiency of inference. As empirical corroboration of this, Figure 3a illustrates that TabPFN with test-time-training allows the model with 200 in-context samples to almost match the performance of the TabPFN without test-time-training with 1000 in-context samples. This means that test-time training reduces the number of required samples for tabular classification by about 5 times. It is worth emphasizing that the sample complexity benefit is more evident near the zero-shot regime as TTT helps the model memorize new task dynamics by improving the accuracy from 0.2 to 0.45. A log-scaled version of Figure 3a is provided in Appendix D (Figure 4) for a clearer view of performance gains across different scales of n .

6.3. Experiments on GPT-2

We demonstrate how the pre-trained model compares against the scratch model under test-time training with distribution shift using GPT-2 architecture (Radford et al., 2019), as shown in Figure 3b. Our results are in agreement with Figure 1b.

Throughout the experiments, we consider the linear model with Gaussian data following Section 3. In order to obtain a pre-trained model, we sample multiple tasks from the distribution with task covariance Σ_β and train the model from scratch on these tasks until convergence. In contrast, in the scratch setting, the model is initialized randomly using the GPT-2 architecture without any pre-training. During test-time training, we evaluate the model on new task parameters sampled from the distribution $\Sigma_{\beta_{\text{TT}}}$, where the covariance structure is provided in the caption of Figure 3b. We then apply test-time training to both models by varying the test-time training set size and updating all parameters of the GPT-2 model. The results are averaged over all the new tasks sampled from $\Sigma_{\beta_{\text{TT}}}$. Finally, we repeat the entire procedure

three times using different random seeds and report the averaged results.

The results in Figure 3b show that using a pre-trained model with TTT is more advantageous when the test-time training set size is small. However, as the training set size increases, training from scratch outperforms the pre-trained model, aligning with Figure 1b.

7. Discussion

We have developed a theoretical framework to characterize how a single-step gradient update at test time enhances in-context learning. In the case of an isotropic covariance matrix, we analyzed the improvement in loss under test-time training as a function of the number of in-context samples (n), embedding dimension (d), and test-time training set size (k). We then extended our analysis to the non-isotropic covariance case, examining how the alignment between the new task parameter (β_{TT}) and the task covariance matrix (Σ_β) affects performance. Our results reveal a phase-transition threshold on k in both isotropic and non-isotropic settings, beyond which zero initialization (training from scratch) can outperform pre-trained initialization. Empirical results align well with our theoretical findings, and together they underscore that single-step test-time training offers significant performance improvements at low computational costs. Overall, our findings highlight test-time training as a lightweight, supervised enhancement to standard in-context learning.

As a future direction, one can investigate the analysis of TTT under a more general architecture and data model. Possible extensions to the architecture model might be to analyze the TTT with softmax attention or multilayer attention, whereas possible extensions to the data model might be to analyze K-Fold cross-validation instead of the current 1-Fold validation approach. Additionally, investigating TTT in unsupervised or semi-supervised in-context learning (ICL) settings presents another promising direction for extending this work. We hope that our results will serve as a foundation for future research on the interaction of TTT and ICL methods.

Impact Statement

This paper presents work on test-time training to provide a stronger understanding of the principles of modern AI models. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

Acknowledgements

H.A.G., M.E.I., X.Z., and S.O. were supported in part by the NSF grants CCF2046816, CCF-2403075, CCF-2008020, and the Office of Naval Research grant N000142412289. M. M. is funded by the European Union (ERC, INF², project number 101161364). Views and opinions expressed are, however, those of the author(s) only and do not necessarily reflect those of the European Union or the European Research Council Executive Agency. Neither the European Union nor the granting authority can be held responsible for them. M.S. is supported by the Packard Fellowship in Science and Engineering, a Sloan Research Fellowship in Mathematics, an NSF-CAREER under award #1846369, DARPA FastNICS program, and NSF-CIF awards #1813877 and #2008443, and NIH DP2LM014564-01. The authors also acknowledge further support from Open Philanthropy, OpenAI, Amazon Research, Google Research, and Microsoft Research.

References

- Agarwal, R., Singh, A., Zhang, L. M., Bohnet, B., Rosias, L., Chan, S., Zhang, B., Anand, A., Abbas, Z., Nova, A., et al. Many-shot in-context learning. *arXiv preprint arXiv:2404.11018*, 2024.
- Ahn, K., Cheng, X., Daneshmand, H., and Sra, S. Transformers learn to implement preconditioned gradient descent for in-context learning, 2023a. URL <https://arxiv.org/abs/2306.00297>.
- Ahn, K., Cheng, X., Daneshmand, H., and Sra, S. Transformers learn to implement preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing Systems*, 36, 2023b.
- Akyürek, E., Damani, M., Qiu, L., Guo, H., Kim, Y., and Andreas, J. The surprising effectiveness of test-time training for abstract reasoning. *arXiv preprint arXiv:2411.07279*, 2024.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.
- Chandrasekaran, V., Recht, B., Parrilo, P. A., and Willsky, A. S. The convex geometry of linear inverse problems. *Foundations of Computational Mathematics*, 12 (6):805–849, October 2012. ISSN 1615-3383. doi: 10.1007/s10208-012-9135-7. URL <http://dx.doi.org/10.1007/s10208-012-9135-7>.
- Finn, C., Abbeel, P., and Levine, S. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pp. 1126–1135. PMLR, 2017.
- Gandelsman, Y., Sun, Y., Chen, X., and Efros, A. Test-time training with masked autoencoders. *Advances in Neural Information Processing Systems*, 35:29374–29385, 2022.
- Gardner, J., Perdomo, J. C., and Schmidt, L. Large scale transfer learning for tabular data via language modeling. *arXiv preprint arXiv:2406.12031*, 2024.
- Garg, S., Tsipras, D., Liang, P. S., and Valiant, G. What can transformers learn in-context? a case study of simple function classes. *Advances in Neural Information Processing Systems*, 35:30583–30598, 2022.
- Hardt, M. and Sun, Y. Test-time training on nearest neighbors for large language models, 2024. URL <https://arxiv.org/abs/2305.18466>.
- Hollmann, N., Müller, S., Eggensperger, K., and Hutter, F. TabPFN: A transformer that solves small tabular classification problems in a second, 2023. URL <https://arxiv.org/abs/2207.01848>.
- Hollmann, N., Müller, S., Purucker, L., Krishnakumar, A., Körfer, M., Hoo, S. B., Schirrmeyer, R. T., and Hutter, F. Accurate predictions on small data with a tabular foundation model. *Nature*, 637(8045):319–326, 2025. doi: 10.1038/s41586-024-08328-6. URL <https://doi.org/10.1038/s41586-024-08328-6>.
- Hübner, J., Bongni, S., Hakimi, I., and Krause, A. Efficiently learning at test-time: Active fine-tuning of llms, 2025. URL <https://arxiv.org/abs/2410.08020>.
- Jaech, A., Kalai, A., Lerer, A., Richardson, A., El-Kishky, A., Low, A., Helyar, A., Madry, A., Beutel, A., Carney, A., et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa, Y. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35: 22199–22213, 2022.
- Li, Y., Ildiz, M. E., Papailiopoulos, D., and Oymak, S. Transformers as algorithms: Generalization and stability in in-context learning. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 19565–19594. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/li231.html>.
- Li, Y., Rawat, A. S., and Oymak, S. Fine-grained analysis of in-context linear estimation: Data, architecture, and beyond. In Globerson, A., Mackey, L., Belgrave, D., Fan, A., Paquet, U., Tomczak, J., and Zhang, C. (eds.),

- Advances in Neural Information Processing Systems*, volume 37, pp. 138324–138364. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/f9dc462382fef56d58279e75de2438f3-Paper-Conference.pdf.
- Li, Y., Tarzanagh, D. A., Rawat, A. S., Fazel, M., and Oymak, S. Gating is weighting: Understanding gated linear attention through in-context learning. *arXiv preprint arXiv:2504.04308*, 2025.
- Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., and Neubig, G. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35, 2023.
- Liu, Y., Kothari, P., Van Delft, B., Bellot-Gurlet, B., Mordan, T., and Alahi, A. Ttt++: When does self-supervised test-time training fail or thrive? *Advances in Neural Information Processing Systems*, 34:21808–21820, 2021.
- Lu, Y., Letey, M., Zavatone-Veth, J. A., Maiti, A., and Pehlevan, C. In-context learning by linear attention: Exact asymptotics and experiments. In *NeurIPS 2024 Workshop on Mathematics of Modern Machine Learning*, 2024.
- Mahankali, A., Hashimoto, T. B., and Ma, T. One step of gradient descent is provably the optimal in-context learner with one layer of linear self-attention, 2023. URL <https://arxiv.org/abs/2307.03576>.
- Mahankali, A. V., Hashimoto, T., and Ma, T. One step of gradient descent is provably the optimal in-context learner with one layer of linear self-attention. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=8p3fu56lKc>.
- Min, S., Lyu, X., Holtzman, A., Artetxe, M., Lewis, M., Hajishirzi, H., and Zettlemoyer, L. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv preprint arXiv:2202.12837*, 2022.
- Niu, S., Wu, J., Zhang, Y., Chen, Y., Zheng, S., Zhao, P., and Tan, M. Efficient test-time model adaptation without forgetting. In *International conference on machine learning*, pp. 16888–16905. PMLR, 2022.
- Petersen, K. B. and Pedersen, M. S. The matrix cookbook, nov 2012. URL <http://www2.compute.dtu.dk/pubdb/pubs/3274-full.html>. Version 20121115.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. Language models are unsupervised multitask learners. *OpenAI Technical Report*, 2019. URL https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf.
- Snell, C., Lee, J., Xu, K., and Kumar, A. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- Sun, Y., Wang, X., Liu, Z., Miller, J., Efros, A., and Hardt, M. Test-time training with self-supervision for generalization under distribution shifts. In *International conference on machine learning*, pp. 9229–9248. PMLR, 2020.
- Sun, Y., Li, X., Dalal, K., Xu, J., Vikram, A., Zhang, G., Dubois, Y., Chen, X., Wang, X., Koyejo, S., et al. Learning to (learn at test time): Rnns with expressive hidden states. *arXiv preprint arXiv:2407.04620*, 2024.
- Thomas, V., Ma, J., Hosseinzadeh, R., Golestan, K., Yu, G., Volkovs, M., and Caterini, A. Retrieval & fine-tuning for in-context tabular models, 2024. URL <https://arxiv.org/abs/2406.05207>.
- Wang, D., Shelhamer, E., Liu, S., Olshausen, B., and Darrell, T. tent: fully test-time adaptation by entropy minimization. 2021.
- Wei, J., Bosma, M., Zhao, V., Guu, K., Yu, A. W., Lester, B., Du, N., Dai, A. M., and Le, Q. V. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*, 2022.
- Wu, J., Zou, D., Chen, Z., Braverman, V., Gu, Q., and Bartlett, P. L. How many pretraining tasks are needed for in-context learning of linear regression? *arXiv preprint arXiv:2310.08391*, 2023.
- Xie, S. M., Raghunathan, A., Liang, P., and Ma, T. An explanation of in-context learning as implicit bayesian inference. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=RdJVFCHJUMI>.
- Yuan, L., Xie, B., and Li, S. Robust test-time adaptation in dynamic scenarios. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15922–15932, 2023.
- Zhang, R., Frei, S., and Bartlett, P. L. Trained transformers learn linear models in-context, 2023. URL <https://arxiv.org/abs/2306.09927>.
- Zhang, R., Frei, S., and Bartlett, P. L. Trained transformers learn linear models in-context. *Journal of Machine Learning Research*, 25(49):1–55, 2024.

A. Proofs for Section 3

Proposition 3.1. *Consider the linear attention model with parameters $\mathbf{W} \in \mathbb{R}^{d \times d}$. Suppose the test-time training loss function is defined as in (2) and define $\mathbf{u}_{\text{context}} := \mathbf{X}_{\text{context}}^\top \mathbf{y}_{\text{context}} \in \mathbb{R}^d$. Then, for any step size $\eta > 0$, the new parameter \mathbf{W}_{TT} after one gradient-descent step from \mathbf{W} is given by the rank-1 update*

$$\mathbf{W}_{\text{TT}} = \mathbf{W} + 2\eta \mathbf{X}_{\text{train}}^\top (\mathbf{y}_{\text{train}} - \mathbf{X}_{\text{train}} \mathbf{W} \mathbf{u}_{\text{context}}) \mathbf{u}_{\text{context}}^\top.$$

Proof. We derive the update on the weight matrix \mathbf{W} via a single gradient step using the loss function defined over the set $\mathbf{Z}_{\text{train}} = [\mathbf{X}_{\text{train}} \ \mathbf{y}_{\text{train}}] \in \mathbb{R}^{k \times (d+1)}$. Note that $\text{SM}(\mathbf{Z}_i, \mathbf{W}) = \bar{\mathbf{x}}_{n+i}^\top \mathbf{W} \mathbf{X}_{\text{context}}^\top \mathbf{y}_{\text{context}} = \bar{\mathbf{x}}_{n+i}^\top \mathbf{W} \mathbf{u}_{\text{context}}$. Then, for the training set $\mathbf{Z}_{\text{train}}$, the vector of predicted values under \mathbf{W} is:

$$\hat{\mathbf{y}}_{\text{train}}(\mathbf{W}) = \mathbf{X}_{\text{train}} \mathbf{W} \mathbf{u}_{\text{context}} \in \mathbb{R}^k.$$

Recall the objective given by (2):

$$\hat{\mathcal{L}}(\mathbf{W}) = \sum_{i=n+1}^{n+k} (\bar{y}_i - \text{SM}(\mathbf{Z}_i, \mathbf{W}))^2 = \sum_{i=n+1}^{n+k} (\bar{y}_i - \bar{\mathbf{x}}_i^\top \mathbf{W} \mathbf{u}_{\text{context}})^2 = \|\mathbf{y}_{\text{train}} - \mathbf{X}_{\text{train}} \mathbf{W} \mathbf{u}_{\text{context}}\|_2^2.$$

The gradient of $\hat{\mathcal{L}}$ w.r.t. \mathbf{W} is then obtained by differentiating the summation:

$$\nabla_{\mathbf{W}} \hat{\mathcal{L}} = -2 \sum_{i=n+1}^{n+k} (\bar{y}_i - \bar{\mathbf{x}}_i^\top \mathbf{W} \mathbf{u}_{\text{context}}) \bar{\mathbf{x}}_i \mathbf{u}_{\text{context}}^\top = -2 \mathbf{X}_{\text{train}}^\top (\mathbf{y}_{\text{train}} - \mathbf{X}_{\text{train}} \mathbf{W} \mathbf{u}_{\text{context}}) \mathbf{u}_{\text{context}}^\top.$$

This is a rank-1 update as $(\mathbf{X}_{\text{train}}^\top (\mathbf{y}_{\text{train}} - \mathbf{X}_{\text{train}} \mathbf{W} \mathbf{u}_{\text{context}})) \mathbf{u}_{\text{context}}^\top = \mathbf{p} \mathbf{q}^\top$ is a rank-1 matrix. Evaluated at \mathbf{W} , the update after one gradient step with step size $\eta > 0$ becomes:

$$\begin{aligned} \mathbf{W}_{\text{TT}} &= \mathbf{W} - \eta \nabla_{\mathbf{W}} \hat{\mathcal{L}}(\mathbf{W}) \\ &= \mathbf{W} + 2\eta \mathbf{X}_{\text{train}}^\top (\mathbf{y}_{\text{train}} - \mathbf{X}_{\text{train}} \mathbf{W} \mathbf{u}_{\text{context}}) \mathbf{u}_{\text{context}}^\top. \end{aligned}$$

This completes the proof. \square

B. Proofs for Section 4

Lemma B.1 (Validity of Gaussian Approximation – Isotropic Case). *Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ have i.i.d $\mathcal{N}(0, 1)$ entries and let $\|\mathbf{w}\| = 1$. Define $\mathbf{q} = \frac{1}{n} \mathbf{X}^\top \mathbf{X} \mathbf{w}$. \mathbf{q} can be written as $\mathbf{q} = \mathbf{w} + \mathbf{g} + \mathbf{e}$ such that $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_d/n)$ and \mathbf{e} is a residual random variable that obeys*

$$\mathbb{E}[\|\mathbf{e}\|^2] \leq 9 \cdot \frac{n+d}{n^2}.$$

Note that \mathbf{e} represents a lower order term as we have $\|\mathbf{w}\| = 1$ and $\mathbb{E}[\|\mathbf{g}\|^2] = d/n$.

Proof. Set $\mathbf{h} = \mathbf{X} \mathbf{w} \sim \mathcal{N}(0, \mathbf{I}_n)$. Let $\mathbf{P} = \mathbf{I} - \mathbf{w} \mathbf{w}^\top$. Decompose $\mathbf{X}^\top = \mathbf{w} \mathbf{w}^\top \mathbf{X}^\top + \mathbf{P} \mathbf{X}^\top$ and observe that $\mathbf{P} \mathbf{X}^\top$ is independent of $\mathbf{w} \mathbf{w}^\top \mathbf{X}^\top = \mathbf{w} \mathbf{h}^\top$. Additionally, introduce independent $\mathbf{v} \sim \mathcal{N}(0, \mathbf{I}_n)$ and set $\mathbf{X}'^\top = \mathbf{P} \mathbf{X}^\top + \mathbf{w} \mathbf{v}^\top$. Finally, set $\mathbf{g} = \sqrt{n} \mathbf{X}'^\top \mathbf{h} / \|\mathbf{h}\|$ and

$$\mathbf{r} = \mathbf{P} \mathbf{X}^\top \mathbf{h} - \mathbf{g} = \underbrace{\mathbf{P} \mathbf{X}^\top (\mathbf{h} - \sqrt{n} \mathbf{h} / \|\mathbf{h}\|)}_{\mathbf{r}_1} - \underbrace{\sqrt{n} \mathbf{w} \mathbf{v}^\top \mathbf{h} / \|\mathbf{h}\|}_{\mathbf{r}_2}.$$

To proceed, observe that \mathbf{X}' has i.i.d $\mathcal{N}(0, 1)$ entries and hence $\mathbf{g} \sim \mathcal{N}(0, n \mathbf{I}_d)$. We can write

$$\mathbf{X}^\top \mathbf{X} \mathbf{w} = \mathbf{X}^\top \mathbf{h} = \mathbf{w} \|\mathbf{h}\|^2 + \mathbf{P} \mathbf{X}^\top \mathbf{h} = \mathbf{w} \|\mathbf{h}\|^2 + \mathbf{g} + \mathbf{r} \quad (7)$$

$$= n \mathbf{w} + \mathbf{g} + \mathbf{r} + (\|\mathbf{h}\|^2 - n) \mathbf{w}. \quad (8)$$

Let us now set $\mathbf{e}' = \mathbf{r} + (\|\mathbf{h}\|^2 - n) \mathbf{w}$ and investigate $\mathbb{E}[\|\mathbf{e}'\|^2]$. Recalling $\mathbf{r} = \mathbf{r}_1 - \mathbf{r}_2$, we can apply the standard upper bounds

$$\mathbb{E}[\|\mathbf{e}'\|^2] \leq 3 \mathbb{E}[(\|\mathbf{h}\|^2 - n)^2] + 3 \mathbb{E}[\|\mathbf{r}_1\|^2] + 3 \mathbb{E}[\|\mathbf{r}_2\|^2]. \quad (9)$$

Next, from standard facts about chi-square random variable with n degrees of freedom, we have $\mathbb{E}[(\|\mathbf{h}\|^2 - n)^2] = 2n$. Similarly, set $\tilde{\mathbf{h}} = \mathbf{h} - \sqrt{n}\mathbf{h}/\|\mathbf{h}\|$. Note that $\mathbb{E}[\|\tilde{\mathbf{h}}\|^2] = \mathbb{E}[(\|\mathbf{h}\| - \sqrt{n})^2] = 2n - 2\sqrt{n}\mathbb{E}[\|\mathbf{h}\|] \leq 2$. This step uses the fact that $\mathbb{E}[\|\mathbf{h}\|] = \frac{\sqrt{2}\Gamma(\frac{n+1}{2})}{\Gamma(\frac{n}{2})} \geq \frac{n}{\sqrt{n+1}}$, which can be obtained by induction as argued in Section 3.1 of (Chandrasekaran et al., 2012). Furthermore, considering that \mathbf{X}' is the sum of two orthogonal vectors, we obtain

$$\mathbb{E}[\|\mathbf{r}_1\|^2] = \mathbb{E}[\|\mathbf{P}\mathbf{X}^\top(\mathbf{h} - \sqrt{n}\mathbf{h}/\|\mathbf{h}\|)\|^2] \leq \mathbb{E}[\|\mathbf{X}'^\top(\mathbf{h} - \sqrt{n}\mathbf{h}/\|\mathbf{h}\|)\|^2] = d\mathbb{E}[\|\tilde{\mathbf{h}}\|^2] \leq 2d.$$

Finally, set $\mathbf{z} = \mathbf{v}^\top \mathbf{h}/\|\mathbf{h}\| \sim \mathcal{N}(0, 1)$. We bound

$$\mathbb{E}[\|\mathbf{r}_2\|^2] \leq \mathbb{E}[\|\sqrt{n}\mathbf{w}\mathbf{z}\|^2] = n.$$

Aggregating these, we obtain

$$\mathbb{E}[\|\mathbf{e}'\|^2] \leq 3(2n + 2d + n) = 9n + 6d.$$

To conclude, set $\mathbf{e} = \mathbf{e}'/n$. After normalizing $\mathbf{X}^\top \mathbf{X}\mathbf{w}$ by n and the change of variable $\mathbf{g} \leftarrow \mathbf{g}/n$ with $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_d/n)$, we obtain

$$n^{-1}\mathbf{X}^\top \mathbf{X}\mathbf{w} = \mathbf{w} + \mathbf{g} + \mathbf{e},$$

such that $\mathbb{E}[\|\mathbf{e}\|^2] = \mathbb{E}[\|\mathbf{e}'\|^2]/n^2 \leq 9/n + 6d/n^2$. \square

Lemma B.2 (Population Loss Expression). *Consider the linear attention model as the sequence model described in Section 3 with weight matrix \mathbf{W} . Suppose that the new task during the test-time training is β_{TT} . In that case, the population loss of the model in (3) with respect to this task is given by:*

$$\begin{aligned} \mathcal{L}(\mathbf{W}) &= \mathbb{E}_{(\mathbf{x}, \mathbf{y}), (\mathbf{x}_i, \mathbf{y}_i) \sim \mathcal{D}_z^{\text{TT}}(\beta_{\text{TT}})} \left[(\mathbf{y} - \mathbf{x}^\top \mathbf{W} \mathbf{X}^\top \mathbf{Y})^2 \right] \\ &= \beta_{\text{TT}}^\top \left[\Sigma_{\mathbf{x}} - n\Sigma_{\mathbf{x}} \mathbf{W} \Sigma_{\mathbf{x}} - n\Sigma_{\mathbf{x}} \mathbf{W}^\top \Sigma_{\mathbf{x}} + n(n+1)\Sigma_{\mathbf{x}} \mathbf{W}^\top \Sigma_{\mathbf{x}} \mathbf{W} \Sigma_{\mathbf{x}} + n \text{tr}(\mathbf{W}^\top \Sigma_{\mathbf{x}} \mathbf{W} \Sigma_{\mathbf{x}}) \Sigma_{\mathbf{x}} \right] \beta_{\text{TT}} + \sigma^2 n \text{tr}(\mathbf{W}^\top \Sigma_{\mathbf{x}} \mathbf{W} \Sigma_{\mathbf{x}}) + \sigma^2. \end{aligned}$$

Proof. Since the target function f is linear, the labels are given by:

$$y_i = f(\mathbf{x}_i) + \xi_i = \beta_{\text{TT}}^\top \mathbf{x}_i + \xi_i.$$

We collect the samples into $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{Y} \in \mathbb{R}^n$ as before, so that:

$$\mathbf{Y} = \mathbf{X}\beta_{\text{TT}} + \boldsymbol{\xi},$$

where $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n)^\top$ represents noise on the context labels. Recall the prediction by the model:

$$\text{SM}(\mathbf{Z}, \mathbf{W}) = \mathbf{x}^\top \mathbf{W} \mathbf{X}^\top \mathbf{Y}.$$

Then, we have the following population loss for the new task β_{TT} :

$$\mathcal{L}(\mathbf{W}) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}), (\mathbf{x}_i, \mathbf{y}_i) \sim \mathcal{D}_z^{\text{TT}}(\beta_{\text{TT}})} \left[(\mathbf{x}^\top \beta_{\text{TT}} - \mathbf{x}^\top \mathbf{W} \mathbf{X}^\top \mathbf{Y})^2 \right] + \sigma^2.$$

Since $\mathbf{Y} = \mathbf{X}\beta_{\text{TT}} + \boldsymbol{\xi}$, we have:

$$\text{SM}(\mathbf{Z}, \mathbf{W}) = \mathbf{x}^\top \mathbf{W} \mathbf{X}^\top \mathbf{X} \beta_{\text{TT}} + \mathbf{x}^\top \mathbf{W} \mathbf{X}^\top \boldsymbol{\xi}.$$

The error is:

$$\begin{aligned} \mathbf{x}^\top \beta_{\text{TT}} - \mathbf{x}^\top \mathbf{W} \mathbf{X}^\top \mathbf{X} \beta_{\text{TT}} - \mathbf{x}^\top \mathbf{W} \mathbf{X}^\top \boldsymbol{\xi} &= \mathbf{x}^\top (\beta_{\text{TT}} - \mathbf{W} \mathbf{X}^\top \mathbf{X} \beta_{\text{TT}}) - \mathbf{x}^\top \mathbf{W} \mathbf{X}^\top \boldsymbol{\xi} \\ &= \mathbf{x}^\top (\mathbf{I} - \mathbf{W} \mathbf{X}^\top \mathbf{X}) \beta_{\text{TT}} - \mathbf{x}^\top \mathbf{W} \mathbf{X}^\top \boldsymbol{\xi}. \end{aligned}$$

Therefore, the population loss is equal to:

$$\begin{aligned} \mathcal{L}(\mathbf{W}) &= \mathbb{E}_{(\mathbf{x}, \mathbf{y}), (\mathbf{x}_i, \mathbf{y}_i) \sim \mathcal{D}_z^{\text{TT}}(\beta_{\text{TT}})} \left[(\mathbf{x}^\top (\mathbf{I} - \mathbf{W} \mathbf{X}^\top \mathbf{X}) \beta_{\text{TT}} - \mathbf{x}^\top \mathbf{W} \mathbf{X}^\top \boldsymbol{\xi})^2 \right] + \sigma^2 \\ &= \mathbb{E}_{\mathbf{X}} \left[\mathbb{E}_{\mathbf{x}, \boldsymbol{\xi}} \left[(\mathbf{x}^\top (\mathbf{I} - \mathbf{W} \mathbf{X}^\top \mathbf{X}) \beta_{\text{TT}} - \mathbf{x}^\top \mathbf{W} \mathbf{X}^\top \boldsymbol{\xi})^2 \right] \right] + \sigma^2 \end{aligned}$$

We will first take the expectation with respect to query \mathbf{x} and noise vector ξ . Then, using independence, we will take the expectation of the resulting expression with respect to the context matrix \mathbf{X} . Expanding the square gives

$$(\mathbf{x}^\top (\mathbf{I} - \mathbf{W}\mathbf{X}^\top \mathbf{X}) \beta_{\text{TT}} - \mathbf{x}^\top \mathbf{W}\mathbf{X}^\top \xi)^2 = (\mathbf{x}^\top (\mathbf{I} - \mathbf{W}\mathbf{X}^\top \mathbf{X}) \beta_{\text{TT}})^2 - 2(\mathbf{x}^\top (\mathbf{I} - \mathbf{W}\mathbf{X}^\top \mathbf{X}) \beta_{\text{TT}})(\mathbf{x}^\top \mathbf{W}\mathbf{X}^\top \xi) + (\mathbf{x}^\top \mathbf{W}\mathbf{X}^\top \xi)^2.$$

Now, let's consider each of these three terms individually. First of all, the expectation of cross-term is:

$$\mathbb{E}_{\mathbf{x}, \xi}[-2(\mathbf{x}^\top (\mathbf{I} - \mathbf{W}\mathbf{X}^\top \mathbf{X}) \beta_{\text{TT}})(\mathbf{x}^\top \mathbf{W}\mathbf{X}^\top \xi)] = 0.$$

Besides, since the first term does not depend on ξ , we have:

$$\mathbb{E}_{\mathbf{x}, \xi}[(\mathbf{x}^\top (\mathbf{I} - \mathbf{W}\mathbf{X}^\top \mathbf{X}) \beta_{\text{TT}})^2] = \beta_{\text{TT}}^\top (\mathbf{I} - \mathbf{W}\mathbf{X}^\top \mathbf{X})^\top \Sigma_{\mathbf{x}} (\mathbf{I} - \mathbf{W}\mathbf{X}^\top \mathbf{X}) \beta_{\text{TT}}.$$

Consider now the last term $(\mathbf{x}^\top \mathbf{W}\mathbf{X}^\top \xi)^2$. Recall that ξ is zero-mean noise with $\mathbb{E}[\xi \xi^\top] = \sigma^2 \mathbf{I}_n$, and $\mathbb{E}[\mathbf{x} \mathbf{x}^\top] = \Sigma_{\mathbf{x}}$. First, condition on \mathbf{x} :

$$(\mathbf{x}^\top \mathbf{W}\mathbf{X}^\top \xi)^2 = \xi^\top (\mathbf{X}\mathbf{W}^\top \mathbf{x} \mathbf{x}^\top \mathbf{W}\mathbf{X}^\top) \xi.$$

Taking expectation over ξ :

$$\mathbb{E}_{\xi}[(\mathbf{x}^\top \mathbf{W}\mathbf{X}^\top \xi)^2 | \mathbf{x}] = \sigma^2 \text{tr}(\mathbf{X}\mathbf{W}^\top \mathbf{x} \mathbf{x}^\top \mathbf{W}\mathbf{X}^\top).$$

Use the cyclic property of trace and the fact that $\mathbf{x} \mathbf{x}^\top$ is rank one:

$$\text{tr}(\mathbf{X}\mathbf{W}^\top \mathbf{x} \mathbf{x}^\top \mathbf{W}\mathbf{X}^\top) = \text{tr}(\mathbf{x} \mathbf{x}^\top \mathbf{W}\mathbf{X}^\top \mathbf{X}\mathbf{W}^\top) = \mathbf{x}^\top (\mathbf{W}\mathbf{X}^\top \mathbf{X}\mathbf{W}^\top) \mathbf{x}.$$

Thus, we obtain:

$$\mathbb{E}_{\xi}[(\mathbf{x}^\top \mathbf{W}\mathbf{X}^\top \xi)^2 | \mathbf{x}] = \sigma^2 \mathbf{x}^\top (\mathbf{W}\mathbf{X}^\top \mathbf{X}\mathbf{W}^\top) \mathbf{x}.$$

Taking expectation over \mathbf{x} yields:

$$\mathbb{E}_{\mathbf{x}}[\mathbf{x}^\top (\mathbf{W}\mathbf{X}^\top \mathbf{X}\mathbf{W}^\top) \mathbf{x}] = \text{tr}(\mathbf{W}\mathbf{X}^\top \mathbf{X}\mathbf{W}^\top \Sigma_{\mathbf{x}}),$$

since $\mathbb{E}[\mathbf{x} \mathbf{x}^\top] = \Sigma_{\mathbf{x}}$. Combining both expectations:

$$\mathbb{E}_{\mathbf{x}, \xi}[(\mathbf{x}^\top \mathbf{W}\mathbf{X}^\top \xi)^2] = \sigma^2 \text{tr}(\mathbf{W}\mathbf{X}^\top \mathbf{X}\mathbf{W}^\top \Sigma_{\mathbf{x}}).$$

Thus, combining the three terms yields the following expression:

$$\beta_{\text{TT}}^\top (\mathbf{I} - \mathbf{W}\mathbf{X}^\top \mathbf{X})^\top \Sigma_{\mathbf{x}} (\mathbf{I} - \mathbf{W}\mathbf{X}^\top \mathbf{X}) \beta_{\text{TT}} + \sigma^2 \text{tr}(\mathbf{W}\mathbf{X}^\top \mathbf{X}\mathbf{W}^\top \Sigma_{\mathbf{x}})$$

Now, we aim to compute the loss function averaged over all possible in-context example sets \mathbf{X} :

$$\mathcal{L}(\mathbf{W}) = \mathbb{E}_{\mathbf{X}} \left[\beta_{\text{TT}}^\top (\mathbf{I} - \mathbf{W}\mathbf{X}^\top \mathbf{X})^\top \Sigma_{\mathbf{x}} (\mathbf{I} - \mathbf{W}\mathbf{X}^\top \mathbf{X}) \beta_{\text{TT}} + \sigma^2 \text{tr}(\mathbf{W}\mathbf{X}^\top \mathbf{X}\mathbf{W}^\top \Sigma_{\mathbf{x}}) \right] + \sigma^2.$$

Rewrite the first expression inside the expectation as

$$\beta_{\text{TT}}^\top \mathbf{A}^\top \Sigma_{\mathbf{x}} \mathbf{A} \beta_{\text{TT}} \quad \text{where} \quad \mathbf{A} = \mathbf{I} - \mathbf{W}\mathbf{X}^\top \mathbf{X}.$$

Then

$$\mathbb{E}_{\mathbf{X}} [\beta_{\text{TT}}^\top \mathbf{A}^\top \Sigma_{\mathbf{x}} \mathbf{A} \beta_{\text{TT}}] = \beta_{\text{TT}}^\top (\mathbb{E}_{\mathbf{X}} [\mathbf{A}^\top \Sigma_{\mathbf{x}} \mathbf{A}]) \beta_{\text{TT}}.$$

We can write

$$\mathbf{A}^\top \Sigma_{\mathbf{x}} \mathbf{A} = (\mathbf{I} - \mathbf{X}^\top \mathbf{X}\mathbf{W}^\top) \Sigma_{\mathbf{x}} (\mathbf{I} - \mathbf{W}\mathbf{X}^\top \mathbf{X}).$$

Hence

$$\mathbb{E}_{\mathbf{X}} [\mathbf{A}^\top \Sigma_{\mathbf{x}} \mathbf{A}] = \Sigma_{\mathbf{x}} - \Sigma_{\mathbf{x}} \mathbf{W} \mathbb{E}[\mathbf{X}^\top \mathbf{X}] - \mathbb{E}[\mathbf{X}^\top \mathbf{X}] \mathbf{W}^\top \Sigma_{\mathbf{x}} + \mathbb{E}[(\mathbf{X}^\top \mathbf{X}) \mathbf{W}^\top \Sigma_{\mathbf{x}} \mathbf{W} (\mathbf{X}^\top \mathbf{X})].$$

We know that $\mathbb{E}[\mathbf{X}^\top \mathbf{X}] = n \Sigma_{\mathbf{x}}$, thus, our expression becomes:

$$\mathbb{E}_{\mathbf{X}} [\mathbf{A}^\top \Sigma_{\mathbf{x}} \mathbf{A}] = \Sigma_{\mathbf{x}} - n \Sigma_{\mathbf{x}} \mathbf{W} \Sigma_{\mathbf{x}} - n \Sigma_{\mathbf{x}} \mathbf{W}^\top \Sigma_{\mathbf{x}} + \mathbb{E}[(\mathbf{X}^\top \mathbf{X}) \mathbf{W}^\top \Sigma_{\mathbf{x}} \mathbf{W} (\mathbf{X}^\top \mathbf{X})].$$

Using Lemma B.3, we know the following fact:

$$\mathbb{E}_{\mathbf{X}} [(\mathbf{X}^\top \mathbf{X}) \mathbf{M} (\mathbf{X}^\top \mathbf{X})] = n(n+1) \Sigma_{\mathbf{x}} \mathbf{M} \Sigma_{\mathbf{x}} + n \text{tr}(\mathbf{M} \Sigma_{\mathbf{x}}) \Sigma_{\mathbf{x}}.$$

Finally, the overall expression becomes:

$$\mathbb{E}_X[A \Sigma_x A] = \Sigma_x - n \Sigma_x W \Sigma_x - n \Sigma_x W^\top \Sigma_x + n(n+1) \Sigma_x W^\top \Sigma_x W \Sigma_x + n \text{tr}(W^\top \Sigma_x W \Sigma_x) \Sigma_x.$$

Adding the constant terms to the expectation gives:

$$\beta_{\text{TT}}^\top \left[\Sigma_x - n \Sigma_x W \Sigma_x - n \Sigma_x W^\top \Sigma_x + n(n+1) \Sigma_x W^\top \Sigma_x W \Sigma_x + n \text{tr}(W^\top \Sigma_x W \Sigma_x) \Sigma_x \right] \beta_{\text{TT}}.$$

Now that we have calculated the expectation of the first term, let's focus on the noise term involving trace inside the expectation. We have:

$$\mathbb{E}_X \left[\text{tr}(W X^\top X W^\top \Sigma_x) \right] = n \text{tr}(W^\top \Sigma_x W \Sigma_x).$$

As a result, we arrive at the final expression:

$$\begin{aligned} \mathcal{L}(W) &= \mathbb{E}_{X, (x, y) \sim \mathcal{D}} \left[\left((x^\top \beta_{\text{TT}} + \xi - x^\top W X^\top Y)^2 \right) \right] \\ &= \beta_{\text{TT}}^\top \left[\Sigma_x - n \Sigma_x W \Sigma_x - n \Sigma_x W^\top \Sigma_x + n(n+1) \Sigma_x W^\top \Sigma_x W \Sigma_x + n \text{tr}(W^\top \Sigma_x W \Sigma_x) \Sigma_x \right] \beta_{\text{TT}} \\ &\quad + \sigma^2 n \text{tr}(W^\top \Sigma_x W \Sigma_x) + \sigma^2. \end{aligned}$$

□

Lemma B.3. Let $X \in \mathbb{R}^{n \times d}$ be a random matrix whose rows $\mathbf{x}_i \in \mathbb{R}^d$, for $i = 1, \dots, n$, are i.i.d. drawn from $\mathcal{N}(\mathbf{0}, \Sigma)$. Let $M \in \mathbb{R}^{d \times d}$ be any fixed symmetric matrix. Then, under these assumptions,

$$\mathbb{E}_X \left[(X^\top X) M (X^\top X) \right] = n(n+1) \Sigma M \Sigma + n \text{tr}(M \Sigma) \Sigma.$$

Proof. We can write the matrix $X^\top X$ as:

$$X^\top X = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top.$$

Therefore,

$$(X^\top X) M (X^\top X) = \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right) M \left(\sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^\top \right) = \sum_{i=1}^n \sum_{j=1}^n \mathbf{x}_i (\mathbf{x}_i^\top M \mathbf{x}_j) \mathbf{x}_j^\top.$$

Taking expectation over X , we split the sum into the diagonal part ($i = j$) and the off-diagonal part ($i \neq j$):

$$\mathbb{E}_X \left[(X^\top X) M (X^\top X) \right] = \underbrace{\sum_{i=1}^n \mathbb{E} \left[\mathbf{x}_i (\mathbf{x}_i^\top M \mathbf{x}_i) \mathbf{x}_i^\top \right]}_{\text{(i) diagonal terms}} + \underbrace{\sum_{i \neq j} \mathbb{E} \left[\mathbf{x}_i (\mathbf{x}_i^\top M \mathbf{x}_j) \mathbf{x}_j^\top \right]}_{\text{(ii) off-diagonal terms}}.$$

(i) Diagonal Terms ($i = j$). Each $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \Sigma)$, so we set

$$\mathbb{E} \left[\mathbf{x}_i (\mathbf{x}_i^\top M \mathbf{x}_i) \mathbf{x}_i^\top \right] = \mathbb{E} \left[(\mathbf{x}_i^\top M \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top \right].$$

Recall the standard result, which follows from Wick's Theorem (Petersen & Pedersen, 2012)[Section 8.2.4] for a zero-mean Gaussian $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma)$ and any symmetric matrix A :

$$\mathbb{E} \left[\mathbf{x} (\mathbf{x}^\top A \mathbf{x}) \mathbf{x}^\top \right] = 2 \Sigma A \Sigma + \text{tr}(A \Sigma) \Sigma.$$

Substituting $A = M$, it follows that

$$\mathbb{E} \left[\mathbf{x}_i (\mathbf{x}_i^\top M \mathbf{x}_i) \mathbf{x}_i^\top \right] = 2 \Sigma M \Sigma + \text{tr}(M \Sigma) \Sigma.$$

Since we have n diagonal terms, summing over $i = 1, \dots, n$ yields

$$\sum_{i=1}^n \mathbb{E} \left[\mathbf{x}_i (\mathbf{x}_i^\top M \mathbf{x}_i) \mathbf{x}_i^\top \right] = n [2 \Sigma M \Sigma + \text{tr}(M \Sigma) \Sigma].$$

(ii) **Off-Diagonal Terms** ($i \neq j$). For $i \neq j$, the vectors \mathbf{x}_i and \mathbf{x}_j are independent. Then,

$$\mathbb{E}[\mathbf{x}_i (\mathbf{x}_i^\top \mathbf{M} \mathbf{x}_j) \mathbf{x}_j^\top] = \mathbb{E}[\mathbf{x}_i \mathbf{x}_i^\top] \mathbf{M} \mathbb{E}[\mathbf{x}_j \mathbf{x}_j^\top] = \mathbf{\Sigma} \mathbf{M} \mathbf{\Sigma}.$$

Hence, considering that we have n diagonal terms and $n(n-1)$ off-diagonal terms, combining (i) + (ii) gives:

$$\mathbb{E}_X[(\mathbf{X}^\top \mathbf{X}) \mathbf{M} (\mathbf{X}^\top \mathbf{X})] = n[2\mathbf{\Sigma} \mathbf{M} \mathbf{\Sigma} + \text{tr}(\mathbf{M} \mathbf{\Sigma}) \mathbf{\Sigma}] + n(n-1)\mathbf{\Sigma} \mathbf{M} \mathbf{\Sigma} = n(n+1)\mathbf{\Sigma} \mathbf{M} \mathbf{\Sigma} + n\text{tr}(\mathbf{M} \mathbf{\Sigma}) \mathbf{\Sigma}.$$

This completes the proof. \square

Proposition 4.1. Let $\mathbf{\Sigma}_\beta, \mathbf{\Sigma}_x = \mathbf{I}, \sigma^2 = 0$ and let β_{TT} be an arbitrary new task parameter. Then, the optimal pre-trained model's parameter \mathbf{W}^* and its population loss are given by

$$\mathbf{W}^* = \frac{1}{n+d+1} \mathbf{I}, \quad \mathcal{L}(\mathbf{W}^*) = \|\beta_{TT}\|^2 \frac{d+1}{n+d+1}.$$

Proof. Recall the loss expression from Lemma B.2:

$$\mathcal{L}(\mathbf{W}) = \beta_{TT}^\top [\mathbf{\Sigma}_x - n\mathbf{\Sigma}_x \mathbf{W} \mathbf{\Sigma}_x - n\mathbf{\Sigma}_x \mathbf{W}^\top \mathbf{\Sigma}_x + n(n+1)\mathbf{\Sigma}_x \mathbf{W}^\top \mathbf{\Sigma}_x \mathbf{W} \mathbf{\Sigma}_x + n\text{tr}(\mathbf{W}^\top \mathbf{\Sigma}_x \mathbf{W} \mathbf{\Sigma}_x) \mathbf{\Sigma}_x] \beta_{TT} + \sigma^2 n\text{tr}(\mathbf{W}^\top \mathbf{\Sigma}_x \mathbf{W} \mathbf{\Sigma}_x) + \sigma^2.$$

We will simplify the expression under the assumptions $\mathbf{\Sigma}_x, \mathbf{\Sigma}_\beta = \mathbf{I}$ and $\mathbf{W}^* = \frac{1}{n+d+1+\sigma^2} \mathbf{I}$, where the optimal pre-trained \mathbf{W}^* follows from Theorem 1 of Li et al. (2024). The first expression of interest is:

$$\mathbf{\Sigma}_x - n\mathbf{\Sigma}_x \mathbf{W} \mathbf{\Sigma}_x - n\mathbf{\Sigma}_x \mathbf{W}^\top \mathbf{\Sigma}_x + n(n+1)\mathbf{\Sigma}_x \mathbf{W}^\top \mathbf{\Sigma}_x \mathbf{W} \mathbf{\Sigma}_x + n\text{tr}(\mathbf{W}^\top \mathbf{\Sigma}_x \mathbf{W} \mathbf{\Sigma}_x) \mathbf{\Sigma}_x.$$

Substituting $\mathbf{\Sigma}_x$ and \mathbf{W}^* yields:

$$\begin{aligned} \mathbf{I} - \frac{n}{n+d+1+\sigma^2} \mathbf{I} - \frac{n}{n+d+1+\sigma^2} \mathbf{I} + \frac{n(n+1)}{(n+d+1+\sigma^2)^2} \mathbf{I} + \frac{nd}{(n+d+1+\sigma^2)^2} \mathbf{I} \\ = \frac{(d+1+\sigma^2-n)(n+d+1+\sigma^2) + n(n+d+1)}{(n+d+1+\sigma^2)^2} \mathbf{I} \\ = \frac{(d+1+\sigma^2)(n+d+1+\sigma^2) - \sigma^2 n}{(n+d+1+\sigma^2)^2} \mathbf{I}. \end{aligned}$$

In addition, the noise term is:

$$\sigma^2 n\text{tr}(\mathbf{W}^{*\top} \mathbf{\Sigma}_x \mathbf{W}^* \mathbf{\Sigma}_x) = \sigma^2 \frac{nd}{(n+d+1+\sigma^2)^2}.$$

Combining these, the final closed-form loss is

$$\mathcal{L}(\mathbf{W}^*) = \|\beta_{TT}\|^2 \frac{(d+1+\sigma^2)(n+d+1+\sigma^2) - \sigma^2 n}{(n+d+1+\sigma^2)^2} + \sigma^2 \frac{nd}{(n+d+1+\sigma^2)^2} + \sigma^2.$$

Finally, setting $\sigma^2 = 0$ in this expression yields the desired result. \square

Theorem 4.2. Let $n/d = \Theta(1)$. Recall the definition of the expected loss $\mathcal{L}_{TT}(\mathbf{W}_{TT})$ with respect to \mathcal{S}_{TT} in (4). In the isotropic covariance and noiseless setting ($\sigma^2 = 0$), the optimal step-size that minimizes $\mathcal{L}_{TT}(\mathbf{W}_{TT})$ is

$$\eta^* \approx \frac{d}{2(k+d)n^2(n+d)\|\beta_{TT}\|_2^2}.$$

With this optimal step-size η^* , the improvement in the loss due to the test-time training is

$$\mathcal{L}(\mathbf{W}^*) - \mathcal{L}_{TT}(\mathbf{W}_{TT}) \approx \frac{k}{k+d} \frac{d^3}{(n+d)^3} \|\beta_{TT}\|_2^2.$$

Proof. We consider the noiseless ($\sigma^2 = 0$) and *isotropic* case where $\Sigma_x = \Sigma_\beta = \mathbf{I}$. Throughout the proof, we omit lower order terms like $+1$ in $(d+1)$ or $2n$ in $(n^2 + 2n)$ as they're negligible compared to the higher order of the same variable (see Remark 4.3 for error discussion). This convention allows us to simplify the subsequent expressions.

Now, we aim to calculate the expected gain in the loss with respect to $(n+k)$ samples used during the test-time training process. That is, we will take the expectation of the loss function given by Lemma B.2 and compute $\mathcal{L}(W^*) - \mathcal{L}_{TT}(W_{TT}) = \mathbb{E}_{X_{\text{train}}, X_{\text{context}}} [\Delta \mathcal{L}]$ where $\Delta \mathcal{L}$ is given by:

$$\begin{aligned} \Delta \mathcal{L} &= \mathcal{L}(W^*) - \mathcal{L}(W_{TT}) \\ &= \beta_{TT}^\top \left[\Sigma_x - n \Sigma_x W^* \Sigma_x - n \Sigma_x W^{*\top} \Sigma_x + n(n+1) \Sigma_x W^{*\top} \Sigma_x W^* \Sigma_x + n \text{tr}(W^{*\top} \Sigma_x W^* \Sigma_x) \Sigma_x \right] \beta_{TT} \\ &\quad - \beta_{TT}^\top \left[\Sigma_x - n \Sigma_x W_{TT} \Sigma_x - n \Sigma_x W_{TT}^\top \Sigma_x + n(n+1) \Sigma_x W_{TT}^\top \Sigma_x W_{TT} \Sigma_x + n \text{tr}(W_{TT}^\top \Sigma_x W_{TT} \Sigma_x) \Sigma_x \right] \beta_{TT} \\ &= \beta_{TT}^\top \left[n \Sigma_x (W_{TT} - W^*) \Sigma_x + n \Sigma_x (W_{TT}^\top - W^{*\top}) \Sigma_x + n(n+1) \Sigma_x (W^{*\top} \Sigma_x W^* - W_{TT}^\top \Sigma_x W_{TT}) \Sigma_x \right. \\ &\quad \left. + n (\text{tr}(W^{*\top} \Sigma_x W^* \Sigma_x) - \text{tr}(W_{TT}^\top \Sigma_x W_{TT} \Sigma_x)) \Sigma_x \right] \beta_{TT}. \end{aligned} \quad (10)$$

In the noiseless and $\Sigma_x = \Sigma_\beta = \mathbf{I}$ case, the optimal pre-trained W^* becomes $W^* = \frac{1}{n+d+1} \mathbf{I}$ as shown in Theorem 1 of Li et al. (2024). Recall that when there's no noise and the linear task is β_{TT} , the gradient update given by Proposition 3.1 becomes:

$$W_{TT} = W^* + 2\eta X_{\text{train}}^\top X_{\text{train}} (\beta_{TT} - W^* X_{\text{context}}^\top X_{\text{context}} \beta_{TT}) \beta_{TT}^\top X_{\text{context}}^\top X_{\text{context}}.$$

A key technical challenge while using the above update is that the expectation of $\Delta \mathcal{L}$ in (10) involves 8th moments of X_{context} . To overcome this problem, we will use the Gaussian approximation following Lemma B.1 by approximating $X_{\text{context}}^\top X_{\text{context}} \beta_{TT} \approx n \Sigma_x \beta_{TT} + \sqrt{n} \Sigma_x^{1/2} \mathbf{g}$ where $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \|\beta_{TT}\|^2 \mathbf{I})$. This non-isotropic form can be obtained by factoring out $\Sigma_x^{1/2}$ and reducing the problem to the isotropic case from Lemma B.1. Also define $\Delta W := W_{TT} - W^*$. Then, using the approximation and the fact that $\Sigma_x = \mathbf{I}$, we can write:

$$\begin{aligned} \Delta W &= 2\eta X_{\text{train}}^\top X_{\text{train}} \left(\beta_{TT} - \frac{n}{n+d+1} \Sigma_x \beta_{TT} - \frac{\sqrt{n}}{n+d+1} \Sigma_x^{1/2} \mathbf{g} \right) (n \beta_{TT}^\top \Sigma_x + \sqrt{n} \mathbf{g}^\top \Sigma_x^{1/2}) \\ &= 2\eta X_{\text{train}}^\top X_{\text{train}} \left(n \beta_{TT} \beta_{TT}^\top \Sigma_x + \sqrt{n} \beta_{TT} \mathbf{g}^\top \Sigma_x^{1/2} - \frac{n^2}{n+d+1} \Sigma_x \beta_{TT} \beta_{TT}^\top \Sigma_x - \frac{n \sqrt{n}}{n+d+1} \Sigma_x \beta_{TT} \mathbf{g}^\top \Sigma_x^{1/2} \right. \\ &\quad \left. - \frac{n \sqrt{n}}{n+d+1} \Sigma_x^{1/2} \mathbf{g} \beta_{TT}^\top \Sigma_x - \frac{n}{n+d+1} \Sigma_x^{1/2} \mathbf{g} \mathbf{g}^\top \Sigma_x^{1/2} \right) \\ &= 2\eta X_{\text{train}}^\top X_{\text{train}} \left(\frac{n(d+1)}{n+d+1} \beta_{TT} \beta_{TT}^\top + \frac{\sqrt{n}(d+1)}{n+d+1} \beta_{TT} \mathbf{g}^\top - \frac{n \sqrt{n}}{n+d+1} \mathbf{g} \beta_{TT}^\top - \frac{n}{n+d+1} \mathbf{g} \mathbf{g}^\top \right) \\ \Delta W^\top &= 2\eta \left(\frac{n(d+1)}{n+d+1} \beta_{TT} \beta_{TT}^\top + \frac{\sqrt{n}(d+1)}{n+d+1} \mathbf{g} \beta_{TT}^\top - \frac{n \sqrt{n}}{n+d+1} \beta_{TT} \mathbf{g}^\top - \frac{n}{n+d+1} \mathbf{g} \mathbf{g}^\top \right) X_{\text{train}}^\top X_{\text{train}}. \end{aligned}$$

Notice that plugging $\Sigma_x = \mathbf{I}$ in Equation (10), we observe the following difference term in two of the expressions:

$$\begin{aligned} W^{*\top} W^* - W_{TT}^\top W_{TT} &= W^{*\top} \Sigma_x W^* - (W^{*\top} + \Delta W^\top) \Sigma_x (W^* + \Delta W) \\ &= -W^{*\top} \Delta W - \Delta W^\top W^* - \Delta W^\top \Delta W. \end{aligned}$$

Hence, we need to calculate both first-order expectations $\mathbb{E}_{X_{\text{train}}, \mathbf{g}} [\Delta W]$, $\mathbb{E}_{X_{\text{train}}, \mathbf{g}} [\Delta W^\top]$ and second-order expectation $\mathbb{E}_{X_{\text{train}}, \mathbf{g}} [\Delta W^\top \Delta W]$. By direct calculation, we obtain:

$$\mathbb{E}_{X_{\text{train}}, \mathbf{g}} [\Delta W] = 2\eta k \left(\frac{n(d+1)}{n+d+1} \beta_{TT} \beta_{TT}^\top - \frac{n}{n+d+1} \|\beta_{TT}\|_2^2 \mathbf{I} \right) = \mathbb{E}_{X_{\text{train}}, \mathbf{g}} [\Delta W^\top]. \quad (11)$$

For the second-order expectation $\mathbb{E}_{X_{\text{train}}} [X_{\text{train}}^\top X_{\text{train}} X_{\text{train}}^\top X_{\text{train}}]$, we will benefit from the fact that $\mathbb{E}_X [(X^\top X)(X^\top X)] = k(k+1)\Sigma^2 + k \text{tr}(\Sigma)\Sigma$ where $X \in \mathbb{R}^{k \times d}$ and its rows are drawn i.i.d from $\mathcal{N}(0, \Sigma)$. This result follows by plugging $M = \mathbf{I}$ into Lemma B.3. In our setting, $\Sigma_x = \mathbf{I}$, which gives

$$\mathbb{E}_{X_{\text{train}}} [X_{\text{train}}^\top X_{\text{train}} X_{\text{train}}^\top X_{\text{train}}] = k(k+d+1)\mathbf{I}.$$

Using the above result, we obtain:

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{x}_{\text{train}}, \mathbf{g}} [\Delta \mathbf{W}^\top \Delta \mathbf{W}] \\
 &= 4\eta^2 \mathbb{E}_{\mathbf{g}} \left[\left(\frac{n(d+1)}{n+d+1} \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top + \frac{\sqrt{n}(d+1)}{n+d+1} \mathbf{g} \boldsymbol{\beta}_{\text{TT}}^\top - \frac{n\sqrt{n}}{n+d+1} \boldsymbol{\beta}_{\text{TT}} \mathbf{g}^\top - \frac{n}{n+d+1} \mathbf{g} \mathbf{g}^\top \right) k(k+d+1) \mathbf{I} \right. \\
 &\quad \left. \left(\frac{n(d+1)}{n+d+1} \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top + \frac{\sqrt{n}(d+1)}{n+d+1} \boldsymbol{\beta}_{\text{TT}} \mathbf{g}^\top - \frac{n\sqrt{n}}{n+d+1} \mathbf{g} \boldsymbol{\beta}_{\text{TT}}^\top - \frac{n}{n+d+1} \mathbf{g} \mathbf{g}^\top \right) \right] \\
 &= 4\eta^2 k(k+d+1) \mathbb{E}_{\mathbf{g}} \left[\left(\frac{n(d+1)}{n+d+1} \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top + \frac{\sqrt{n}(d+1)}{n+d+1} \mathbf{g} \boldsymbol{\beta}_{\text{TT}}^\top - \frac{n\sqrt{n}}{n+d+1} \boldsymbol{\beta}_{\text{TT}} \mathbf{g}^\top - \frac{n}{n+d+1} \mathbf{g} \mathbf{g}^\top \right) \right. \\
 &\quad \left. \left(\frac{n(d+1)}{n+d+1} \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top + \frac{\sqrt{n}(d+1)}{n+d+1} \boldsymbol{\beta}_{\text{TT}} \mathbf{g}^\top - \frac{n\sqrt{n}}{n+d+1} \mathbf{g} \boldsymbol{\beta}_{\text{TT}}^\top - \frac{n}{n+d+1} \mathbf{g} \mathbf{g}^\top \right) \right] \\
 &= 4\eta^2 \frac{k(k+d+1)n}{(n+d+1)^2} \mathbb{E}_{\mathbf{g}} \left[\left(\sqrt{n}(d+1) \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top + (d+1) \mathbf{g} \boldsymbol{\beta}_{\text{TT}}^\top - n \boldsymbol{\beta}_{\text{TT}} \mathbf{g}^\top - \sqrt{n} \mathbf{g} \mathbf{g}^\top \right) \left(\sqrt{n}(d+1) \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top \right. \right. \\
 &\quad \left. \left. + (d+1) \boldsymbol{\beta}_{\text{TT}} \mathbf{g}^\top - n \mathbf{g} \boldsymbol{\beta}_{\text{TT}}^\top - \sqrt{n} \mathbf{g} \mathbf{g}^\top \right) \right] \\
 &= 4\eta^2 \frac{k(k+d+1)n}{(n+d+1)^2} \mathbb{E}_{\mathbf{g}} \left[n(d+1)^2 \|\boldsymbol{\beta}_{\text{TT}}\|_2^2 \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top - n(d+1) \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top \mathbf{g} \mathbf{g}^\top + (d+1)^2 \|\boldsymbol{\beta}_{\text{TT}}\|_2^2 \mathbf{g} \mathbf{g}^\top \right. \\
 &\quad \left. + n^2 \|\mathbf{g}\|_2^2 \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top - n(d+1) \mathbf{g} \mathbf{g}^\top \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top - n(d+1) \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top \mathbf{g} \mathbf{g}^\top - n(d+1) \mathbf{g} \mathbf{g}^\top \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top + n \mathbf{g} \mathbf{g}^\top \mathbf{g} \mathbf{g}^\top \right] \\
 &= 4\eta^2 \frac{k(k+d+1)n}{(n+d+1)^2} \left[n(d+1)^2 \|\boldsymbol{\beta}_{\text{TT}}\|_2^2 \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top - n(d+1) \|\boldsymbol{\beta}_{\text{TT}}\|_2^2 \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top + (d+1)^2 \|\boldsymbol{\beta}_{\text{TT}}\|_2^4 \mathbf{I} \right. \\
 &\quad \left. + n^2 d \|\boldsymbol{\beta}_{\text{TT}}\|_2^2 \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top - 3n(d+1) \|\boldsymbol{\beta}_{\text{TT}}\|_2^2 \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top + n(d+2) \|\boldsymbol{\beta}_{\text{TT}}\|_2^4 \mathbf{I} \right] \\
 &= 4\eta^2 \frac{k(k+d+1)n}{(n+d+1)^2} \left[n((d+1)(d-3) + nd) \|\boldsymbol{\beta}_{\text{TT}}\|_2^2 \boldsymbol{\beta}_{\text{TT}} \boldsymbol{\beta}_{\text{TT}}^\top + ((d+1)^2 + n(d+2)) \|\boldsymbol{\beta}_{\text{TT}}\|_2^4 \mathbf{I} \right].
 \end{aligned}$$

Applying the first-order expectation results from (11), we obtain that the first two terms of the loss difference are:

$$\begin{aligned}
 \mathbb{E}_{\mathbf{x}_{\text{train}}, \mathbf{g}} [\boldsymbol{\beta}_{\text{TT}}^\top n \boldsymbol{\Sigma}_{\mathbf{x}} (\mathbf{W}_{\text{TT}} - \mathbf{W}^*) \boldsymbol{\Sigma}_{\mathbf{x}} \boldsymbol{\beta}_{\text{TT}}] &= 2\eta k n \left(\frac{n(d+1)}{n+d+1} \|\boldsymbol{\beta}_{\text{TT}}\|_2^4 - \frac{n}{n+d+1} \|\boldsymbol{\beta}_{\text{TT}}\|_2^4 \right) \\
 &= 2\eta k n \frac{nd}{n+d+1} \|\boldsymbol{\beta}_{\text{TT}}\|_2^4,
 \end{aligned} \tag{12}$$

$$\mathbb{E}_{\mathbf{x}_{\text{train}}, \mathbf{g}} [\boldsymbol{\beta}_{\text{TT}}^\top n \boldsymbol{\Sigma}_{\mathbf{x}} (\mathbf{W}_{\text{TT}}^\top - \mathbf{W}^{*\top}) \boldsymbol{\Sigma}_{\mathbf{x}} \boldsymbol{\beta}_{\text{TT}}] = 2\eta k n \frac{nd}{n+d+1} \|\boldsymbol{\beta}_{\text{TT}}\|_2^4. \tag{13}$$

The other two terms in Equation (10) are:

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{x}_{\text{train}}, \mathbf{g}} [\boldsymbol{\beta}_{\text{TT}}^\top n(n+1) \boldsymbol{\Sigma}_{\mathbf{x}} (\mathbf{W}^{*\top} \boldsymbol{\Sigma}_{\mathbf{x}} \mathbf{W}^* - \mathbf{W}_{\text{TT}}^\top \boldsymbol{\Sigma}_{\mathbf{x}} \mathbf{W}_{\text{TT}}) \boldsymbol{\Sigma}_{\mathbf{x}} \boldsymbol{\beta}_{\text{TT}}] \\
 &= n(n+1) \boldsymbol{\beta}_{\text{TT}}^\top \mathbb{E}_{\mathbf{x}_{\text{train}}, \mathbf{g}} [-\mathbf{W}^{*\top} \Delta \mathbf{W} - \Delta \mathbf{W}^\top \mathbf{W}^* - \Delta \mathbf{W}^\top \Delta \mathbf{W}] \boldsymbol{\beta}_{\text{TT}} \\
 &= -4\eta k n(n+1) \frac{nd}{(n+d+1)^2} \|\boldsymbol{\beta}_{\text{TT}}\|_2^4 - 4\eta^2 \frac{k(k+d+1)n^2(n+1)}{(n+d+1)^2} \left[n(d^2 - 2d - 3 + nd) \|\boldsymbol{\beta}_{\text{TT}}\|_2^6 + ((d+1)^2 + n(d+2)) \|\boldsymbol{\beta}_{\text{TT}}\|_2^6 \right],
 \end{aligned} \tag{14}$$

and

$$\begin{aligned}
 & \mathbb{E}_{\mathbf{x}_{\text{train}}, \mathbf{g}} [\boldsymbol{\beta}_{\text{TT}}^\top n (\text{tr}(\mathbf{W}^{*\top} \boldsymbol{\Sigma}_{\mathbf{x}} \mathbf{W}^* \boldsymbol{\Sigma}_{\mathbf{x}}) - \text{tr}(\mathbf{W}_{\text{TT}}^\top \boldsymbol{\Sigma}_{\mathbf{x}} \mathbf{W}_{\text{TT}} \boldsymbol{\Sigma}_{\mathbf{x}})) \boldsymbol{\Sigma}_{\mathbf{x}} \boldsymbol{\beta}_{\text{TT}}] \\
 &= n \|\boldsymbol{\beta}_{\text{TT}}\|_2^2 \mathbb{E}_{\mathbf{x}_{\text{train}}, \mathbf{g}} [\text{tr}(-\mathbf{W}^{*\top} \Delta \mathbf{W} - \Delta \mathbf{W}^\top \mathbf{W}^* - \Delta \mathbf{W}^\top \Delta \mathbf{W})] \\
 &= -4\eta^2 \frac{k(k+d+1)n^2}{(n+d+1)^2} \left[n(d^2 - 2d - 3 + nd) \|\boldsymbol{\beta}_{\text{TT}}\|_2^6 + d((d+1)^2 + n(d+2)) \|\boldsymbol{\beta}_{\text{TT}}\|_2^6 \right] - 4\eta k \frac{n^2}{(n+d+1)^2} \|\boldsymbol{\beta}_{\text{TT}}\|_2^4.
 \end{aligned} \tag{15}$$

Combining (12), (13), (14), (15) in Equation (10), the expected improvement in the loss is the following:

$$\begin{aligned}
 \mathcal{L}(\mathbf{W}^*) - \mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}}) &= 4\eta k \frac{n^2 d}{n+d+1} \|\boldsymbol{\beta}_{\text{TT}}\|_2^4 - 4\eta k \frac{n^2}{(n+d+1)^2} \|\boldsymbol{\beta}_{\text{TT}}\|_2^4 - 4\eta k \frac{n^2(n+1)d}{(n+d+1)^2} \|\boldsymbol{\beta}_{\text{TT}}\|_2^4 \\
 &\quad - 4\eta^2 \frac{k(k+d+1)n^2}{(n+d+1)^2} \left[n(n+2)(d^2 - 2d - 3 + nd) \|\boldsymbol{\beta}_{\text{TT}}\|_2^6 + (n+d+1)((d+1)^2 + n(d+2)) \|\boldsymbol{\beta}_{\text{TT}}\|_2^6 \right].
 \end{aligned}$$

Rearranging and approximating by omitting lower-order terms, we obtain:

$$4\eta k \frac{n^2 d^2}{(n+d+1)^2} \|\beta_{\text{TT}}\|_2^4 - 4\eta^2 \frac{k(k+d+1)n^2 d(n+d)(n^2+d)}{(n+d+1)^2} \|\beta_{\text{TT}}\|_2^6. \quad (16)$$

Taking the derivative of this expression and setting to 0 yields that the optimal η^* solution is approximately:

$$\eta^* \approx \frac{d}{2(k+d+1)(n^2+d)(n+d) \|\beta_{\text{TT}}\|_2^2}. \quad (17)$$

Plugging the optimal η^* value in (17) into Equation (16), we obtain that the loss improvement is approximately:

$$\begin{aligned} \Delta \mathcal{L} &\approx \frac{2kd^3}{(n+d+1)^2(k+d+1)(n+d)(1+\frac{d}{n^2})} \|\beta_{\text{TT}}\|_2^2 - \frac{kd^3}{(n+d+1)^2(k+d+1)(n+d)(1+\frac{d}{n^2})} \|\beta_{\text{TT}}\|_2^2 \\ &= \frac{k}{k+d+1} \frac{d^3}{(n+d+1)^2(n+d)(1+\frac{d}{n^2})} \|\beta_{\text{TT}}\|_2^2. \end{aligned}$$

Considering that $n/d = \Theta(1)$, we conclude that

$$\Delta \mathcal{L} \approx \frac{k}{k+d} \frac{d^3}{(n+d)^3} \|\beta_{\text{TT}}\|_2^2.$$

□

Corollary 4.4. Recall the definitions of $\gamma = k/d$, $\alpha = n/d$, and consider the loss $\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}})$ as a function of α in the isotropic covariance and noiseless setting. If $\gamma > \frac{1}{2}$, then the loss $\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}})$ is non-monotonic in α . Specifically, the loss $\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}})$ is increasing for $\alpha < \sqrt{\frac{3\gamma}{\gamma+1}} - 1$ and decreasing for $\alpha > \sqrt{\frac{3\gamma}{\gamma+1}} - 1$. Conversely, if $\gamma < \frac{1}{2}$, then $\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}})$ is monotonic in α .

Proof. Using Theorem 4.2, the new loss $\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}})$ is approximately:

$$\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}}) \approx \|\beta_{\text{TT}}\|_2^2 \left(\frac{d}{n+d} - \frac{k}{k+d} \left(\frac{d}{n+d} \right)^3 \right) = \|\beta_{\text{TT}}\|_2^2 \left(a - \frac{k}{k+d} a^3 \right) \text{ where } a := \frac{d}{n+d} \in (0, 1).$$

Since $\|\beta_{\text{TT}}\|_2$ is fixed, we focus on the second multiplier and define $f(a) := a - \frac{k}{k+d} a^3$. Computing the derivative of $f(a)$ gives:

$$f'(a) = 1 - 3 \frac{k}{k+d} a^2.$$

Solving $f'(a) = 0$ yields $a_{\text{crit}} = \sqrt{\frac{k+d}{3k}}$. This lies in $(0, 1)$ when $\sqrt{\frac{k+d}{3k}} < 1$, i.e. $k > \frac{d}{2}$. In that case, $f(a)$ has exactly one turning point in $(0, 1)$, so f increases on $(0, a_{\text{crit}})$ and decreases on $(a_{\text{crit}}, 1)$, making it non-monotonic. On the other hand, in the $k < \frac{d}{2}$ regime, one finds

$$\sqrt{\frac{k+d}{3k}} > 1 \implies a_{\text{crit}} \notin (0, 1),$$

hence $f'(a)$ cannot change sign over $(0, 1)$. Therefore, $f(a)$ is monotonic in $a = \frac{d}{n+d}$, which implies that $\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}})$ is also monotonic in the ratio n/d . Finally, since $a = \frac{1}{\alpha+1}$, the threshold until which $\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}})$ increases translates to $\sqrt{\frac{3\gamma}{\gamma+1}} - 1$ in terms of α . Consequently, for $\gamma > 1/2$, the new loss is non-monotonic with respect to $\alpha = n/d$, whereas for $\gamma < 1/2$, it remains monotonic. This completes the argument. □

Theorem 4.5. Consider the isotropic covariance and noiseless setting ($\sigma^2 = 0$). Suppose the initial weight matrix is $\mathbf{W}^* = \mathbf{0}_{d \times d}$. Then, the optimal step-size that minimizes $\mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}})$ is

$$\eta^* = \frac{1}{2(k+d+1)(n^2+4n+3+d) \|\beta_{\text{TT}}\|_2^2}.$$

With this optimal step-size η^* , the improvement in the loss due to test-time training is

$$\mathcal{L}(\mathbf{W}^*) - \mathcal{L}_{TT}(\mathbf{W}_{TT}) = \frac{k}{k+d+1} \frac{n^2}{n^2+4n+3+d} \|\boldsymbol{\beta}_{TT}\|_2^2,$$

where $\mathcal{L}(\mathbf{W}^*) = \|\boldsymbol{\beta}_{TT}\|_2^2$.

Proof. In the proof, we will handle the general noisy scenario, and setting $\sigma^2 = 0$ in the final expression will recover the result for the noiseless setting. Similar to the proof of Theorem 4.2, we aim to calculate the expected gain in the loss with respect to $(n+k)$ samples used during the test-time training process. During test time, we draw $(n+k)$ total samples $\{(\bar{\mathbf{x}}_i, \bar{\mathbf{y}}_i)\}_{i=1}^{n+k}$ drawn at test time, each with an associated noise variable $\{\xi_i\}_{i=1}^{n+k}$. Let's denote $\boldsymbol{\xi}_{\text{train}} = [\xi_1 \dots \xi_n]^\top$ and $\boldsymbol{\xi}_{\text{context}} = [\xi_{n+1} \dots \xi_{n+k}]^\top$. We will then evaluate the expectation of the loss function given by Lemma B.2 and compute $\mathcal{L}(\mathbf{W}^*) - \mathcal{L}_{TT}(\mathbf{W}_{TT}) = \mathbb{E}_{\mathbf{X}_{\text{train}}, \boldsymbol{\xi}_{\text{train}}, \mathbf{X}_{\text{context}}, \boldsymbol{\xi}_{\text{context}}} [\Delta \mathcal{L}]$ where $\Delta \mathcal{L}$ is given by:

$$\begin{aligned} \Delta \mathcal{L} &= \mathcal{L}(\mathbf{W}^*) - \mathcal{L}(\mathbf{W}_{TT}) \\ &= \boldsymbol{\beta}_{TT}^\top \left[\boldsymbol{\Sigma}_x - n \boldsymbol{\Sigma}_x \mathbf{W}^* \boldsymbol{\Sigma}_x - n \boldsymbol{\Sigma}_x \mathbf{W}^{*\top} \boldsymbol{\Sigma}_x + n(n+1) \boldsymbol{\Sigma}_x \mathbf{W}^{*\top} \boldsymbol{\Sigma}_x \mathbf{W}^* \boldsymbol{\Sigma}_x + n \text{tr}(\mathbf{W}^{*\top} \boldsymbol{\Sigma}_x \mathbf{W}^* \boldsymbol{\Sigma}_x) \boldsymbol{\Sigma}_x \right] \boldsymbol{\beta}_{TT} \\ &\quad + \sigma^2 n \text{tr}(\mathbf{W}^{*\top} \boldsymbol{\Sigma}_x \mathbf{W}^* \boldsymbol{\Sigma}_x) \\ &\quad - \boldsymbol{\beta}_{TT}^\top \left[\boldsymbol{\Sigma}_x - n \boldsymbol{\Sigma}_x \mathbf{W}_{TT} \boldsymbol{\Sigma}_x - n \boldsymbol{\Sigma}_x \mathbf{W}_{TT}^\top \boldsymbol{\Sigma}_x + n(n+1) \boldsymbol{\Sigma}_x \mathbf{W}_{TT}^\top \boldsymbol{\Sigma}_x \mathbf{W}_{TT} \boldsymbol{\Sigma}_x + n \text{tr}(\mathbf{W}_{TT}^\top \boldsymbol{\Sigma}_x \mathbf{W}_{TT} \boldsymbol{\Sigma}_x) \boldsymbol{\Sigma}_x \right] \boldsymbol{\beta}_{TT} \\ &\quad - \sigma^2 n \text{tr}(\mathbf{W}_{TT}^\top \boldsymbol{\Sigma}_x \mathbf{W}_{TT} \boldsymbol{\Sigma}_x) \\ &= \boldsymbol{\beta}_{TT}^\top \left[n \boldsymbol{\Sigma}_x (\mathbf{W}_{TT} - \mathbf{W}^*) \boldsymbol{\Sigma}_x + n \boldsymbol{\Sigma}_x (\mathbf{W}_{TT}^\top - \mathbf{W}^{*\top}) \boldsymbol{\Sigma}_x + n(n+1) \boldsymbol{\Sigma}_x (\mathbf{W}^{*\top} \boldsymbol{\Sigma}_x \mathbf{W}^* - \mathbf{W}_{TT}^\top \boldsymbol{\Sigma}_x \mathbf{W}_{TT}) \boldsymbol{\Sigma}_x \right. \\ &\quad \left. + n (\text{tr}(\mathbf{W}^{*\top} \boldsymbol{\Sigma}_x \mathbf{W}^* \boldsymbol{\Sigma}_x) - \text{tr}(\mathbf{W}_{TT}^\top \boldsymbol{\Sigma}_x \mathbf{W}_{TT} \boldsymbol{\Sigma}_x)) \boldsymbol{\Sigma}_x \right] \boldsymbol{\beta}_{TT} + \sigma^2 n (\text{tr}(\mathbf{W}^{*\top} \boldsymbol{\Sigma}_x \mathbf{W}^* \boldsymbol{\Sigma}_x) - \text{tr}(\mathbf{W}_{TT}^\top \boldsymbol{\Sigma}_x \mathbf{W}_{TT} \boldsymbol{\Sigma}_x)). \quad (18) \end{aligned}$$

Recall that when the linear task is $\boldsymbol{\beta}_{TT}$ and $\mathbf{W}^* = \mathbf{0}_{d \times d}$, the gradient update given by Proposition 3.1 becomes:

$$\begin{aligned} \mathbf{W}_{TT} &= \mathbf{W}^* + 2\eta \left[\mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} (\boldsymbol{\beta}_{TT} - \mathbf{W}^* \mathbf{X}_{\text{context}}^\top \mathbf{y}_{\text{context}}) + \mathbf{X}_{\text{train}}^\top \boldsymbol{\xi}_{\text{train}} \right] (\mathbf{X}_{\text{context}}^\top \mathbf{y}_{\text{context}})^\top \\ &= \mathbf{W}^* + 2\eta \left[\mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} (\boldsymbol{\beta}_{TT} - \mathbf{W}^* \mathbf{X}_{\text{context}}^\top (\mathbf{X}_{\text{context}} \boldsymbol{\beta}_{TT} + \boldsymbol{\xi}_{\text{context}})) + \mathbf{X}_{\text{train}}^\top \boldsymbol{\xi}_{\text{train}} \right] (\mathbf{X}_{\text{context}}^\top (\mathbf{X}_{\text{context}} \boldsymbol{\beta}_{TT} + \boldsymbol{\xi}_{\text{context}}))^\top \\ &= \mathbf{W}^* + 2\eta \left[\mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} \boldsymbol{\beta}_{TT} + \mathbf{X}_{\text{train}}^\top \boldsymbol{\xi}_{\text{train}} \right] (\mathbf{X}_{\text{context}}^\top (\mathbf{X}_{\text{context}} \boldsymbol{\beta}_{TT} + \boldsymbol{\xi}_{\text{context}}))^\top \\ &= 2\eta \left[\mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} \boldsymbol{\beta}_{TT} + \mathbf{X}_{\text{train}}^\top \boldsymbol{\xi}_{\text{train}} \right] \left[\boldsymbol{\beta}_{TT}^\top \mathbf{X}_{\text{context}}^\top \mathbf{X}_{\text{context}} + \boldsymbol{\xi}_{\text{context}}^\top \mathbf{X}_{\text{context}} \right]. \end{aligned}$$

Then, we obtain its transpose as

$$\mathbf{W}_{TT}^\top = 2\eta \left[\mathbf{X}_{\text{context}}^\top \mathbf{X}_{\text{context}} \boldsymbol{\beta}_{TT} + \mathbf{X}_{\text{context}}^\top \boldsymbol{\xi}_{\text{context}} \right] \left[\boldsymbol{\beta}_{TT}^\top \mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} + \boldsymbol{\xi}_{\text{train}}^\top \mathbf{X}_{\text{train}} \right].$$

Thus, the expression $\mathbf{W}_{TT}^\top \mathbf{W}_{TT}$ becomes:

$$\begin{aligned} \mathbf{W}_{TT}^\top \mathbf{W}_{TT} &= 4\eta^2 \left[\mathbf{X}_{\text{context}}^\top \mathbf{X}_{\text{context}} \boldsymbol{\beta}_{TT} + \mathbf{X}_{\text{context}}^\top \boldsymbol{\xi}_{\text{context}} \right] \left[\boldsymbol{\beta}_{TT}^\top \mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} + \boldsymbol{\xi}_{\text{train}}^\top \mathbf{X}_{\text{train}} \right] \\ &\quad \left[\mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} \boldsymbol{\beta}_{TT} + \mathbf{X}_{\text{train}}^\top \boldsymbol{\xi}_{\text{train}} \right] \left[\boldsymbol{\beta}_{TT}^\top \mathbf{X}_{\text{context}}^\top \mathbf{X}_{\text{context}} + \boldsymbol{\xi}_{\text{context}}^\top \mathbf{X}_{\text{context}} \right]. \end{aligned}$$

Therefore, we can write:

$$\begin{aligned} \mathbb{E}_{\mathbf{X}_{\text{train}}, \boldsymbol{\xi}_{\text{train}}, \mathbf{X}_{\text{context}}, \boldsymbol{\xi}_{\text{context}}} [\mathbf{W}_{TT}^\top \mathbf{W}_{TT}] &= 4\eta^2 \mathbb{E}_{\mathbf{X}_{\text{train}}, \boldsymbol{\xi}_{\text{train}}, \mathbf{X}_{\text{context}}, \boldsymbol{\xi}_{\text{context}}} \left[\mathbf{X}_{\text{context}}^\top \boldsymbol{\xi}_{\text{context}} \boldsymbol{\xi}_{\text{train}}^\top \mathbf{X}_{\text{train}} \mathbf{X}_{\text{train}}^\top \boldsymbol{\xi}_{\text{train}} \boldsymbol{\xi}_{\text{context}}^\top \mathbf{X}_{\text{context}} \right. \\ &\quad + \mathbf{X}_{\text{context}}^\top \boldsymbol{\xi}_{\text{context}} \boldsymbol{\beta}_{TT}^\top \mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} \mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} \boldsymbol{\beta}_{TT} \boldsymbol{\xi}_{\text{context}}^\top \mathbf{X}_{\text{context}} \\ &\quad + \mathbf{X}_{\text{context}}^\top \mathbf{X}_{\text{context}} \boldsymbol{\beta}_{TT} \boldsymbol{\xi}_{\text{train}}^\top \mathbf{X}_{\text{train}} \mathbf{X}_{\text{train}}^\top \boldsymbol{\xi}_{\text{train}} \boldsymbol{\beta}_{TT}^\top \mathbf{X}_{\text{context}}^\top \mathbf{X}_{\text{context}} \\ &\quad \left. + \mathbf{X}_{\text{context}}^\top \mathbf{X}_{\text{context}} \boldsymbol{\beta}_{TT} \boldsymbol{\beta}_{TT}^\top \mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} \mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} \boldsymbol{\beta}_{TT} \boldsymbol{\beta}_{TT}^\top \mathbf{X}_{\text{context}}^\top \mathbf{X}_{\text{context}} \right]. \end{aligned}$$

By plugging $\mathbf{M} = \mathbf{I}$ in Lemma B.3 and considering that $\boldsymbol{\Sigma}_x = \mathbf{I}$, we have the following fact:

$$\mathbb{E}_{\mathbf{X}_{\text{train}}} [\mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} \mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}}] = k(k+d+1)\mathbf{I}.$$

Applying this identity twice and using the fact that $\mathbb{E}_{X_{\text{train}}} [X_{\text{train}} X_{\text{train}}^\top] = \text{tr}(\Sigma_x) \mathbf{I} = d\mathbf{I}$, we obtain:

$$\begin{aligned}
 & \mathbb{E}_{X_{\text{train}}, \xi_{\text{train}}, X_{\text{context}}, \xi_{\text{context}}} [\mathbf{W}_{\text{TT}}^\top \mathbf{W}_{\text{TT}}] \\
 &= 4\eta^2 \left(\sigma^4 k n d \mathbf{I} + \sigma^2 k (k + d + 1) n \|\beta_{\text{TT}}\|_2^2 \mathbf{I} + \sigma^2 k d (n(n + 1) \beta_{\text{TT}} \beta_{\text{TT}}^\top + n \|\beta_{\text{TT}}\|_2^2 \mathbf{I}) \right. \\
 &\quad \left. + k (k + d + 1) \|\beta_{\text{TT}}\|_2^2 (n(n + 1) \beta_{\text{TT}} \beta_{\text{TT}}^\top + n \|\beta_{\text{TT}}\|_2^2 \mathbf{I}) \right) \\
 &= 4\eta^2 \left(\sigma^4 k n (\sigma^2 d + (k + d + 1) \|\beta_{\text{TT}}\|_2^2) \mathbf{I} + k (n(n + 1) \beta_{\text{TT}} \beta_{\text{TT}}^\top + n \|\beta_{\text{TT}}\|_2^2 \mathbf{I}) (\sigma^2 d + (k + d + 1) \|\beta_{\text{TT}}\|_2^2) \right) \\
 &= 4\eta^2 k n (\sigma^2 d + (k + d + 1) \|\beta_{\text{TT}}\|_2^2) (\sigma^2 \mathbf{I} + (n + 1) \beta_{\text{TT}} \beta_{\text{TT}}^\top + \|\beta_{\text{TT}}\|_2^2 \mathbf{I}). \tag{19}
 \end{aligned}$$

Besides that, we calculate the first order expectations of \mathbf{W}_{TT} as:

$$\begin{aligned}
 \mathbb{E}_{X_{\text{train}}, \xi_{\text{train}}, X_{\text{context}}, \xi_{\text{context}}} [\beta_{\text{TT}}^\top n \Sigma_x (\mathbf{W}_{\text{TT}} - \mathbf{W}^*) \Sigma_x \beta_{\text{TT}}] &= 2\eta n \beta_{\text{TT}}^\top \mathbb{E}_{X_{\text{train}}, X_{\text{context}}} [X_{\text{train}}^\top X_{\text{train}} \beta_{\text{TT}} \beta_{\text{TT}}^\top X_{\text{context}}^\top X_{\text{context}}] \beta_{\text{TT}} \\
 &= 2\eta k n^2 \|\beta_{\text{TT}}\|_2^4, \tag{20}
 \end{aligned}$$

and similarly,

$$\mathbb{E}_{X_{\text{train}}, \xi_{\text{train}}, X_{\text{context}}, \xi_{\text{context}}} [\beta_{\text{TT}}^\top n \Sigma_x (\mathbf{W}_{\text{TT}}^\top - \mathbf{W}^{*\top}) \Sigma_x \beta_{\text{TT}}] = 2\eta k n^2 \|\beta_{\text{TT}}\|_2^4. \tag{21}$$

The next term in the loss difference is:

$$\begin{aligned}
 & \mathbb{E}_{X_{\text{train}}, \xi_{\text{train}}, X_{\text{context}}, \xi_{\text{context}}} [\beta_{\text{TT}}^\top n (n + 1) \Sigma_x (\mathbf{W}^{*\top} \Sigma_x \mathbf{W}^* - \mathbf{W}_{\text{TT}}^\top \Sigma_x \mathbf{W}_{\text{TT}}) \Sigma_x \beta_{\text{TT}}] \\
 &= -4\eta^2 k n^2 (n + 1) ((n + 2) \|\beta_{\text{TT}}\|_2^2 + \sigma^2) \|\beta_{\text{TT}}\|_2^2 (\sigma^2 d + (k + d + 1) \|\beta_{\text{TT}}\|_2^2). \tag{22}
 \end{aligned}$$

Using (19), the last term is:

$$\begin{aligned}
 & \mathbb{E}_{X_{\text{train}}, \xi_{\text{train}}, X_{\text{context}}, \xi_{\text{context}}} [\beta_{\text{TT}}^\top n (\text{tr}(\mathbf{W}^{*\top} \Sigma_x \mathbf{W}^* \Sigma_x) - \text{tr}(\mathbf{W}_{\text{TT}}^\top \Sigma_x \mathbf{W}_{\text{TT}} \Sigma_x)) \Sigma_x \beta_{\text{TT}}] \\
 &= -4\eta^2 n \beta_{\text{TT}}^\top \text{tr}(\mathbb{E}_{X_{\text{train}}, \xi_{\text{train}}, X_{\text{context}}, \xi_{\text{context}}} [\mathbf{W}_{\text{TT}}^\top \mathbf{W}_{\text{TT}}]) \beta_{\text{TT}}. \\
 &= -4\eta^2 k n^2 (\sigma^2 d + (k + d + 1) \|\beta_{\text{TT}}\|_2^2) (\sigma^2 d + (n + d + 1) \|\beta_{\text{TT}}\|_2^2) \|\beta_{\text{TT}}\|_2^2. \tag{23}
 \end{aligned}$$

Finally, the noise term is equal to:

$$\begin{aligned}
 & \mathbb{E}_{X_{\text{train}}, \xi_{\text{train}}, X_{\text{context}}, \xi_{\text{context}}} [\sigma^2 n (\text{tr}(\mathbf{W}^{*\top} \Sigma_x \mathbf{W}^* \Sigma_x) - \text{tr}(\mathbf{W}_{\text{TT}}^\top \Sigma_x \mathbf{W}_{\text{TT}} \Sigma_x))] \\
 &= -4\eta^2 k n^2 \sigma^2 (\sigma^2 d + (k + d + 1) \|\beta_{\text{TT}}\|_2^2) (\sigma^2 d + (n + d + 1) \|\beta_{\text{TT}}\|_2^2). \tag{24}
 \end{aligned}$$

Combining the Equations (20), (21), (22), (23), (24) in the expectation of the loss difference (18), we obtain the following quadratic expression of η :

$$\begin{aligned}
 & \mathcal{L}(\mathbf{W}^*) - \mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}}) \\
 &= 4\eta k n^2 \|\beta_{\text{TT}}\|_2^4 - 4\eta^2 k n^2 (n + 1) ((n + 2) \|\beta_{\text{TT}}\|_2^2 + \sigma^2) \|\beta_{\text{TT}}\|_2^2 (\sigma^2 d + (k + d + 1) \|\beta_{\text{TT}}\|_2^2) \\
 &\quad - 4\eta^2 k n^2 (\sigma^2 d + (k + d + 1) \|\beta_{\text{TT}}\|_2^2) (\sigma^2 d + (n + d + 1) \|\beta_{\text{TT}}\|_2^2) \|\beta_{\text{TT}}\|_2^2 \\
 &\quad - 4\eta^2 k n^2 \sigma^2 (\sigma^2 d + (k + d + 1) \|\beta_{\text{TT}}\|_2^2) (\sigma^2 d + (n + d + 1) \|\beta_{\text{TT}}\|_2^2) \\
 &= 4k n^2 \left[\eta \|\beta_{\text{TT}}\|_2^4 - \eta^2 (\sigma^2 d + (k + d + 1) \|\beta_{\text{TT}}\|_2^2) [(\sigma^2 + \|\beta_{\text{TT}}\|_2^2) (\sigma^2 d + (n + d + 1) \|\beta_{\text{TT}}\|_2^2) \right. \\
 &\quad \left. + (n + 1) ((n + 2) \|\beta_{\text{TT}}\|_2^2 + \sigma^2) \|\beta_{\text{TT}}\|_2^2] \right] \\
 &= 4k n^2 \left[\eta \|\beta_{\text{TT}}\|_2^4 - \eta^2 (\sigma^2 d + (k + d + 1) \|\beta_{\text{TT}}\|_2^2) (\sigma^4 d + \|\beta_{\text{TT}}\|_2^2 ((n^2 + 4n + 3 + d) \|\beta_{\text{TT}}\|_2^2 + 2\sigma^2 (n + d + 1))) \right].
 \end{aligned}$$

Setting the derivative to 0 and solving for η gives:

$$\eta^* = \frac{\|\beta_{\text{TT}}\|_2^4}{2 (\sigma^2 d + (k + d + 1) \|\beta_{\text{TT}}\|_2^2) (\sigma^4 d + \|\beta_{\text{TT}}\|_2^2 (n^2 + 4n + 3 + d) + 2\sigma^2 (n + d + 1) \|\beta_{\text{TT}}\|_2^2)}.$$

At this optimal value η^* , the improvement on the loss is:

$$\frac{k \|\beta_{\text{TT}}\|_2^2}{\sigma^2 d + (k + d + 1) \|\beta_{\text{TT}}\|_2^2} \frac{n^2 \|\beta_{\text{TT}}\|_2^4}{\left(\sigma^4 d + \|\beta_{\text{TT}}\|_2^4 (n^2 + 4n + 3 + d) + 2\sigma^2 (n + d + 1) \|\beta_{\text{TT}}\|_2^2 \right)} \|\beta_{\text{TT}}\|_2^2.$$

In the noiseless setting, the optimal η^* corresponds to:

$$\eta^* = \frac{1}{2(k + d + 1)(n^2 + 4n + 3 + d) \|\beta_{\text{TT}}\|_2^2}.$$

At that optimal value η^* , the improvement is:

$$\frac{k}{k + d + 1} \frac{n^2}{n^2 + 4n + 3 + d} \|\beta_{\text{TT}}\|_2^2.$$

This completes the proof. As a final note, recall that the initial loss is $\|\beta_{\text{TT}}\|_2^2$. Thus, the new loss of the system after the test-time-training update is:

$$\left(1 - \frac{k}{k + d + 1} \frac{n^2}{n^2 + 4n + 3 + d} \right) \|\beta_{\text{TT}}\|_2^2.$$

In the proportional n, d regime and when $k/d \rightarrow \infty$, this gives an error of order $O(n^{-1})$ -matching that of the optimal \mathbf{W}^{opt} from Lemma C.1:

$$\frac{4n + d + 3}{n^2 + 4n + 3 + d} \|\beta_{\text{TT}}\|_2^2.$$

□

Corollary 4.6. Recall the definitions of $\alpha = n/d$ and $\gamma = k/d$. Consider the setting in Theorem 4.2 and test-time training described in Proposition 3.1 with both pre-trained and zero initializations. Then, under the isotropic covariance and noiseless setting, there exists a threshold γ^* given by $\gamma^* \approx \frac{(\alpha + 1)^2}{(\alpha + 2)}$ such that $\gamma < \gamma^*$ if and only if it is better to utilize the pre-trained initialization over the zero initialization $\mathbf{0}_{d \times d}$.

Proof. Considering that the initial loss of the weight matrix $\mathbf{W} = \mathbf{0}_{d \times d}$ is $\|\beta_{\text{TT}}\|_2^2$ and the improvement given by Theorem 4.5, the new loss after the test-time-training update is approximately:

$$\left(1 - \frac{k}{k + d} \right) \|\beta_{\text{TT}}\|_2^2.$$

On the other hand, recall that by Proposition 4.1, the loss of \mathbf{W}^* before the update is $\frac{d+1}{n+d+1} \|\beta_{\text{TT}}\|_2^2$. Combining this with the improvement given by Theorem 4.2, the new loss after the test-time-training update is approximately:

$$\left(\frac{d}{n + d} - \frac{k}{k + d} \frac{d^3}{(n + d)^3} \right) \|\beta_{\text{TT}}\|_2^2.$$

Let us define $\beta := \frac{k}{k+d}$ and $\theta := \frac{d}{n+d}$. Then, we check when it's better (yields smaller loss) to use the pre-trained matrix \mathbf{W}^* over zero initialization with the below inequality:

$$\begin{aligned} 1 - \beta > \theta - \beta\theta^3 &\iff 1 - \theta > \beta(1 - \theta^3) \\ &\iff \frac{1}{\beta} > \theta^2 + \theta + 1 \\ &\iff \frac{d}{k} > \theta^2 + \theta \\ &\iff \frac{d}{\frac{d}{n+d} \left(1 + \frac{d}{n+d} \right)} > k \\ &\iff \frac{(n + d)^2}{n + 2d} > k \iff \frac{(\alpha + 1)^2}{\alpha + 2} > \gamma = k/d \text{ where } \alpha = \frac{n}{d}. \end{aligned}$$

Thus, we conclude our argument. □

C. Proofs for Section 5

Lemma C.1 (Optimal \mathbf{W}). *When the system is noiseless ($\sigma^2 = 0$), the optimal \mathbf{W} which minimizes the population risk with respect to the new task β_{TT} given by Lemma B.2 and its corresponding loss are:*

$$\mathbf{W}^{\text{opt}} = \frac{\beta_{\text{TT}} \beta_{\text{TT}}^\top}{(n+2) \|\beta_{\text{TT}}\|_2^2} \quad \mathcal{L}(\mathbf{W}^{\text{opt}}) = \frac{2}{n+2} \|\beta_{\text{TT}}\|_2^2.$$

Proof. Recall that when $\Sigma_{\mathbf{x}} = \mathbf{I}$ and $\sigma^2 = 0$, the loss formula given by Lemma B.2 becomes:

$$\mathcal{L}(\mathbf{W}) = \beta_{\text{TT}}^\top \left[\mathbf{I} - n\mathbf{W} - n\mathbf{W}^\top + n(n+1)\mathbf{W}^\top \mathbf{W} + n \text{tr}(\mathbf{W}^\top \mathbf{W}) \right] \beta_{\text{TT}}. \quad (25)$$

The gradient of this expression with respect to \mathbf{W} is:

$$\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}) = -2n \beta_{\text{TT}} \beta_{\text{TT}}^\top + 2n(n+1) \mathbf{W} \beta_{\text{TT}} \beta_{\text{TT}}^\top + 2n \beta_{\text{TT}}^\top \beta_{\text{TT}} \mathbf{W} = 0.$$

Simplifying, this gives:

$$(n+1) \mathbf{W} \beta_{\text{TT}} \beta_{\text{TT}}^\top + \beta_{\text{TT}}^\top \beta_{\text{TT}} \mathbf{W} = \beta_{\text{TT}} \beta_{\text{TT}}^\top. \quad (26)$$

Multiply (26) on the right by β_{TT} . This yields:

$$\mathbf{W} \beta_{\text{TT}} = \frac{1}{n+2} \beta_{\text{TT}}.$$

So β_{TT} is an eigenvector of \mathbf{W} corresponding to the eigenvalue $1/(n+2)$. Next, consider any vector \mathbf{v} that is orthonormal (or simply orthogonal) to β_{TT} , that is $\mathbf{v}^\top \beta_{\text{TT}} = 0$. Multiplying (26) on the right by \mathbf{v} :

$$\|\beta_{\text{TT}}\|^2 \mathbf{W} \mathbf{v} = 0 \implies \mathbf{W} \mathbf{v} = 0 \quad \text{whenever} \quad \mathbf{v}^\top \beta_{\text{TT}} = 0.$$

Hence, all other eigenvalues of \mathbf{W} are 0. Using the eigendecomposition of \mathbf{W} , this uniquely specifies \mathbf{W} as $c \cdot \beta_{\text{TT}} \beta_{\text{TT}}^\top$. Solving for c yields the unique solution to $\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}) = 0$:

$$\mathbf{W}^{\text{opt}} = \frac{\beta_{\text{TT}} \beta_{\text{TT}}^\top}{(n+2) \|\beta_{\text{TT}}\|_2^2}$$

Finally, plugging this \mathbf{W}^{opt} to (25) yields the following error:

$$\mathcal{L}(\mathbf{W}^{\text{opt}}) = \frac{2}{n+2} \|\beta_{\text{TT}}\|_2^2.$$

□

Lemma C.2 (Eigenvalues). *Consider the optimal pre-trained \mathbf{W}^* matrix, which minimizes the population loss over all possible tasks sampled from Σ_{β} in (1). Then, all eigenvalues of $\bar{\mathbf{W}}^* = \Sigma_{\mathbf{x}}^{1/2} \mathbf{W}^* \Sigma_{\mathbf{x}}^{1/2}$ lie in the interval $[0, \frac{1}{n+1}]$.*

Proof. By Theorem 1 of Li et al. (2024), we write the following:

$$\mathbf{W}^* = \Sigma_{\mathbf{x}}^{-1/2} \bar{\mathbf{W}}^* \Sigma_{\mathbf{x}}^{-1/2} \quad \text{where} \quad \bar{\mathbf{W}}^* = \left((n+1) \mathbf{I}_d + M \mathbf{A}^{-1} \right)^{-1} \quad \text{where} \quad M = \text{tr}(\mathbf{A}) + \sigma^2 \quad \text{and} \quad \mathbf{A} = \Sigma_{\mathbf{x}}^{1/2} \Sigma_{\beta} \Sigma_{\mathbf{x}}^{1/2}$$

We know that the matrices $\Sigma_{\mathbf{x}}$ and Σ_{β} are jointly diagonalizable with $\Sigma_{\mathbf{x}} = \mathbf{Q} \Lambda_{\mathbf{x}} \mathbf{Q}^\top$ and $\Sigma_{\beta} = \mathbf{Q} \Lambda_{\beta} \mathbf{Q}^\top$ where \mathbf{Q} is an orthonormal matrix $\mathbf{Q} \in \mathbb{R}^{d \times d}$ and $\Lambda_{\mathbf{x}}, \Lambda_{\beta} \in \mathbb{R}^{d \times d}$ are diagonal matrices with entries $\{\lambda_i\}$ and $\{\beta_i\}$, respectively. Then,

$$\mathbf{A} = \mathbf{Q} \left(\Lambda_{\mathbf{x}}^{1/2} \Lambda_{\beta} \Lambda_{\mathbf{x}}^{1/2} \right) \mathbf{Q}^\top.$$

Define $\Lambda_{\mathbf{A}} := \Lambda_{\mathbf{x}}^{1/2} \Lambda_{\beta} \Lambda_{\mathbf{x}}^{1/2}$, which is also diagonal. Hence

$$\mathbf{A}^{-1} = \mathbf{Q} \Lambda_{\mathbf{A}}^{-1} \mathbf{Q}^\top,$$

where $\Lambda_A^{-1} = \text{diag}(1/(\lambda_i \beta_i))$. Inside \bar{W}^* , we have

$$(n+1)\mathbf{I} + M\mathbf{A}^{-1} = \mathbf{Q}(n+1)\mathbf{I}\mathbf{Q}^\top + M(\mathbf{Q}\Lambda_A^{-1}\mathbf{Q}^\top) = \mathbf{Q}[(n+1)\mathbf{I} + M\Lambda_A^{-1}]\mathbf{Q}^\top.$$

We know that $\Lambda_{\text{diag}} := (n+1)\mathbf{I} + M\Lambda_A^{-1}$ is also diagonal with diagonal entries $(n+1) + \frac{M}{\lambda_i \beta_i}$. It follows that

$$\bar{W}^* = ((n+1)\mathbf{I} + M\mathbf{A}^{-1})^{-1} = \mathbf{Q} \text{diag}\left(\frac{\lambda_i \beta_i}{(n+1)\lambda_i \beta_i + M}\right) \mathbf{Q}^\top,$$

where $M = \text{tr}(\mathbf{A}) + \sigma^2 = \text{tr}(\Lambda_\beta \Lambda_x) = \sigma^2 + \sum_{i=1}^d \lambda_i \beta_i$. This concludes the proof. \square

Lemma C.3 (Covariance Shift). *Consider the setting in Section 3 where*

- The true labels are generated by $y = \mathbf{x}^\top \boldsymbol{\beta} + \xi$ with $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma_x)$, $\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \Sigma_\beta)$, and noise $\xi \sim \mathcal{N}(0, \sigma^2)$.
- The prediction of a linear attention model is of the form

$$\hat{y} = \mathbf{x}^\top \mathbf{W} \mathbf{X}^\top \mathbf{y},$$

where the context \mathbf{X} in $\mathbb{R}^{n \times d}$ collects previous samples \mathbf{x}_i^\top as rows, \mathbf{y} is the corresponding label vector, and $\mathbf{W} \in \mathbb{R}^{d \times d}$ is the model parameter matrix.

Then, the invertible transformation

$$(\mathbf{x}, \mathbf{X}, \boldsymbol{\beta}, \mathbf{W}) \mapsto (\bar{\mathbf{x}}, \bar{\mathbf{X}}, \bar{\boldsymbol{\beta}}, \bar{\mathbf{W}}) \text{ where } \bar{\mathbf{x}} = \Sigma_x^{-1/2} \mathbf{x}, \quad \bar{\mathbf{X}} = \mathbf{X} \Sigma_x^{-1/2}, \quad \bar{\boldsymbol{\beta}} = \Sigma_x^{1/2} \boldsymbol{\beta}, \quad \bar{\mathbf{W}} = \Sigma_x^{1/2} \mathbf{W} \Sigma_x^{1/2},$$

preserves the population risks in (1) and (3). More precisely, defining $\bar{\mathcal{L}}(\bar{\mathbf{W}})$ to be the population loss (3) evaluated under the transformed setting $(\bar{\mathbf{x}}, \bar{\mathbf{X}}, \bar{\boldsymbol{\beta}}, \bar{\mathbf{W}})$, we have

$$\mathcal{L}(\mathbf{W}) = \bar{\mathcal{L}}(\bar{\mathbf{W}}).$$

Proof. Since $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma_x)$, multiplying by $\Sigma_x^{-1/2}$ yields $\bar{\mathbf{x}} = \Sigma_x^{-1/2} \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. Likewise, because $\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \Sigma_\beta)$, we have $\bar{\boldsymbol{\beta}} = \Sigma_x^{1/2} \boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \Sigma_x^{1/2} \Sigma_\beta \Sigma_x^{1/2})$. Next, we observe

$$\bar{\mathbf{X}} := \mathbf{X} \Sigma_x^{-1/2} \implies \bar{\mathbf{X}} \bar{\boldsymbol{\beta}} = (\mathbf{X} \Sigma_x^{-1/2}) (\Sigma_x^{1/2} \boldsymbol{\beta}) = \mathbf{X} \boldsymbol{\beta},$$

ensuring that label vector corresponding to context data $\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \xi = \bar{\mathbf{X}} \bar{\boldsymbol{\beta}} + \xi$ stays the same. Also, note that the label is preserved for the query token $y = \mathbf{x}^\top \boldsymbol{\beta} + \xi = \bar{\mathbf{x}}^\top \bar{\boldsymbol{\beta}} + \xi$, as well. Under the map $\mathbf{x} \mapsto \bar{\mathbf{x}} = \Sigma_x^{-1/2} \mathbf{x}$, $\mathbf{X} \mapsto \bar{\mathbf{X}} = \mathbf{X} \Sigma_x^{-1/2}$, $\boldsymbol{\beta} \mapsto \bar{\boldsymbol{\beta}} = \Sigma_x^{1/2} \boldsymbol{\beta}$, $\mathbf{W} \mapsto \bar{\mathbf{W}} = \Sigma_x^{1/2} \mathbf{W} \Sigma_x^{1/2}$, the prediction of the linear attention model is also preserved:

$$\mathbf{x}^\top \mathbf{W} \mathbf{X}^\top \mathbf{y} = (\Sigma_x^{-1/2} \mathbf{x})^\top \Sigma_x^{1/2} \mathbf{W} \Sigma_x^{1/2} (\Sigma_x^{-1/2} \mathbf{X}^\top \mathbf{y}) = \bar{\mathbf{x}}^\top (\Sigma_x^{1/2} \mathbf{W} \Sigma_x^{1/2}) (\bar{\mathbf{X}}^\top \mathbf{y}) = \bar{\mathbf{x}}^\top \bar{\mathbf{W}} \bar{\mathbf{X}}^\top \mathbf{y}.$$

As a result, the errors in the labels remain numerically unchanged as

$$(y - \mathbf{x}^\top \mathbf{W} \mathbf{X}^\top \mathbf{y})^2 = (y - \bar{\mathbf{x}}^\top \bar{\mathbf{W}} \bar{\mathbf{X}}^\top \mathbf{y})^2.$$

Hence, this implies that the population losses described in (1), (3) are preserved.

$$\mathcal{L}(\mathbf{W}) = \bar{\mathcal{L}}(\bar{\mathbf{W}}).$$

In particular, if \mathbf{W}^* is the unique minimizer of the pre-training loss in (1) in the original coordinates, its counterpart in the transformed system is precisely

$$\bar{\mathbf{W}}^* = \Sigma_x^{1/2} \mathbf{W}^* \Sigma_x^{1/2}.$$

This completes the proof. \square

Theorem C.4. Let $n/d = \Theta(1)$ and $\sigma^2 = 0$. Suppose the covariance matrices are jointly diagonalizable by an orthogonal matrix $\mathbf{Q} \in \mathbb{R}^{d \times d}$, so that $\Sigma_{\mathbf{x}} = \mathbf{Q}\Lambda_{\mathbf{x}}\mathbf{Q}^\top$ and $\Sigma_{\beta} = \mathbf{Q}\Lambda_{\beta}\mathbf{Q}^\top$. Let \mathbf{W} be any matrix that's also jointly diagonalizable by \mathbf{Q} such that if $\Sigma_{\mathbf{x}}^{1/2} \mathbf{W} \Sigma_{\mathbf{x}}^{1/2} = \mathbf{Q}\Lambda\mathbf{Q}^\top$, then all eigenvalues of Λ lie in the interval $[0, \frac{1}{n+1}]$. Define $\tilde{\beta}_{\text{TT}} := \mathbf{Q}^\top \Sigma_{\mathbf{x}}^{1/2} \beta_{\text{TT}}$, $A := \tilde{\beta}_{\text{TT}}^\top (\mathbf{I} - n\Lambda)^2 \tilde{\beta}_{\text{TT}}$, and $B := n\|\tilde{\beta}_{\text{TT}}\|_2^2 \text{tr}(\Lambda^2)$. Then, the optimal step size that minimizes the population loss given in (3) after the test-time training update is

$$\eta^* \approx \frac{A}{2(k+d)n^2\|\tilde{\beta}_{\text{TT}}\|_2^2(A+B)}.$$

With this optimal step-size η^* , the improvement in the loss due to test-time training and the initial loss are approximately

$$\mathcal{L}(\mathbf{W}) - \mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}}) \approx \frac{k}{k+d} \frac{A^2}{A+B}, \quad \mathcal{L}(\mathbf{W}) \approx A+B.$$

Proof. We consider the general non-isotropic covariance scenario where $\Sigma_{\mathbf{x}}$ and Σ_{β} may be arbitrary covariance matrices. Our analysis will hold for any initial weight matrix \mathbf{W} in $\mathbb{R}^{d \times d}$ such that if $\Sigma_{\mathbf{x}}^{1/2} \mathbf{W} \Sigma_{\mathbf{x}}^{1/2} = \mathbf{Q}\Lambda\mathbf{Q}^\top$, all eigenvalues of Λ are in $[0, \frac{1}{n+1}]$. Now, we aim to calculate the expected gain in the loss with respect to $(n+k)$ samples used during the test-time training process. That is, we will take the expectation of the loss function given by Lemma B.2 and compute $\mathcal{L}(\mathbf{W}) - \mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}}) = \mathbb{E}_{\mathbf{X}_{\text{train}}, \mathbf{X}_{\text{context}}} [\Delta \mathcal{L}]$ where $\Delta \mathcal{L}$ is given by:

$$\begin{aligned} \Delta \mathcal{L} &= \mathcal{L}(\mathbf{W}) - \mathcal{L}(\mathbf{W}_{\text{TT}}) \\ &= \beta_{\text{TT}}^\top \left[n \Sigma_{\mathbf{x}} (\mathbf{W}_{\text{TT}} - \mathbf{W}) \Sigma_{\mathbf{x}} + n \Sigma_{\mathbf{x}} (\mathbf{W}_{\text{TT}}^\top - \mathbf{W}^\top) \Sigma_{\mathbf{x}} + n(n+1) \Sigma_{\mathbf{x}} (\mathbf{W}^\top \Sigma_{\mathbf{x}} \mathbf{W} - \mathbf{W}_{\text{TT}}^\top \Sigma_{\mathbf{x}} \mathbf{W}_{\text{TT}}) \Sigma_{\mathbf{x}} \right. \\ &\quad \left. + n (\text{tr}(\mathbf{W}^\top \Sigma_{\mathbf{x}} \mathbf{W} \Sigma_{\mathbf{x}}) - \text{tr}(\mathbf{W}_{\text{TT}}^\top \Sigma_{\mathbf{x}} \mathbf{W}_{\text{TT}} \Sigma_{\mathbf{x}})) \Sigma_{\mathbf{x}} \right] \beta_{\text{TT}}. \end{aligned} \quad (27)$$

Now, recall the test-time training update from Proposition 3.1 when the new task is β_{TT} :

$$\begin{aligned} \mathbf{W}_{\text{TT}} - \mathbf{W} &= 2\eta \mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}} (\mathbf{I} - \mathbf{W} \mathbf{X}_{\text{context}}^\top \mathbf{X}_{\text{context}}) \beta_{\text{TT}} \beta_{\text{TT}}^\top \mathbf{X}_{\text{context}}^\top \mathbf{X}_{\text{context}}, \\ \mathbf{W}_{\text{TT}}^\top - \mathbf{W}^\top &= 2\eta \mathbf{X}_{\text{context}}^\top \mathbf{X}_{\text{context}} \beta_{\text{TT}} \beta_{\text{TT}}^\top (\mathbf{I} - \mathbf{X}_{\text{context}}^\top \mathbf{X}_{\text{context}} \mathbf{W}^\top) \mathbf{X}_{\text{train}}^\top \mathbf{X}_{\text{train}}. \end{aligned}$$

As in previous theorems, we will use the Gaussian approximation following Lemma B.1 by approximating $\mathbf{X}_{\text{context}}^\top \mathbf{X}_{\text{context}} \beta_{\text{TT}} \approx n \Sigma_{\mathbf{x}} \beta_{\text{TT}} + \sqrt{n} \Sigma_{\mathbf{x}}^{1/2} \mathbf{g}$ where $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \|\beta_{\text{TT}}\|^2 \mathbf{I})$. This non-isotropic form can be obtained by factoring out $\Sigma_{\mathbf{x}}^{1/2}$ and reducing the problem to the isotropic case from Lemma B.1. Before going forward, let's apply the covariance shift discussed in Lemma C.3, which transforms:

$$\bar{\mathbf{x}} = \Sigma_{\mathbf{x}}^{-1/2} \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d), \quad \bar{\beta}_{\text{TT}} = \Sigma_{\mathbf{x}}^{1/2} \beta_{\text{TT}} \sim \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{x}}^{1/2} \Sigma_{\beta} \Sigma_{\mathbf{x}}^{1/2}), \quad \text{and} \quad \bar{\mathbf{W}} = \Sigma_{\mathbf{x}}^{1/2} \mathbf{W} \Sigma_{\mathbf{x}}^{1/2}.$$

Likewise, we define the transformed versions of training and context data as $\bar{\mathbf{X}}_{\text{train}} := \mathbf{X}_{\text{train}} \Sigma_{\mathbf{x}}^{-1/2}$ and $\bar{\mathbf{X}}_{\text{context}} := \mathbf{X}_{\text{context}} \Sigma_{\mathbf{x}}^{-1/2}$. Now, suppose that $\bar{\mathbf{W}} = \mathbf{Q}\Lambda\mathbf{Q}^\top$ and let $\Delta \mathbf{W} := \bar{\mathbf{W}}_{\text{TT}} - \bar{\mathbf{W}} = \Sigma_{\mathbf{x}}^{1/2} \mathbf{W}_{\text{TT}} \Sigma_{\mathbf{x}}^{1/2} - \Sigma_{\mathbf{x}}^{1/2} \mathbf{W} \Sigma_{\mathbf{x}}^{1/2}$. Then,

$$\begin{aligned} \Delta \mathbf{W} &= 2\eta \bar{\mathbf{X}}_{\text{train}}^\top \bar{\mathbf{X}}_{\text{train}} (\bar{\beta}_{\text{TT}} - \mathbf{Q}\Lambda\mathbf{Q}^\top \bar{\mathbf{X}}_{\text{context}}^\top \bar{\mathbf{X}}_{\text{context}} \bar{\beta}_{\text{TT}}) \bar{\beta}_{\text{TT}}^\top \bar{\mathbf{X}}_{\text{context}}^\top \bar{\mathbf{X}}_{\text{context}} \\ &\approx 2\eta \bar{\mathbf{X}}_{\text{train}}^\top \bar{\mathbf{X}}_{\text{train}} (\bar{\beta}_{\text{TT}} - n\mathbf{Q}\Lambda\mathbf{Q}^\top \bar{\beta}_{\text{TT}} - \sqrt{n}\mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g}) (n\bar{\beta}_{\text{TT}}^\top + \sqrt{n}\mathbf{g}^\top) \\ &= 2\eta \bar{\mathbf{X}}_{\text{train}}^\top \bar{\mathbf{X}}_{\text{train}} (n\bar{\beta}_{\text{TT}} \bar{\beta}_{\text{TT}}^\top + \sqrt{n}\bar{\beta}_{\text{TT}} \mathbf{g}^\top - n^2\mathbf{Q}\Lambda\mathbf{Q}^\top \bar{\beta}_{\text{TT}} \bar{\beta}_{\text{TT}}^\top - n\sqrt{n}\mathbf{Q}\Lambda\mathbf{Q}^\top \bar{\beta}_{\text{TT}} \mathbf{g}^\top - n\sqrt{n}\mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g} \bar{\beta}_{\text{TT}}^\top - n\mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g} \mathbf{g}^\top) \\ &= 2\eta \bar{\mathbf{X}}_{\text{train}}^\top \bar{\mathbf{X}}_{\text{train}} (n(\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\beta}_{\text{TT}} \bar{\beta}_{\text{TT}}^\top + \sqrt{n}(\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\beta}_{\text{TT}} \mathbf{g}^\top - n\sqrt{n}\mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g} \bar{\beta}_{\text{TT}}^\top - n\mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g} \mathbf{g}^\top) \\ \Delta \mathbf{W}^\top &= 2\eta (n\bar{\beta}_{\text{TT}} \bar{\beta}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) + \sqrt{n}\mathbf{g} \bar{\beta}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) - n\sqrt{n}\bar{\beta}_{\text{TT}} \mathbf{g}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top - n\mathbf{g} \mathbf{g}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\mathbf{X}}_{\text{train}}^\top \bar{\mathbf{X}}_{\text{train}} \end{aligned}$$

By plugging the above definitions of $\bar{\mathbf{W}}_{\text{TT}}$, $\bar{\mathbf{W}}$ into Equation (27), we encounter the following difference term in two of the expressions:

$$\begin{aligned} \bar{\mathbf{W}}^\top \bar{\mathbf{W}} - \bar{\mathbf{W}}_{\text{TT}}^\top \bar{\mathbf{W}}_{\text{TT}} &= \bar{\mathbf{W}}^\top \bar{\mathbf{W}} - (\bar{\mathbf{W}}^\top + \Delta \mathbf{W}^\top)(\bar{\mathbf{W}} + \Delta \mathbf{W}) \\ &= -\bar{\mathbf{W}}^\top \Delta \mathbf{W} - \Delta \mathbf{W}^\top \bar{\mathbf{W}} - \Delta \mathbf{W}^\top \Delta \mathbf{W}. \end{aligned} \quad (28)$$

Hence, we need to calculate both the first-order expectations $\mathbb{E}_{\bar{\mathbf{X}}_{\text{train},g}}[\Delta \mathbf{W}]$, $\mathbb{E}_{\bar{\mathbf{X}}_{\text{train},g}}[\Delta \mathbf{W}^\top]$ and the second-order expectation $\mathbb{E}_{\bar{\mathbf{X}}_{\text{train},g}}[\Delta \mathbf{W}^\top \Delta \mathbf{W}]$. Calculating the first-order expectations gives:

$$\mathbb{E}_{\bar{\mathbf{X}}_{\text{train},g}}[\Delta \mathbf{W}] = 2\eta kn \left((\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top - \|\bar{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \mathbf{Q}\Lambda\mathbf{Q}^\top \right), \quad (29)$$

$$\mathbb{E}_{\bar{\mathbf{X}}_{\text{train},g}}[\Delta \mathbf{W}^\top] = 2\eta kn \left(\bar{\boldsymbol{\beta}}_{\text{TT}} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) - \|\bar{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \mathbf{Q}\Lambda\mathbf{Q}^\top \right). \quad (30)$$

Using (29) and (30), the first two terms of the loss difference are:

$$\mathbb{E}_{\bar{\mathbf{X}}_{\text{train},g}} \left[\bar{\boldsymbol{\beta}}_{\text{TT}}^\top n (\bar{\mathbf{W}}_{\text{TT}} - \bar{\mathbf{W}}) \bar{\boldsymbol{\beta}}_{\text{TT}} \right] = 2\eta kn \left(n\|\bar{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} - n\|\bar{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \bar{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top \bar{\boldsymbol{\beta}}_{\text{TT}} \right), \quad (31)$$

$$\mathbb{E}_{\bar{\mathbf{X}}_{\text{train},g}} \left[\bar{\boldsymbol{\beta}}_{\text{TT}}^\top n (\bar{\mathbf{W}}_{\text{TT}}^\top - \bar{\mathbf{W}}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \right] = 2\eta kn \left(n\|\bar{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} - n\|\bar{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \bar{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top \bar{\boldsymbol{\beta}}_{\text{TT}} \right). \quad (32)$$

Similar to proofs of previous theorems, by plugging $\mathbf{M} = \mathbf{I}$ in Lemma B.3, we know that $\mathbb{E}_{\mathbf{X}}[(\mathbf{X}^\top \mathbf{X})(\mathbf{X}^\top \mathbf{X})] = k(k+1)\boldsymbol{\Sigma}^2 + k\text{tr}(\boldsymbol{\Sigma})\boldsymbol{\Sigma}$ where $\mathbf{X} \in \mathbb{R}^{k \times d}$ has its rows drawn i.i.d from $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$. Therefore,

$$\mathbb{E}_{\bar{\mathbf{X}}_{\text{train}}} \left[\bar{\mathbf{X}}_{\text{train}}^\top \bar{\mathbf{X}}_{\text{train}} \bar{\mathbf{X}}_{\text{train}}^\top \bar{\mathbf{X}}_{\text{train}} \right] = k(k+d+1)\mathbf{I}.$$

Utilizing the above fact, we compute:

$$\begin{aligned} & \mathbb{E}_{\bar{\mathbf{X}}_{\text{train},g}} [\Delta \mathbf{W}^\top \Delta \mathbf{W}] \\ &= 4\eta^2 \mathbb{E}_g \left[\left(n\bar{\boldsymbol{\beta}}_{\text{TT}} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) + \sqrt{n} \mathbf{g} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) - n\sqrt{n} \bar{\boldsymbol{\beta}}_{\text{TT}} \mathbf{g}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top - n\mathbf{g} \mathbf{g}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top \right) \right. \\ & \quad \left. k(k+d+1)\mathbf{I} \left(n(\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top + \sqrt{n} (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \mathbf{g}^\top - n\sqrt{n} \mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top - n\mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g} \mathbf{g}^\top \right) \right] \\ &= 4\eta^2 k(k+d+1) \mathbb{E}_g \left[n^2 \bar{\boldsymbol{\beta}}_{\text{TT}} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top - n^2 \bar{\boldsymbol{\beta}}_{\text{TT}} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g} \mathbf{g}^\top \right. \\ & \quad \left. - n^2 \mathbf{g} \mathbf{g}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top + n^2 \mathbf{g} \mathbf{g}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g} \mathbf{g}^\top + n\mathbf{g} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \mathbf{g}^\top \right. \\ & \quad \left. + n^3 \bar{\boldsymbol{\beta}}_{\text{TT}} \mathbf{g}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top - n^2 \mathbf{g} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top - n^2 \bar{\boldsymbol{\beta}}_{\text{TT}} \mathbf{g}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \mathbf{g}^\top \right]. \end{aligned}$$

Let's inspect each term inside the expectation in the above sum:

$$\begin{aligned} \mathbb{E}_g \left[n^3 \bar{\boldsymbol{\beta}}_{\text{TT}} \mathbf{g}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top \right] &= n^3 \bar{\boldsymbol{\beta}}_{\text{TT}} \mathbb{E}_g \left[\mathbf{g}^\top \mathbf{Q}\Lambda^2 \mathbf{Q}^\top \mathbf{g} \right] \bar{\boldsymbol{\beta}}_{\text{TT}}^\top \\ &= n^3 \bar{\boldsymbol{\beta}}_{\text{TT}} \mathbb{E}_z \left[\mathbf{z}^\top \Lambda^2 \mathbf{z} \right] \bar{\boldsymbol{\beta}}_{\text{TT}}^\top \text{ where } \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \|\bar{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \mathbf{I}) \\ &= n^3 \|\bar{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \text{tr}(\Lambda^2) \bar{\boldsymbol{\beta}}_{\text{TT}} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top \\ \mathbb{E}_g \left[n^2 \mathbf{g} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top \right] &= n^2 \mathbb{E}_g \left[\mathbf{g} \left(\bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g} \right) \right] \bar{\boldsymbol{\beta}}_{\text{TT}}^\top \\ &= n^2 \mathbb{E}_g \left[\mathbf{g} \left(\mathbf{g}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \right) \right] \bar{\boldsymbol{\beta}}_{\text{TT}}^\top \\ &= n^2 \|\bar{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \mathbf{Q}\Lambda\mathbf{Q}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top \\ \mathbb{E}_g \left[n^2 \bar{\boldsymbol{\beta}}_{\text{TT}} \mathbf{g}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \mathbf{g}^\top \right] &= n^2 \bar{\boldsymbol{\beta}}_{\text{TT}} \mathbb{E}_g \left[\left(\mathbf{g}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \right) \mathbf{g}^\top \right] \\ &= n^2 \bar{\boldsymbol{\beta}}_{\text{TT}} \mathbb{E}_g \left[\left(\bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g} \right) \mathbf{g}^\top \right] \\ &= n^2 \|\bar{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \bar{\boldsymbol{\beta}}_{\text{TT}} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \mathbf{Q}\Lambda\mathbf{Q}^\top \\ \mathbb{E}_g \left[n\mathbf{g} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \mathbf{g}^\top \right] &= n \mathbb{E}_g \left[\mathbf{g} \left(\bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \right) \mathbf{g}^\top \right] \\ &= n \bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \mathbb{E}_g \left[\mathbf{g} \mathbf{g}^\top \right] \\ &= n \|\bar{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \left(\bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \right) \mathbf{I} \\ \mathbb{E}_g \left[n^2 \mathbf{g} \mathbf{g}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g} \mathbf{g}^\top \right] &= n^2 \mathbb{E}_g \left[\mathbf{g} \mathbf{g}^\top \mathbf{Q}\Lambda^2 \mathbf{Q}^\top \mathbf{g} \mathbf{g}^\top \right] \\ &= n^2 \|\bar{\boldsymbol{\beta}}_{\text{TT}}\|_2^4 \left(\text{tr}(\Lambda^2) \mathbf{I} + 2\mathbf{Q}\Lambda^2 \mathbf{Q}^\top \right) \\ \mathbb{E}_g \left[n^2 \mathbf{g} \mathbf{g}^\top \mathbf{Q}\Lambda\mathbf{Q}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top \right] &= n^2 \|\bar{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \mathbf{Q}\Lambda\mathbf{Q}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \bar{\boldsymbol{\beta}}_{\text{TT}} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top \\ \mathbb{E}_g \left[n^2 \bar{\boldsymbol{\beta}}_{\text{TT}} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \mathbf{Q}\Lambda\mathbf{Q}^\top \mathbf{g} \mathbf{g}^\top \right] &= n^2 \|\bar{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \bar{\boldsymbol{\beta}}_{\text{TT}} \bar{\boldsymbol{\beta}}_{\text{TT}}^\top (\mathbf{I} - n\mathbf{Q}\Lambda\mathbf{Q}^\top) \mathbf{Q}\Lambda\mathbf{Q}^\top. \end{aligned}$$

Hence, considering that the eigenvalues of Λ are smaller than $\frac{1}{n+1}$, we drop lower order terms after combining the results above:

$$\begin{aligned}\mathbb{E}_{\tilde{\mathbf{X}}_{\text{train},g}}[\Delta \mathbf{W}^\top \Delta \mathbf{W}] &= 4\eta^2 k(k+d+1) \left[n^3 \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \text{tr}(\Lambda^2) \tilde{\boldsymbol{\beta}}_{\text{TT}} \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top + n^2 \tilde{\boldsymbol{\beta}}_{\text{TT}} \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\mathbf{I} - n\Lambda)^2 \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}} \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \right. \\ &\quad - 2n^2 \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \mathbf{Q} (\Lambda - n\Lambda^2) \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}} \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top - 2n^2 \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \tilde{\boldsymbol{\beta}}_{\text{TT}} \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\Lambda - n\Lambda^2) \mathbf{Q}^\top \\ &\quad \left. + n \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 (\tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\mathbf{I} - n\Lambda)^2 \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}}) \mathbf{I} + n^2 \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^4 (\text{tr}(\Lambda^2) \mathbf{I} + 2\mathbf{Q} \Lambda^2 \mathbf{Q}^\top) \right] \\ &\approx 4\eta^2 k(k+d+1) \left[n^3 \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \text{tr}(\Lambda^2) \tilde{\boldsymbol{\beta}}_{\text{TT}} \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top + n^2 \tilde{\boldsymbol{\beta}}_{\text{TT}} \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\mathbf{I} - n\Lambda)^2 \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}} \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \right. \\ &\quad \left. + n \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 (\tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\mathbf{I} - n\Lambda)^2 \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}}) \mathbf{I} \right].\end{aligned}$$

Thus, we reach the following result:

$$\mathbb{E}_{\tilde{\mathbf{X}}_{\text{train},g}}[\tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \Delta \mathbf{W}^\top \Delta \mathbf{W} \tilde{\boldsymbol{\beta}}_{\text{TT}}] \approx 4\eta^2 k(k+d+1) \left[n^3 \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^6 \text{tr}(\Lambda^2) + (n^2 + n) \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^4 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\mathbf{I} - n\Lambda)^2 \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}} \right]. \quad (33)$$

Now, recall from Equation (28) that the second-order difference is in the form:

$$\bar{\mathbf{W}}^\top \bar{\mathbf{W}} - \bar{\mathbf{W}}_{\text{TT}}^\top \bar{\mathbf{W}}_{\text{TT}} = -\bar{\mathbf{W}}^\top \Delta \mathbf{W} - \Delta \mathbf{W}^\top \bar{\mathbf{W}} - \Delta \mathbf{W}^\top \Delta \mathbf{W}.$$

Therefore, we also need to calculate the expectations of the terms $\bar{\mathbf{W}}^\top \Delta \mathbf{W}$, $\Delta \mathbf{W}^\top \bar{\mathbf{W}}$. Using previous results (29) and (30), we get:

$$\begin{aligned}\mathbb{E}_{\tilde{\mathbf{X}}_{\text{train},g}}[-\bar{\mathbf{W}}^\top \Delta \mathbf{W}] &= -2\eta k \mathbf{Q} \Lambda \mathbf{Q}^\top \left(n (\mathbf{I} - n\mathbf{Q} \Lambda \mathbf{Q}^\top) \tilde{\boldsymbol{\beta}}_{\text{TT}} \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top - n \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \mathbf{Q} \Lambda \mathbf{Q}^\top \right) \\ &= -2\eta k n (\mathbf{Q} (\Lambda - n\Lambda^2) \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}} \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top - \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \mathbf{Q} \Lambda^2 \mathbf{Q}^\top) \\ \mathbb{E}_{\tilde{\mathbf{X}}_{\text{train},g}}[-\Delta \mathbf{W}^\top \bar{\mathbf{W}}] &= -2\eta k n (\tilde{\boldsymbol{\beta}}_{\text{TT}} \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\Lambda - n\Lambda^2) \mathbf{Q}^\top - \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \mathbf{Q} \Lambda^2 \mathbf{Q}^\top).\end{aligned}$$

This gives us:

$$\mathbb{E}_{\tilde{\mathbf{X}}_{\text{train},g}}[-\tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \bar{\mathbf{W}}^\top \Delta \mathbf{W} \tilde{\boldsymbol{\beta}}_{\text{TT}}] = -2\eta k n (\|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\Lambda - n\Lambda^2) \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}} - \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} \Lambda^2 \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}}) \quad (34)$$

$$\mathbb{E}_{\tilde{\mathbf{X}}_{\text{train},g}}[-\tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \Delta \mathbf{W}^\top \bar{\mathbf{W}} \tilde{\boldsymbol{\beta}}_{\text{TT}}] = -2\eta k n (\|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\Lambda - n\Lambda^2) \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}} - \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} \Lambda^2 \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}}). \quad (35)$$

Exploiting (33), (34) and (35), we get

$$\begin{aligned}\mathbb{E}_{\tilde{\mathbf{X}}_{\text{train},g}}[\tilde{\boldsymbol{\beta}}_{\text{TT}}^\top n(n+1) (\bar{\mathbf{W}}^\top \bar{\mathbf{W}} - \bar{\mathbf{W}}_{\text{TT}}^\top \bar{\mathbf{W}}_{\text{TT}}) \tilde{\boldsymbol{\beta}}_{\text{TT}}] \\ \approx -4\eta^2 k(k+d+1) n(n+1) \left[n^3 \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^6 \text{tr}(\Lambda^2) + (n^2 + n) \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^4 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\mathbf{I} - n\Lambda)^2 \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}} \right. \\ \left. - 4\eta k n^2 (n+1) (\|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\Lambda - n\Lambda^2) \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}} - \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} \Lambda^2 \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}}) \right] \quad (36)\end{aligned}$$

$$\begin{aligned}\mathbb{E}_{\tilde{\mathbf{X}}_{\text{train},g}}[\tilde{\boldsymbol{\beta}}_{\text{TT}}^\top n (\text{tr}(\bar{\mathbf{W}}^\top \bar{\mathbf{W}}) - \text{tr}(\bar{\mathbf{W}}_{\text{TT}}^\top \bar{\mathbf{W}}_{\text{TT}})) \tilde{\boldsymbol{\beta}}_{\text{TT}}] \\ \approx -4\eta^2 k(k+d+1) \left[n^4 \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^6 \text{tr}(\Lambda^2) + n^2 (n+d) \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^4 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\mathbf{I} - n\Lambda)^2 \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}} \right. \\ \left. - 4\eta k n^2 [\|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\Lambda - n\Lambda^2) \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}} - \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^4 \text{tr}(\Lambda^2)] \right]. \quad (37)\end{aligned}$$

Combining (31), (32), (36) and (37), and plugging into the Equation (27) gives the expected loss improvement:

$$\mathcal{L}(\mathbf{W}) - \mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}}) \quad (38)$$

$$\begin{aligned}&\approx 4\eta k n^2 (\|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\mathbf{I} - n\Lambda) \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}} - \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} \Lambda \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}}) \\ &\quad - 4\eta k n^2 (n+1) (\|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\Lambda - n\Lambda^2) \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}} - \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} \Lambda^2 \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}}) \\ &\quad - 4\eta k n^2 [\|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\Lambda - n\Lambda^2) \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}} - \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^4 \text{tr}(\Lambda^2)] \\ &\quad - 4\eta^2 k(k+d+1) n [n^3 (n+2) \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^6 \text{tr}(\Lambda^2) + n((n+1)^2 + n+d) \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^4 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\mathbf{I} - n\Lambda)^2 \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}}] \\ &\approx 4\eta k n^2 \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^2 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\mathbf{I} - n\Lambda)^2 \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}} - 4\eta^2 k(k+d+1) n^2 [n^3 \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^6 \text{tr}(\Lambda^2) + (n^2 + d) \|\tilde{\boldsymbol{\beta}}_{\text{TT}}\|_2^4 \tilde{\boldsymbol{\beta}}_{\text{TT}}^\top \mathbf{Q} (\mathbf{I} - n\Lambda)^2 \mathbf{Q}^\top \tilde{\boldsymbol{\beta}}_{\text{TT}}]. \quad (39)\end{aligned}$$

This is a quadratic expression in η , and we can solve for the optimal η^* by setting the derivative to 0. This way, the optimal step size is approximately:

$$\eta^* \approx \frac{\bar{\beta}_{\text{TT}}^\top \mathbf{Q} (\mathbf{I} - n\Lambda)^2 \mathbf{Q}^\top \bar{\beta}_{\text{TT}}}{2(k+d+1) \|\bar{\beta}_{\text{TT}}\|_2^2 \left[n^3 \|\bar{\beta}_{\text{TT}}\|_2^2 \text{tr}(\Lambda^2) + (n^2 + d) \bar{\beta}_{\text{TT}}^\top \mathbf{Q} (\mathbf{I} - n\Lambda)^2 \mathbf{Q}^\top \bar{\beta}_{\text{TT}} \right]}.$$

Plugging this optimal η^* value into Equation (39), the population loss gain becomes:

$$\frac{k}{k+d+1} \frac{(\bar{\beta}_{\text{TT}}^\top \mathbf{Q} (\mathbf{I} - n\Lambda)^2 \mathbf{Q}^\top \bar{\beta}_{\text{TT}})^2}{n \|\bar{\beta}_{\text{TT}}\|_2^2 \text{tr}(\Lambda^2) + \left(1 + \frac{d}{n^2}\right) \bar{\beta}_{\text{TT}}^\top \mathbf{Q} (\mathbf{I} - n\Lambda)^2 \mathbf{Q}^\top \bar{\beta}_{\text{TT}}}.$$

Thus, considering the definitions $\tilde{\beta}_{\text{TT}} = \mathbf{Q}^\top \bar{\beta}_{\text{TT}}$ and $A = \tilde{\beta}_{\text{TT}}^\top (\mathbf{I} - n\Lambda)^2 \tilde{\beta}_{\text{TT}}$, $B = n \text{tr}(\Lambda^2) \|\tilde{\beta}_{\text{TT}}\|_2^2$ (notice that ℓ_2 norm is unitarily-invariant, so that $\|\tilde{\beta}_{\text{TT}}\| = \|\bar{\beta}_{\text{TT}}\|$) and recalling $n/d = \Theta(1)$ recovers the desired final expression. Also, by viewing the initial task as $\tilde{\beta}_{\text{TT}}$, we match the diagonal covariance form described in Section 5. Therefore, we conclude our argument

$$\mathcal{L}(\mathbf{W}) - \mathcal{L}_{\text{TT}}(\mathbf{W}_{\text{TT}}) \approx \frac{k}{k+d} \frac{A^2}{A+B}.$$

For the second part of the proposition, let's write the initial loss by applying the covariance shift from Lemma C.3 so that the transformed features $\tilde{\mathbf{x}} = \Sigma_{\mathbf{x}}^{-1/2} \mathbf{x}$ have identity covariance and $\tilde{\mathbf{W}} = \Sigma_{\mathbf{x}}^{1/2} \mathbf{W}^* \Sigma_{\mathbf{x}}^{1/2}$ is the corresponding minimizer of the transformed system. Again, we set $\tilde{\beta}_{\text{TT}} = \mathbf{Q}^\top \bar{\beta}_{\text{TT}}$ in order to have diagonal covariances. Assuming $\sigma = 0$ and plugging these in the population loss formula from Lemma B.2 with respect to the task $\tilde{\beta}_{\text{TT}}$ gives:

$$\begin{aligned} \mathcal{L}(\mathbf{W}) &= \tilde{\beta}_{\text{TT}}^\top \left[\mathbf{I} - n\mathbf{I} \tilde{\mathbf{W}} \mathbf{I} - n\mathbf{I} \tilde{\mathbf{W}}^\top \mathbf{I} + n(n+1) \mathbf{I} \tilde{\mathbf{W}}^\top \mathbf{I} \tilde{\mathbf{W}} \mathbf{I} + n \text{tr}(\tilde{\mathbf{W}}^\top \mathbf{I} \tilde{\mathbf{W}} \mathbf{I}) \right] \tilde{\beta}_{\text{TT}} + \sigma^2 n \text{tr}(\tilde{\mathbf{W}}^\top \mathbf{I} \tilde{\mathbf{W}} \mathbf{I}) + \sigma^2 \\ &= \|\tilde{\beta}_{\text{TT}}\|_2^2 - 2n \tilde{\beta}_{\text{TT}}^\top \mathbf{Q} \Lambda \mathbf{Q}^\top \tilde{\beta}_{\text{TT}} + n(n+1) \tilde{\beta}_{\text{TT}}^\top \mathbf{Q} \Lambda^2 \mathbf{Q}^\top \tilde{\beta}_{\text{TT}} + n \text{tr}(\Lambda^2) \|\tilde{\beta}_{\text{TT}}\|_2^2 \\ &\approx \|\tilde{\beta}_{\text{TT}}\|_2^2 - 2n \tilde{\beta}_{\text{TT}}^\top \mathbf{Q} \Lambda \mathbf{Q}^\top \tilde{\beta}_{\text{TT}} + n^2 \tilde{\beta}_{\text{TT}}^\top \mathbf{Q} \Lambda^2 \mathbf{Q}^\top \tilde{\beta}_{\text{TT}} + n \text{tr}(\Lambda^2) \|\tilde{\beta}_{\text{TT}}\|_2^2 \\ &= \tilde{\beta}_{\text{TT}}^\top \mathbf{Q} (\mathbf{I} - n\Lambda)^2 \mathbf{Q}^\top \tilde{\beta}_{\text{TT}} + n \|\tilde{\beta}_{\text{TT}}\|_2^2 \text{tr}(\Lambda^2) \\ &= \tilde{\beta}_{\text{TT}}^\top (\mathbf{I} - n\Lambda)^2 \tilde{\beta}_{\text{TT}} + n \|\tilde{\beta}_{\text{TT}}\|_2^2 \text{tr}(\Lambda^2) \\ &= A + B. \end{aligned}$$

Hence, $\mathcal{L}(\mathbf{W}) \approx A + B$. The proof is complete. \square

Proposition 5.2. Let $\Sigma_{\mathbf{x}} = \mathbf{I}$, $\sigma^2 = 0$ and suppose Σ_{β} is an arbitrary diagonal covariance matrix with the first r diagonal entries being non-zero. Then, the minimizer \mathbf{W}^* of the pre-training loss (1) that has minimal Frobenius norm and its population loss are given by

$$\mathbf{W}^* = \left((n+1) \mathbf{I}' + \text{tr}(\Sigma_{\beta}) \Sigma_{\beta}^{\dagger} \right)^{\dagger}, \quad \mathcal{L}(\mathbf{W}^*) \approx A + B,$$

where $\mathbf{I}' = \text{diag}(\mathbf{1}_r, \mathbf{0}_{d-r})$.

Proof. While Li et al. (2024) derives the solution form in the full-rank case, the same closed-form expression can be adapted when Σ_{β} is low-rank. Indeed, if Σ_{β} has zeros in its last $d-r$ diagonal entries, those coordinates do not appear in the pre-training loss. Consequently, there is a set of minimizers \mathbf{W} which differ only in those degenerate directions. Among these, the minimal-Frobenius-norm criterion forces all degenerate coordinates to be zero as any nonzero component in that subspace would increase $\|\mathbf{W}\|_F$. Consequently,

$$\mathbf{W}^* = \left((n+1) \mathbf{I}' + \text{tr}(\Sigma_{\beta}) \Sigma_{\beta}^{\dagger} \right)^{-1}$$

zeroes out precisely the degenerate directions and is the unique Frobenius-minimal solution.

For the second part of the proposition, the argument directly follows from the proof of Theorem 5.3. \square

Corollary 5.4. Consider the setting in Theorem 5.3. Recall the definitions of $\alpha = n/d$ and $\gamma = k/d$ and define $c_1 = \frac{A}{\|\beta_{\text{TT}}\|_2^2}$,

$c_2 = \frac{B}{\|\beta_{\text{TT}}\|_2^2}$. Then, there exists a phase transition point $\gamma^* \approx \frac{(c_1 + c_2) - (c_1 + c_2)^2}{c_2 (2c_1 + c_2)}$ such that $\gamma < \gamma^*$ if and only if it is better to utilize the pre-trained \mathbf{W}^* over the null initialization $\mathbf{0}_{d \times d}$.

Proof. Recall the population loss for \mathbf{W} given by Lemma B.2:

$$\mathcal{L}(\mathbf{W}) = \beta_{\text{TT}}^\top \left[\Sigma_x - n \Sigma_x \mathbf{W} \Sigma_x - n \Sigma_x \mathbf{W}^\top \Sigma_x + n(n+1) \Sigma_x \mathbf{W}^\top \Sigma_x \mathbf{W} \Sigma_x + n \text{tr}(\mathbf{W}^\top \Sigma_x \mathbf{W} \Sigma_x) \Sigma_x \right] \beta_{\text{TT}} + \sigma^2 n \text{tr}(\mathbf{W}^\top \Sigma_x \mathbf{W} \Sigma_x) + \sigma^2.$$

First, plugging $\mathbf{W} = \mathbf{0}_{d \times d}$ yields the initial loss of $\|\tilde{\beta}_{\text{TT}}\|_2^2$. Also, setting $\bar{\mathbf{W}}^* = \mathbf{0}_{d \times d}$ in Theorem 5.3, we know that the improvement by test-time-training is approximately:

$$\frac{k}{k+d} \|\tilde{\beta}_{\text{TT}}\|_2^2.$$

At the same time, still by Theorem 5.3, we know that if our initial weight matrix is the pre-trained \mathbf{W}^* the improvement is approximately:

$$\frac{k}{k+d} \frac{A^2}{A+B}.$$

Whereas, the corresponding initial loss is approximately $\mathcal{L}(\mathbf{W}^*) \approx A+B$. Therefore, we check when it's better to use pre-trained matrix \mathbf{W}^* over $\mathbf{0}_{d \times d}$ (null) matrix with the below inequality, which compares the losses after the test-time-training update:

$$A+B - \frac{k}{k+d} \frac{A^2}{A+B} < \|\tilde{\beta}_{\text{TT}}\|_2^2 - \frac{k}{k+d} \|\tilde{\beta}_{\text{TT}}\|_2^2.$$

In lines with the proof of Corollary 4.6, denote $\beta := \frac{k}{k+d}$. Then, a series of algebraic manipulations give:

$$\begin{aligned} A+B - \beta \frac{A^2}{A+B} < \|\tilde{\beta}_{\text{TT}}\|_2^2 - \beta \|\tilde{\beta}_{\text{TT}}\|_2^2 &\iff c_1 + c_2 - \beta \frac{c_1^2}{c_1 + c_2} < 1 - \beta \\ &\iff \beta \left(1 - \frac{c_1^2}{c_1 + c_2} \right) < 1 - (c_1 + c_2) \\ &\iff \frac{1}{\beta} > \frac{c_1 + c_2 - c_1^2}{(c_1 + c_2)(1 - (c_1 + c_2))} \\ &\iff 1 + \frac{d}{k} > 1 + \frac{c_2(2c_1 + c_2)}{c_1 + c_2 - (c_1 + c_2)^2} \\ &\iff \frac{k}{d} = \gamma < \frac{c_1 + c_2 - (c_1 + c_2)^2}{c_2(2c_1 + c_2)}. \end{aligned}$$

This completes our argument. □

D. Further Experimental Results for Section 6

To illustrate the improvement more clearly, we also provide a version of Figure 3a with the x-axis on a log scale, shown in Figure 4.

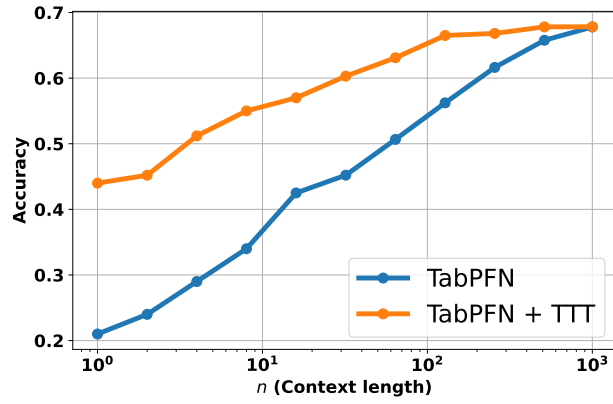


Figure 4. Accuracy of TabPFN model with and without test-time-training as a function of number of in-context samples n with the x-axis in log scale.