
Towards Understanding Out-of-Distribution Generalization for In-Context Learning via Low-Dimensional Subspaces

Soo Min Kwon*
University of Michigan

Alec S. Xu*
University of Michigan

Can Yaras
University of Michigan

Laura Balzano
University of Michigan

Qing Qu
University of Michigan

Abstract

The transformer’s remarkable ability to perform in-context learning (ICL) has sparked a wide range of studies designed to understand its strengths and limitations. However, a theoretical understanding of when ICL can and cannot generalize beyond its pre-training data still remains unclear. This paper puts forth a minimal mathematical model that provably identifies when ICL can generalize out-of-distribution (OOD). By studying linear regression tasks parameterized with low-rank covariance matrices, we model distribution shifts as varying angles between subspaces and derive conditions under which a single-layer linear attention model interpolates across all angles. We show that if pre-training task vectors are drawn from a union of subspaces, transformers can generalize to all angle shifts, enabling ICL even in regions with zero probability mass in the training distribution. On the other hand, if the pre-training tasks are drawn from a single Gaussian, the test risk shows a non-negligible dependence on the angle, implying that ICL cannot generalize OOD. We empirically show that our results also hold for models such as GPT-2, and present experiments on how our results extend to nonlinear function classes.

Lopez-Paz, 2023; Zhang et al., 2024; Li et al., 2024a; Pan et al., 2023; Kossen et al., 2024; Zhang et al., 2024; Huang et al., 2024; Li et al., 2024a; Akyürek et al., 2023; Von Oswald et al., 2023; Ahn et al., 2023; Li et al., 2024b). However, its generalization capabilities, particularly whether ICL can generalize beyond its pre-training data, remain unclear. There are conflicting views in the literature: Garg et al. (2022); Zhang et al. (2024) showed that ICL is relatively robust to distribution shifts in several settings, while Wang et al. (2025) empirically demonstrated that ICL can only solve in-distribution (language) tasks in general. To resolve these contrasting perspectives, Goddard et al. (2025) presented yet another setting in which ICL can generalize OOD, attributing this capability to pre-training task diversity. However, the complex nature of their definition of the task vector made theoretical analysis difficult, restricting their study to empirical results. Overall, these observations highlight the lack of a theoretical framework that clearly explains when ICL can and cannot generalize OOD.

In this work, we present a mathematical model to demystify and quantify the OOD generalization capabilities of ICL. We primarily focus on task distribution shifts by studying ICL in a single-layer linear attention model performing linear regression, where the weight (or task) vectors are sampled from low-dimensional subspaces. This allows us to quantify the distribution shift in the task vector via the principal angles between subspaces, and characterize the OOD test risk as a function of these angles. Unlike Goddard et al. (2025), we then provide theoretical guarantees by proving conditions on the pre-training task vectors under which OOD generalization is possible. Furthermore, we empirically show that our findings extend beyond linear settings, in that (i) our results hold for nonlinear transformers such as GPT-2; and (ii) our results apply to nonlinear function classes.

1 INTRODUCTION

The remarkable capability of ICL in transformer-based large language models (LLMs) (Vaswani et al., 2017) has sparked a wide range of both empirical and theoretical research dedicated to understanding its foundations (Garg et al., 2022; Raventós et al., 2023; Yadlowsky et al., 2023; Wang et al., 2025; Ahuja and

2 PROBLEM SETUP

2.1 Preliminaries

Given a sequence of n input-output example pairs $\{\mathbf{x}_i, y_i\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$, the objective of ICL is to predict the output $y_{n+1} \in \mathbb{R}$ corresponding to an unseen query $\mathbf{x}_{n+1} \in \mathbb{R}^d$. Following prior works (Garg et al., 2022), we assume each output is generated via $y_i = f(\mathbf{x}_i)$ for some function $f(\cdot)$, where $f \in \mathcal{F}$ is sampled from a distribution over a function class \mathcal{F} . By convention, a transformer takes in these $n+1$ pairs via the following input prompt $\mathbf{Z} \in \mathbb{R}^{(n+1) \times (d+1)}$:

$$\mathbf{Z} = [\mathbf{z}_1 \ \dots \ \mathbf{z}_n \ \mathbf{z}_{n+1}]^\top = \begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_n & \mathbf{x}_{n+1} \\ y_1 & \dots & y_n & 0 \end{bmatrix}^\top,$$

where $\mathbf{z}_i := [\mathbf{x}_i^\top \ y_i]^\top$ and $\mathbf{z}_{n+1} := [\mathbf{x}_{n+1}^\top \ 0]^\top$. Then, the transformer g_{ATT} , parameterized by weights \mathcal{W} , is trained by minimizing the following expected squared loss with respect to \mathcal{W} :

$$\min_{\mathcal{W}} \mathcal{L}_{\text{ATT}}(\mathcal{W}) := \mathbb{E} \left[(y_{n+1} - g_{\text{ATT}}(\mathbf{z}_{n+1}, \mathbf{Z}))^2 \right]. \quad (1)$$

At inference, we test the trained model g_{ATT}^* using $m+1$ paired examples $\{\mathbf{x}_j, \tilde{y}_j\}_{j=1}^{m+1}$, where the prompts are constructed in the same manner:

$$\tilde{\mathbf{Z}} = \begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_n & \mathbf{x}_{n+1} \\ \tilde{y}_1 & \dots & \tilde{y}_n & 0 \end{bmatrix}^\top \quad \text{and} \quad \tilde{\mathbf{z}}_{m+1} = \begin{bmatrix} \mathbf{x}_{m+1} \\ 0 \end{bmatrix}.$$

Since we are interested in the OOD capabilities of transformers, the labels are generated via $\tilde{y}_j = \tilde{f}(\mathbf{x}_j)$ for some function $\tilde{f} \neq f$, i.e., the task function between training and testing are different.

2.2 Single-Layer Linear Attention

We now introduce the linear attention architecture. First, define the masked input prompts as

$$\mathbf{Z}_{\mathcal{M}} = [\mathbf{z}_1 \ \dots \ \mathbf{z}_n \ \mathbf{0}]^\top \quad \text{and} \quad \mathbf{z}_{n+1} = \begin{bmatrix} \mathbf{x}_{n+1} \\ 0 \end{bmatrix}.$$

Given the input sequence $\mathbf{Z}_{\mathcal{M}}$ and query \mathbf{z}_{n+1} , the predicted label \hat{y}_{m+1} is the output of a linear attention model g_{ATT} , where

$$g_{\text{ATT}}(\mathbf{z}_{n+1}, \mathbf{Z}_{\mathcal{M}}) = \frac{(\mathbf{z}_{n+1}^\top \mathbf{W}_Q \mathbf{W}_K^\top \mathbf{Z}_{\mathcal{M}}^\top) \mathbf{Z}_{\mathcal{M}} \mathbf{W}_V \mathbf{p}}{n}, \quad (2)$$

$\mathbf{p} = [\mathbf{0}_d \ 1]^\top$, and $\mathbf{W}_K, \mathbf{W}_Q, \mathbf{W}_V \in \mathbb{R}^{(d+1) \times (d+1)}$ are the key, query, and value weight matrices, respectively, meaning $\mathcal{W} = \{\mathbf{W}_K, \mathbf{W}_Q, \mathbf{W}_V\}$. Then, letting $\mathcal{W}^* = \{\mathbf{W}_K^*, \mathbf{W}_Q^*, \mathbf{W}_V^*\}$ denote the optimal weights obtained via Equation (1). At inference, the predicted test-time

label \hat{y}_{m+1} is the output of the trained linear attention model g_{att}^* ($\tilde{\mathbf{z}}_{m+1}, \tilde{\mathbf{Z}}_{\mathcal{M}}$), where

$$g_{\text{att}}^* \left(\tilde{\mathbf{z}}_{m+1}, \tilde{\mathbf{Z}}_{\mathcal{M}} \right) = \frac{(\mathbf{z}_{n+1}^\top \mathbf{W}_Q^* \mathbf{W}_K^{*\top} \mathbf{Z}_{\mathcal{M}}^\top) \mathbf{Z}_{\mathcal{M}} \mathbf{W}_V^* \mathbf{p}}{m}. \quad (3)$$

2.3 ICL with Linear Regression

We consider linear regression tasks, i.e., $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$ for some input $\mathbf{x} \in \mathbb{R}^d$ and task vector $\mathbf{w} \in \mathbb{R}^d$.

Training Distribution. Suppose that $d \geq 2r$ and let $\mathbf{U}_s \in \mathbb{R}^{d \times r}$ and $\mathbf{U}_{s,\perp} \in \mathbb{R}^{d \times r}$ be two r -dimensional orthonormal bases in \mathbb{R}^d . Consider the following two covariance matrices:

$$\begin{aligned} \Sigma_s &= \mathbf{U}_s \mathbf{U}_s^\top + \epsilon \cdot \mathbf{I}_d \\ \Sigma_{s,\perp} &= \mathbf{U}_{s,\perp} \mathbf{U}_{s,\perp}^\top + \epsilon \cdot \mathbf{I}_d, \end{aligned}$$

For training, we consider two separate models trained on two different task vector distributions. Specifically, each feature and label pair (\mathbf{x}_i, y_i) is generated as follows. For all $i \in [n+1]$, let $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and

$$y_i = \mathbf{w}^\top \mathbf{x}_i + \eta_i, \quad (4)$$

where $\eta_i \sim \mathcal{N}(0, \sigma^2)$ is iid noise with variance σ^2 . We consider the following distributions for the training task vector \mathbf{w} :

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_s), \quad (\text{Single Subspace})$$

or

$$\mathbf{w} \sim \begin{cases} \mathcal{N}(\mathbf{0}, \Sigma_s) & \text{w.p. } \gamma, \\ \mathcal{N}(\mathbf{0}, \Sigma_{s,\perp}) & \text{w.p. } 1 - \gamma, \end{cases} \quad (\text{Union of Subs.})$$

where $0 < \gamma < 1$ denotes the mixture probability. Under these two distributions, we show that the trained models exhibit different OOD generalization behaviors with respect to the testing distribution defined in the following section. Finally, in the union of subspaces distribution, while we focus on $K = 2$ bases for ease of exposition, we also generalize our results to consider a union of $K > 2$ bases in Section A.

Testing Distribution. We define a testing subspace $\mathbf{U}_t \in \mathbb{R}^{d \times r}$ such that it represents a region between the training subspaces. To this end, we parameterize \mathbf{U}_t as such (Absil et al., 2004, Section 3.8):

$$\mathbf{U}_t = \mathbf{U}_s \cdot \cos(\theta) + \mathbf{U}_{s,\perp} \cdot \sin(\theta), \quad (5)$$

where $\theta \in [0, \frac{\pi}{2}]$ denotes the (r identical) principal angles between \mathbf{U}_s and \mathbf{U}_t . Finally, we define the testing covariance Σ_t as

$$\Sigma_t = \mathbf{U}_t \mathbf{U}_t^\top + \epsilon \cdot \mathbf{I}_d. \quad (6)$$

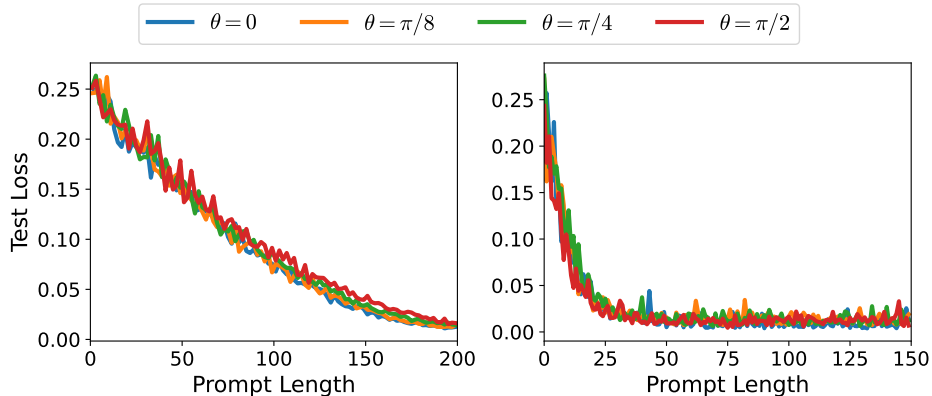


Figure 1: Plot of the test risk as a function of the prompt length for a linear transformer (left) and a GPT-2 model (right) under the data setting of Theorem 1. When the prompt length at test time is large enough, the test risk goes nearly to zero for all $\theta \in [0, \frac{\pi}{2}]$, corroborating Theorem 1.

We generate each testing pair $(\mathbf{x}_j, \tilde{y}_j)$ independent of the training data in a similar fashion: for all $j \in [m+1]$, let $\mathbf{x}_j \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and

$$\tilde{y}_j = \tilde{\mathbf{w}}^\top \mathbf{x}_j + \eta_j, \quad \text{where } \tilde{\mathbf{w}} \sim \mathcal{N}(\mathbf{0}, \Sigma_t) \quad (7)$$

and $\eta_j \sim \mathcal{N}(0, \sigma^2)$ is again iid noise.

3 MAIN RESULTS

3.1 Transformers Can Generalize to the Span When Trained on a Union of Subspaces

In this section, we show that when the linear attention model is trained on prompts whose task vectors are sampled via Equation (Union of Subs.), and tested on a task vector sampled via Equation (7) at any $\theta \in [0, \frac{\pi}{2}]$, the test risk can be arbitrarily close to the optimal test risk (the label noise variance). Notably, this test risk is *independent* of θ , implying linear attention is robust to such subspace shifts in this setting.

Theorem 1 (Test Risk Under a Union of Subspaces). *Let g_{ATT}^* denote the optimal linear attention model corresponding to the independent data setting in Equation (4), where the task vector is drawn from Equation (Union of Subs.) with $\gamma = 0.5$. For all $j \in [m+1]$, suppose that the prompts at test time are constructed with features $\mathbf{x}_j \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and labels whose task vectors are drawn from Equation (6). For any $\theta \in [0, \frac{\pi}{2}]$ and $\delta \in (0, r)$, if*

$$m \geq n > \frac{(2(r + \sigma^2) + 1)r}{\delta} - (2(r + \sigma^2) + 1),$$

then we have

$$\lim_{\epsilon \rightarrow 0} \mathbb{E} \left[\left(\tilde{y}_{m+1} - g_{ATT}^*(\tilde{\mathbf{z}}_{m+1}, \tilde{\mathbf{Z}}_{\mathcal{M}}) \right)^2 \right] < \sigma^2 + \delta.$$

Task vectors sampled from the test subspace \mathbf{U}_t have zero probability density in the training task distribution. Hence, Theorem 1 implies ICL generalizes well to this OOD task. We hypothesize this can explain why ICL achieves OOD generalization: the test data actually lies within the span of the training data. We empirically corroborate Theorem 1 in Figure 1 on both linear attention and GPT-2, deferring experimental details to Section B, and the proof to Section C.2.

3.2 Transformers Cannot Generalize When Trained on a Single Subspace

We now aim to identify a setting in which generalization beyond the pre-training data is not possible. The following result shows that if the task vectors are drawn from a Gaussian whose covariance matrix spans only a single subspace $\mathbf{U}_s \in \mathbb{R}^{d \times r}$, then testing a linear attention model on a task vector shifted away from the training subspace by an angle θ yields a test risk with a non-negligible dependence on θ , even as the prompt length tends to infinity.

Proposition 1 (Test Risk Under a Single Subspace). *Let g_{ATT}^* denote the optimal linear attention model corresponding to the independent data setting in Equation (4), where the task vector is drawn from Equation (Single Subspace). For all $j \in [m+1]$, suppose that the prompts at test time are constructed with features $\mathbf{x}_j \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and labels whose task vec-*

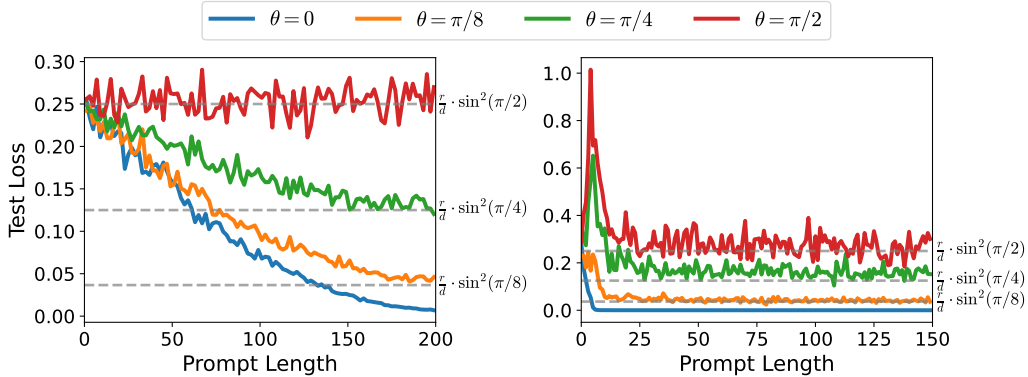


Figure 2: Plot of the normalized test risk as a function of the prompt length for a linear transformer (left) and a GPT-2 model (right) under the data setting of Proposition 1. The test risk exactly matches the predicted risk from Proposition 1.

tors are drawn from Equation (6). Then, we have

$$\lim_{m \rightarrow \infty} \lim_{n \rightarrow \infty} \lim_{\epsilon \rightarrow 0} \mathbb{E} \left[\left(\tilde{y}_{m+1} - g_{ATT}^*(\tilde{\mathbf{z}}_{m+1}, \tilde{\mathbf{Z}}_{\mathcal{M}}) \right)^2 \right] = r \sin^2(\theta) + \sigma^2,$$

where $\theta \in [0, \frac{\pi}{2}]$ are the r principal angles between $\mathbf{U}_s \in \mathbb{R}^{d \times r}$ and $\mathbf{U}_t \in \mathbb{R}^{d \times r}$.

When $\theta = 0$, the $\sin(\cdot)$ term vanishes, allowing perfect recovery up to the label noise variance. However, as θ increases from 0 to $\frac{\pi}{2}$, the test risk increases with respect to θ . At $\theta = \frac{\pi}{2}$, the test risk becomes exactly the rank of the covariance matrix, which is the largest possible error in this setting (Garg et al., 2022; Oko et al., 2025). We empirically corroborate Proposition 1 in Figure 2, again deferring experimental details to Section B, and the proof to Section C.3.

4 EXPERIMENTS

In this section, we provide experimental results showing our theoretical results hold beyond linear function classes. We defer specific experimental details to Section B. We look at square-integrable functions under the uniform measure. We construct an orthonormal basis via cosines, i.e., $\psi_n(x) = (1/\sqrt{2}) \cos(n\pi x)$ for $n \in \mathbb{N}$. As in previous sections, we consider two settings: observing instances of a single (one-dimensional) subspace, as well as for a union of three (one-dimensional) subspaces. We draw the inputs x via $x \sim \mathcal{U}([0, 1])$. The results are shown in Figure 3: as seen on the left, transformers are not robust to subspace shifts, as the test risk increases with the subspace angle from the training subspace, in accordance with Proposition 1. However, in Figure 3 (right), we have

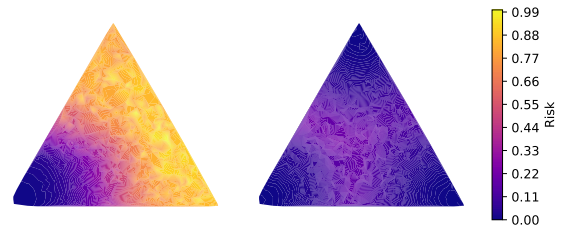


Figure 3: Visualization of the generalization behavior of transformers for learning cosines in-context. Each corner of a triangle represents a one-dimensional subspace spanned by ψ_1 (bottom left), ψ_2 (bottom right), or ψ_3 (top), with all possible convex combinations given by the interior. We show the risk when evaluated at different points in $\text{span}(\{\psi_1, \psi_2, \psi_3\})$ for the appropriate function space. **Left:** train on prompts drawn from $\text{span}(\{\psi_1\})$. **Right:** train on prompts drawn from $\text{span}(\{\psi_1\}) \cup \text{span}(\{\psi_2\}) \cup \text{span}(\{\psi_3\})$.

the generalization behavior described by Theorem 1, where training on the mixture of subspaces results in low risk in the space spanned by the basis vectors.

5 CONCLUSION

In this work, we proposed a mathematical framework to analyze the OOD capabilities of ICL. Unlike prior studies, our framework provides theoretical guarantees for when linear transformers can or cannot generalize OOD, based on their pre-training data, which we parameterize using low-dimensional subspaces. We show that pre-training task diversity, defined as a union of subspaces, enables generalization to regions with zero probability density in the training distribution, whereas tasks drawn from a single subspace do not have this property.

References

- Absil, P.-A., Mahony, R., and Sepulchre, R. (2004). Riemannian geometry of grassmann manifolds with a view on algorithmic computation. *Acta Applicandae Mathematica*, 80(2):199–220.
- Ahn, K., Cheng, X., Daneshmand, H., and Sra, S. (2023). Transformers learn to implement preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing Systems*, 36:45614–45650.
- Ahuja, K. and Lopez-Paz, D. (2023). A closer look at in-context learning under distribution shifts. *arXiv preprint arXiv:2305.16704*.
- Akyürek, E., Schuurmans, D., Andreas, J., Ma, T., and Zhou, D. (2023). What learning algorithm is in-context learning? investigations with linear models. *The Eleventh International Conference on Learning Representations*.
- Garg, S., Tsipras, D., Liang, P., and Valiant, G. (2022). What can transformers learn in-context? A case study of simple function classes. In *Advances in Neural Information Processing Systems*.
- Goddard, C., Smith, L. M., Ngampruetikorn, V., and Schwab, D. J. (2025). When can in-context learning generalize out of task distribution? In *Forty-second International Conference on Machine Learning*.
- Huang, Y., Cheng, Y., and Liang, Y. (2024). In-context convergence of transformers. In *International Conference on Machine Learning*, pages 19660–19722. PMLR.
- Kossen, J., Gal, Y., and Rainforth, T. (2024). In-context learning learns label relationships but is not conventional learning. In *The Twelfth International Conference on Learning Representations*.
- Li, H., Wang, M., Lu, S., Cui, X., and Chen, P.-Y. (2024a). How do nonlinear transformers learn and generalize in in-context learning? In *International Conference on Machine Learning*, pages 28734–28783. PMLR.
- Li, Y., Rawat, A. S., and Oymak, S. (2024b). Fine-grained analysis of in-context linear estimation: Data, architecture, and beyond. *arXiv preprint arXiv:2407.10005*.
- Mattei, P.-A. (2017). Multiplying a gaussian matrix by a gaussian vector. *Statistics & Probability Letters*, 128:67–70.
- Oko, K., Song, Y., Suzuki, T., and Wu, D. (2025). Pretrained transformer efficiently learns low-dimensional target functions in-context. *Advances in Neural Information Processing Systems*, 37:77316–77365.
- Pan, J., Gao, T., Chen, H., and Chen, D. (2023). What in-context learning ”learns” in-context: Disentangling task recognition and task learning. In *The 61st Annual Meeting Of The Association For Computational Linguistics*.
- Petersen, K. B., Pedersen, M. S., et al. (2008). The matrix cookbook. *Technical University of Denmark*, 7(15):510.
- Raventós, A., Paul, M., Chen, F., and Ganguli, S. (2023). Pretraining task diversity and the emergence of non-bayesian in-context learning for regression. *Advances in neural information processing systems*, 36:14228–14246.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.
- Von Oswald, J., Niklasson, E., Randazzo, E., Sacramento, J., Mordvintsev, A., Zhmoginov, A., and Vladymyrov, M. (2023). Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pages 35151–35174. PMLR.
- Wang, Q., Wang, Y., Wang, Y., and Ying, X. (2025). Can in-context learning really generalize to out-of-distribution tasks? In *The Thirteenth International Conference on Learning Representations*.
- Yadlowsky, S., Doshi, L., and Tripuraneni, N. (2023). Pretraining data mixtures enable narrow model selection capabilities in transformer models. *arXiv preprint arXiv:2311.00871*.
- Zhang, R., Frei, S., and Bartlett, P. L. (2024). Trained transformers learn linear models in-context. *Journal of Machine Learning Research*, 25(49):1–55.

Supplementary Materials

A Beyond a Mixture of Two Gaussians

We now generalize Theorem 1 to consider a mixture of $K > 2$ Gaussians. Assume $d \geq Kr$, and let $\mathbf{U} = [\mathbf{u}_1 \dots \mathbf{u}_d] \in \mathbb{R}^{d \times d}$ be an orthonormal basis for \mathbb{R}^d . For all $k \in [K]$, we define $\mathbf{U}_{s,k} = [\mathbf{u}_{(k-1) \cdot r+1} \dots \mathbf{u}_{kr}] \in \mathbb{R}^{d \times r}$. Note that $\mathbf{U}_{s,k}^\top \mathbf{U}_{s,l} = \mathbf{0}_{r \times r}$ for all $k \neq l$. Then, we assume the training task $\mathbf{w} \in \mathbb{R}^d$ is sampled as such:

$$\mathbf{w} \sim \sum_{k=1}^K \gamma_k \cdot \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{s,k}), \quad (8)$$

where $\boldsymbol{\Sigma}_{s,k} = \mathbf{U}_{s,k} \mathbf{U}_{s,k}^\top + \epsilon \cdot \mathbf{I}_d$ and $\sum_{k=1}^K \gamma_k = 1$. At inference time, we define an orthonormal basis $\bar{\mathbf{U}}_t \in \mathbb{R}^{d \times r}$ that lies within the span of $\{\mathbf{U}_{s,k}\}_{k=1}^K$:

$$\bar{\mathbf{U}}_t = \sum_{k=1}^K \alpha_k \mathbf{U}_{s,k}, \text{ for } \{\alpha_k\}_{k=1}^K \text{ s.t. } \sum_{k=1}^K \alpha_k^2 = 1, \quad (9)$$

where the constraint on $\{\alpha_k\}_{k=1}^K$ ensures $\bar{\mathbf{U}}_t \in \mathbb{R}^{d \times r}$ is an orthonormal basis. Then, similar to Theorem 1, we consider testing on task vectors $\tilde{\mathbf{w}} \sim \mathcal{N}(\mathbf{0}, \bar{\boldsymbol{\Sigma}}_t)$ with $\bar{\boldsymbol{\Sigma}}_t = \bar{\mathbf{U}}_t \bar{\mathbf{U}}_t^\top + \epsilon \cdot \mathbf{I}_d$. We again emphasize $\bar{\mathbf{U}}_t$ is unseen during training, but lies within the span of the training subspaces.

Theorem 2. Let g_{ATT}^* denote the optimal linear attention model corresponding to the independent data setting in Equation (4), where the task vector is drawn from Equation (8) with $\gamma_k = \frac{1}{K}$ for all $k \in [K]$. For all $j \in [m+1]$, suppose that the test prompts are constructed with features $\mathbf{x}_j \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and labels whose task vectors are constructed with subspaces defined in Equation (9). For any $\{\alpha_k\}_{k=1}^K$ s.t. $\sum_{k=1}^K \alpha_k^2 = 1$ and $\delta \in (0, r)$, if

$$m \geq n > \frac{(K(r + \sigma^2) + 1)r}{\delta} - (K(r + \sigma^2) + 1),$$

then we have

$$\lim_{\epsilon \rightarrow 0} \mathbb{E} \left[\left(\tilde{y} - g_{ATT}^*(\tilde{\mathbf{z}}_{m+1}, \tilde{\mathbf{Z}}_{\mathcal{M}}) \right)^2 \right] < \sigma^2 + \delta.$$

The proof is deferred to Section C.2. Similar to Theorem 1, if the linear attention model is trained on task vectors that lie in a union of K subspaces, it can generalize well to any region within the span of the K subspaces, even if those regions have zero probability density during training. Lastly, note that by setting $K = 2$, $\alpha_1 = \cos(\theta)$, and $\alpha_2 = \sin(\theta)$, we exactly recover Theorem 1.

B Experimental Details

In this section, we provide specific experimental details for Figures 1 to 3. To generate linear regression data, we set $d = 20$, $r = 5$, $\sigma^2 = 0$, and $\epsilon = 10^{-5}$. To construct the train and test subspaces, we sample an orthogonal matrix $\mathbf{U} \in \mathbb{R}^{d \times d}$ uniformly at random, set \mathbf{U}_s to be the first r columns of \mathbf{U} , and set $\mathbf{U}_{s,\perp}$ to be the second r columns. Given this setup, we typically consider a mixture of $K = 2$ subspaces for the experiments.

For the experiments with the linear transformer, we plug in the optimal weights according to their respective settings (e.g., optimal weights using a single subspace or a mixture of subspaces) and set $m = n = 200$. For the nonlinear transformer, following Garg et al. (2022), we use a small GPT-2 model with 6 layers, 4 heads, and a 128-dimensional embedding space. We append a learnable linear transformation to map the vector predicted by the model to a scalar. We use a learning rate of $\eta = 10^{-4}$, batch size 128, prompt lengths $m = n = 150$, and train for 100K iterations. We run all experiments on a single A100 GPU.

C Deferred Proofs

This section presents all deferred proofs. First, Section C.1 provides auxiliary results used to support both the task and feature shift proofs. Next, Sections C.2 and C.3 provide proofs of Theorems 1 and 2 and proposition 1.

C.1 Supporting Results

We first derive an expression for the test risk under a general distribution shift for the task vector.

Lemma 1 (Test Risk under General Task Distribution Shift). *Let g_{ATT}^* denote the optimal linear attention model corresponding to the independent data setting in Equation (4), where \mathbf{w} follows a (mixture of) Gaussian distribution(s) with zero mean and covariance Σ_s . For all $j \in [m+1]$, suppose that the prompts at test time are constructed with features $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and labels*

$$\tilde{y}_j = \tilde{\mathbf{w}}^\top \mathbf{x}_j + \eta_j, \quad \text{where } \tilde{\mathbf{w}} \sim \mathcal{N}(\mathbf{0}, \Sigma_t) \quad \text{and} \quad \eta_j \sim \mathcal{N}(0, \sigma^2).$$

Then,

$$\mathbb{E} \left[\left(\tilde{y}_{m+1} - g_{ATT}^*(\tilde{\mathbf{z}}_q, \tilde{\mathbf{Z}}_{\mathcal{M}}) \right)^2 \right] = M_t - \text{Tr}(\Sigma_t \mathbf{A}) + \frac{M_t}{m} \text{Tr}(\mathbf{A}^\top \mathbf{A}) - \text{Tr}(\Sigma_t \mathbf{A}) + \frac{m+1}{m} \text{Tr}(\mathbf{A} \Sigma_t \mathbf{A}),$$

where $M_t = \text{Tr}(\Sigma_t) + \sigma^2$.

Proof. Recall at inference time,

$$\tilde{\mathbf{Z}}_{\mathcal{M}} = [\tilde{\mathbf{z}}_1 \quad \dots \quad \tilde{\mathbf{z}}_m \quad \mathbf{0}]^\top = \begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_m & \mathbf{0} \\ \tilde{y}_1 & \dots & \tilde{y}_m & 0 \end{bmatrix}^\top \quad \text{and} \quad \tilde{\mathbf{z}}_q = \begin{bmatrix} \mathbf{x}_{m+1} \\ 0 \end{bmatrix} := \begin{bmatrix} \mathbf{x}_q \\ 0 \end{bmatrix}. \quad (10)$$

Then, let us define

$$\mathbf{X}_{te} := [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_m]^\top, \quad \mathbf{y}_{te} := [\tilde{y}_1 \quad \tilde{y}_2 \quad \dots \quad \tilde{y}_m]^\top, \quad \boldsymbol{\eta}_{te} := [\eta_1 \quad \eta_2 \quad \dots \quad \eta_m]^\top,$$

and $\eta_q := \eta_{m+1}$. Note $\mathbf{y}_{te} = \mathbf{X}_{te} \tilde{\mathbf{w}} + \boldsymbol{\eta}_{te}$. By Lemma 2, we have

$$g_{ATT}^*(\tilde{\mathbf{z}}_q, \tilde{\mathbf{Z}}_{\mathcal{M}}) = \frac{1}{m} \mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \mathbf{y}_{te} = \mathbf{x}_q^\top \underbrace{\left(\frac{1}{m} \mathbf{A} \mathbf{X}_{te}^\top \mathbf{y}_{te} \right)}_{:= \hat{\mathbf{w}}},$$

where $\mathbf{A} = \left(\frac{n+1}{n} \mathbf{I}_d + \frac{M_s}{n} \Sigma_s \right)^{-1}$. By plugging this into the risk and linearity of expectation,

$$\mathbb{E} \left[\left(\tilde{\mathbf{w}}^\top \mathbf{x}_q + \eta_q - \hat{\mathbf{w}}^\top \mathbf{x}_q \right)^2 \right] = \underbrace{\mathbb{E} \left[\left(\tilde{\mathbf{w}}^\top \mathbf{x}_q + \eta_q \right)^2 \right]}_{(a)} - 2 \underbrace{\mathbb{E} \left[\left(\tilde{\mathbf{w}}^\top \mathbf{x}_q + \eta_q \right) \left(\mathbf{x}_q^\top \hat{\mathbf{w}} \right) \right]}_{(b)} + \underbrace{\mathbb{E} \left[\left(\hat{\mathbf{w}}^\top \mathbf{x}_q \right)^2 \right]}_{(c)}. \quad (11)$$

It suffices to analyze each individual term.

Analyzing (a). We first evaluate $\mathbb{E} \left[\left(\tilde{\mathbf{w}}^\top \mathbf{x}_q + \eta_q \right)^2 \right]$. First, we note

$$\mathbb{E} \left[\left(\tilde{\mathbf{w}}^\top \mathbf{x}_q + \eta_q \right)^2 \right] = \mathbb{E} \left[\left(\tilde{\mathbf{w}}^\top \mathbf{x}_q \right)^2 \right] + 2 \mathbb{E} \left[\eta_q \tilde{\mathbf{w}}^\top \mathbf{x}_q \right] + \mathbb{E} \left[\eta_q^2 \right] = \mathbb{E} \left[\left(\tilde{\mathbf{w}}^\top \mathbf{x}_q \right)^2 \right] + \sigma^2,$$

so it suffices to analyze $\mathbb{E} [(\tilde{\mathbf{w}}^\top \mathbf{x}_q)^2]$. By law of total expectation and the fact that $\tilde{\mathbf{w}}, \mathbf{x}_q$ are independent,

$$\mathbb{E} [(\tilde{\mathbf{w}}^\top \mathbf{x}_q)^2] = \mathbb{E}_{\tilde{\mathbf{w}}} \left[\mathbb{E}_{\mathbf{x}_q} [(\tilde{\mathbf{w}}^\top \mathbf{x}_q)^2 \mid \tilde{\mathbf{w}}] \right].$$

Conditioned on $\tilde{\mathbf{w}}$, $\tilde{\mathbf{w}}^\top \mathbf{x}_q \sim \mathcal{N}(0, \|\tilde{\mathbf{w}}\|^2)$, so $\mathbb{E} [(\tilde{\mathbf{w}}^\top \mathbf{x}_q)^2 \mid \tilde{\mathbf{w}}] = \text{Var}(\tilde{\mathbf{w}}^\top \mathbf{x}_q \mid \tilde{\mathbf{w}}) = \|\tilde{\mathbf{w}}\|^2$. Therefore,

$$\mathbb{E}_{\tilde{\mathbf{w}}} \left[\mathbb{E}_{\mathbf{x}_q} [(\tilde{\mathbf{w}}^\top \mathbf{x}_q)^2 \mid \tilde{\mathbf{w}}] \right] = \mathbb{E} [\|\tilde{\mathbf{w}}\|^2] = \text{Tr}(\mathbb{E}[\tilde{\mathbf{w}}\tilde{\mathbf{w}}^\top]) = \text{Tr}(\boldsymbol{\Sigma}_t).$$

Therefore,

$$\mathbb{E} [(\tilde{\mathbf{w}}^\top \mathbf{x}_q + \eta_q)^2] = \text{Tr}(\boldsymbol{\Sigma}_t) + \sigma^2.$$

Analyzing (b). Next, we analyze $\mathbb{E} [(\tilde{\mathbf{w}}^\top \mathbf{x}_q + \eta_q)(\mathbf{x}_q^\top \hat{\mathbf{w}})]$. We first note

$$\mathbb{E} [(\tilde{\mathbf{w}}^\top \mathbf{x}_q + \eta_q)(\mathbf{x}_q^\top \hat{\mathbf{w}})] = \mathbb{E} [(\tilde{\mathbf{w}}^\top \mathbf{x}_q)(\mathbf{x}_q^\top \hat{\mathbf{w}})] + \underbrace{\mathbb{E} [\eta_q \mathbf{x}_q^\top \hat{\mathbf{w}}]}_{=0} = \mathbb{E} [(\tilde{\mathbf{w}}^\top \mathbf{x}_q)(\mathbf{x}_q^\top \hat{\mathbf{w}})],$$

so it suffices to analyze $\mathbb{E} [(\tilde{\mathbf{w}}^\top \mathbf{x}_q)(\mathbf{x}_q^\top \hat{\mathbf{w}})]$. Substituting $\hat{\mathbf{w}} := \frac{1}{m} \mathbf{A} \mathbf{X}_{te}^\top \mathbf{y}_{te} = \frac{1}{m} \mathbf{A} \mathbf{X}_{te}^\top (\mathbf{X}_{te} \tilde{\mathbf{w}} + \boldsymbol{\eta}_{te})$ yields

$$\begin{aligned} \mathbb{E} [(\tilde{\mathbf{w}}^\top \mathbf{x}_q)(\mathbf{x}_q^\top \hat{\mathbf{w}})] &= \frac{1}{m} \mathbb{E} [\tilde{\mathbf{w}}^\top \mathbf{x}_q \mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top (\mathbf{X}_{te} \tilde{\mathbf{w}} + \boldsymbol{\eta}_{te})] \\ &= \frac{1}{m} \left(\mathbb{E} [\tilde{\mathbf{w}}^\top \mathbf{x}_q \mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \mathbf{X}_{te} \tilde{\mathbf{w}}] + \mathbb{E} [\tilde{\mathbf{w}}^\top \mathbf{x}_q \mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \boldsymbol{\eta}_{te}] \right) \\ &= \frac{1}{m} \left(\mathbb{E} [\tilde{\mathbf{w}}^\top \mathbf{x}_q \mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \mathbf{X}_{te} \tilde{\mathbf{w}}] + \underbrace{\mathbb{E} [\tilde{\mathbf{w}}^\top \mathbf{x}_q \mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top]}_{=0} \mathbb{E} [\boldsymbol{\eta}_{te}] \right) \\ &= \frac{1}{m} \mathbb{E} [\tilde{\mathbf{w}}^\top \mathbf{x}_q \mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \mathbf{X}_{te} \tilde{\mathbf{w}}] = \frac{1}{m} \mathbb{E} [\text{Tr}(\tilde{\mathbf{w}} \tilde{\mathbf{w}}^\top \mathbf{x}_q \mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \mathbf{X}_{te})] \\ &= \frac{1}{m} \text{Tr} \left(\mathbb{E} [\tilde{\mathbf{w}} \tilde{\mathbf{w}}^\top \mathbf{x}_q \mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \mathbf{X}_{te}] \right) \\ &= \frac{1}{m} \text{Tr} \left(\underbrace{\mathbb{E} [\tilde{\mathbf{w}} \tilde{\mathbf{w}}^\top]}_{\boldsymbol{\Sigma}_t} \underbrace{\mathbb{E} [\mathbf{x}_q \mathbf{x}_q^\top]}_{\mathbf{I}_d} \underbrace{\mathbf{A} \mathbb{E} [\mathbf{X}_{te}^\top \mathbf{X}_{te}]}_{m \cdot \mathbf{I}_d} \right) = \text{Tr}(\boldsymbol{\Sigma}_t \mathbf{A}). \end{aligned}$$

Analyzing (c). Finally, we analyze $\mathbb{E} [(\mathbf{x}_q^\top \hat{\mathbf{w}})^2]$:

$$\begin{aligned} \mathbb{E} [(\mathbf{x}_q^\top \hat{\mathbf{w}})^2] &= \frac{1}{m^2} \mathbb{E} [(\mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top (\mathbf{X}_{te} \tilde{\mathbf{w}} + \boldsymbol{\eta}_{te}))^2] = \frac{1}{m^2} \mathbb{E} [(\mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \mathbf{X}_{te} \tilde{\mathbf{w}} + \mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \boldsymbol{\eta}_{te})^2] \\ &= \frac{1}{m^2} \left(\mathbb{E} [(\mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \mathbf{X}_{te} \tilde{\mathbf{w}})^2] + 2 \underbrace{\mathbb{E} [(\mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \mathbf{X}_{te} \tilde{\mathbf{w}})(\mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \boldsymbol{\eta}_{te})]}_{=0} + \mathbb{E} [(\mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \boldsymbol{\eta}_{te})^2] \right) \\ &= \frac{1}{m^2} \left(\underbrace{\mathbb{E} [\mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \mathbf{X}_{te} \tilde{\mathbf{w}} \tilde{\mathbf{w}}^\top \mathbf{X}_{te}^\top \mathbf{X}_{te} \mathbf{A}^\top \mathbf{x}_q]}_{(d)} + \underbrace{\mathbb{E} [\mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \boldsymbol{\eta}_{te} \boldsymbol{\eta}_{te}^\top \mathbf{X}_{te} \mathbf{A}^\top \mathbf{x}_q]}_{(e)} \right). \end{aligned}$$

We first focus on (d):

$$\begin{aligned} \mathbb{E} [\mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \mathbf{X}_{te} \tilde{\mathbf{w}} \tilde{\mathbf{w}}^\top \mathbf{X}_{te}^\top \mathbf{X}_{te} \mathbf{A}^\top \mathbf{x}_q] &= \mathbb{E} \left[\text{Tr}(\mathbf{x}_q \mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \mathbf{X}_{te} \tilde{\mathbf{w}} \tilde{\mathbf{w}}^\top \mathbf{X}_{te}^\top \mathbf{X}_{te} \mathbf{A}^\top) \right] \\ &= \text{Tr} \left(\underbrace{\mathbb{E} [\mathbf{x}_q \mathbf{x}_q^\top]}_{\mathbf{I}_d} \mathbf{A} \mathbb{E} [\mathbf{X}_{te}^\top \mathbf{X}_{te} \tilde{\mathbf{w}} \tilde{\mathbf{w}}^\top \mathbf{X}_{te}^\top \mathbf{X}_{te}] \mathbf{A}^\top \right) \\ &= \mathbb{E} \left[\text{Tr}(\mathbf{A} \mathbf{X}_{te}^\top \mathbf{X}_{te} \tilde{\mathbf{w}} \tilde{\mathbf{w}}^\top \mathbf{X}_{te}^\top \mathbf{X}_{te} \mathbf{A}^\top) \right] = \mathbb{E} \left[\text{Tr}(\tilde{\mathbf{w}} \tilde{\mathbf{w}}^\top \mathbf{X}_{te}^\top \mathbf{X}_{te} \mathbf{A}^\top \mathbf{A} \mathbf{X}_{te}^\top \mathbf{X}_{te}) \right] \\ &= \text{Tr} \left(\underbrace{\mathbb{E} [\tilde{\mathbf{w}} \tilde{\mathbf{w}}^\top]}_{\boldsymbol{\Sigma}_t} \mathbb{E} [\mathbf{X}_{te}^\top \mathbf{X}_{te} \mathbf{A}^\top \mathbf{A} \mathbf{X}_{te}^\top \mathbf{X}_{te}] \right) = \mathbb{E} \left[\text{Tr}(\boldsymbol{\Sigma}_t \mathbf{X}_{te}^\top \mathbf{X}_{te} \mathbf{A}^\top \mathbf{A} \mathbf{X}_{te}^\top \mathbf{X}_{te}) \right] \\ &= \mathbb{E} \left[\text{Tr}(\mathbf{A} \mathbf{X}_{te}^\top \mathbf{X}_{te} \boldsymbol{\Sigma}_t^{1/2} \boldsymbol{\Sigma}_t^{1/2} \mathbf{X}_{te}^\top \mathbf{X}_{te} \mathbf{A}^\top) \right] := \mathbb{E} \left[\text{Tr}(\tilde{\mathbf{X}}_{te}^\top \bar{\mathbf{X}}_{te} \bar{\mathbf{X}}_{te}^\top \tilde{\mathbf{X}}_{te}) \right], \end{aligned}$$

where $\tilde{\mathbf{X}}_{te}^\top := \mathbf{A}\mathbf{X}_{te}^\top$ and $\bar{\mathbf{X}}_{te}^\top := \Sigma_t^{1/2}\mathbf{X}_{te}^\top$. Note $\tilde{\mathbf{X}}_{te}^\top = [\mathbf{A}\mathbf{x}_1 \ \dots \ \mathbf{A}\mathbf{x}_m] := [\tilde{\mathbf{x}}_1 \ \dots \ \tilde{\mathbf{x}}_m]$ where $\tilde{\mathbf{x}}_i := \mathbf{A}\mathbf{x}_i \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}_d, \mathbf{A}\mathbf{A}^\top)$, and $\bar{\mathbf{X}}_{te}^\top = [\Sigma_t^{1/2}\mathbf{x}_1 \ \dots \ \Sigma_t^{1/2}\mathbf{x}_m] := [\bar{\mathbf{x}}_1 \ \dots \ \bar{\mathbf{x}}_m]$ where $\bar{\mathbf{x}}_i := \Sigma_t^{1/2}\mathbf{x}_i \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}_d, \Sigma_t)$. We can express $\tilde{\mathbf{X}}_{te}^\top \bar{\mathbf{X}}_{te}$ and $\bar{\mathbf{X}}_{te}^\top \tilde{\mathbf{X}}_{te}$ as such:

$$\tilde{\mathbf{X}}_{te}^\top \bar{\mathbf{X}}_{te} = \sum_{i=1}^m \tilde{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \quad \text{and} \quad \bar{\mathbf{X}}_{te}^\top \tilde{\mathbf{X}}_{te} = \sum_{j=1}^m \bar{\mathbf{x}}_j \tilde{\mathbf{x}}_j^\top.$$

Therefore,

$$\begin{aligned} \mathbb{E} \left[\text{Tr} \left(\tilde{\mathbf{X}}_{te}^\top \bar{\mathbf{X}}_{te} \bar{\mathbf{X}}_{te}^\top \tilde{\mathbf{X}}_{te} \right) \right] &= \text{Tr} \left(\mathbb{E} \left[\tilde{\mathbf{X}}_{te}^\top \bar{\mathbf{X}}_{te} \bar{\mathbf{X}}_{te}^\top \tilde{\mathbf{X}}_{te} \right] \right) \\ &= \text{Tr} \left(\sum_{i=1}^m \sum_{j=1}^m \mathbb{E} [\tilde{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \bar{\mathbf{x}}_j \tilde{\mathbf{x}}_j^\top] \right) \\ &= \text{Tr} \left(\sum_{i=1}^m \sum_{j \neq i} \mathbb{E} [\tilde{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \bar{\mathbf{x}}_j \tilde{\mathbf{x}}_j^\top] \right) + \text{Tr} \left(\sum_{i=1}^m \mathbb{E} [\tilde{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \bar{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top] \right) \end{aligned}$$

We first consider the case when $i \neq j$. In this setting, \mathbf{x}_i and \mathbf{x}_j are independent, so

$$\mathbb{E} [\tilde{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \bar{\mathbf{x}}_j \tilde{\mathbf{x}}_j^\top] = \mathbb{E} [\tilde{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top] \mathbb{E} [\bar{\mathbf{x}}_j \tilde{\mathbf{x}}_j^\top] = \mathbf{A} \underbrace{\mathbb{E} [\mathbf{x}_i \mathbf{x}_i^\top]}_{\mathbf{I}_d} \Sigma_t \underbrace{\mathbb{E} [\mathbf{x}_j \mathbf{x}_j^\top]}_{\mathbf{I}_d} \mathbf{A}^\top = \mathbf{A} \Sigma_t \mathbf{A}^\top.$$

Therefore,

$$\text{Tr} \left(\sum_{i=1}^m \sum_{j \neq i} \mathbb{E} [\tilde{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \bar{\mathbf{x}}_j \tilde{\mathbf{x}}_j^\top] \right) = m \cdot (m-1) \cdot \text{Tr} (\mathbf{A} \Sigma_t \mathbf{A}^\top).$$

We now consider the case where $i = j$:

$$\begin{aligned} \text{Tr} \left(\sum_{i=1}^m \mathbb{E} [\tilde{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \bar{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top] \right) &= \sum_{i=1}^m \mathbb{E} \left[\text{Tr} (\tilde{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \bar{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top) \right] \\ &= \sum_{i=1}^m \mathbb{E} [\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \bar{\mathbf{x}}_i] = \sum_{i=1}^m \mathbb{E} [(\mathbf{x}_i^\top \mathbf{A}^\top \mathbf{A} \mathbf{x}_i)(\mathbf{x}_i^\top \Sigma_t \mathbf{x}_i)] \\ &\stackrel{(i)}{=} m \cdot \left(2 \text{Tr} (\mathbf{A} \Sigma_t \mathbf{A}^\top) + \text{Tr} (\mathbf{A}^\top \mathbf{A}) \text{Tr} (\Sigma_t) \right), \end{aligned}$$

where (i) is because for $\mathbf{a} \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ and fixed $\mathbf{Q}, \mathbf{R} \in \mathbb{R}^{d \times d}$, $\mathbb{E} [(\mathbf{a}^\top \mathbf{Q} \mathbf{a})(\mathbf{a}^\top \mathbf{R} \mathbf{a})] = \text{Tr} (\mathbf{Q}(\mathbf{R} + \mathbf{R}^\top)) + \text{Tr} (\mathbf{Q}) \text{Tr} (\mathbf{R})$ (see Section 8.2.4 in [Petersen et al. \(2008\)](#)).

We now focus on (e):

$$\begin{aligned} \mathbb{E} [\mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \boldsymbol{\eta}_{te} \boldsymbol{\eta}_{te}^\top \mathbf{X}_{te} \mathbf{A}^\top \mathbf{x}_q] &= \mathbb{E} \left[\text{Tr} (\mathbf{x}_q \mathbf{x}_q^\top \mathbf{A} \mathbf{X}_{te}^\top \boldsymbol{\eta}_{te} \boldsymbol{\eta}_{te}^\top \mathbf{X}_{te} \mathbf{A}^\top) \right] \\ &= \text{Tr} \left(\underbrace{\mathbb{E} [\mathbf{x}_q \mathbf{x}_q^\top]}_{\mathbf{I}_d} \mathbf{A} \mathbb{E} [\mathbf{X}_{te}^\top \boldsymbol{\eta}_{te} \boldsymbol{\eta}_{te}^\top \mathbf{X}_{te}] \mathbf{A}^\top \right) = \text{Tr} \left(\mathbb{E} [\mathbf{A} \mathbf{X}_{te}^\top \boldsymbol{\eta}_{te} \boldsymbol{\eta}_{te}^\top \mathbf{X}_{te} \mathbf{A}^\top] \right) \\ &:= \text{Tr} \left(\mathbb{E} [\tilde{\boldsymbol{\eta}}_{te} \tilde{\boldsymbol{\eta}}_{te}^\top] \right), \end{aligned}$$

where $\tilde{\boldsymbol{\eta}}_{te} := \mathbf{A} \mathbf{X}_{te}^\top \boldsymbol{\eta}_{te} = \tilde{\mathbf{X}}_{te}^\top \boldsymbol{\eta}_{te}$. Note the columns of $\tilde{\mathbf{X}}_{te}^\top$ are iid Gaussian with covariance $\mathbf{A} \mathbf{A}^\top$. By Corollary 6 in [Mattei \(2017\)](#), $\tilde{\boldsymbol{\eta}}_{te} \sim \text{GAL}_d(2\sigma^2 \mathbf{A} \mathbf{A}^\top, \mathbf{0}_d, m/2)$, where $\text{GAL}_p(\boldsymbol{\Sigma}, \boldsymbol{\mu}, s)$ denotes a p -dimensional *multivariate generalized asymmetric Laplace distribution* with mean $s\boldsymbol{\mu}$ and covariance $s(\boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}^\top)$ (Definition 1 and Proposition 2 in [Mattei \(2017\)](#)). Therefore,

$$\text{Tr} \left(\mathbb{E} [\tilde{\boldsymbol{\eta}}_{te} \tilde{\boldsymbol{\eta}}_{te}^\top] \right) = \text{Tr} \left(\text{Cov}(\tilde{\boldsymbol{\eta}}_{te}) \right) = m\sigma^2 \text{Tr} (\mathbf{A} \mathbf{A}^\top).$$

Adding (a), (b), and (c). Adding the expressions for (a), (b), and (c), where (c) = (d) + (e), yields and combining like terms yields the following expression:

$$\begin{aligned} \mathbb{E} \left[\left(\tilde{\mathbf{w}}^\top \mathbf{x}_q + \eta_q - \hat{\mathbf{w}}^\top \mathbf{x}_q \right)^2 \right] &= \underbrace{\text{Tr}(\boldsymbol{\Sigma}_t) + \sigma^2}_{=(a)} - 2 \underbrace{\text{Tr}(\boldsymbol{\Sigma}_t \mathbf{A})}_{=(b)} \\ &+ \frac{1}{m^2} \left(\underbrace{m(m-1) \text{Tr}(\mathbf{A} \boldsymbol{\Sigma}_t \mathbf{A}^\top) + 2m \text{Tr}(\mathbf{A} \boldsymbol{\Sigma}_t \mathbf{A}^\top)}_{=(d)} + \underbrace{m \text{Tr}(\boldsymbol{\Sigma}_t) \text{Tr}(\mathbf{A}^\top \mathbf{A}) + m \sigma^2 \text{Tr}(\mathbf{A}^\top \mathbf{A})}_{=(e)} \right). \end{aligned}$$

Combining like terms yields

$$\begin{aligned} \mathbb{E} \left[\left(\tilde{\mathbf{w}}^\top \mathbf{x}_q + \eta_q - \hat{\mathbf{w}}^\top \mathbf{x}_q \right)^2 \right] &= \left(\frac{1}{m} \text{Tr}(\mathbf{A}^\top \mathbf{A}) + 1 \right) \left(\text{Tr}(\boldsymbol{\Sigma}_t) + \sigma^2 \right) - 2 \text{Tr}(\boldsymbol{\Sigma}_t \mathbf{A}) + \frac{m+1}{m} \text{Tr}(\mathbf{A} \boldsymbol{\Sigma}_t \mathbf{A}^\top) \\ &= M_t - \text{Tr}(\boldsymbol{\Sigma}_t \mathbf{A}) + \frac{M_t}{m} \text{Tr}(\mathbf{A}^\top \mathbf{A}) - \text{Tr}(\boldsymbol{\Sigma}_t \mathbf{A}) + \frac{m+1}{m} \text{Tr}(\mathbf{A} \boldsymbol{\Sigma}_t \mathbf{A}^\top), \end{aligned}$$

which is exactly Equation (11). This completes the proof. \square

C.2 Proof of Theorems 1 and 2

We now provide a proof of Theorems 1 and 2. We specifically focus on the setting of Theorem 2. However, we also emphasize that the proof reduces down to a proof of Theorem 1 when $K = 2$, $\alpha_1 = \sin(\theta)$, and $\alpha_2 = \cos(\theta)$ for any $\theta \in [0, \pi/2]$.

Proof. Let $\tilde{y} := \tilde{y}_{m+1}$. By Lemmas 1 and 3, we have

$$\mathbb{E} \left[\left(\tilde{y} - g_{\text{ATT}}^*(\tilde{\mathbf{z}}_q, \tilde{\mathbf{Z}}_{\mathcal{M}}) \right)^2 \right] = \left(\frac{1}{m} \text{Tr}(\mathbf{A}^\top \mathbf{A}) + 1 \right) \left(\text{Tr}(\bar{\boldsymbol{\Sigma}}_t) + \sigma^2 \right) - 2 \text{Tr}(\bar{\boldsymbol{\Sigma}}_t \mathbf{A}) + \frac{m+1}{m} \text{Tr}(\mathbf{A} \bar{\boldsymbol{\Sigma}}_t \mathbf{A}^\top), \quad (12)$$

where $\mathbf{A} = \left(\frac{n+1}{n} \mathbf{I}_d + \frac{M_s}{n} \boldsymbol{\Sigma}^{-1} \right)^{-1}$, $M_s = \text{Tr}(\boldsymbol{\Sigma}) + \sigma^2$, and $\boldsymbol{\Sigma} = \sum_{k=1}^K \gamma_k \cdot \boldsymbol{\Sigma}_{s,k}$.

Let $\mathbf{U} := [\mathbf{U}_{s,1} \quad \mathbf{U}_{s,2} \quad \dots \quad \mathbf{U}_{s,K} \quad \mathbf{U}_\perp]$, where $\mathbf{U}_\perp \in \mathbb{R}^{d \times (d-Kr)}$ completes the orthonormal basis for \mathbb{R}^d . By Lemma 4,

$$\mathbf{A} = \left(\frac{n+1}{n} \mathbf{I}_d + \frac{M_s}{n} \boldsymbol{\Sigma}^{-1} \right)^{-1} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top,$$

where

$$\boldsymbol{\Lambda} = \begin{bmatrix} \nu_1 \mathbf{I}_r & & & \\ & \ddots & & \\ & & \nu_K \mathbf{I}_r & \\ & & & \nu_{K+1} \mathbf{I}_{d-Kr} \end{bmatrix}$$

with $\nu_k = \frac{n(\gamma_k + \epsilon)}{(n+1)(\gamma_k + \epsilon) + M_s}$ for all $k \in [K]$, and $\nu_{K+1} = \frac{n\epsilon}{(n+1)\epsilon + r + \epsilon d + \sigma^2}$.

Simplifying $\text{Tr}(\bar{\boldsymbol{\Sigma}}_t)$. We can write $\text{Tr}(\bar{\boldsymbol{\Sigma}}_t)$ as such:

$$\text{Tr}(\bar{\boldsymbol{\Sigma}}_t) = \text{Tr}(\bar{\mathbf{U}}_t \bar{\mathbf{U}}_t^\top) + \epsilon \text{Tr}(\mathbf{I}_d) = r + \epsilon d.$$

Simplifying $\text{Tr}(\mathbf{A})$ and $\text{Tr}(\mathbf{A}^\top \mathbf{A})$. We can write $\text{Tr}(\mathbf{A})$ and $\text{Tr}(\mathbf{A}^\top \mathbf{A})$ as such:

$$\text{Tr}(\mathbf{A}) = r \sum_{k=1}^K \nu_k + (d - Kr) \nu_{K+1} \quad \text{and} \quad \text{Tr}(\mathbf{A}^\top \mathbf{A}) = \text{Tr}(\mathbf{A}^2) = r \sum_{k=1}^K \nu_k^2 + (d - Kr) \nu_{K+1}^2.$$

Simplifying $\text{Tr}(\bar{\Sigma}_t \mathbf{A})$ and $\text{Tr}(\mathbf{A} \bar{\Sigma}_t \mathbf{A}^\top)$. Note $\text{Tr}(\mathbf{A} \bar{\Sigma}_t \mathbf{A}^\top) = \text{Tr}(\bar{\Sigma}_t \mathbf{A}^2)$. We first focus on $\text{Tr}(\bar{\Sigma}_t \mathbf{A})$:

$$\begin{aligned} \bar{\Sigma}_t \mathbf{A} &= \left(\bar{\mathbf{U}}_t \bar{\mathbf{U}}_t^\top + \epsilon \mathbf{I}_d \right) \mathbf{U} \Lambda \mathbf{U}^\top = \bar{\mathbf{U}}_t \bar{\mathbf{U}}_t^\top \mathbf{U} \Lambda \mathbf{U}^\top + \epsilon \mathbf{U} \Lambda \mathbf{U}^\top \\ \implies \text{Tr}(\bar{\Sigma}_t \mathbf{A}) &= \text{Tr} \left(\bar{\mathbf{U}}_t^\top \mathbf{U} \Lambda \mathbf{U}^\top \bar{\mathbf{U}}_t \right) + \epsilon \text{Tr}(\mathbf{A}). \end{aligned}$$

Recall $\bar{\mathbf{U}}_t = \sum_{k=1}^K \alpha_k \mathbf{U}_{s,k}$ where $\sum_{k=1}^K \alpha_k^2 = 1$, and so we have

$$\bar{\mathbf{U}}_t^\top \mathbf{U} = \left(\sum_{k=1}^K \alpha_k \mathbf{U}_k \right)^\top \begin{bmatrix} \mathbf{U}_{s,1} & \dots & \mathbf{U}_{s,K} & \mathbf{U}_\perp \end{bmatrix} = \begin{bmatrix} \alpha_1 \mathbf{I}_r & & & \\ & \ddots & & \\ & & \alpha_K \mathbf{I}_r & \\ & & & \mathbf{0}_{(d-Kr) \times (d-Kr)} \end{bmatrix}$$

Thus,

$$\text{Tr} \left(\bar{\mathbf{U}}_t^\top \mathbf{U} \Lambda \mathbf{U}^\top \bar{\mathbf{U}}_t \right) = \text{Tr} \left(\begin{bmatrix} \alpha_1^2 \nu_1 \mathbf{I}_r & & & \\ & \ddots & & \\ & & \alpha_K^2 \nu_K \mathbf{I}_r & \\ & & & \mathbf{0}_{(d-Kr) \times (d-Kr)} \end{bmatrix} \right) = r \sum_{k=1}^K \alpha_k^2 \nu_k$$

Using a similar argument,

$$\text{Tr} \left(\bar{\Sigma}_t^\top \mathbf{A}^2 \right) = r \sum_{k=1}^K \alpha_k^2 \nu_k^2 + \epsilon \text{Tr}(\mathbf{A}^2).$$

Simplifying the test risk. Substituting the expressions for the $\text{Tr}(\cdot)$ terms into Equation (12) yields

$$\begin{aligned} \mathbb{E} \left[\left(\tilde{y} - g_{\text{ATT}}^*(\tilde{\mathbf{z}}_q, \tilde{\mathbf{Z}}_{\mathcal{M}}) \right)^2 \right] &= \left(\frac{1}{m} \left(r \sum_{k=1}^K \nu_k^2 + (d-Kr) \nu_{K+1}^2 \right) + 1 \right) (r + \epsilon d + \sigma^2) \\ &\quad - 2 \left(r \sum_{k=1}^K \alpha_k^2 \nu_k + \left(r \sum_{k=1}^K \nu_k + (d-Kr) \nu_{K+1} \right) \epsilon \right) \\ &\quad + \frac{m+1}{m} \left(r \sum_{k=1}^K \alpha_k^2 \nu_k^2 + \left(r \sum_{k=1}^K \nu_k^2 + (d-Kr) \nu_{K+1}^2 \right) \epsilon \right). \end{aligned}$$

Taking $\epsilon \rightarrow 0$ results in the following expression for the test risk:

$$\begin{aligned} \lim_{\epsilon \rightarrow 0} \mathbb{E} \left[\left(\tilde{y} - g_{\text{ATT}}^*(\tilde{\mathbf{z}}_q, \tilde{\mathbf{Z}}_{\mathcal{M}}) \right)^2 \right] &= r + \sigma^2 + \frac{(r + \sigma^2)r}{m} \sum_{k=1}^K \left(\frac{\gamma_k n}{\gamma_k(n+1) + r + \sigma^2} \right)^2 \\ &\quad - 2r \sum_{k=1}^K \frac{\alpha_k^2 \gamma_k n}{\gamma_k(n+1) + r + \sigma^2} + \frac{(m+1)r}{m} \sum_{k=1}^K \left(\frac{\alpha_k \gamma_k n}{\gamma_k(n+1) + r + \sigma^2} \right)^2 \end{aligned}$$

Substituting $\gamma_k = \frac{1}{K}$ for all $k \in [K]$ and combining like terms yields

$$\lim_{\epsilon \rightarrow 0} \mathbb{E} \left[\left(\tilde{y} - g_{\text{ATT}}^*(\tilde{\mathbf{z}}_q, \tilde{\mathbf{Z}}_{\mathcal{M}}) \right)^2 \right] = r + \sigma^2 + \frac{m+1+K(r+\sigma^2)}{m} \cdot \frac{rn^2}{(n+1+K(r+\sigma^2))^2} - \frac{2rn}{n+1+K(r+\sigma^2)}.$$

Now suppose $n \leq m$. Then, we have

$$\lim_{\epsilon \rightarrow 0} \mathbb{E} \left[\left(\tilde{y} - g_{\text{ATT}}^*(\tilde{\mathbf{z}}_q, \tilde{\mathbf{Z}}_{\mathcal{M}}) \right)^2 \right] \leq r + \sigma^2 - \frac{rn^2}{n+1+K(r+\sigma^2)}.$$

Upper bounding this by $\sigma^2 + \delta$ for some $\delta \in (0, r)$, then solving for n , yields the following result. For any $\delta \in (0, r)$, if

$$m \geq n > \frac{(K(r+\sigma^2)+1)r}{\delta} - (K(r+\sigma^2)+1),$$

then $\lim_{\epsilon \rightarrow 0} \mathbb{E} \left[\left(\tilde{y} - g_{\text{ATT}}^*(\tilde{\mathbf{z}}_q, \tilde{\mathbf{Z}}_{\mathcal{M}}) \right)^2 \right] < \sigma^2 + \delta$, which completes the proof. \square

C.3 Proof of Proposition 1

We now provide the proof of Proposition 1.

Proof. For simplicity, we denote $\tilde{y} := \tilde{y}_{m+1}$. Recall $\mathbf{U} := [\mathbf{U}_s \quad \mathbf{U}_{s,\perp} \quad \mathbf{U}_{2r,\perp}] \in \mathbb{R}^{d \times d}$, where $\mathbf{U}_s, \mathbf{U}_{s,\perp} \in \mathbb{R}^{d \times r}$ and $\mathbf{U}_{2r,\perp} \in \mathbb{R}^{d \times (d-2r)}$ all have orthonormal columns, while $\mathbf{U}_s^\top \mathbf{U}_{\perp,s} = \mathbf{0}_{r \times r}$ and $\mathbf{U}_s^\top \mathbf{U}_\perp = \mathbf{U}_{s,\perp}^\top \mathbf{U}_{2r,\perp} = \mathbf{0}_{r \times (d-2r)}$. We re-write Σ_s as such:

$$\Sigma_s = \mathbf{U}_s \mathbf{U}_s^\top + \epsilon \mathbf{I}_d = \mathbf{U} \begin{bmatrix} \mathbf{I}_r & & \\ & \mathbf{0}_{(d-r) \times (d-r)} & \\ & & \end{bmatrix} \mathbf{U}^\top + \epsilon \mathbf{I} = \mathbf{U} \begin{bmatrix} (1+\epsilon)\mathbf{I}_r & & \\ & & \\ & & \epsilon \mathbf{I}_{d-r} \end{bmatrix} \mathbf{U}^\top.$$

Note this is a valid eigendecomposition of Σ_s . Thus, by Lemma 4, we have

$$\mathbf{A} = \left(\frac{n+1}{n} \mathbf{I}_d + \frac{M_s}{n} \Sigma_s^{-1} \right)^{-1} = \mathbf{U} \Lambda \mathbf{U}^\top, \quad (13)$$

where

$$\Lambda = \begin{bmatrix} \frac{n(1+\epsilon)}{(n+1)\epsilon + M_s} \cdot \mathbf{I}_r & & \\ & \frac{n\epsilon}{(n+1)\epsilon + M_s} \cdot \mathbf{I}_{d-r} & \\ & & \end{bmatrix} := \begin{bmatrix} \nu_1 \mathbf{I}_r & & \\ & & \\ & & \nu_2 \mathbf{I}_{d-r} \end{bmatrix}.$$

and $M_s = \text{Tr}(\Sigma_s) + \sigma^2$.

By Lemma 1 (and omitting the subscripts in the expectation),

$$\mathbb{E} \left[\left(\tilde{y} - g_{\text{ATT}}^*(\tilde{\mathbf{z}}_q, \tilde{\mathbf{Z}}_{\mathcal{M}}) \right)^2 \right] = \left(\frac{1}{m} \text{Tr}(\mathbf{A}^\top \mathbf{A}) + 1 \right) \left(\text{Tr}(\Sigma_t) + \sigma^2 \right) - 2 \text{Tr}(\Sigma_t \mathbf{A}) + \frac{m+1}{m} \text{Tr}(\mathbf{A} \Sigma_t \mathbf{A}^\top). \quad (14)$$

We simplify the remaining $\text{Tr}(\cdot)$ terms using Equation (13).

Simplifying $\text{Tr}(\mathbf{A})$ and $\text{Tr}(\mathbf{A}^\top \mathbf{A})$. Directly from Equation (13):

$$\text{Tr}(\mathbf{A}) = r \cdot \nu_1 + (d-r) \cdot \nu_2 \quad \text{and} \quad \text{Tr}(\mathbf{A}^\top \mathbf{A}) = \text{Tr}(\mathbf{A}^2) = r \cdot \nu_1^2 + (d-r) \cdot \nu_2^2,$$

where $\mathbf{A}^2 = \mathbf{U} \Lambda^2 \mathbf{U}^\top$.

Simplifying $\text{Tr}(\Sigma_t \mathbf{A})$ and $\text{Tr}(\mathbf{A} \Sigma_t \mathbf{A}^\top)$. First note $\text{Tr}(\mathbf{A} \Sigma_t \mathbf{A}^\top) = \text{Tr}(\Sigma_t \mathbf{A}^2)$. We first focus on $\text{Tr}(\Sigma_t \mathbf{A})$:

$$\begin{aligned} \Sigma_t \mathbf{A} &= (\mathbf{U}_t \mathbf{U}_t^\top + \epsilon \mathbf{I}_d) \mathbf{U} \Lambda \mathbf{U}^\top = \mathbf{U}_t \mathbf{U}_t^\top \mathbf{U} \Lambda \mathbf{U}^\top + \epsilon \mathbf{U} \Lambda \mathbf{U}^\top \\ &\implies \text{Tr}(\Sigma_t \mathbf{A}) = \text{Tr}(\mathbf{U}_t^\top \mathbf{U} \Lambda \mathbf{U}^\top \mathbf{U}_t) + \epsilon \text{Tr}(\mathbf{A}). \end{aligned}$$

Recall we defined \mathbf{U}_t in Equation (5) as follows:

$$\mathbf{U}_t = \mathbf{U}_s \cos(\Theta) + \mathbf{U}_{s,\perp} \sin(\Theta).$$

Therefore:

$$\mathbf{U}_t^\top \mathbf{U} = (\mathbf{U}_s \cos(\Theta) + \mathbf{U}_{s,\perp} \sin(\Theta))^\top [\mathbf{U}_s \quad \mathbf{U}_{s,\perp} \quad \mathbf{U}_\perp] = [\cos(\Theta) \quad \sin(\Theta) \quad \mathbf{0}_{d \times (d-2r)}],$$

and thus,

$$\begin{aligned} \text{Tr}(\mathbf{U}_t^\top \mathbf{U} \Lambda \mathbf{U}^\top \mathbf{U}_t) &= \text{Tr} \left([\cos(\Theta) \quad \sin(\Theta) \quad \mathbf{0}_{d \times (d-2r)}] \begin{bmatrix} \nu_1 \mathbf{I}_r & & \\ & \nu_2 \mathbf{I}_r & \\ & & \nu_2 \mathbf{I}_{d-2r} \end{bmatrix} \begin{bmatrix} \cos(\Theta) \\ \sin(\Theta) \\ \mathbf{0}_{(d-2r) \times d} \end{bmatrix} \right) \\ &= \text{Tr} \left(\begin{bmatrix} \nu_1 \cos^2(\Theta) & & \\ & \nu_2 \sin^2(\Theta) & \\ & & \mathbf{0}_{(d-2r) \times (d-2r)} \end{bmatrix} \right) = r \cdot \nu_1 \cdot \cos^2(\theta) + r \cdot \nu_2 \cdot \sin^2(\theta), \end{aligned}$$

where we used the fact that the principal angles are all equal to θ . Using a similar argument,

$$\text{Tr}(\Sigma_t^\top \mathbf{A}^2) = r \cdot \nu_1^2 \cdot \cos^2(\theta) + r \cdot \nu_2^2 \cdot \sin^2(\theta) + \epsilon \text{Tr}(\mathbf{A}^2)$$

Simplifying the Test Risk. Substituting the expressions for the $\text{Tr}(\cdot)$ terms into Equation (14) yields

$$\begin{aligned} \mathbb{E} \left[\left(\tilde{y} - g_{\text{ATT}}^*(\tilde{\mathbf{z}}_q, \tilde{\mathbf{Z}}_{\mathcal{M}}) \right)^2 \right] &= \left(\frac{1}{m} (r\nu_1^2 + (d-r)\nu_2^2) + 1 \right) (r + \epsilon d + \sigma^2) \\ &\quad - 2 (r\nu_1 \cos^2(\theta) + r\nu_2 \sin^2(\theta) + (r\nu_1 + (d-r)\nu_2)\epsilon) \\ &\quad + \frac{m+1}{m} (r\nu_1^2 \cos^2(\theta) + r\nu_2^2 \sin^2(\theta) + (r\nu_1^2 + (d-r)\nu_2^2)\epsilon) \end{aligned}$$

Substituting the expressions for ν_1 and ν_2 and taking $\epsilon \rightarrow 0$ results in the following:

$$\begin{aligned} \lim_{\epsilon \rightarrow 0} \mathbb{E} \left[\left(\tilde{y} - g_{\text{ATT}}^*(\tilde{\mathbf{z}}_q, \tilde{\mathbf{Z}}_{\mathcal{M}}) \right)^2 \right] &= \left(\frac{rn^2}{m(n+1+r+\sigma^2)^2} + 1 \right) (r + \sigma^2) \\ &\quad - \frac{2rn \cos^2(\theta)}{n+1+r+\sigma^2} + \frac{(m+1)rn^2 \cos^2(\theta)}{m(n+1+r+\sigma^2)^2} \end{aligned}$$

Subsequently taking $m, n \rightarrow \infty$ yields

$$\lim_{m \rightarrow \infty} \lim_{n \rightarrow \infty} \lim_{\epsilon \rightarrow 0} \mathbb{E} \left[\left(\tilde{y} - g_{\text{ATT}}^*(\tilde{\mathbf{z}}_q, \tilde{\mathbf{Z}}_{\mathcal{M}}) \right)^2 \right] = r + \sigma^2 - r \cos^2(\theta) = r \sin^2(\theta) + \sigma^2,$$

which completes the proof. \square

C.4 Auxiliary Results

Here, we provide auxiliary results to support the proofs in Sections C.1 to C.3.

C.4.1 Optimal Linear Attention Weights

We first provide results on the form of the weights matrices after training a single-layer linear attention model on the objective Equation (1). The following results are largely inspired by Theorem 1 in Li et al. (2024b), but are slightly different since we consider a normalization factor of $1/n$ in our linear attention model.

Lemma 2 (Optimal Attention Weights (Li et al., 2024b)). *Consider the independent data model in Equation (4) with $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_s)$, and let $n \in \mathbb{N}$ denote the in-context prompt length used at training. Then, the optimal linear attention weights obtained by minimizing the loss in Equation (1) are given by*

$$\mathbf{W}_K^* = \mathbf{W}_V^* = \mathbf{I}_{d+1}, \quad \text{and} \quad \mathbf{W}_Q^* = \begin{bmatrix} \mathbf{A} & \mathbf{0}_d \\ \mathbf{0}_d^\top & 0 \end{bmatrix} \quad (15)$$

where $\mathbf{A} = \left(\frac{n+1}{n} \mathbf{I}_d + \frac{M_s}{n} \Sigma_s^{-1} \right)^{-1}$ and $M_s = \text{Tr}(\Sigma_s) + \sigma^2$, with empirical risk $\mathcal{L}_s^* = M_s - \text{Tr}(\Sigma_s \mathbf{A})$.

Proof. The proof is the same as that of Theorem 1 in Li et al. (2024b) by absorbing the $1/n$ factor into \mathbf{W}_Q . \square

Lemma 3 (Optimal Attention Weights for Mixture of K Gaussians). *Consider the independent data model in Equation (4) with $\mathbf{w} \sim \sum_{k=1}^K \gamma_k \cdot \mathcal{N}(\mathbf{0}, \Sigma_{s,k})$ for $\gamma_k \in (0, 1)$ for all $k \in [K]$ and $\sum_{k=1}^K \gamma_k = 1$. Let $n \in \mathbb{N}$ denote the in-context prompt length used at training. Define $\Sigma = \sum_{k=1}^K \gamma_k \cdot \Sigma_{s,k}$. Then, the optimal linear attention weights obtained by minimizing the loss in Equation (1) are given by*

$$\mathbf{W}_K^* = \mathbf{W}_V^* = \mathbf{I}_{d+1}, \quad \text{and} \quad \mathbf{W}_Q^* = \begin{bmatrix} \mathbf{A} & \mathbf{0}_d \\ \mathbf{0}_d^\top & 0 \end{bmatrix}, \quad (16)$$

where $\mathbf{A} = \left(\frac{n+1}{n} \mathbf{I}_d + \frac{M_s}{n} \Sigma^{-1} \right)^{-1}$ and $M_s = \text{Tr}(\Sigma) + \sigma^2$, with empirical risk $\mathcal{L}_s^* = M_s - \text{Tr}(\Sigma \mathbf{A})$.

Proof. It is straightforward to see that if $\mathbf{w} \sim \sum_{k=1}^K \gamma_k \gamma \cdot \mathcal{N}(\mathbf{0}, \Sigma_{s,k})$, then

$$\Sigma_s := \text{Cov}(\mathbf{w}) = \sum_{k=1}^K \gamma_k \cdot \Sigma_{s,k}.$$

Then, the proof is equivalent to that of Lemma 2 under the new form of Σ_s . \square

C.4.2 Miscellaneous Results

Lemma 4. Let $0 \prec \Sigma \in \mathbb{R}^{d \times d}$ and $c, k > 0$ be constants. Then,

$$(c \cdot \mathbf{I}_d + k \cdot \Sigma^{-1})^{-1} = \mathbf{V} \begin{bmatrix} \frac{\lambda_1}{c \cdot \lambda_1 + k} & 0 & \dots & 0 \\ 0 & \frac{\lambda_2}{c \cdot \lambda_2 + k} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{\lambda_d}{c \cdot \lambda_d + k} \end{bmatrix} \mathbf{V}^\top, \quad (17)$$

where $\mathbf{V} \in \mathbb{R}^{d \times d}$ is an orthonormal matrix whose columns are eigenvectors of Σ , and λ_i is the i^{th} largest eigenvalue of Σ .

Proof. Since $\Sigma \succ 0$, there exists an eigendecomposition $\Sigma = \mathbf{V} \Lambda \mathbf{V}^\top$ such that \mathbf{V} is an orthonormal matrix and Λ is a diagonal matrix consisting of the real, positive eigenvalues of Σ , denoted as $\lambda_1, \lambda_2, \dots, \lambda_d$. Thus,

$$\begin{aligned} \Sigma^{-1} = \mathbf{V} \Lambda^{-1} \mathbf{V}^\top &\implies c \cdot \mathbf{I}_d + k \cdot \Sigma^{-1} = \mathbf{V} \underbrace{\begin{bmatrix} c + \frac{k}{\lambda_1} & 0 & \dots & 0 \\ 0 & c + \frac{k}{\lambda_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & c + \frac{k}{\lambda_d} \end{bmatrix}}_{\tilde{\Lambda}} \mathbf{V}^\top \\ &\implies (c \cdot \mathbf{I}_d + k \cdot \Sigma^{-1})^{-1} = \mathbf{V} \tilde{\Lambda}^{-1} \mathbf{V}^\top, \end{aligned}$$

which completes the proof. □