
How Important is the Train-Validation Split in Meta-Learning?

Anonymous Author(s)

Affiliation

Address

email

Abstract

Meta-learning aims to perform fast adaptation on a new task through learning a “prior” from multiple existing tasks. A common practice in meta-learning is to perform a *train-validation split* where the prior adapts to the task on one split of the data, and the resulting predictor is evaluated on another split. Despite its prevalence, the importance of the train-validation split is not well understood either in theory or in practice, particularly in comparison to the more direct *non-splitting* method, which uses all the per-task data for both training and evaluation.

We provide a detailed theoretical study on whether and when the train-validation split is helpful on the linear centroid meta-learning problem, in the asymptotic setting where the number of tasks goes to infinity. We show that the splitting method converges to the optimal prior as expected, whereas the non-splitting method does not in general without structural assumptions on the data. In contrast, if the data are generated from linear models (the realizable regime), we show that both the splitting and non-splitting methods converge to the optimal prior. Further, perhaps surprisingly, our main result shows that the non-splitting method achieves a *strictly better* asymptotic excess risk under this data distribution, even when the regularization parameter and split ratio are optimally tuned for both methods. Our results highlight that data splitting may not always be preferable, especially when the data is realizable by the model. We validate our theories by experimentally showing that the non-splitting method can indeed outperform the splitting method, on both simulations and real meta-learning tasks.

1 Introduction

Meta-learning, also known as “learning to learn”, has recently emerged as a powerful paradigm for learning to adapt to unseen tasks (Schmidhuber, 1987). The high-level methodology in meta-learning is akin to how human beings learn new skills, which is typically done by relating to certain prior experience that makes the learning process easier. More concretely, meta-learning does not train one model for each individual task, but rather learns a “prior” model from multiple existing tasks so that it is able to quickly adapt to unseen new tasks. Meta-learning has been successfully applied to many real problems, including few-shot image classification (Finn et al., 2017; Snell et al., 2017), hyper-parameter optimization (Franceschi et al., 2018), low-resource machine translation (Gu et al., 2018) and short event sequence modeling (Xie et al., 2019).

A common practice in meta-learning algorithms is to perform a *sample splitting*, where the data within each task is divided into a *training split* which the prior uses to adapt to a task-specific predictor, and a *validation split* on which we evaluate the performance of the task-specific predictor (Nichol et al., 2018; Rajeswaran et al., 2019; Fallah et al., 2020; Wang et al., 2020a). This sample splitting is believed to be crucial, as it matches the evaluation criterion at meta-test time,

where we perform adaptation on training data from a new task but evaluate its performance on unseen data from the same task. However, despite the aforementioned importance, performing the train-validation split has a potential drawback from the data efficiency perspective — Because of the split, neither the training nor the evaluation stage is able to use all the available per-task data. Further, performing the train-validation split is also not the only option in practice: there exist algorithms such as Reptile (Nichol & Schulman, 2018) and Meta-MinibatchProx (Zhou et al., 2019) that can instead use all the per-task data for training the task-specific predictor and also perform well empirically on benchmark tasks. These algorithms modify the loss function in the outer loop so that the training loss no longer matches the meta-test loss, but may have the advantage in terms of data efficiency for the overall problem of learning the best prior. So far it is theoretically unclear how these two approaches (with/without train-validation split) compare with each other, which motivates us to ask the following

Question: Is the train-validation split *necessary* and *optimal* in meta-learning?

In this paper, we perform a detailed theoretical study on the importance of the train-validation split. We consider the linear centroid meta-learning problem (Denevi et al., 2018b), where for each task we learn a linear predictor that is close to a common centroid in the inner loop, and find the best centroid in the outer loop (see Section 2 for the detailed problem setup). We compare two outer-loop algorithms of either performing the train-validation split (the *train-val method*) or using all the per-task data for both training and evaluation (the *train-train method*). Specifically, we compare the two methods when the number of tasks T is large, and examine if and how fast they converge to the (properly defined) best centroid at meta-test time. We summarize our contributions as follows:

- On the linear centroid meta-learning problem, we show that the train-validation split is necessary in the general agnostic setting: As $T \rightarrow \infty$, the train-val method converges to the optimal centroid for test-time adaptation, whereas the train-train method does not without further assumptions on the tasks (Section 3). The convergence of the train-val method is expected since its (population) training loss is equivalent to the meta-test time loss, whereas the non-convergence of the train-train method is because these two losses are not equivalent in general.
- Our main theoretical contribution is to show that the train-validation split is not necessary and even non-optimal, in the perhaps more interesting regime when there are structural assumptions on the tasks: When the data are generated from noiseless linear models, both the train-val and train-train methods converge to the common best centroid, and the train-train method achieves strictly better (asymptotic) estimation error and test loss than the train-val method (Section 4). This is in stark contrast with the agnostic case, and suggests that data efficiency may indeed be more important when the tasks have a nice structure.
- We perform meta-learning experiments on simulations and benchmark few-shot image classification tasks, showing that the train-train method can consistently outperform the train-val method (Section 5 & Appendix A). This validates our theories and presents empirical evidence that sample splitting may not be crucial; methods that utilize the per-task data more efficiently may be preferred.

1.1 Related work

Meta-learning and representation learning theory Baxter (2000) provided the first theoretical analysis of meta-learning via covering numbers, and Maurer et al. (2016) improved the analysis via Gaussian complexity techniques. Another recent line of theoretical work analyzed gradient-based meta-learning methods (Denevi et al., 2018a; Finn et al., 2019; Khodak et al., 2019) and showed guarantees for convex losses by using tools from online convex optimization. Saunshi et al. (2020) proved the success of Reptile in a one-dimensional subspace setting. Wang et al. (2020b) compared the performance of train-train and train-val methods for learning the learning rate. Denevi et al. (2018b) proposed the linear centroid model studied in this paper, and provided generalization error bounds for train-val method; the bounds proved also hold for train-train, so are not sharp enough to compare the two algorithms. On the representation learning end, Du et al. (2020); Tripuraneni et al. (2020) showed that ERM can successfully pool data across tasks to learn the representation. Yet the focus is on the accurate estimation of the common representation, not on the fast adaptation of the learned prior. Lastly, we remark that there are analyses for other representation learning schemes (McNamara & Balcan, 2017; Galanti et al., 2016; Alquier et al., 2016). **Multi-task learning** Multi-task learning also exploits structures and similarities across multiple tasks. The earliest idea dates back to Caruana (1997); Thrun & Pratt (1998); Baxter (2000), initially in connections to

neural network models. They further motivated other approaches using kernel methods (Evgeniou et al., 2005; Argyriou et al., 2007) and multivariate linear regression models with structured sparsity (Liu et al., 2009, 2015). More recent advances on deep multi-task learning focus on learning shared intermediate representations across tasks Ruder (2017). These multi-task learning approaches usually minimize the joint empirical risk over all tasks, and the models for different tasks are enforced to share a large amount of parameters. In contrast, meta-learning only requires the models to share the same “prior”, which is more flexible than multi-task learning.

2 Preliminaries

Linear centroid meta-learning We instantiate our study on the *linear centroid meta-learning problem* (also known as learning to learn around a common mean (Denevi et al., 2018b)), where we wish to learn a task-specific linear predictor $\mathbf{w}_t \in \mathbb{R}^d$ in the inner loop for each task t , and learn a “centroid” \mathbf{w}_0 in the outer loop that enables fast adaptation to \mathbf{w}_t within each task:

Find the best centroid $\mathbf{w}_0 \in \mathbb{R}^d$ for adapting to a linear predictor \mathbf{w}_t on each task t .

Formally, we assume that we observe training data from $T \geq 1$ tasks, where for each task index t we sample a task p_t (a distribution over $\mathbb{R}^d \times \mathbb{R}$) from some distribution of tasks Π , and observe n examples $(\mathbf{X}_t, \mathbf{y}_t) \in \mathbb{R}^{n \times d} \times \mathbb{R}^n$ that are drawn i.i.d. from p_t :

$$p_t \sim \Pi, \quad (\mathbf{X}_t, \mathbf{y}_t) = \{(\mathbf{x}_{t,i}, y_{t,i})\}_{i=1}^n \quad \text{where } (\mathbf{x}_{t,i}, y_{t,i}) \stackrel{\text{iid}}{\sim} p_t. \quad (1)$$

We do not make further assumptions on (n, d) ; in particular, we allow the underdetermined setting $n \leq d$, in which there exists (one or many) interpolators $\tilde{\mathbf{w}}_t$ that perfectly fit the data: $\mathbf{X}_t \tilde{\mathbf{w}}_t = \mathbf{y}_t$.

Inner loop: Ridge solver with biased regularization towards the centroid Our goal in the inner loop is to find a linear predictor \mathbf{w}_t that fits the data in task t while being close to the given “centroid” $\mathbf{w}_0 \in \mathbb{R}^d$. We instantiate this through ridge regression (i.e. linear regression with L_2 regularization) where the regularization biases \mathbf{w}_t towards the centroid. Formally, for any $\mathbf{w}_0 \in \mathbb{R}^d$ and any dataset (\mathbf{X}, \mathbf{y}) , we consider the algorithm

$$\mathcal{A}_\lambda(\mathbf{w}_0; \mathbf{X}, \mathbf{y}) := \arg \min_{\mathbf{w}} \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w} - \mathbf{w}_0\|^2 = \mathbf{w}_0 + (\mathbf{X}^\top \mathbf{X} + n\lambda \mathbf{I}_d)^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}_0),$$

where $\lambda > 0$ is the regularization strength (typically a tunable hyper-parameter).

Outer loop: finding the best centroid In the outer loop, our goal is to find the best centroid \mathbf{w}_0 . The standard approach in meta-learning is to perform a *train-validation split*, that is, (1) execute the inner solver on a first split of the task-specific data, and (2) evaluate the loss on a second split, yielding a function of \mathbf{w}_0 that we can optimize. This two-stage procedure can be written as

Compute $\mathbf{w}_t(\mathbf{w}_0) = \mathcal{A}_\lambda(\mathbf{w}_0; \mathbf{X}_t^{\text{train}}, \mathbf{y}_t^{\text{train}})$, and evaluate $\|\mathbf{y}_t^{\text{val}} - \mathbf{X}_t^{\text{val}} \mathbf{w}_t(\mathbf{w}_0)\|^2$.
 where $(\mathbf{X}_t^{\text{train}}, \mathbf{y}_t^{\text{train}}) = \{(\mathbf{x}_{t,i}, y_{t,i})\}_{i=1}^{n_1}$ and $(\mathbf{X}_t^{\text{val}}, \mathbf{y}_t^{\text{val}}) = \{(\mathbf{x}_{t,i}, y_{t,i})\}_{i=n_1+1}^n$ are two disjoint splits of the per-task data $(\mathbf{X}_t, \mathbf{y}_t)$ of size (n_1, n_2) , with $n_1 + n_2 = n$. Written concisely, this is to consider the “split loss”

$$\ell_t^{\text{tr-val}}(\mathbf{w}_0) := \frac{1}{2n_2} \|\mathbf{y}_t^{\text{val}} - \mathbf{X}_t^{\text{val}} \mathcal{A}_\lambda(\mathbf{w}_0; \mathbf{X}_t^{\text{train}}, \mathbf{y}_t^{\text{train}})\|^2. \quad (2)$$

In this paper, we will also consider an alternative version, where we do not perform the train-validation split, but instead use *all the per-task data for both training and evaluation*. Mathematically, this is to look at the “non-split loss”

$$\ell_t^{\text{tr-tr}}(\mathbf{w}_0) := \frac{1}{2n} \|\mathbf{y}_t - \mathbf{X}_t \mathcal{A}_\lambda(\mathbf{w}_0; \mathbf{X}_t, \mathbf{y}_t)\|^2. \quad (3)$$

Our overall algorithm is to solve the empirical risk minimization (ERM) problem on the T observed tasks, using either one of the two losses above:

$$\begin{aligned} \hat{L}_T^{\text{tr-val}}(\mathbf{w}_0) &:= \frac{1}{T} \sum_{t=1}^T \ell_t^{\text{tr-val}}(\mathbf{w}_0) \quad \text{and} \quad \hat{L}_T^{\text{tr-tr}}(\mathbf{w}_0) := \frac{1}{T} \sum_{t=1}^T \ell_t^{\text{tr-tr}}(\mathbf{w}_0), \\ \hat{\mathbf{w}}_{0,T}^{\{\text{tr-val}, \text{tr-tr}\}} &:= \arg \min_{\mathbf{w}_0} \hat{L}_T^{\{\text{tr-val}, \text{tr-tr}\}}(\mathbf{w}_0). \end{aligned} \quad (4)$$

Let $L^{\{\text{tr-val}, \text{tr-tr}\}}(\mathbf{w}_0) := \mathbb{E}_{p_t \sim \Pi, (\mathbf{X}_t, \mathbf{y}_t) \sim p_t} [\ell_t^{\{\text{tr-val}, \text{tr-tr}\}}(\mathbf{w}_0)]$ be the population risks.

(Meta-)Test time The meta-test time performance of any meta-learning algorithm is a joint function of the (learned) centroid \mathbf{w}_0 and the inner algorithm Alg. Upon receiving a new task $p_{T+1} \sim \Pi$ and training data $(\mathbf{X}_{T+1}, \mathbf{y}_{T+1}) \in \mathbb{R}^{n \times d} \times \mathbb{R}^n$, we run the inner loop Alg with prior \mathbf{w}_0 on the training data, and evaluate it on an (unseen) test example $(\mathbf{x}', y') \sim p_{T+1}$:

$$L^{\text{test}}(\mathbf{w}_0; \text{Alg}) := \mathbb{E}_{p_{T+1} \sim \Pi} \mathbb{E}_{(\mathbf{X}_{T+1}, \mathbf{y}_{T+1}), (\mathbf{x}', y') \sim p_{T+1}} \left[\frac{1}{2} (\mathbf{x}'^\top \text{Alg}(\mathbf{w}_0; \mathbf{X}_{T+1}, \mathbf{y}_{T+1}) - y')^2 \right].$$

Additionally, for both train-val and train-train methods, we need to ensure that the inner loop used for meta-testing is exactly the same as that used in meta-training. Therefore, the meta-test performance for the train-val and train-train methods above should be evaluated as $L_{\lambda, n_1}^{\text{test}}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}) := L^{\text{test}}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}; \mathcal{A}_{\lambda, n_1})$, and $L_{\lambda, n}^{\text{test}}(\hat{\mathbf{w}}_{0,T}^{\text{tr-tr}}) := L^{\text{test}}(\hat{\mathbf{w}}_{0,T}^{\text{tr-tr}}; \mathcal{A}_{\lambda, n})$, where $\mathcal{A}_{\lambda, m}$ denotes the ridge solver with regularization strength $\lambda > 0$ on $m \leq n$ data points. Finally, we let

$$\mathbf{w}_{0,*}(\lambda; n) = \arg \min_{\mathbf{w}_0} L_{\lambda, n}^{\text{test}}(\mathbf{w}_0) \quad (5)$$

denote the best centroid if the inner loop uses $\mathcal{A}_{\lambda, n}$. The performance of the train-val algorithm $\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}$ should be compared against $\mathbf{w}_{0,*}(\lambda, n_1)$, whereas the train-train algorithm $\hat{\mathbf{w}}_{0,T}^{\text{tr-tr}}$ should be compared against $\mathbf{w}_{0,*}(\lambda, n)$.

2.1 Task-abundant setting through asymptotic analysis

In this paper we are interested in the *task-abundant* setting where we fix some finite (d, n) and let T be very large. We analyze such a task-abundant setting through the asymptotic analysis framework, that is, examine the limiting properties of the algorithm (e.g. $\hat{\mathbf{w}}_{0,T}^{\{\text{tr-val}, \text{tr-tr}\}}$) as $T \rightarrow \infty$. Here we set up the basic notation of asymptotic analysis required in this paper, and refer the readers to (Van der Vaart, 2000) for a detailed treatment.

Asymptotic rate of estimation & excess risk Let L be any population risk with minimizer $\mathbf{w}_{0,*}$ (which we assume is unique), \hat{L}_T be the empirical risk on the observed data from T tasks, and $\hat{\mathbf{w}}_{0,T}$ be the minimizer of \hat{L}_T (i.e. the ERM). We say that $\hat{\mathbf{w}}_{0,T}$ is **consistent** if $\hat{\mathbf{w}}_{0,T} \rightarrow \mathbf{w}_{0,*}$ in probability as $T \rightarrow \infty$. In typical scenarios, for consistent ERMs, the limiting distribution of $\hat{\mathbf{w}}_{0,T}$ is asymptotically normal with a known covariance matrix, as is characterized in the following classical result (see, e.g. (Van der Vaart, 2000, Section 5.3) and also (Liang, 2016)).

Proposition 1 (Asymptotic normality and excess risk of ERMs). *Suppose the ERM $\hat{\mathbf{w}}_{0,T}$ is consistent and certain regularity conditions hold, then we have*

$$\begin{aligned} \sqrt{T} \cdot (\hat{\mathbf{w}}_{0,T} - \mathbf{w}_{0,*}) &\xrightarrow{d} \mathcal{N}(\mathbf{0}, \nabla^2 L(\mathbf{w}_{0,*})^{-1} \text{Cov}(\nabla \ell_t(\mathbf{w}_{0,*})) \nabla^2 L(\mathbf{w}_{0,*})^{-1}) =: P_{\mathbf{w}}, \\ T \cdot (L(\hat{\mathbf{w}}_{0,T}) - L(\mathbf{w}_{0,*})) &\xrightarrow{d} \Delta^\top \nabla^2 L(\mathbf{w}_{0,*}) \Delta \quad \text{where } \Delta \sim P_{\mathbf{w}}. \end{aligned}$$

where \xrightarrow{d} denotes convergence in distribution and $\ell_t : \mathbb{R}^d \rightarrow \mathbb{R}$ is the loss function on a single task.

When this happens, we define the asymptotic rate of estimation (in MSE loss) and asymptotic excess risk $\hat{\mathbf{w}}_{0,T}$ as those of its limiting distribution:

$$\begin{aligned} \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}) &:= \mathbb{E}_{\Delta \sim P_{\mathbf{w}}} [\|\Delta\|^2] = \text{tr}(\nabla^2 L(\mathbf{w}_{0,*})^{-1} \text{Cov}(\nabla \ell_t(\mathbf{w}_{0,*})) \nabla^2 L(\mathbf{w}_{0,*})^{-1}) \\ \text{AsymExcessRisk}(\hat{\mathbf{w}}_{0,T}) &:= \mathbb{E}_{\Delta \sim P_{\mathbf{w}}} [\Delta^\top \nabla^2 L(\mathbf{w}_{0,*}) \Delta] = \text{tr}(\nabla^2 L(\mathbf{w}_{0,*})^{-1} \text{Cov}(\nabla \ell_t(\mathbf{w}_{0,*}))). \end{aligned}$$

3 The importance of sample splitting

We begin by analyzing whether the algorithms $\hat{\mathbf{w}}_{0,T}^{\{\text{tr-val}, \text{tr-tr}\}}$ defined in (4) converge to the best test-time centroid $\mathbf{w}_{0,*}(\lambda; n_1)$ or $\mathbf{w}_{0,*}(\lambda; n)$ (defined (5)) respectively as $T \rightarrow \infty$, in the general situation where we do not make structural assumptions on the data distribution p_t .

Proposition 2 (Consistency and asymptotics of train-val method). *Suppose $\mathbb{E}_{\mathbf{x} \sim p_t} [\mathbf{x} \mathbf{x}^\top] \succ \mathbf{0}$, $\mathbb{E}_{\mathbf{x} \sim p_t} [\|\mathbf{x}\|^4] < \infty$ and $\mathbb{E}_{(\mathbf{x}, y) \sim p_t} [\|\mathbf{x} y\|] < \infty$ for almost surely all $p_t \sim \Pi$. Then for any $\lambda > 0$*

and any (n_1, n_2) such that $n_1 + n_2 = n$, the sample splitting algorithm $\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}$ converges to the best test-time centroid: $\hat{\mathbf{w}}_{0,T}^{\text{tr-val}} \rightarrow \mathbf{w}_{0,*}(\lambda, n_1)$ almost surely as $T \rightarrow \infty$. Further, we have

$$\begin{aligned} \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}) &= \text{tr}(\nabla^{-2} L_{\lambda, n_1}^{\text{test}}(\mathbf{w}_{0,*}(\lambda, n_1)) \cdot \text{Cov}(\nabla \ell_t^{\text{tr-val}}(\mathbf{w}_{0,*}(\lambda, n_1))) \cdot \nabla^{-2} L_{\lambda, n_1}^{\text{test}}(\mathbf{w}_{0,*}(\lambda, n_1))), \\ \text{AsymExcessRisk}_{L_{\lambda, n_1}^{\text{test}}}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}) &= \text{tr}(\nabla^{-2} L_{\lambda, n_1}^{\text{test}}(\mathbf{w}_{0,*}(\lambda, n_1)) \cdot \text{Cov}(\nabla \ell_t^{\text{tr-val}}(\mathbf{w}_{0,*}(\lambda, n_1)))). \end{aligned}$$

Proposition 3 (Inconsistency of train-train method). *There exists a distribution of tasks Π on $d = 1$ satisfying the conditions in Proposition 2 on which the train-train method does not converge to the best test-time centroid: for any $n \geq 1$ and any $\lambda > 0$, the estimation error $\|\hat{\mathbf{w}}_{0,T}^{\text{tr-tr}} - \mathbf{w}_{0,*}(\lambda, n)\|$ and the excess risk $L_{\lambda, n}^{\text{test}}(\hat{\mathbf{w}}_{0,T}^{\text{tr-tr}}) - L_{\lambda, n}^{\text{test}}(\mathbf{w}_{0,*}(\lambda, n))$ are both bounded away from 0 almost surely as $T \rightarrow \infty$.*

Propositions 2 and 3 justify the importance of sample splitting: the train-val method converges the best test-time centroid, whereas the train-train method does not converge to the best centroid in general. Appendix B gives the proofs of Proposition 2 and 3.

4 Is sample splitting always optimal?

Proposition 3 states a negative result for the train-train method, showing that it does not converge to the best test-time centroid without further assumptions on the data distribution. However, such a negative result is inherently *worst-case*, and does not preclude the possibility that there exists a data distribution on which the train-train method can also work well. In this section, we construct a simple data distribution in which we can analyze the performance of both the train-val and the train-train methods more explicitly, showing that sample splitting is indeed not optimal, and the train-train method can work better.

Realizable linear model We consider the following instantiation of the (generic) data distribution assumption in (1): We assume that each task p_t is specified by a $\mathbf{w}_t \in \mathbb{R}^d$ sampled from some distribution Π (overloading notation), and the observed data follows the noiseless linear model with ground truth parameter \mathbf{w}_t :

$$\mathbf{y}_t = \mathbf{X}_t \mathbf{w}_t, \quad (6)$$

where the inputs $\mathbf{x}_{t,i} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \mathbf{I}_d)$ and are independent of \mathbf{w}_t . We assume that Π has a finite second moment (i.e. $\mathbb{E}_{\mathbf{w}_t \sim \Pi}[\|\mathbf{w}_t\|^2] < \infty$).

4.1 Consistency and asymptotic rates

We begin by showing that on this task and data distribution, the population best centroids $\mathbf{w}_{0,*}(\lambda, n) = \arg \min_{\mathbf{w}_0} L_{\lambda, n}^{\text{test}}(\mathbf{w}_0)$ is the same for any (λ, n) , and both the train-val and train-train methods are asymptotically consistent and converge to same best centroid.

Theorem 4 (Consistency of both train-val and train-train methods). *On the realizable linear model (6), the test-time meta loss for all $\lambda > 0$ and all n is minimized at the same point, that is, the mean of the ground truth parameters:*

$$\mathbf{w}_{0,*}(\lambda, n) = \arg \min_{\mathbf{w}_0} L_{\lambda, n}^{\text{test}}(\mathbf{w}_0) = \mathbf{w}_{0,*} := \mathbb{E}_{\mathbf{w}_t \sim \Pi}[\mathbf{w}_t], \quad \text{for all } \lambda > 0, n.$$

Further, both the train-val method and the train-train method is asymptotically consistent: for any $\lambda > 0, n$, and (n_1, n_2) , we have

$$\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}(n_1, n_2; \lambda) \rightarrow \mathbf{w}_{0,*} \quad \text{and} \quad \hat{\mathbf{w}}_{0,T}^{\text{tr-tr}}(n; \lambda) \rightarrow \mathbf{w}_{0,*} \quad \text{almost surely as } T \rightarrow \infty.$$

Theorem 4 suggests that we are now able to compare performance of the two methods based on their asymptotic parameter estimation error (for estimating $\mathbf{w}_{0,*}$), which we state in the following result. Throughout the rest of this section, let

$$R^2 := \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}_{0,*}\|^2]$$

denote the variance of \mathbf{w}_t .

Theorem 5 (Exact asymptotic rates of the train-val and train-train methods). *Let $\rho_{\text{tr-tr}} = \frac{\mathbb{E}[\sum_{i=1}^d (\sigma_i^{(n)})^2 / (\sigma_i^{(n)} + \lambda)^4]}{(\mathbb{E}[\sum_{i=1}^d \sigma_i^{(n)} / (\sigma_i^{(n)} + \lambda)^2])^2}$ and $\rho_{\text{tr-val}} = \frac{\mathbb{E}[(\sum_{i=1}^d \lambda^2 / (\sigma_i^{(n_1)} + \lambda)^2)^2 + (n_2 + 1) \sum_{i=1}^d \lambda^4 / (\sigma_i^{(n_1)} + \lambda)^4]}{(\mathbb{E}[\sum_{i=1}^d \lambda^2 / (\sigma_i^{(n_1)} + \lambda)^2])^2}$, where for any n , $\sigma_1^{(n)} \geq \dots \geq \sigma_d^{(n)}$ denotes the eigenvalues of the matrix $\frac{1}{n} \mathbf{X}_t^\top \mathbf{X}_t \in \mathbb{R}^{d \times d}$, where we recall $\mathbf{X}_t \in \mathbb{R}^{n \times d}$ is a random matrix with i.i.d. standard Gaussian entries. For any (n, d) , we have on the realizable linear model (6) that*

$$\text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-tr}}(n; \lambda)) = dR^2 \rho_{\text{tr-tr}}, \quad \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}(n_1, n_2; \lambda)) = \frac{dR^2 \rho_{\text{tr-val}}}{n_2}.$$

See its proof in Appendix C.2. Theorem 5 follows straightforwardly from the classical asymptotic result for empirical risk minimization (Van der Vaart, 2000) and simplifications of certain matrix traces in terms of the spectrum of the empirical covariance matrix $\frac{1}{n} \sum_{t=1}^n \mathbf{X}_t \mathbf{X}_t^\top$.

4.2 Comparison of train-val vs. train-train

We now present our main theoretical result, which shows that the train-train method achieves a strictly better asymptotic MSE than the train-val method, in the proportional limit of $d, n \rightarrow \infty$ and $d/n \rightarrow \gamma \in (0, 1)$.

Theorem 6 (Optimal rates of the train-val and train-train method in the proportional limit). *In the high-dimensional limiting regime $d, n \rightarrow \infty$, $d/n \rightarrow \gamma \in (0, \infty)$, the optimal rate of the train-train method obtained by tuning the regularization $\lambda \in (0, \infty)$ satisfies*

$$\inf_{\lambda > 0} \lim_{d, n \rightarrow \infty, d/n = \gamma} \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-tr}}(n; \lambda)) = \inf_{\lambda > 0} \rho_{\lambda, \gamma} R^2 \stackrel{(*)}{\leq} \max \left\{ 1 + \frac{5}{27} \gamma, \frac{5}{27} + \gamma \right\} \cdot R^2,$$

and the inequality becomes equality at $\gamma = 1$. In contrast, the optimal rate of the train-val method by tuning the regularization $\lambda \in (0, \infty)$ and split ratio $s \in (0, 1)$ is

$$\inf_{\lambda > 0, s \in (0, 1)} \lim_{d, n \rightarrow \infty, d/n = \gamma} \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}(ns, n(1-s); \lambda)) = (1 + \gamma) R^2.$$

As $\max\{1 + 5\gamma/27, 5/27 + \gamma\} < 1 + \gamma$ ($\forall \gamma > 0$), the train-train method achieves a strictly better asymptotic rate than the train-val method when λ and s are optimally tuned in both methods.

Implications Comparison between the analytical upper bound $\max\{1 + 5\gamma/27, 5/27 + \gamma\} R^2$ for train-train $(1 + \gamma) R^2$ for train-val in Theorem 6 shows that the train-train method achieves a strictly better asymptotic MSE (and thus also asymptotic excess risk) than the train-val method, for any $\gamma > 0$. (See Figure 1(a) for a visualization of the optimal rates.) Perhaps surprisingly, this suggests that the train-train method is not only “correct” (converging to the best centroid), but can be even better than the train-val method, when the data is structured. While we reached such a conclusion on this particular problem of linear centroid meta-learning, we suspect that this phenomenon to be not restricted to this problem, and may hold in more generality when data is structured or when the signal-to-noise ratio is high. The proof of Theorem 6 is deferred to Appendix C.7.

5 Experiments

Simulation. We experiment on the realizable linear model studied in Section 4 and show that the train-train method indeed achieves better performance than the train-val method. Both performances agree well with the asymptotic theoretical prediction, even at a fairly small $(n, d) = (20, 60)$ (Figure 1, more details in Appendix A.1).

Real data. We additionally find that the train-train method consistently outperforms the train-val method on few-shot image classification benchmarks (Table 1, more details in Appendix A.2).

6 Conclusion

We study the importance of train-validation split on the linear-centroid meta-learning problem, and show that the necessity and optimality of train-validation split depends greatly on whether the tasks are structured: the sample splitting is necessary in general situations, and not necessary and non-optimal when the tasks are nicely structured. It would be of interest to study whether a similar conclusion holds on other meta-learning problems such as learning a representation, or whether our insights can be used towards designing meta-learning algorithms with better empirical performance.

References

- Pierre Alquier, The Tien Mai, and Massimiliano Pontil. Regret bounds for lifelong learning. *arXiv preprint arXiv:1610.08628*, 2016.
- Greg W Anderson, Alice Guionnet, and Ofer Zeitouni. *An introduction to random matrices*, volume 118. Cambridge university press, 2010.
- Andreas Argyriou, Theodoros Evgeniou, and Massimiliano Pontil. Multi-task feature learning. In *Advances in neural information processing systems*, pp. 41–48, 2007.
- Zhidong Bai and Jack W Silverstein. *Spectral analysis of large dimensional random matrices*, volume 20. Springer, 2010.
- Jonathan Baxter. A model of inductive bias learning. *J. Artif. Int. Res.*, 2000.
- Rich Caruana. Multitask learning. *Machine Learning*, 28(1):41–75, Jul 1997. ISSN 1573-0565. doi: 10.1023/A:1007379606734. URL <https://doi.org/10.1023/A:1007379606734>.
- Giulia Denevi, Carlo Ciliberto, Dimitris Stamos, and Massimiliano Pontil. Incremental learning-to-learn with statistical guarantees. *arXiv preprint arXiv:1803.08089*, 2018a.
- Giulia Denevi, Carlo Ciliberto, Dimitris Stamos, and Massimiliano Pontil. Learning to learn around a common mean. In *Advances in Neural Information Processing Systems*, pp. 10169–10179, 2018b.
- Edgar Dobriban, Stefan Wager, et al. High-dimensional asymptotics of prediction: Ridge regression and classification. *The Annals of Statistics*, 46(1):247–279, 2018.
- Simon S Du, Wei Hu, Sham M Kakade, Jason D Lee, and Qi Lei. Few-shot learning via learning the representation, provably. *arXiv preprint arXiv:2002.09434*, 2020.
- Theodoros Evgeniou, Charles A Micchelli, and Massimiliano Pontil. Learning multiple tasks with kernel methods. *Journal of machine learning research*, 6(Apr):615–637, 2005.
- Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. On the convergence theory of gradient-based model-agnostic meta-learning algorithms. In *International Conference on Artificial Intelligence and Statistics*, pp. 1082–1092, 2020.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 1126–1135, 2017.
- Chelsea Finn, Aravind Rajeswaran, Sham Kakade, and Sergey Levine. Online meta-learning. In *Proceedings of the 36th International Conference on Machine Learning*, 2019.
- Luca Franceschi, Paolo Frasconi, Saverio Salzo, Riccardo Grazi, and Massimiliano Pontil. Bilevel programming for hyperparameter optimization and meta-learning. *arXiv preprint arXiv:1806.04910*, 2018.
- Tomer Galanti, Lior Wolf, and Tamir Hazan. A theoretical framework for deep transfer learning. *Information and Inference: A Journal of the IMA*, 5(2):159–209, 2016.
- Jiatao Gu, Yong Wang, Yun Chen, Kyunghyun Cho, and Victor OK Li. Meta-learning for low-resource neural machine translation. *arXiv preprint arXiv:1808.08437*, 2018.
- Mikhail Khodak, Maria-Florina Balcan, and Ameet Talwalkar. Adaptive gradient-based meta-learning methods. *arXiv preprint arXiv:1906.02717*, 2019.
- A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional neural networks. pp. 1097–1105, 2012.
- Kwonjoon Lee, Subhansu Maji, Avinash Ravichandran, and Stefano Soatto. Meta-learning with differentiable convex optimization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 10657–10665, 2019.

287 Percy Liang. Cs229t/stat231: Statistical learning theory (winter 2016), 2016.

288 Han Liu, Mark Palatucci, and Jian Zhang. Blockwise coordinate descent procedures for the multi-
 289 task lasso, with applications to neural semantic basis discovery. In *Proceedings of the 26th Annual*
 290 *International Conference on Machine Learning*, pp. 649–656, 2009.

291 Han Liu, Lie Wang, and Tuo Zhao. Calibrated multivariate regression with application to neural
 292 semantic basis discovery. *Journal of machine learning research: JMLR*, 16:1579, 2015.

293 Andreas Maurer, Massimiliano Pontil, and Bernardino Romera-Paredes. The benefit of multitask
 294 representation learning. *The Journal of Machine Learning Research*, 17(1):2853–2884, 2016.

295 Daniel McNamara and Maria-Florina Balcan. Risk bounds for transferring representations with and
 296 without fine-tuning. In *Proceedings of the 34th International Conference on Machine Learning-*
 297 *Volume 70*, pp. 2373–2381. JMLR. org, 2017.

298 A. Nichol and J. Schulman. Reptile: a scalable metalearning algorithm. *arXiv preprint*
 299 *arXiv:1803.02999*, 2, 2018.

300 Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. *arXiv*
 301 *preprint arXiv:1803.02999*, 2018.

302 Aravind Rajeswaran, Chelsea Finn, Sham M Kakade, and Sergey Levine. Meta-learning with im-
 303 plicit gradients. In *Advances in Neural Information Processing Systems*, pp. 113–124, 2019.

304 S. Ravi and H. Larochelle. Optimization as a model for few-shot learning. 2017.

305 M. Ren, E. Triantafillou, S. Ravi, J. Snell, K. Swersky, J. Tenenbaum, H. Larochelle, and R. Zemel.
 306 Meta-learning for semi-supervised few-shot classification. *arXiv preprint arXiv:1803.00676*,
 307 2018.

308 Sebastian Ruder. An overview of multi-task learning in deep neural networks. *arXiv preprint*
 309 *arXiv:1706.05098*, 2017.

310 O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla,
 311 and M. Bernstein. Imagenet large scale visual recognition challenge. 115(3):211–252, 2015.

312 Nikunj Saunshi, Yi Zhang, Mikhail Khodak, and Sanjeev Arora. A sample complexity separation
 313 between non-convex and convex meta-learning. *arXiv preprint arXiv:2002.11172*, 2020.

314 Jürgen Schmidhuber. *Evolutionary principles in self-referential learning, or on learning how to*
 315 *learn: the meta-meta-... hook*. PhD thesis, Technische Universität München, 1987.

316 J. Snell, K. Swersky, and R. Zemel. Prototypical networks for few-shot learning. pp. 4077–4087,
 317 2017.

318 Sebastian Thrun and Lorien Pratt. *Learning to Learn: Introduction and Overview*, pp. 3–17.
 319 Springer US, Boston, MA, 1998. ISBN 978-1-4615-5529-2. doi: 10.1007/978-1-4615-5529-2_1.
 320 URL https://doi.org/10.1007/978-1-4615-5529-2_1.

321 Nilesh Tripuraneni, Chi Jin, and Michael I Jordan. Provable meta-learning of linear representations.
 322 *arXiv preprint arXiv:2002.11684*, 2020.

323 Aad W Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.

324 Lingxiao Wang, Qi Cai, Zhuoran Yang, and Zhaoran Wang. On the global optimality of model-
 325 agnostic meta-learning. *arXiv preprint arXiv:2006.13182*, 2020a.

326 Xiang Wang, Shuai Yuan, Chenwei Wu, and Rong Ge. Guarantees for tuning the step size using a
 327 learning-to-learn approach. *arXiv preprint arXiv:2006.16495*, 2020b.

328 Yujia Xie, Haoming Jiang, Feng Liu, Tuo Zhao, and Hongyuan Zha. Meta learning with relational
 329 information for short sequences. In *Advances in Neural Information Processing Systems*, pp.
 330 9904–9915, 2019.

331 Pan Zhou, Xiaotong Yuan, Huan Xu, Shuicheng Yan, and Jiashi Feng. Efficient meta learning via
 332 minibatch proximal update. In *Advances in Neural Information Processing Systems*, pp. 1534–
 333 1544, 2019.

A Experiments

Here we investigate our theory via simulations and benchmark few-shot classification tasks.

A.1 Simulations

We experiment on the realizable linear model studied in Section 4. Recall that the observed data of the t -th task are generated as

$$\mathbf{y}_t = \mathbf{X}_t \mathbf{w}_t, \quad \text{with } \mathbf{x}_{t,i} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \mathbf{I}_d).$$

We independently generate $\mathbf{w}_t \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{w}_{0,*}, \mathbf{I}_d/\sqrt{d})$, where $\mathbf{w}_{0,*}$ is the linear centroid and the corresponding $R^2 = 1$ here. The goal is to learn the linear centroid $\mathbf{w}_{0,*}$ using the train-train method and train-val method, i.e., minimizing $\hat{L}_T^{\text{tr-tr}}$ and $\hat{L}_T^{\text{tr-val}}$, respectively. Note that both $\hat{L}_T^{\text{tr-tr}}$ and $\hat{L}_T^{\text{tr-val}}$ are quadratic in \mathbf{w}_0 , therefore, we can find the close-form solutions $\hat{\mathbf{w}}_{0,T}^{\{\text{tr-tr}, \text{tr-val}\}}$. We measure the performance of train-train and train-val methods using the ℓ_2 -error $\|\mathbf{w}_{0,*} - \hat{\mathbf{w}}_{0,T}^{\{\text{tr-tr}, \text{tr-val}\}}\|^2$.

We present the comparison among train-train and train-val methods in Figure 1 with scatter plots representing the simulation outputs under different settings. Across all the simulations, we well-tune the regularization coefficient λ in the train-train method, and use a sufficiently large $\lambda = 10000$ in the train-val method according to Corollary 9. The simulated results concentrate around the reference curves corresponding to our theoretical findings. This corroborates our analyses and demonstrates the better performance of train-train method on the realizable linear model.

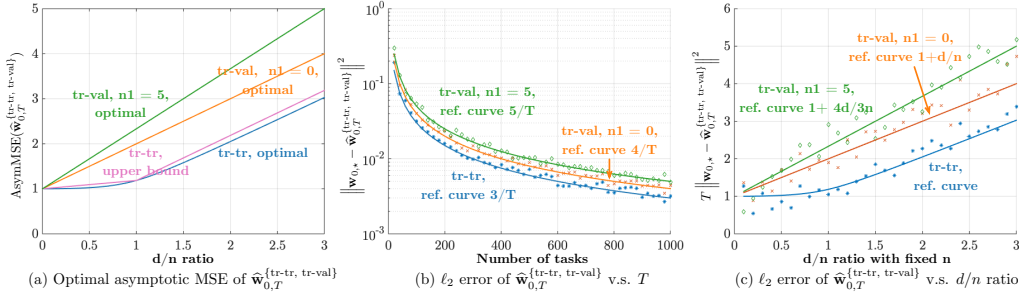


Figure 1: Panel (a) presents the optimal $\text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-tr}})$ (blue) in Theorem 10 via grid search, and the optimal $\text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}})$ in Corollary 9 with $n_1 = 0$ (orange) and $n_1 = 5$ (green), as well as the upper bound of $\text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\{\text{tr-tr}, \text{tr-val}\}})$ (magenta) in Corollary 6. The optimal $\text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\{\text{tr-tr}, \text{tr-val}\}})$ are used as reference curves in plots (b) and (c). Panel (b) plots the ℓ_2 -error of $\hat{\mathbf{w}}_{0,T}^{\{\text{tr-tr}, \text{tr-val}\}}$ as the total number of tasks increases from 20 to 1000 with an increment of 20. We fix data dimension $d = 60$ and per-task sample size $n = 20$. For the train-val method, we experiment on $n_1 = 0$ and $n_1 = 5$. Panel (c) shows the scaled (by T) ℓ_2 -error of $\hat{\mathbf{w}}_{0,T}^{\{\text{tr-tr}, \text{tr-val}\}}$ as the ratio d/n varies from 0 to 3 ($n = 20$ and $T = 1000$ are fixed).

A.2 Deep meta-learning

Baselines. We further testify our theoretical results in real meta-learning experiments. We compare train-val method (iMAML Rajeswaran et al. (2019)) and train-train method (Meta-MinibatchProx Zhou et al. (2019)) on benchmark few-shot classification tasks. iMAML splits the per-task data into training and validation sets. Using the training set, an approximate task-specific optimal \mathbf{w}_t is obtained by optimizing the standard classification loss with a regularization towards the centroid \mathbf{w}_0 . Next, iMAML evaluates \mathbf{w}_t on the validation set and optimizes \mathbf{w}_0 . In contrast, Meta-MinibatchProx does not split the per-task data, and pools all samples to optimize the same classification loss as in iMAML. In meta-test, we receive n_1 training and n_2 test samples from a test task. Both iMAML and Meta-MinibatchProx adapt the learned \mathbf{w}_0 to a test-task specific parameter \mathbf{w}_{ts} by running a few gradient updates to optimize the standard classification loss with n_1 training samples. Then they evaluate \mathbf{w}_{ts} on n_2 test samples.

Table 1: Few-shot classification accuracy (%) on the miniImageNet and tieredImageNet datasets.

miniImageNet	method	1-shot 5-way	5-shot 5-way	1-shot 20-way	5-shot 20-way
	train-val method	48.76 \pm 0.87	63.56 \pm 0.95	17.52 \pm 0.49	21.32 \pm 0.54
	train-train method	50.77 \pm 0.90	67.43 \pm 0.89	21.17 \pm 0.38	34.30 \pm 0.41
tieredImageNet	method	1-shot 5-way	5-shot 5-way	1-shot 10-way	5-shot 10-way
	train-val method	50.61 \pm 1.12	67.30 \pm 0.98	29.18 \pm 0.57	43.15 \pm 0.72
	train-train method	54.37 \pm 0.93	71.45 \pm 0.94	35.56 \pm 0.60	54.50 \pm 0.71

Datasets. We experiment on miniImageNet (Ravi & Larochelle, 2017) and tieredImageNet (Ren et al., 2018). MiniImageNet consists of 100 classes of images from ImageNet (Krizhevsky et al., 2012) and each class has 600 images of resolution $84 \times 84 \times 3$. We use 64 classes for training, 16 classes for validation, and the remaining 20 classes for testing (Ravi & Larochelle, 2017). TieredImageNet consists of 608 classes from the ILSVRC-12 data set (Russakovsky et al., 2015) and each image is also of resolution $84 \times 84 \times 3$. TieredImageNet groups classes into broader hierarchy categories corresponding to higher-level nodes in the ImageNet. Specifically, its top hierarchy has 20 training categories (351 classes), 6 validation categories (97 classes) and 8 test categories (160 classes). This structure ensures that all training classes are distinct from the testing classes, providing a more realistic few-shot learning scenario.

Experimental settings. We adopt the episodic training procedure (Finn et al., 2017; Zhou et al., 2019; Rajeswaran et al., 2019). In meta-test, we sample a set of N -way ($K + 1$)-shot test tasks. The first K instances are for training and the remaining one is for testing. In meta-training, we use the “higher way” training strategy. We set the per-task sample size to be 30 following Zhou et al. (2019); Rajeswaran et al. (2019). For iMAML, training and validation samples are evenly split, i.e., both being 15. We evaluate baselines under the transduction setting where the information is shared between the test data via batch normalization. We use the standard 4-layer convolution network in (Finn et al., 2017; Zhou et al., 2019) as the backbone. We report the average accuracy over 2,000 random test episodes with 95% confidence interval.

Results. Table 1 presents the percent classification accuracy on miniImageNet and tieredImageNet. Train-train method (Meta-MinibatchProx), consistently outperforms the train-val method (iMAML): On miniImageNet, train-train method respectively makes about 2.01%, 3.87% improvements on the 1-shot 5-way and 5-shot 5-way tasks; On tieredImageNet, train-train method averagely improves by about 6.40% on the four testing cases. These results show the advantages of train-train method over train-val and well support our theoretical findings in Corollary 6.

We further tune the split (n_1, n_2) in iMAML and report the results in Table 2. As can be seen, as the number of test samples n_2 increases, the percent classification accuracy on both the miniImageNet and tieredImageNet datasets becomes higher. This testifies our theoretical affirmation in Corollary 9.

Table 2: Investigation of the effects of training/validation splitting ratio in the train-val method (iMAML) to the few-shot classification accuracy (%) on miniImageNet and tieredImageNet.

datasets	$n_1 = 25, n_2 = 5$	$n_1 = 15, n_2 = 15$	$n_1 = 5, n_2 = 25$
miniImageNet	62.09 \pm 0.97	63.56 \pm 0.95	63.92 \pm 1.04
tieredImageNet	66.45 \pm 1.05	67.30 \pm 0.98	67.50 \pm 0.94

B Proofs for Section 3

B.1 Generic result on asymptotic normality

We first present the full version of Proposition 1, taken from (Van der Vaart, 2000, Theorem 5.21).

Proposition 7 (Asymptotic normality and excess risk of ERM; formal version of Proposition 1). *Assume the population minimizer $\mathbf{w}_{0,*}$ is unique and the ERM $\hat{\mathbf{w}}_{0,T}$ is consistent (i.e. it converges to $\mathbf{w}_{0,*}$ in probability as $T \rightarrow \infty$). Further assume the following regularity conditions:*

396 (a) There exists some random variable $A_t = A(p_t, \mathbf{X}_t, \mathbf{y}_t)$ such that $\mathbb{E}[A_t^2] < \infty$ and

$$\|\nabla \ell_t(\mathbf{w}_1) - \nabla \ell_t(\mathbf{w}_2)\| \leq A_t \|\mathbf{w}_1 - \mathbf{w}_2\|$$

397 for all $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^d$;

398 (b) $\mathbb{E}[\|\nabla \ell_t(\mathbf{w}_{0,*})\|^2] < \infty$;

399 (c) L is twice-differentiable with $\nabla^2 L(\mathbf{w}_{0,*}) \succ \mathbf{0}$,

400 then the ERM $\widehat{\mathbf{w}}_{0,T}$ is asymptotically normally distributed, with

$$\begin{aligned} \sqrt{T} \cdot (\widehat{\mathbf{w}}_{0,T} - \mathbf{w}_{0,*}) &\xrightarrow{d} \mathbf{N}(\mathbf{0}, \nabla^2 L(\mathbf{w}_{0,*})^{-1} \text{Cov}(\nabla \ell_t(\mathbf{w}_{0,*})) \nabla^2 L(\mathbf{w}_{0,*})^{-1}) =: P_{\mathbf{w}}, \\ T \cdot (L(\widehat{\mathbf{w}}_{0,T}) - L(\mathbf{w}_{0,*})) &\xrightarrow{d} \mathbf{\Delta}^\top \nabla^2 L(\mathbf{w}_{0,*}) \mathbf{\Delta} \quad \text{where } \mathbf{\Delta} \sim P_{\mathbf{w}}. \end{aligned}$$

401 where \xrightarrow{d} denotes convergence in distribution and $\ell_t : \mathbb{R}^d \rightarrow \mathbb{R}$ is the loss function on a single task.

402 B.2 Proof of Proposition 2

403 **Equivalence of test-time risk and training loss for train-val method** We first show that

$$L^{\text{tr-val}}(\mathbf{w}_0) = \mathbb{E}[\ell_t^{\text{tr-val}}(\mathbf{w}_0)] = L_{\lambda, n_1}^{\text{test}}(\mathbf{w}_0)$$

404 for all \mathbf{w}_0 , that is, the population meta-test loss is exactly the same as the population risk of the train-
405 val method. This is straightforward: as the tasks are i.i.d. and $\mathcal{A}_\lambda(\mathbf{w}_0; \mathbf{X}_t^{\text{train}}, \mathbf{y}_t^{\text{train}})$ is independent
406 of the test points $(\mathbf{X}_t^{\text{val}}, \mathbf{y}_t^{\text{val}})$, we have for any \mathbf{w}_0 that

$$\begin{aligned} \mathbb{E}[\ell_t^{\text{tr-val}}(\mathbf{w}_0)] &= \mathbb{E}_{p_t \sim \Pi, (\mathbf{X}_t, \mathbf{y}_t) \sim p_t} \left[\frac{1}{2n_2} \|\mathbf{y}_t^{\text{val}} - \mathbf{X}_t^{\text{val}} \mathcal{A}_\lambda(\mathbf{w}_0; \mathbf{X}_t^{\text{train}}, \mathbf{y}_t^{\text{train}})\|^2 \right] \\ &= \mathbb{E}_{p_t \sim \Pi, (\mathbf{X}_t, \mathbf{y}_t) \sim p_t} \left[\frac{1}{2} (y_{t,1}^{\text{val}} - \mathbf{x}_{t,1}^{\text{val}\top} \mathcal{A}_\lambda(\mathbf{w}_0; \mathbf{X}_t^{\text{train}}, \mathbf{y}_t^{\text{train}}))^2 \right] \\ &= \mathbb{E}_{p_{T+1} \sim \Pi, (\mathbf{X}_{T+1}, \mathbf{y}_{T+1}), (\mathbf{x}', \mathbf{y}') \stackrel{\text{iid}}{\sim} p_t} \left[\frac{1}{2} (y' - \mathbf{x}'^\top \mathcal{A}_{\lambda, n_1}(\mathbf{w}_0; \mathbf{X}_{T+1}, \mathbf{y}_{T+1}))^2 \right] \\ &= L_{\lambda, n_1}^{\text{test}}(\mathbf{w}_0). \end{aligned}$$

407 Therefore the train-val method is acutally a valid ERM for the test loss $L_{\lambda, n_1}^{\text{test}}$, and it remains to show
408 that the train-val method is (itself) consistent.

409 **Consistency** We expand the empirical risk of the train-val method as

$$\begin{aligned} \widehat{L}_T^{\text{tr-val}}(\mathbf{w}_0) &= \frac{1}{T} \sum_{t=1}^T \frac{1}{2n_2} \|\mathbf{y}_t^{\text{val}} - \mathbf{X}_t^{\text{val}} \mathcal{A}_\lambda(\mathbf{w}_0; \mathbf{X}_t^{\text{train}}, \mathbf{y}_t^{\text{train}})\|^2 \\ &= \frac{1}{T} \sum_{t=1}^T \frac{1}{2n_2} \|\mathbf{y}_t^{\text{val}} - \mathbf{X}_t^{\text{val}} [\mathbf{w}_0 + (\mathbf{X}_t^{\text{train}\top} \mathbf{X}_t^{\text{train}} + n_1 \lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^{\text{train}\top} (\mathbf{y}_t^{\text{train}} - \mathbf{X}_t^{\text{train}} \mathbf{w}_0)]\|^2 \\ &= \frac{1}{T} \sum_{t=1}^T \frac{1}{2n_2} \|\mathbf{y}_t^{\text{val}} - \mathbf{X}_t^{\text{val}} (\mathbf{X}_t^{\text{train}\top} \mathbf{X}_t^{\text{train}} + n_1 \lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^{\text{train}\top} \mathbf{y}_t^{\text{train}} - \mathbf{X}_t^{\text{val}} n_1 \lambda (\mathbf{X}_t^{\text{train}\top} \mathbf{X}_t^{\text{train}} + n_1 \lambda \mathbf{I}_d)^{-1} \mathbf{w}_0\|^2 \\ &= \frac{1}{2} \mathbf{w}_0^\top \mathbf{M}_T \mathbf{w}_0 - \mathbf{w}_0^\top \mathbf{b}_T + \text{const}, \end{aligned}$$

410 where

$$\begin{aligned} \mathbf{M}_T &:= \frac{1}{T} \sum_{t=1}^T \lambda^2 (\mathbf{X}_t^{\text{train}\top} \mathbf{X}_t^{\text{train}} / n_1 + \lambda \mathbf{I}_d)^{-1} \frac{\mathbf{X}_t^{\text{val}\top} \mathbf{X}_t^{\text{val}}}{n_2} (\mathbf{X}_t^{\text{train}\top} \mathbf{X}_t^{\text{train}} / n_1 + \lambda \mathbf{I}_d)^{-1}, \\ \mathbf{b}_T &:= \frac{1}{T} \sum_{t=1}^T \lambda (\mathbf{X}_t^{\text{train}\top} \mathbf{X}_t^{\text{train}} / n_1 + \lambda \mathbf{I}_d)^{-1} \cdot \frac{1}{n_2} \mathbf{X}_t^{\text{val}\top} (\mathbf{y}_t^{\text{val}} - \mathbf{X}_t^{\text{val}} (\mathbf{X}_t^{\text{train}\top} \mathbf{X}_t^{\text{train}} + n_1 \lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^{\text{train}\top} \mathbf{y}_t^{\text{train}}). \end{aligned}$$

411 Noticing that $(\mathbf{X}_t^{\text{train}\top} \mathbf{X}_t^{\text{train}} / n_1 + \lambda \mathbf{I}_d)^{-1} \preceq \lambda^{-1} \mathbf{I}_d$ and by the assumption that $\mathbb{E}_{(\mathbf{x}, y) \sim p_t} [\mathbf{x} \mathbf{x}^\top] \prec$
 412 ∞ , $\mathbb{E}_{(\mathbf{x}, y) \sim p_t} [\mathbf{x} \mathbf{y}] < \infty$, we have $\mathbb{E}[\|\mathbf{M}_T\|] < \infty$ and $\mathbb{E}[\|\mathbf{b}_T\|] < \infty$. Since the task p_t 's are i.i.d.,
 413 by the law of large numbers, we have with probability one that

$$\begin{aligned} \mathbf{M}_T &\rightarrow \mathbb{E}[\mathbf{M}_T] \\ &= \mathbb{E}_{p_t, (\mathbf{X}_t, \mathbf{y}_t)} \left[\lambda^2 (\mathbf{X}_t^{\text{train}\top} \mathbf{X}_t^{\text{train}} / n_1 + \lambda \mathbf{I}_d)^{-1} \frac{\mathbf{X}_t^{\text{val}\top} \mathbf{X}_t^{\text{val}}}{n_2} (\mathbf{X}_t^{\text{train}\top} \mathbf{X}_t^{\text{train}} / n_1 + \lambda \mathbf{I}_d)^{-1} \right] \quad (7) \\ &= \mathbb{E}_{p_t, (\mathbf{X}_t, \mathbf{y}_t)} \left[\lambda^2 (\mathbf{X}_t^{\text{train}\top} \mathbf{X}_t^{\text{train}} / n_1 + \lambda \mathbf{I}_d)^{-1} \boldsymbol{\Sigma}_t (\mathbf{X}_t^{\text{train}\top} \mathbf{X}_t^{\text{train}} / n_1 + \lambda \mathbf{I}_d)^{-1} \right] \succ \mathbf{0}, \end{aligned}$$

414 (where $\boldsymbol{\Sigma}_t = \mathbb{E}_{\mathbf{x} \sim p_t} [\mathbf{x} \mathbf{x}^\top] \succ \mathbf{0}$) and

$$\begin{aligned} \mathbf{b}_T &\rightarrow \mathbb{E}[\mathbf{b}_T] \\ &= \mathbb{E}_{p_t, (\mathbf{X}_t, \mathbf{y}_t)} \left[\lambda (\mathbf{X}_t^{\text{train}\top} \mathbf{X}_t^{\text{train}} / n_1 + \lambda \mathbf{I}_d)^{-1} \cdot \frac{1}{n_2} \mathbf{X}_t^{\text{val}\top} (\mathbf{y}_t^{\text{val}} - \mathbf{X}_t^{\text{val}} (\mathbf{X}_t^{\text{train}\top} \mathbf{X}_t^{\text{train}} / n_1 + \lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^{\text{train}\top} \mathbf{y}_t^{\text{train}}) \right] \\ &< \infty \end{aligned} \quad (8)$$

415 as $T \rightarrow \infty$. Therefore, by Slutsky's Theorem, we have

$$\hat{\mathbf{w}}_{0,T} = \mathbf{M}_T^{-1} \mathbf{b}_T \rightarrow \mathbb{E}[\mathbf{M}_T]^{-1} \mathbb{E}[\mathbf{b}_T] = \arg \min_{\mathbf{w}_0} L^{\text{tr-val}}(\mathbf{w}_0) = \arg \min_{\mathbf{w}_0} L_{\lambda, n_1}^{\text{test}}(\mathbf{w}_0) = \mathbf{w}_{0,*}(\lambda, n_1)$$

416 as $T \rightarrow \infty$. This proves the consistency of the train-val method.

417 **Asymptotic normality** Similar as above, we can write the per-task loss as

$$\ell_t(\mathbf{w}_0) = \frac{1}{2} \|\mathbf{A}_t \mathbf{w}_0 - \mathbf{c}_t\|^2,$$

418 where

$$\begin{aligned} \mathbf{A}_t &= \frac{\lambda}{\sqrt{n_2}} \mathbf{X}_t^{\text{val}} (\mathbf{X}_t^{\text{train}\top} \mathbf{X}_t^{\text{train}} / n_1 + \lambda \mathbf{I}_d)^{-1}, \\ \mathbf{c}_t &= \frac{1}{\sqrt{n_2}} \left(\mathbf{y}_t^{\text{val}} - \mathbf{X}_t^{\text{val}} (\mathbf{X}_t^{\text{train}\top} \mathbf{X}_t^{\text{train}} / n_1 + \lambda \mathbf{I}_d)^{-1} \frac{1}{n_1} \mathbf{X}_t^{\text{train}\top} \mathbf{y}_t^{\text{train}} \right). \end{aligned}$$

419 In order to show the desired asymptotic normality result, it suffices to check the conditions in Propo-
 420 sition 7. First, we have

$$\nabla \ell_t(\mathbf{w}_0) = \mathbf{A}_t^\top (\mathbf{A}_t \mathbf{w}_0 - \mathbf{c}_t).$$

421 This is Lipschitz in \mathbf{w}_0 with Lipschitz constant

$$\|\mathbf{A}_t^\top \mathbf{A}_t\|_{\text{op}} \leq \|\mathbf{A}_t\|_{\text{Fr}}^2 < \frac{1}{n_2} \|\mathbf{X}_t^{\text{val}}\|_{\text{Fr}}^2.$$

422 As $\mathbb{E}_{\mathbf{x} \sim p_t} [\|\mathbf{x}\|^4] < \infty$, the above quantity is clearly square integrable, therefore verifying (a). As
 423 $\mathbf{w}_{0,*} = \mathbf{w}_{0,*}(\lambda, n_1)$ is finite, we can use similar arguments as above to show (b) holds. Finally, we
 424 have already seen L is twice-differentiable (since it is quadratic in \mathbf{w}_0) and $\nabla^2 L(\mathbf{w}_{0,*}) \succ \mathbf{0}$, which
 425 verifies (c). Therefore the conditions of Proposition 7 hold, which yields the desired asymptotic
 426 normality result. \square

427 B.3 Proof of Proposition 3

428 **High-level idea** At a high level, this proof proceeds by showing that the train-train method is also
 429 consistent to the (population) minimizer of $L^{\text{tr-tr}}$, and constructing a simple counter-example on
 430 which the minimizers of $L^{\text{tr-tr}}$ is not equal to that of $L_{\lambda, n}^{\text{test}}$.

431 **Population minimizers of $L^{\text{tr-tr}}$ and $L^{\text{tr-val}}$** We begin by simplifying the non-splitting risk. We
 432 have

$$\begin{aligned} \ell_t^{\text{tr-tr}}(\mathbf{w}_0) &= \frac{1}{2n} \|\mathbf{y}_t - \mathbf{X}_t \mathcal{A}_\lambda(\mathbf{w}_0; \mathbf{X}_t, \mathbf{y}_t)\|^2 \\ &= \frac{1}{2n} \|\mathbf{y}_t - \mathbf{X}_t [\mathbf{w}_0 + (\mathbf{X}_t^\top \mathbf{X}_t + n \lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^\top (\mathbf{y}_t - \mathbf{X}_t \mathbf{w}_0)]\|^2 \\ &= \frac{1}{2} \|\mathbf{A}_t \mathbf{w}_0 - \mathbf{c}_t\|^2, \end{aligned}$$

433 where

$$\mathbf{A}_t = \frac{1}{\sqrt{n}} n \lambda \mathbf{X}_t (\mathbf{X}_t^\top \mathbf{X}_t + n \lambda \mathbf{I}_d)^{-1} \quad \text{and} \quad \mathbf{c}_t = \frac{1}{\sqrt{n}} (\mathbf{I}_n - \mathbf{X}_t (\mathbf{X}_t^\top \mathbf{X}_t + n \lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^\top) \mathbf{y}_t.$$

434 Using similar arguments as in the proof of Proposition 2 (Appendix B.2), we see that the train-
435 train method $\hat{\mathbf{w}}_{0,T}^{\text{tr-tr}}$ converges with probability one to the minimizer of the population risk $L^{\text{tr-tr}}$,
436 which is

$$\begin{aligned} \mathbf{w}_{0,\star}^{\text{tr-tr}} &= \arg \min_{\mathbf{w}_0} L^{\text{tr-tr}}(\mathbf{w}_0) = (\mathbb{E}[\mathbf{A}_t^\top \mathbf{A}_t])^{-1} \mathbb{E}[\mathbf{A}_t^\top \mathbf{c}_t] \\ &= \mathbb{E} \left[\lambda^2 (\mathbf{X}_t^\top \mathbf{X}_t / n + \lambda \mathbf{I}_d)^{-2} \frac{\mathbf{X}_t^\top \mathbf{X}_t}{n} \right]^{-1} \cdot \mathbb{E} \left[\frac{1}{n} \lambda (\mathbf{X}_t^\top \mathbf{X}_t / n + \lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^\top (\mathbf{I}_n - \mathbf{X}_t (\mathbf{X}_t^\top \mathbf{X}_t + n \lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^\top) \mathbf{y}_t \right] \\ &= \mathbb{E} \left[\lambda^2 (\mathbf{X}_t^\top \mathbf{X}_t / n + \lambda \mathbf{I}_d)^{-2} \frac{\mathbf{X}_t^\top \mathbf{X}_t}{n} \right]^{-1} \cdot \mathbb{E} \left[\lambda^2 (\mathbf{X}_t^\top \mathbf{X}_t / n + \lambda \mathbf{I}_d)^{-2} \frac{1}{n} \mathbf{X}_t^\top \mathbf{y}_t \right]. \end{aligned} \tag{9}$$

437 On the other hand, recall from Proposition 2 ((7) and (8)) that the population minimizer of $L_{\lambda,n}^{\text{test}}$ is

$$\begin{aligned} \mathbf{w}_{0,\star}(\lambda, n) &= \arg \min_{\mathbf{w}_0} L_{\lambda,n}^{\text{test}}(\mathbf{w}_0) \\ &= \mathbb{E} \left[\lambda^2 (\mathbf{X}_t^\top \mathbf{X}_t / n + \lambda \mathbf{I}_d)^{-1} \Sigma_t (\mathbf{X}_t^\top \mathbf{X}_t / n + \lambda \mathbf{I}_d)^{-1} \right]^{-1} \cdot \left\{ \lambda \mathbb{E} [(\mathbf{X}_t^\top \mathbf{X}_t / n + \lambda \mathbf{I}_d)^{-1}] \mathbb{E}_{p_t, (\mathbf{x}', \mathbf{y}') \sim p_t} [\mathbf{x}' \mathbf{y}'] \right. \\ &\quad \left. - \lambda \mathbb{E} \left[(\mathbf{X}_t^\top \mathbf{X}_t / n + \lambda \mathbf{I}_d)^{-1} \Sigma_t (\mathbf{X}_t^\top \mathbf{X}_t / n + \lambda \mathbf{I}_d)^{-1} \frac{1}{n} \mathbf{X}_t^\top \mathbf{y}_t \right] \right\}. \end{aligned} \tag{10}$$

438 **Construction of the counter-example** We now construct a distribution for which (9) is not equal
439 to (10). Let $d = 1$ and let all p_t be the following distribution:

$$p_t : (x_{t,i}, y_{t,i}) = \begin{cases} (1, 3) & \text{with probability } 1/2; \\ (3, -1) & \text{with probability } 1/2. \end{cases}$$

440 Clearly, we have $\Sigma_t = 5$, $s_t := \mathbf{X}_t^\top \mathbf{X}_t / n \in [1, 9]$, and $\mathbb{E}_{x', y' \sim p_t} [x' y'] = 0$. Therefore we have

$$\mathbf{w}_{0,\star}^{\text{tr-tr}} = \mathbb{E}[(s_t + \lambda)^{-2} s_t]^{-1} \cdot \mathbb{E} \left[(s_t + \lambda)^{-2} \frac{1}{n} \sum_{i=1}^n x_{t,i} y_{t,i} \right],$$

441 and

$$\begin{aligned} \mathbf{w}_{0,\star}(\lambda, n) &= -\mathbb{E}[\lambda^2 (s_t + \lambda)^{-2}]^{-1} \cdot \mathbb{E} \left[5 \lambda (s_t + \lambda)^{-2} \frac{1}{n} \sum_{i=1}^n x_{t,i} y_{t,i} \right] \\ &= -\mathbb{E}[\lambda (s_t + \lambda)^{-2}]^{-1} \cdot \mathbb{E} \left[(s_t + \lambda)^{-2} \frac{1}{n} \sum_{i=1}^n x_{t,i} y_{t,i} \right]. \end{aligned}$$

442 We now show that $\mathbf{w}_{0,\star}^{\text{tr-tr}} \neq \mathbf{w}_{0,\star}(\lambda, n)$ by showing that

$$\mathbb{E} \left[(s_t + \lambda)^{-2} \frac{1}{n} \sum_{i=1}^n x_{t,i} y_{t,i} \right] = \mathbb{E} \left[\frac{x_{t,1} y_{t,1}}{(s_t + \lambda)^2} \right] \neq 0$$

443 for any $\lambda > 0$. Indeed, conditioning on $(x_{t,1}, y_{t,1}) = (1, 3)$, we know that the sum-of-squares in s_t
444 has one term that equals 1, and all others i.i.d. being 1 or 9 with probability one half. On the other
445 hand, if we condition on $(x_{t,1}, y_{t,1}) = (3, -1)$, then we know the sum in s_t has one term that equals
446 9 and all others i.i.d.. This means that the negative contribution in the expectation is smaller than
447 the positive contribution, in other words

$$\begin{aligned} \mathbb{E} \left[\frac{x_{t,1} y_{t,1}}{(s_t + \lambda)^2} \right] &= \frac{1}{2} \cdot 3 \mathbb{E} \left[\frac{1}{(s_t + \lambda)^2} \middle| (x_{t,1}, y_{t,1}) = (1, 3) \right] \\ &\quad + \frac{1}{2} \cdot -3 \mathbb{E} \left[\frac{1}{(s_t + \lambda)^2} \middle| (x_{t,1}, y_{t,1}) = (3, -1) \right] > 0. \end{aligned}$$

448 This shows $\mathbf{w}_{0,\star}^{\text{tr-tr}} \neq \mathbf{w}_{0,\star}(\lambda, n)$ and consequently the $\widehat{\mathbf{w}}_{0,T}^{\text{tr-tr}}$ does not converge to $\mathbf{w}_{0,\star}(\lambda, n)$ and the
 449 difference is bounded away from zero as $T \rightarrow \infty$.

450 Finally, for this distribution, the risk $L_{\lambda,n}^{\text{test}}(\mathbf{w}_0)$ is strongly convex (since it has a positive second
 451 derivative), this further implies that $L_{\lambda,n}^{\text{test}}(\widehat{\mathbf{w}}_{0,T}^{\text{tr-tr}}) - L_{\lambda,n}^{\text{test}}(\mathbf{w}_{0,\star}(\lambda, n))$ is bounded away from zero
 452 almost surely as $T \rightarrow \infty$.

453 C Proofs for Section 4

454 C.1 Proof of Theorem 4

455 We first show that $\mathbf{w}_{0,\star} = \mathbb{E}_{\mathbf{w}_t \sim \Pi}[\mathbf{w}_t]$ is a global optimizer for $L^{\text{tr-tr}}$ and $L^{\text{tr-val}}$ with any regulariza-
 456 tion coefficient $\lambda > 0$, any n , and any split (n_1, n_2) . To do this, it suffices to check that the gradient
 457 at $\mathbf{w}_{0,\star}$ is zero and the Hessian is positive definite (PD).

458 **Optimality of $\mathbf{w}_{0,\star}$ in both $L^{\text{tr-tr}}$ and $L^{\text{tr-val}}$.** We first look at $L^{\text{tr-tr}}$: for any $\mathbf{w}_0 \in \mathbb{R}^d$ we have

$$\begin{aligned} L^{\text{tr-tr}}(\mathbf{w}_0) &= \mathbb{E}[\ell_t^{\text{tr-tr}}(\mathbf{w}_0)] \\ &= \frac{1}{2n} \mathbb{E} \left[\left\| \mathbf{X}_t \mathbf{w}_t - \mathbf{X}_t \left[(\mathbf{X}_t^\top \mathbf{X}_t + n\lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^\top (\mathbf{X}_t \mathbf{w}_t - \mathbf{X}_t \mathbf{w}_0) + \mathbf{w}_0 \right] \right\|^2 \right] \\ &= \frac{1}{2n} \mathbb{E} \left[\left\| \mathbf{X}_t \left(\mathbf{I}_d - (\mathbf{X}_t^\top \mathbf{X}_t + n\lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^\top \mathbf{X}_t \right) (\mathbf{w}_t - \mathbf{w}_0) \right\|^2 \right]. \end{aligned} \quad (11)$$

459 Similarly, $L^{\text{tr-val}}$ can be written as

$$\begin{aligned} L^{\text{tr-val}}(\mathbf{w}_0) &= \mathbb{E}[\ell_t^{\text{tr-val}}(\mathbf{w}_0)] \\ &= \frac{1}{2n_2} \mathbb{E} \left[\left\| \mathbf{X}_t^{\text{val}} \mathbf{w}_t - \mathbf{X}_t^{\text{val}} \left[((\mathbf{X}_t^{\text{train}})^\top \mathbf{X}_t^{\text{train}} + n_1 \lambda \mathbf{I}_d)^{-1} (\mathbf{X}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{train}} \mathbf{w}_t - \mathbf{X}_t^{\text{train}} \mathbf{w}_0) + \mathbf{w}_0 \right] \right\|^2 \right] \\ &= \frac{1}{2n_2} \mathbb{E} \left[\left\| \mathbf{X}_t^{\text{val}} \left(\mathbf{I}_d - ((\mathbf{X}_t^{\text{train}})^\top \mathbf{X}_t^{\text{train}} + n_1 \lambda \mathbf{I}_d)^{-1} (\mathbf{X}_t^{\text{train}})^\top \mathbf{X}_t^{\text{train}} \right) (\mathbf{w}_t - \mathbf{w}_0) \right\|^2 \right]. \end{aligned} \quad (13)$$

460 We denote

$$\mathbf{M}_t^{\text{tr-tr}} = \mathbf{X}_t \left(\mathbf{I}_d - (\mathbf{X}_t^\top \mathbf{X}_t + n\lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^\top \mathbf{X}_t \right) \quad \text{and} \quad (14)$$

$$\mathbf{M}_t^{\text{tr-val}} = \mathbf{X}_t^{\text{val}} \left(\mathbf{I}_d - ((\mathbf{X}_t^{\text{train}})^\top \mathbf{X}_t^{\text{train}} + n_1 \lambda \mathbf{I}_d)^{-1} (\mathbf{X}_t^{\text{train}})^\top \mathbf{X}_t^{\text{train}} \right) \quad (15)$$

461 to simplify the notations in (11) and (13). We take gradient of $L^{\text{tr-tr}}$ and $L^{\text{tr-val}}$ with respect to \mathbf{w}_0 :

$$\nabla_{\mathbf{w}_0} L^{\text{tr-tr}}(\mathbf{w}_0) = -\frac{1}{n} \mathbb{E} [(\mathbf{M}_t^{\text{tr-tr}})^\top \mathbf{M}_t^{\text{tr-tr}} (\mathbf{w}_t - \mathbf{w}_0)], \quad (16)$$

$$\nabla_{\mathbf{w}_0} L^{\text{tr-val}}(\mathbf{w}_0) = -\frac{1}{n_2} \mathbb{E} [(\mathbf{M}_t^{\text{tr-val}})^\top \mathbf{M}_t^{\text{tr-val}} (\mathbf{w}_t - \mathbf{w}_0)]. \quad (17)$$

462 Substituting $\mathbf{w}_{0,\star}$ into (16) and taking expectation, we deduce

$$\nabla_{\mathbf{w}_0} L^{\text{tr-tr}}(\mathbf{w}_{0,\star}) = -\frac{1}{n} \mathbb{E} [(\mathbf{M}_t^{\text{tr-tr}})^\top \mathbf{M}_t^{\text{tr-tr}} (\mathbf{w}_t - \mathbf{w}_{0,\star})] = \mathbf{0}. \quad (18)$$

463 To see this, observe that by definition $\mathbb{E}[\mathbf{w}_t - \mathbf{w}_{0,\star}] = \mathbf{0}$. Combining with \mathbf{w}_t being generated
 464 independently of \mathbf{X}_t , we have the first term in RHS of (18) vanish. In addition, \mathbf{z}_t is independent
 465 white noise, therefore, the second term in RHS of (18) also vanishes. Following the same argument,
 466 we can show

$$\nabla_{\mathbf{w}_0} L^{\text{tr-val}}(\mathbf{w}_{0,\star}) = \mathbf{0},$$

467 since \mathbf{X}_t' is also independent of \mathbf{w}_t . The above reasonings indicates that $\mathbf{w}_{0,\star}$ is a stationary point
 468 of both $L^{\text{tr-tr}}$ and $L^{\text{tr-val}}$. The remaining step is to check $\nabla_{\mathbf{w}_0} L^{\text{tr-tr}}(\mathbf{w}_{0,\star})$ and $\nabla_{\mathbf{w}_0} L^{\text{tr-val}}(\mathbf{w}_{0,\star})$ are
 469 PD. From (16) and (17), we derive respectively the hessian of $L^{\text{tr-tr}}$ and $L^{\text{tr-val}}$ as

$$\begin{aligned} \nabla_{\mathbf{w}_0}^2 L^{\text{tr-tr}}(\mathbf{w}_{0,\star}) &= \frac{1}{n} \mathbb{E} [(\mathbf{M}_t^{\text{tr-tr}})^\top \mathbf{M}_t^{\text{tr-tr}}] \quad \text{and} \\ \nabla_{\mathbf{w}_0}^2 L^{\text{tr-val}}(\mathbf{w}_{0,\star}) &= \frac{1}{n_2} \mathbb{E} [(\mathbf{M}_t^{\text{tr-val}})^\top \mathbf{M}_t^{\text{tr-val}}]. \end{aligned}$$

Let $\mathbf{v} \in \mathbb{R}^d$ be any nonzero vector, we check $\mathbf{v}^\top \nabla_{\mathbf{w}_0}^2 L^{\text{tr-tr}}(\mathbf{w}_{0,*}) \mathbf{v} > 0$. A key observation is that $(\mathbf{I}_d - (\mathbf{X}_t^\top \mathbf{X}_t + n\lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^\top \mathbf{X}_t)$ is positive definite for any $\lambda \neq 0$. To see this, let $\sigma_1 \geq \dots \geq \sigma_d$ be eigenvalues of $\frac{1}{n} \mathbf{X}_t^\top \mathbf{X}_t$, some algebra yields the eigenvalues of $(\mathbf{I}_d - (\mathbf{X}_t^\top \mathbf{X}_t + n\lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^\top \mathbf{X}_t)$ are $\frac{\lambda}{\lambda + \sigma_i} > 0$ for $\lambda \neq 0$ and $i = 1, \dots, d$. Hence, we deduce

$$\mathbf{v}^\top \nabla_{\mathbf{w}_0}^2 L^{\text{tr-tr}}(\mathbf{w}_{0,*}) \mathbf{v} = \frac{1}{n} \mathbb{E}[\mathbf{v}^\top \mathbf{X}_t^\top (\mathbf{I}_d - (\mathbf{X}_t^\top \mathbf{X}_t + n\lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^\top \mathbf{X}_t)^2 \mathbf{X}_t \mathbf{v}] > 0, \quad (19)$$

since \mathbf{X}_t is isotropic (an explicit computation of the hessian matrix can be found in Appendix C.2). As a consequence, we have shown that $\mathbf{w}_{0,*}$ is a global optimum of $L^{\text{tr-tr}}$. The same argument applies to $L^{\text{tr-val}}$, and the proof is complete.

Consistency of $\hat{\mathbf{w}}_{0,T}^{\{\text{tr-tr}, \text{tr-val}\}}$. To check the consistency, we need to verify the conditions (a) – (c) in Proposition 7.

For condition (a), we derive from (16) and (17) that

$$\left\| \nabla \ell_t^{\{\text{tr-tr}, \text{tr-val}\}}(\mathbf{w}_1) - \nabla \ell_t^{\{\text{tr-tr}, \text{tr-val}\}}(\mathbf{w}_2) \right\| \leq \frac{1}{n} \left\| (\mathbf{M}_t^{\{\text{tr-tr}, \text{tr-val}\}})^\top \mathbf{M}_t^{\{\text{tr-tr}, \text{tr-val}\}} \right\|_{\text{op}} \|\mathbf{w}_1 - \mathbf{w}_2\|,$$

where n should be replaced by n_2 for the split method (we slightly abuse the notation for simplicity).

It suffices to show $\mathbb{E} \left[\left\| (\mathbf{M}_t^{\{\text{tr-tr}, \text{tr-val}\}})^\top \mathbf{M}_t^{\{\text{tr-tr}, \text{tr-val}\}} \right\|_{\text{op}} \right] < \infty$, which follows from the same argument in the proof of Proposition 2. In particular, we know $\mathbf{0} \preceq \mathbf{I}_d - (\mathbf{X}_t^\top \mathbf{X}_t + n\lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^\top \mathbf{X}_t \preceq \mathbf{I}_d$

and $\mathbb{E}[\|\mathbf{X}_t^\top \mathbf{X}_t\|_{\text{op}}] < \infty$ since \mathbf{X}_t is Gaussian. As a consequence, for the no-split method, we have $\mathbb{E}[\|(\mathbf{M}_t^{\text{tr-tr}})^\top \mathbf{M}_t^{\text{tr-tr}}\|_{\text{op}}^2] < \infty$. For the split method, we also have $\mathbf{0} \preceq \mathbf{I}_d - ((\mathbf{X}_t^{\text{train}})^\top \mathbf{X}_t^{\text{train}} + n\lambda \mathbf{I}_d)^{-1} (\mathbf{X}_t^{\text{train}})^\top \mathbf{X}_t^{\text{train}} \preceq \mathbf{I}_d$ and $\mathbb{E}[\|(\mathbf{X}_t^{\text{val}})^\top \mathbf{X}_t^{\text{val}}\|_{\text{op}}] < \infty$, which implies

$$\mathbb{E}[\|(\mathbf{M}_t^{\text{tr-val}})^\top \mathbf{M}_t^{\text{tr-val}}\|_{\text{op}}^2] < \infty.$$

For condition (b), using a similar argument in condition (a) and combining with $R^2 = \mathbb{E}[\|\mathbf{w}_{0,*} - \mathbf{w}_t\|^2]$, we have $\mathbb{E}[\|\nabla \ell_t^{\{\text{tr-tr}, \text{tr-val}\}}(\mathbf{w}_{0,*})\|^2] < \infty$.

For condition (c), using (19), we directly verify that $L^{\{\text{tr-tr}, \text{tr-val}\}}$ is twice-differentiable and $\nabla^2 L^{\{\text{tr-tr}, \text{tr-val}\}} \succ \mathbf{0}$.

C.2 Proof of Theorem 5

Proof. In this section we prove Theorem 5. Using the asymptotic normality in Proposition 7, the asymptotic covariance is $\nabla^{-2} L^{\{\text{tr-tr}, \text{tr-val}\}} \text{Cov}[\nabla \ell_t^{\{\text{tr-tr}, \text{tr-val}\}}] \nabla^{-2} L^{\{\text{tr-tr}, \text{tr-val}\}}$. Therefore, in the following, we only need to find $\nabla^{-2} L^{\{\text{tr-tr}, \text{tr-val}\}}$ and $\text{Cov}[\nabla \ell_t^{\{\text{tr-tr}, \text{tr-val}\}}]$.

• **Asymptotic variance of $\hat{\mathbf{w}}_{0,T}^{\text{tr-tr}}$.** We begin with the computation of the expected Hessian $\frac{1}{n} \mathbb{E}[(\mathbf{M}_t^{\text{tr-tr}})^\top \mathbf{M}_t^{\text{tr-tr}}]$.

$$\begin{aligned} & \mathbb{E}[(\mathbf{M}_t^{\text{tr-tr}})^\top \mathbf{M}_t^{\text{tr-tr}}] \\ &= \mathbb{E} \left[\left(\mathbf{I}_d - (\mathbf{X}_t^\top \mathbf{X}_t + n\lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^\top \mathbf{X}_t \right)^\top \mathbf{X}_t^\top \mathbf{X}_t \left(\mathbf{I}_d - (\mathbf{X}_t^\top \mathbf{X}_t + n\lambda \mathbf{I}_d)^{-1} \mathbf{X}_t^\top \mathbf{X}_t \right) \right] \\ &\stackrel{(i)}{=} \mathbb{E} \left[\mathbf{V}_t (\mathbf{I}_d - (\mathbf{D}_t^\top \mathbf{D}_t + n\lambda \mathbf{I}_d)^{-1} \mathbf{D}_t^\top \mathbf{D}_t)^\top \mathbf{D}_t^\top \mathbf{D}_t (\mathbf{I}_d - (\mathbf{D}_t^\top \mathbf{D}_t + n\lambda \mathbf{I}_d)^{-1} \mathbf{D}_t^\top \mathbf{D}_t) \mathbf{V}_t^\top \right], \end{aligned} \quad (20)$$

where the equality (i) is obtained by plugging in the SVD of $\mathbf{X}_t = \mathbf{U}_t \mathbf{D}_t \mathbf{V}_t^\top$ with $\mathbf{U}_t \in \mathbb{R}^{n \times n}$, $\mathbf{D}_t \in \mathbb{R}^{n \times d}$, and $\mathbf{V}_t \in \mathbb{R}^{d \times d}$. A key observation is that \mathbf{U}_t and \mathbf{V}_t are independent, since \mathbf{X}_t is isotropic, i.e., homogeneous in each orthogonal direction. To see this, for any orthogonal matrices $\mathbf{Q} \in \mathbb{R}^{n \times n}$ and $\mathbf{P} \in \mathbb{R}^{d \times d}$, we know \mathbf{X}_t and $\mathbf{Q} \mathbf{X}_t \mathbf{P}^\top$ share the same distribution. Moreover, we have $\mathbf{Q} \mathbf{X}_t \mathbf{P}^\top = (\mathbf{Q} \mathbf{U}_t) \mathbf{D}_t (\mathbf{P} \mathbf{V}_t)^\top$ as the SVD. This shows that the left and right singular matrices

are independent and both uniformly distributed on all the orthogonal matrices of the corresponding dimensions ($\mathbb{R}^{n \times n}$ and $\mathbb{R}^{d \times d}$, respectively).

Recall that we denote $\sigma_1^{(n)} \geq \dots \geq \sigma_d^{(n)}$ as the eigenvalues of $\frac{1}{n} \mathbf{X}_t^\top \mathbf{X}_t$. Thus, we have $\mathbf{D}_t^\top \mathbf{D}_t = \text{Diag}(n\sigma_1^{(n)}, \dots, n\sigma_d^{(n)})$. We can further simplify (20) as

$$\begin{aligned} & \mathbb{E} \left[\mathbf{V}_t (\mathbf{I}_d - (\mathbf{D}_t^\top \mathbf{D}_t + n\lambda \mathbf{I}_d)^{-1} \mathbf{D}_t^\top \mathbf{D}_t)^\top \mathbf{D}_t^\top \mathbf{D}_t (\mathbf{I}_d - (\mathbf{D}_t^\top \mathbf{D}_t + n\lambda \mathbf{I}_d)^{-1} \mathbf{D}_t^\top \mathbf{D}_t) \mathbf{V}_t^\top \right] \\ &= \mathbb{E} \left[\mathbf{V}_t \text{Diag} \left(\frac{n\lambda^2 \sigma_1^{(n)}}{(\sigma_1^{(n)} + \lambda)^2}, \dots, \frac{n\lambda^2 \sigma_d^{(n)}}{(\sigma_d^{(n)} + \lambda)^2} \right) \mathbf{V}_t^\top \right] \\ &= \mathbb{E} \left[\sum_{i=1}^d \frac{n\lambda^2 \sigma_i^{(n)}}{(\sigma_i^{(n)} + \lambda)^2} \mathbf{v}_{t,i} \mathbf{v}_{t,i}^\top \right]. \end{aligned} \quad (21)$$

We will utilize the isotropicity of \mathbf{X}_t to find (21). Recall that we have shown that \mathbf{V}_t is uniform on all the orthogonal matrices. Let $\mathbf{P} \in \mathbb{R}^{d \times d}$ be any permutation matrix, then $\mathbf{V}_t \mathbf{P}$ has the same distribution as \mathbf{V}_t . For this permuted data matrix $\mathbf{V}_t \mathbf{P}$, (21) becomes

$$\mathbb{E} \left[\sum_{i=1}^d \frac{n\lambda^2 \sigma_i^{(n)}}{(\sigma_i^{(n)} + \lambda)^2} \mathbf{v}_{t,\tau_p(i)} \mathbf{v}_{t,\tau_p(i)}^\top \right] \quad \text{with } \tau_p(i) \text{ denotes the permutation of the } i\text{-th element in } \mathbf{P}.$$

Summing over all the permutations \mathbf{P} (and there are totally $d!$ instances), we deduce

$$\begin{aligned} & d! \mathbb{E}[(\mathbf{M}_t^{\text{tr-tr}})^\top \mathbf{M}_t^{\text{tr-tr}}] \\ &= \sum_{\text{all permutation } \tau_p} \mathbb{E} \left[\sum_{i=1}^d \frac{n\lambda^2 \sigma_i^{(n)}}{(\sigma_i^{(n)} + \lambda)^2} \mathbf{v}_{t,\tau_p(i)} \mathbf{v}_{t,\tau_p(i)}^\top \right] \\ &= (d-1)! \mathbb{E} \left[\sum_{j=1}^d \left[\sum_{i=1}^d \frac{n\lambda^2 \sigma_i^{(n)}}{(\sigma_i^{(n)} + \lambda)^2} \right] \mathbf{v}_{t,j} \mathbf{v}_{t,j}^\top \right] \\ &= (d-1)! \mathbb{E} \left[\mathbf{V}_t \text{Diag} \left(\sum_{i=1}^d \frac{\lambda^2 \sigma_i^{(n)}}{(\lambda + \sigma_i^{(n)})^2}, \dots, \sum_{i=1}^d \frac{\lambda^2 \sigma_i^{(n)}}{(\lambda + \sigma_i^{(n)})^2} \right) \mathbf{V}_t^\top \right] \\ &= (d-1)! \mathbb{E} \left[\sum_{i=1}^d \frac{\lambda^2 \sigma_i^{(n)}}{(\lambda + \sigma_i^{(n)})^2} \mathbf{V}_t \mathbf{I}_d \mathbf{V}_t^\top \right] \end{aligned} \quad (22)$$

Dividing $(d-1)!$ on both sides of (22) yields

$$\mathbb{E}[(\mathbf{M}_t^{\text{tr-tr}})^\top \mathbf{M}_t^{\text{tr-tr}}] = \frac{n}{d} \mathbb{E} \left[\sum_{i=1}^d \frac{\lambda^2 \sigma_i^{(n)}}{(\lambda + \sigma_i^{(n)})^2} \right] \mathbf{I}_d. \quad (23)$$

Next, we find the expected covariance matrix $\frac{1}{n_2} \mathbb{E}[\nabla \ell_t^{\text{tr-tr}}(\mathbf{w}_{0,*})(\nabla \ell_t^{\text{tr-tr}}(\mathbf{w}_{0,*}))^\top]$.

$$\begin{aligned} & \mathbb{E}[\nabla \ell_t^{\text{tr-tr}}(\mathbf{w}_{0,*})(\nabla \ell_t^{\text{tr-tr}}(\mathbf{w}_{0,*}))^\top] \\ &= \mathbb{E}[(\mathbf{M}_t^{\text{tr-tr}})^\top \mathbf{M}_t^{\text{tr-tr}}(\mathbf{w}_{0,*} - \mathbf{w}_t)(\mathbf{w}_{0,*} - \mathbf{w}_t)^\top (\mathbf{M}_t^{\text{tr-tr}})^\top \mathbf{M}_t^{\text{tr-tr}}] \\ &\stackrel{(i)}{=} \mathbb{E} \left[\mathbf{V}_t \text{Diag} \left(\frac{n\lambda^2 \sigma_1^{(n)}}{(\sigma_1^{(n)} + \lambda)^2}, \dots, \frac{n\lambda^2 \sigma_d^{(n)}}{(\sigma_d^{(n)} + \lambda)^2} \right) \mathbf{V}_t^\top (\mathbf{w}_{0,*} - \mathbf{w}_t)(\mathbf{w}_{0,*} - \mathbf{w}_t)^\top \right. \\ &\quad \left. \cdot \mathbf{V}_t \text{Diag} \left(\frac{n\lambda^2 \sigma_1^{(n)}}{(\sigma_1^{(n)} + \lambda)^2}, \dots, \frac{n\lambda^2 \sigma_d^{(n)}}{(\sigma_d^{(n)} + \lambda)^2} \right) \mathbf{V}_t^\top \right]. \end{aligned} \quad (24)$$

512 Here step (i) uses the SVD of \mathbf{X}_t and the computation in (21). Combining (23) and (24), we derive
 513 the asymptotic covariance matrix of using $L^{\text{tr-tr}}$ as

$$\begin{aligned}
 & \text{AsymCov}(\hat{\mathbf{w}}_{0,T}^{\text{tr-tr}}) \\
 &= \mathbb{E}[\nabla^{-2} \ell_t^{\text{tr-tr}}] \text{Cov}[\nabla \ell_t^{\text{tr-tr}}(\mathbf{w}_{0,*})] \mathbb{E}[\nabla^{-2} \ell_t^{\text{tr-tr}}] \\
 &= d^2 \left(\mathbb{E} \left[\sum_{i=1}^d \frac{\lambda^2 \sigma_i^{(n)}}{(\lambda + \sigma_i^{(n)})^2} \right] \right)^{-2} \\
 &\quad \cdot \mathbb{E} \left[\mathbf{V}_t \text{Diag} \left(\frac{n\lambda^2 \sigma_1^{(n)}}{(\sigma_1^{(n)} + \lambda)^2}, \dots, \frac{n\lambda^2 \sigma_d^{(n)}}{(\sigma_d^{(n)} + \lambda)^2} \right) \mathbf{V}_t^\top (\mathbf{w}_{0,*} - \mathbf{w}_t)(\mathbf{w}_{0,*} - \mathbf{w}_t)^\top \right. \\
 &\quad \left. \cdot \mathbf{V}_t \text{Diag} \left(\frac{n\lambda^2 \sigma_1^{(n)}}{(\sigma_1^{(n)} + \lambda)^2}, \dots, \frac{n\lambda^2 \sigma_d^{(n)}}{(\sigma_d^{(n)} + \lambda)^2} \right) \mathbf{V}_t^\top \right]. \tag{25}
 \end{aligned}$$

514 Taking trace in (25), we deduce

$$\begin{aligned}
 & \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-tr}}) \\
 &= \text{tr}(\text{AsymCov}(\hat{\mathbf{w}}_{0,T}^{\text{tr-tr}})) \\
 &= d^2 \left(\mathbb{E} \left[\sum_{i=1}^d \frac{\lambda^2 \sigma_i^{(n)}}{(\lambda + \sigma_i^{(n)})^2} \right] \right)^{-2} \\
 &\quad \cdot \text{tr} \left(\mathbb{E} \left[\mathbf{V}_t \text{Diag} \left(\frac{n^2 \lambda^4 (\sigma_1^{(n)})^2}{(\sigma_1^{(n)} + \lambda)^4}, \dots, \frac{n^2 \lambda^4 (\sigma_d^{(n)})^2}{(\sigma_d^{(n)} + \lambda)^4} \right) \mathbf{V}_t^\top (\mathbf{w}_{0,*} - \mathbf{w}_t)(\mathbf{w}_{0,*} - \mathbf{w}_t)^\top \right] \right) \\
 &\stackrel{(i)}{=} d^2 \left(\mathbb{E} \left[\sum_{i=1}^d \frac{\lambda^2 \sigma_i^{(n)}}{(\lambda + \sigma_i^{(n)})^2} \right] \right)^{-2} \\
 &\quad \cdot \frac{n^2}{d} \text{tr} \left(\mathbb{E} \left[\sum_{i=1}^d \frac{\lambda^4 (\sigma_i^{(n)})^2}{(\lambda + \sigma_i^{(n)})^4} \mathbf{V}_t \mathbf{I}_d \mathbf{V}_t^\top (\mathbf{w}_{0,*} - \mathbf{w}_t)(\mathbf{w}_{0,*} - \mathbf{w}_t)^\top \right] \right) \\
 &= d \frac{\mathbb{E} \left[\sum_{i=1}^d \lambda^4 (\sigma_i^{(n)})^2 / (\lambda + \sigma_i^{(n)})^4 \right]}{\left(\mathbb{E} \left[\sum_{i=1}^d \lambda^2 \sigma_i^{(n)} / (\lambda + \sigma_i^{(n)})^2 \right] \right)^2} \cdot \text{tr}(\mathbb{E}[(\mathbf{w}_{0,*} - \mathbf{w}_t)(\mathbf{w}_{0,*} - \mathbf{w}_t)^\top]) \\
 &= dR^2 \frac{\mathbb{E} \left[\sum_{i=1}^d (\sigma_i^{(n)})^2 / (\lambda + \sigma_i^{(n)})^4 \right]}{\left(\mathbb{E} \left[\sum_{i=1}^d \sigma_i^{(n)} / (\lambda + \sigma_i^{(n)})^2 \right] \right)^2}, \tag{26}
 \end{aligned}$$

515 where step (i) utilizes the independence between \mathbf{w}_t and \mathbf{X}_t and applies the permutation trick in
 516 (22) to find $\mathbb{E} \left[\mathbf{V}_t \text{Diag} \left(\frac{n^2 \lambda^4 (\sigma_1^{(n)})^2}{(\sigma_1^{(n)} + \lambda)^4}, \dots, \frac{n^2 \lambda^4 (\sigma_d^{(n)})^2}{(\sigma_d^{(n)} + \lambda)^4} \right) \mathbf{V}_t^\top \right]$.

517 • **Asymptotic variance of $\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}$.** Similar to the no-split case, we compute the Hessian
 518 $\frac{1}{n_2} \mathbb{E}[\nabla^2 \ell_t^{\text{tr-val}}] = \frac{1}{n_2} \mathbb{E}[(\mathbf{M}_t^{\text{tr-val}})^\top \mathbf{M}_t^{\text{tr-val}}]$ first.

$$\begin{aligned}
 & \mathbb{E}[(\mathbf{M}_t^{\text{tr-val}})^\top \mathbf{M}_t^{\text{tr-val}}] \\
 &= \mathbb{E} \left[\left(\mathbf{I}_d - ((\mathbf{X}_t^{\text{train}})^\top \mathbf{X}_t^{\text{train}} + n_1 \lambda \mathbf{I}_d)^{-1} (\mathbf{X}_t^{\text{train}})^\top \mathbf{X}_t^{\text{train}} \right)^\top (\mathbf{X}_t^{\text{val}})^\top \mathbf{X}_t^{\text{val}} \right. \\
 & \quad \left. \cdot \left(\mathbf{I}_d - ((\mathbf{X}_t^{\text{train}})^\top \mathbf{X}_t^{\text{train}} + n_1 \lambda \mathbf{I}_d)^{-1} (\mathbf{X}_t^{\text{train}})^\top \mathbf{X}_t^{\text{train}} \right) \right] \\
 &\stackrel{(i)}{=} n_2 \mathbb{E} \left[\left(\mathbf{I}_d - ((\mathbf{X}_t^{\text{train}})^\top \mathbf{X}_t^{\text{train}} + n_1 \lambda \mathbf{I}_d)^{-1} ((\mathbf{X}_t^{\text{train}})^\top \mathbf{X}_t^{\text{train}}) \right)^\top \right. \\
 & \quad \left. \cdot \left(\mathbf{I}_d - ((\mathbf{X}_t^{\text{train}})^\top \mathbf{X}_t^{\text{train}} + n_1 \lambda \mathbf{I}_d)^{-1} (\mathbf{X}_t^{\text{train}})^\top \mathbf{X}_t^{\text{train}} \right) \right] \\
 &\stackrel{(ii)}{=} n_2 \mathbb{E} \left[\mathbf{V}_t^{\text{train}} (\mathbf{I}_d - ((\mathbf{D}_t^{\text{train}})^\top \mathbf{D}_t^{\text{train}} + n_1 \lambda \mathbf{I}_d)^{-1} (\mathbf{D}_t^{\text{train}})^\top \mathbf{D}_t^{\text{train}})^2 (\mathbf{V}_t^{\text{train}})^\top \right], \quad (27)
 \end{aligned}$$

519 where (i) uses the data generating assumption $\mathbb{E}[(\mathbf{X}_t^{\text{val}})^\top \mathbf{X}_t^{\text{val}}] = n_2 \mathbf{I}_d$ and the independence between $\mathbf{X}_t^{\text{train}}$ and $\mathbf{X}_t^{\text{val}}$, and (ii) follows from the SVD of $\mathbf{X}_t^{\text{train}} = \mathbf{U}_t^{\text{train}} \mathbf{D}_t^{\text{train}} (\mathbf{V}_t^{\text{train}})^\top$.

521 Here we denote $\sigma_1^{(n_1)} \geq \dots \geq \sigma_d^{(n_1)}$ as the eigenvalues of $\frac{1}{n_1} (\mathbf{X}_t^{\text{train}})^\top \mathbf{X}_t^{\text{train}}$. Thus, we have
 522 $(\mathbf{D}_t^{\text{train}})^\top \mathbf{D}_t^{\text{train}} = \text{Diag}(n_1 \sigma_1^{(n_1)}, \dots, n_1 \sigma_d^{(n_1)})$. We can now further simplify (27) as

$$\begin{aligned}
 & n_2 \mathbb{E} \left[\mathbf{V}_t^{\text{train}} (\mathbf{I}_d - ((\mathbf{D}_t^{\text{train}})^\top \mathbf{D}_t^{\text{train}} + n_1 \lambda \mathbf{I}_d)^{-1} (\mathbf{D}_t^{\text{train}})^\top \mathbf{D}_t^{\text{train}})^2 (\mathbf{V}_t^{\text{train}})^\top \right] \\
 &\stackrel{(i)}{=} n_2 \mathbb{E} \left[\mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda^2}{(\sigma_1^{(n_1)} + \lambda)^2}, \dots, \frac{\lambda^2}{(\sigma_d^{(n_1)} + \lambda)^2} \right) (\mathbf{V}_t^{\text{train}})^\top \right] \\
 &\stackrel{(ii)}{=} \frac{n_2}{d} \mathbb{E} \left[\sum_{i=1}^d \frac{\lambda^2}{(\lambda + \sigma_i^{(n_1)})^2} \right] \mathbf{I}_d. \quad (28)
 \end{aligned}$$

523 Step (i) follows from the same computation in (21), and step (ii) uses the permutation trick in (22).

524 Next, we find the expected covariance matrix $\frac{1}{n_2} \mathbb{E}[\nabla \ell_t^{\text{tr-tr}}(\mathbf{w}_{0,*}) (\nabla \ell_t^{\text{tr-tr}}(\mathbf{w}_{0,*}))^\top]$.

$$\begin{aligned}
 & \mathbb{E}[\nabla \ell_t^{\text{tr-tr}}(\mathbf{w}_{0,*}) (\nabla \ell_t^{\text{tr-tr}}(\mathbf{w}_{0,*}))^\top] \\
 &= \mathbb{E}[(\mathbf{M}_t^{\text{tr-tr}})^\top \mathbf{M}_t^{\text{tr-tr}}(\mathbf{w}_{0,*} - \mathbf{w}_t)(\mathbf{w}_{0,*} - \mathbf{w}_t)^\top (\mathbf{M}_t^{\text{tr-tr}})^\top \mathbf{M}_t^{\text{tr-tr}}] \\
 &= \mathbb{E} \left[\mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \right. \\
 & \quad \cdot \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{w}_{0,*} - \mathbf{w}_t)(\mathbf{w}_{0,*} - \mathbf{w}_t)^\top \\
 & \quad \cdot \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \\
 & \quad \left. \cdot \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top \right]. \quad (29)
 \end{aligned}$$

525 Combining (28) and (29), we derive the asymptotic covariance matrix of using $L^{\text{tr-val}}$ as

$$\begin{aligned}
& \text{AsymCov}(\widehat{\mathbf{w}}_{0,T}^{\text{tr-val}}) \\
&= \mathbb{E}[\nabla^{-2} \ell_t^{\text{tr-val}}] \text{Cov}[\nabla \ell_t^{\text{tr-val}}(\mathbf{w}_{0,*})] \mathbb{E}[\nabla^{-2} \ell_t^{\text{tr-val}}] \\
&= \frac{d^2}{n_2} \left(\mathbb{E} \left[\sum_{i=1}^d \frac{\lambda^2}{(\lambda + \sigma_i^{(n_1)})^2} \right] \right)^{-2} \\
&\quad \cdot \mathbb{E} \left[\mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \right. \\
&\quad \cdot \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{w}_{0,*} - \mathbf{w}_t)(\mathbf{w}_{0,*} - \mathbf{w}_t)^\top \\
&\quad \cdot \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \\
&\quad \left. \cdot \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top \right]. \tag{30}
\end{aligned}$$

526 Taking trace in (30), we deduce

$$\begin{aligned}
& \text{AsymMSE}(\widehat{\mathbf{w}}_{0,T}^{\text{tr-tr}}) \\
&= \text{tr}(\text{AsymCov}(\widehat{\mathbf{w}}_{0,T}^{\text{tr-tr}})) \\
&= \frac{d^2}{n_2} \left(\mathbb{E} \left[\sum_{i=1}^d \frac{\lambda^2}{(\lambda + \sigma_i^{(n_1)})^2} \right] \right)^{-2} \\
&\quad \cdot \text{tr} \left(\mathbb{E} \left[\mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \right. \right. \\
&\quad \cdot \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{w}_{0,*} - \mathbf{w}_t)(\mathbf{w}_{0,*} - \mathbf{w}_t)^\top \\
&\quad \cdot \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \\
&\quad \left. \left. \cdot \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top \right] \right) \\
&\stackrel{(i)}{=} \frac{d^2}{n_2} \left(\mathbb{E} \left[\sum_{i=1}^d \frac{\lambda^2}{(\lambda + \sigma_i^{(n)})^2} \right] \right)^{-2} \\
&\quad \cdot \text{tr} \left(\mathbb{E} \left[\mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \right. \right. \\
&\quad \cdot \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda^2}{(\sigma_1^{(n_1)} + \lambda)^2}, \dots, \frac{\lambda^2}{(\sigma_d^{(n_1)} + \lambda)^2} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \\
&\quad \left. \left. \cdot \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{w}_{0,*} - \mathbf{w}_t)(\mathbf{w}_{0,*} - \mathbf{w}_t)^\top \right] \right), \tag{31}
\end{aligned}$$

527 where (i) follows from the cyclic property of the matrix trace operation. Due to the isotropicity of
 528 $\mathbf{X}_t^{\text{train}}$ and $\mathbf{X}_t^{\text{val}}$, we claim that

$$\begin{aligned} & \mathbb{E} \left[\mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \right. \\ & \quad \cdot \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda^2}{(\sigma_1^{(n_1)} + \lambda)^2}, \dots, \frac{\lambda^2}{(\sigma_d^{(n_1)} + \lambda)^2} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \\ & \quad \left. \cdot \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top \right] \end{aligned} \quad (32)$$

529 is a diagonal matrix $c\mathbf{I}_d$ with all the diagonal elements identical. We can show the claim bying
 530 taking expectation with respect to $\mathbf{X}_t^{\text{val}}$ first. Since $\mathbf{V}_t^{\text{train}}$ is an orthogonal matrix, $\mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}}$ has
 531 the same distribution as $\mathbf{X}_t^{\text{val}}$ and independent of \mathbf{X}_t . We verify that any off-diagonal element of the
 532 matrix expectation

$$\begin{aligned} \mathbf{A} := \mathbb{E}_{\mathbf{X}_t^{\text{val}}} & \left[(\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda^2}{(\sigma_1^{(n_1)} + \lambda)^2}, \dots, \frac{\lambda^2}{(\sigma_d^{(n_1)} + \lambda)^2} \right) \right. \\ & \left. \cdot (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \right] \end{aligned}$$

533 is zero. We denote $\mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n_2 \times d}$ with $\mathbf{x}_i \stackrel{\text{iid}}{\sim} \mathbf{N}(\mathbf{0}, \mathbf{I}_d)$. For $k \neq \ell$, the
 534 (k, ℓ) -th entry $A_{k, \ell}$ of \mathbf{A} is

$$\begin{aligned} A_{k, \ell} &= \mathbb{E} \left[\sum_j \left(\frac{\lambda^2}{(\sigma_j^{(n_1)} + \lambda)^2} \left(\sum_i x_{k, i} x_{j, i} \right) \left(\sum_i x_{j, i} x_{\ell, i} \right) \right) \right] \\ &= \mathbb{E} \left[\sum_j \frac{\lambda^2}{(\sigma_j^{(n_1)} + \lambda)^2} \left(\sum_{m, n} x_{k, m} x_{j, m} x_{j, n} x_{\ell, n} \right) \right] \\ &\stackrel{(i)}{=} 0, \end{aligned}$$

535 where $x_{i, j}$ denotes the j -th element of \mathbf{x}_i . Equality (i) holds, since either $x_{k, m}$ or $x_{\ell, n}$ only appears
 536 once in each summand. Therefore, we can write $\mathbf{A} = \text{Diag}(A_{1,1}, \dots, A_{d,d})$ with $A_{k,k}$ being

$$\begin{aligned} A_{k,k} &= \mathbb{E} \left[\sum_j \frac{\lambda^2}{(\sigma_j^{(n_1)} + \lambda)^2} \left(\sum_{m, n} x_{k, m} x_{j, m} x_{j, n} x_{k, n} \right) \right] \\ &= \mathbb{E} \left[\frac{\lambda^2}{(\sigma_k^{(n_1)} + \lambda)^2} \left(\sum_{m, n} x_{k, m} x_{k, m} x_{k, n} x_{k, n} \right) \right]. \end{aligned}$$

537 Observe that $A_{k,k}$ only depends on $\sigma_k^{(n_1)}$. Plugging back into (32), we have

$$\begin{aligned}
& \mathbb{E} \left[\mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \right. \\
& \quad \cdot \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda^2}{(\sigma_1^{(n_1)} + \lambda)^2}, \dots, \frac{\lambda^2}{(\sigma_d^{(n_1)} + \lambda)^2} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \\
& \quad \cdot \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top \Big] \\
&= \mathbb{E} \left[\mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) \text{Diag}(A_{1,1}, \dots, A_{d,d}) \right. \\
& \quad \cdot \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top \Big] \\
&= \mathbb{E} \left[\mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda^2 A_{1,1}}{(\sigma_1^{(n_1)} + \lambda)^2}, \dots, \frac{\lambda^2 A_{d,d}}{(\sigma_d^{(n_1)} + \lambda)^2} \right) (\mathbf{V}_t^{\text{train}})^\top \right] \\
&\stackrel{(i)}{=} c \mathbf{I}_d,
\end{aligned}$$

538 where equality (i) utilizes the permutation trick in (23). To this end, it is sufficient to find c as

$$\begin{aligned}
c &= \frac{1}{d} \text{tr} \left(\mathbb{E} \left[\mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \right. \right. \\
& \quad \cdot \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda^2}{(\sigma_1^{(n_1)} + \lambda)^2}, \dots, \frac{\lambda^2}{(\sigma_d^{(n_1)} + \lambda)^2} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \\
& \quad \cdot \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda}{\sigma_1^{(n_1)} + \lambda}, \dots, \frac{\lambda}{\sigma_d^{(n_1)} + \lambda} \right) (\mathbf{V}_t^{\text{train}})^\top \Big] \Big) \\
&= \frac{1}{d} \text{tr} \left(\mathbb{E} \left[\mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda^2}{(\sigma_1^{(n_1)} + \lambda)^2}, \dots, \frac{\lambda^2}{(\sigma_d^{(n_1)} + \lambda)^2} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \right. \right. \\
& \quad \cdot \mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \text{Diag} \left(\frac{\lambda^2}{(\sigma_1^{(n_1)} + \lambda)^2}, \dots, \frac{\lambda^2}{(\sigma_d^{(n_1)} + \lambda)^2} \right) (\mathbf{V}_t^{\text{train}})^\top (\mathbf{X}_t^{\text{val}})^\top \Big] \Big). \quad (33)
\end{aligned}$$

539 Observe again that $\mathbf{X}_t^{\text{val}} \mathbf{V}_t^{\text{train}} \in \mathbb{R}^{n_2 \times d}$ is a Gaussian random matrix. We rewrite (33) as

$$c = \frac{1}{d} \mathbb{E} \left[\left(\sum_{i,j=1}^{n_2} \mathbf{v}_i^\top \text{Diag} \left(\frac{\lambda^2}{(\sigma_1^{(n_1)} + \lambda)^2}, \dots, \frac{\lambda^2}{(\sigma_d^{(n_1)} + \lambda)^2} \right) \mathbf{v}_j \right)^2 \right], \quad (34)$$

540 where $\mathbf{v}_i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ is i.i.d. Gaussian random vectors for $i = 1, \dots, n_2$. To compute (34), we need
541 the following result.

542 **Claim 8.** Given any symmetric matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ and i.i.d. standard Gaussian random vectors
543 $\mathbf{v}, \mathbf{u} \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$, we have

$$\mathbb{E} [(\mathbf{v}^\top \mathbf{A} \mathbf{v})^2] = 2 \|\mathbf{A}\|_{\text{Fr}}^2 + \text{tr}^2(\mathbf{A}) \quad \text{and} \quad (35)$$

$$\mathbb{E} [(\mathbf{v}^\top \mathbf{A} \mathbf{u})^2] = \|\mathbf{A}\|_{\text{Fr}}^2. \quad (36)$$

544 *Proof of Claim 8.* We show (35) first. We denote $A_{i,j}$ as the (i,j) -th element of \mathbf{A} and v_i as the i -th
 545 element of \mathbf{v} . Expanding the quadratic form, we have

$$\begin{aligned}
 \mathbb{E} [(\mathbf{v}^\top \mathbf{A} \mathbf{v})^2] &= \mathbb{E} \left[\sum_{i,j,k,\ell \leq d} v_i v_j v_k v_\ell A_{i,j} A_{k,\ell} \right] \\
 &= \mathbb{E} \left[\sum_{i \leq d} v_i^4 A_{i,i}^2 \right] + \mathbb{E} \left[\sum_{i \neq j} v_i^2 v_j^2 (A_{i,j}^2 + A_{i,i} A_{j,j} + A_{i,j} A_{j,i}) \right] \\
 &= 3 \sum_{i \leq d} A_{i,i}^2 + \sum_{i \neq j} (A_{i,j}^2 + A_{i,i} A_{j,j} + A_{i,j} A_{j,i}) \\
 &= \text{tr}^2(\mathbf{A}) + 2 \sum_{i \leq d} A_{i,i}^2 + \sum_{i \neq j} (A_{i,j}^2 + A_{i,j} A_{j,i}) \\
 &= \text{tr}^2(\mathbf{A}) + 2 \|\mathbf{A}\|_{\text{Fr}}^2.
 \end{aligned}$$

546 Next, we show (36) by the cyclic property of trace.

$$\mathbb{E} [(\mathbf{v}^\top \mathbf{A} \mathbf{u})^2] = \text{tr} (\mathbb{E} [\mathbf{u} \mathbf{u}^\top \mathbf{A} \mathbf{v} \mathbf{v}^\top \mathbf{A}]) = \text{tr}(\mathbf{A}^2) = \|\mathbf{A}\|_{\text{Fr}}^2.$$

547

□

548 We back to the computation of (34) using Claim 8.

$$\begin{aligned}
 c &= \frac{1}{d} \mathbb{E} \left[\sum_{i,j=1}^{n_2} \left(\mathbf{v}_i^\top \text{Diag} \left(\frac{\lambda^2}{(\sigma_1^{(n_1)} + \lambda)^2}, \dots, \frac{\lambda^2}{(\sigma_d^{(n_1)} + \lambda)^2} \right) \mathbf{v}_j \right)^2 \right] \\
 &= \frac{1}{d} \mathbb{E} \left[\sum_{i=1}^{n_2} \left(\mathbf{v}_i^\top \text{Diag} \left(\frac{\lambda^2}{(\sigma_1^{(n_1)} + \lambda)^2}, \dots, \frac{\lambda^2}{(\sigma_d^{(n_1)} + \lambda)^2} \right) \mathbf{v}_i \right)^2 \right] \\
 &\quad + \frac{1}{d} \mathbb{E} \left[\sum_{i \neq j} \left(\mathbf{v}_i^\top \text{Diag} \left(\frac{\lambda^2}{(\sigma_1^{(n_1)} + \lambda)^2}, \dots, \frac{\lambda^2}{(\sigma_d^{(n_1)} + \lambda)^2} \right) \mathbf{v}_j \right)^2 \right] \\
 &= \frac{n_2}{d} \mathbb{E} \left[\text{tr}^2 \left(\text{Diag} \left(\frac{\lambda^2}{(\sigma_1^{(n_1)} + \lambda)^2}, \dots, \frac{\lambda^2}{(\sigma_d^{(n_1)} + \lambda)^2} \right) \right) \right] \\
 &\quad + 2 \frac{n_2}{d} \mathbb{E} \left[\left\| \text{Diag} \left(\frac{\lambda^2}{(\sigma_1^{(n_1)} + \lambda)^2}, \dots, \frac{\lambda^2}{(\sigma_d^{(n_1)} + \lambda)^2} \right) \right\|_{\text{Fr}}^2 \right] \\
 &\quad + \frac{n_2(n_2 - 1)}{d} \mathbb{E} \left[\left\| \text{Diag} \left(\frac{\lambda^2}{(\sigma_1^{(n_1)} + \lambda)^2}, \dots, \frac{\lambda^2}{(\sigma_d^{(n_1)} + \lambda)^2} \right) \right\|_{\text{Fr}}^2 \right] \\
 &= \frac{n_2}{d} \left(\mathbb{E} \left[\sum_{i=1}^d \frac{\lambda^2}{(\sigma_i^{(n_1)} + \lambda)^2} \right]^2 + (n_2 + 1) \mathbb{E} \left[\sum_{i=1}^d \frac{\lambda^4}{(\sigma_i^{(n_1)} + \lambda)^4} \right] \right). \tag{37}
 \end{aligned}$$

Combining (37) and (32), by the independence between \mathbf{w}_t and $\mathbf{X}_t^{\text{train}}, \mathbf{X}_t^{\text{val}}$, we compute (31) as

$$\begin{aligned} & \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-tr}}) \\ &= \frac{d^2}{n_2} \left(\mathbb{E} \left[\sum_{i=1}^d \frac{\lambda^2}{(\lambda + \sigma_i^{(n_1)})^2} \right] \right)^{-2} \\ & \quad \cdot \frac{n_2}{d} \left(\mathbb{E} \left[\sum_{i=1}^d \frac{\lambda^2}{(\sigma_i^{(n_1)} + \lambda)^2} \right]^2 + (n_2 + 1) \mathbb{E} \left[\sum_{i=1}^d \frac{\lambda^4}{(\sigma_i^{(n_1)} + \lambda)^4} \right] \right) \\ & \quad \cdot \mathbb{E} [(\mathbf{w}_{0,*} - \mathbf{w}_t)(\mathbf{w}_{0,*} - \mathbf{w}_t)^\top] \\ &= dR^2 \frac{\mathbb{E} \left[\sum_{i=1}^d \lambda^2 / (\sigma_i^{(n_1)} + \lambda)^2 \right]^2 + (n_2 + 1) \mathbb{E} \left[\sum_{i=1}^d \lambda^4 / (\sigma_i^{(n_1)} + \lambda)^4 \right]}{\left(\mathbb{E} \left[\sum_{i=1}^d \lambda^2 / (\lambda + \sigma_i^{(n_1)})^2 \right] \right)^2}. \end{aligned}$$

The proof is complete. \square

C.3 Optimally tuned rate of the train-val method at finite (n, d)

Theorem 5 provides explicit expressions for the asymptotic rates of both the train-val and the train-train methods in terms of expectations that involve the spectrum of $\frac{1}{n} \mathbf{X}_t^\top \mathbf{X}_t$. It may be difficult to compare the rates without further simplifying these expectations. However, if we only care about the *optimal* rates obtained by tuning λ and n_1, n_2 , then we show that the asymptotic rate of the train-val method can be dramatically simplified at any finite (n, d) :

Corollary 9 (Optimal rate of the train-val method at finite (n, d)). *For any (n, d) and any split ratio $(n_1, n_2) = (n_1, n - n_1)$, the optimal rate (by tuning the regularization $\lambda > 0$) of the train-val method is achieved at*

$$\inf_{\lambda > 0} \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}(n_1, n_2; \lambda)) = \lim_{\lambda \rightarrow \infty} \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}(n_1, n_2; \lambda)) = \frac{(d + n_2 + 1)R^2}{n_2}.$$

Further optimizing the rate over n_2 , the best rate is taken at $(n_1, n_2) = (0, n)$, in which the rate is

$$\inf_{\lambda > 0, n_2 \in [n]} \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}(n_1, n_2; \lambda)) = \frac{(d + n + 1)R^2}{n}.$$

Using all data as validation The proof of Corollary 9 can be found in Appendix C.4. Corollary 9 suggests that the optimal asymptotic rate of the train-val method is obtained at $\lambda = \infty$ and $(n_1, n_2) = (0, n)$. In other words, the optimal choice for the train-val method is to *use all the data as validation*. In this case, since there is no training data, the inner solver reduces to the identity map: $\mathcal{A}_{\infty,0}(\mathbf{w}_0; \mathbf{X}_t, \mathbf{y}_t) = \mathbf{w}_0$, and the outer loop reduces to learning a single linear model \mathbf{w}_0 on all the tasks combined. We remark that while the optimality of such a split ratio is likely an artifact of the data distribution we assumed (noiseless realizable linear model) and may not generalize to other meta-learning problems, we do find experimentally that using more data as validation (than training) can also improve the performance on real meta-learning tasks (see Table 2).

C.4 Proof of Corollary 9

Fix $n_1 \in [n]$ and $n_2 = n - n_1$. Recall from Theorem 5 that

$$\text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}(n_1, n_2; \lambda)) = \frac{dR^2}{n_2} \cdot \frac{\mathbb{E} \left[\left(\sum_{i=1}^d \lambda^2 / (\sigma_i^{(n_1)} + \lambda)^2 \right)^2 + (n_2 + 1) \sum_{i=1}^d \lambda^4 / (\sigma_i^{(n_1)} + \lambda)^4 \right]}{\left(\mathbb{E} \left[\sum_{i=1}^d \lambda^2 / (\sigma_i^{(n_1)} + \lambda)^2 \right] \right)^2}.$$

Clearly, as $\lambda \rightarrow \infty$, we have

$$\lim_{\lambda \rightarrow \infty} \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}(n_1, n_2; \lambda)) = \frac{dR^2}{n_2} \cdot \frac{d^2 + (n_2 + 1)d}{d^2} = \frac{(d + n_2 + 1)R^2}{n_2}.$$

573 It remains to show that the above quantity is a lower bound for $\text{AsymMSE}(\widehat{\mathbf{w}}_{0,T}^{\text{tr-val}}(n_1, n_2; \lambda))$ for
 574 any $\lambda > 0$, which is equivalent to

$$\frac{\mathbb{E}\left[\left(\sum_{i=1}^d \lambda^2/(\sigma_i^{(n_1)} + \lambda)^2\right)^2 + (n_2 + 1) \sum_{i=1}^d \lambda^4/(\sigma_i^{(n_1)} + \lambda)^4\right]}{\left(\mathbb{E}\left[\sum_{i=1}^d \lambda^2/(\sigma_i^{(n_1)} + \lambda)^2\right]\right)^2} \geq \frac{d + n_2 + 1}{d}, \quad \text{for all } \lambda > 0. \quad (38)$$

575 We now prove (38). For $i \in [n_1]$, define random variables

$$X_i := \frac{\lambda^2}{(\sigma_i^{(n_1)} + \lambda)^2} \in [0, 1] \quad \text{and} \quad Y_i := 1 - X_i \in [0, 1].$$

576 Then the left-hand side of (38) can be rewritten as

$$\begin{aligned} & \frac{\mathbb{E}\left[(d - n_1 + \sum_{i=1}^{n_1} X_i)^2 + (n_2 + 1)(d - n_1 + \sum_{i=1}^{n_1} X_i^2)\right]}{(\mathbb{E}[d - n_1 + \sum_{i=1}^{n_1} X_i])^2} \\ &= \frac{\mathbb{E}\left[(d - \sum_{i=1}^{n_1} Y_i)^2 + (n_2 + 1)(d - 2 \sum_{i=1}^{n_1} Y_i + \sum_{i=1}^{n_1} Y_i^2)\right]}{(\mathbb{E}[d - \sum_{i=1}^{n_1} Y_i])^2} \\ &= \frac{d^2 + (n_2 + 1)d - 2(d + n_2 + 1)\mathbb{E}[\sum Y_i] + \mathbb{E}[(\sum Y_i)^2] + (n_2 + 1)\mathbb{E}[\sum Y_i^2]}{d^2 - 2d\mathbb{E}[\sum Y_i] + (\mathbb{E}[\sum Y_i])^2} \end{aligned}$$

577 By algebraic manipulation, inequality (38) is equivalent to showing that

$$\frac{\mathbb{E}[(\sum Y_i)^2] + (n_2 + 1)\mathbb{E}[\sum Y_i^2]}{(\mathbb{E}[\sum Y_i])^2} \geq \frac{d + n_2 + 1}{d}. \quad (39)$$

578 Clearly, $\mathbb{E}[(\sum Y_i)^2] \geq (\mathbb{E}[\sum Y_i])^2$. By Cauchy-Schwarz we also have

$$\mathbb{E}[\sum Y_i^2] \geq \frac{1}{n_1} \mathbb{E}\left[\left(\sum Y_i\right)^2\right] \geq \frac{1}{n_1} \left(\mathbb{E}[\sum Y_i]\right)^2.$$

579 Therefore we have

$$\frac{\mathbb{E}[(\sum Y_i)^2] + (n_2 + 1)\mathbb{E}[\sum Y_i^2]}{(\mathbb{E}[\sum Y_i])^2} \geq 1 + \frac{n_2 + 1}{n_1} \geq 1 + \frac{n_2 + 1}{d} = \frac{d + n_2 + 1}{d},$$

580 where we have used that $n_1 \leq n \leq d$. This shows (39) and consequently (38). \square

581 C.5 Proof of Theorem 10

582 Let $\widehat{\Sigma}_n := \frac{1}{n} \mathbf{X}_t \mathbf{X}_t^\top$ denote the sample covariance matrix of the inputs in a single task (t). By
 583 Theorem 5, we have

$$\begin{aligned} \text{AsymMSE}(\widehat{\mathbf{w}}_{0,T}^{\text{tr-tr}}(n; \lambda)) &= R^2 \cdot \frac{\frac{1}{d} \mathbb{E}\left[\sum_{i=1}^d \sigma_i(\widehat{\Sigma}_n)^2/(\sigma_i(\widehat{\Sigma}_n) + \lambda)^4\right]}{\left(\frac{1}{d} \mathbb{E}\left[\sum_{i=1}^d \sigma_i(\widehat{\Sigma}_n)/(\sigma_i(\widehat{\Sigma}_n) + \lambda)^2\right]\right)^2} \\ &= R^2 \cdot \underbrace{\frac{1}{d} \mathbb{E}\left[\text{tr}\left((\widehat{\Sigma}_n + \lambda \mathbf{I}_d)^{-4} \widehat{\Sigma}_n^2\right)\right]}_{\text{I}_{n,d}} \bigg/ \underbrace{\left\{ \frac{1}{d} \mathbb{E}\left[\text{tr}\left((\widehat{\Sigma}_n + \lambda \mathbf{I}_d)^{-2} \widehat{\Sigma}_n\right)\right] \right\}^2}_{\text{II}_{n,d}}. \end{aligned} \quad (40)$$

584 We now evaluate quantities $\text{I}_{n,d}$ and $\text{II}_{n,d}$ in the high-dimensional limit of $d, n \rightarrow \infty$, $d/n \rightarrow \gamma \in$
 585 $(0, \infty)$. Consider the (slightly generalized) Stieltjes transform of $\widehat{\Sigma}_n$ defined for all $\lambda_1, \lambda_2 > 0$:

$$s(\lambda_1, \lambda_2) := \lim_{d, n \rightarrow \infty, d/n \rightarrow \gamma} \frac{1}{d} \mathbb{E}\left[\text{tr}\left((\lambda_1 \mathbf{I}_d + \lambda_2 \widehat{\Sigma}_n)^{-1}\right)\right].$$

As the entries of \mathbf{X}_t are i.i.d. $N(0, 1)$, the above limiting Stieltjes transform is the Stieltjes form of the Marchenko-Pastur law, which has a closed form (see, e.g. (Dobriban et al., 2018, Equation (7)))

$$\begin{aligned} s(\lambda_1, \lambda_2) &= \lambda_2^{-1} s(\lambda_1/\lambda_2, 1) = \frac{1}{\lambda_2} \cdot \frac{\gamma - 1 - \lambda_1/\lambda_2 + \sqrt{(\lambda_1/\lambda_2 + 1 + \gamma)^2 - 4\gamma}}{2\gamma\lambda_1/\lambda_2} \\ &= \frac{\gamma - 1 - \lambda_1/\lambda_2 + \sqrt{(\lambda_1/\lambda_2 + 1 + \gamma)^2 - 4\gamma}}{2\gamma\lambda_1}. \end{aligned} \quad (41)$$

Now observe that differentiating $s(\lambda_1, \lambda_2)$ yields quantity Π (known as the derivative trick of Stieltjes transforms). Indeed, we have

$$\begin{aligned} -\frac{d}{d\lambda_2} s(\lambda_1, \lambda_2) &= -\frac{d}{d\lambda_2} \lim_{d, n \rightarrow \infty, d/n \rightarrow \gamma} \frac{1}{d} \mathbb{E} \left[\text{tr} \left((\lambda_1 \mathbf{I}_d + \lambda_2 \widehat{\Sigma}_n)^{-1} \right) \right] \\ &= \lim_{d, n \rightarrow \infty, d/n \rightarrow \gamma} \frac{1}{d} \mathbb{E} \left[-\frac{d}{d\lambda_2} \text{tr} \left((\lambda_1 \mathbf{I}_d + \lambda_2 \widehat{\Sigma}_n)^{-1} \right) \right] \\ &= \lim_{d, n \rightarrow \infty, d/n \rightarrow \gamma} \frac{1}{d} \mathbb{E} \left[\text{tr} \left((\lambda_1 \mathbf{I}_d + \lambda_2 \widehat{\Sigma}_n)^{-2} \widehat{\Sigma}_n \right) \right]. \end{aligned}$$

(Above, the exchange of differentiation and limit is due to the uniform convergence of the derivatives, which holds at any $\lambda_1, \lambda_2 > 0$.) Taking $\lambda_1 = \lambda$ and $\lambda_2 = 1$, we get

$$\lim_{d, n \rightarrow \infty, d/n \rightarrow \gamma} \Pi_{n,d} = \lim_{d, n \rightarrow \infty, d/n \rightarrow \gamma} \frac{1}{d} \mathbb{E} \left[\text{tr} \left((\lambda \mathbf{I}_d + \widehat{\Sigma}_n)^{-2} \widehat{\Sigma}_n \right) \right] = -\frac{d}{d\lambda_2} s(\lambda_1, \lambda_2) |_{\lambda_1=\lambda, \lambda_2=1}.$$

Similarly we have

$$\lim_{d, n \rightarrow \infty, d/n \rightarrow \gamma} \mathbf{I}_{n,d} = \lim_{d, n \rightarrow \infty, d/n \rightarrow \gamma} \frac{1}{d} \mathbb{E} \left[\text{tr} \left((\lambda \mathbf{I}_d + \widehat{\Sigma}_n)^{-4} \widehat{\Sigma}_n^2 \right) \right] = -\frac{1}{6} \frac{d}{d\lambda_1} \frac{d^2}{d\lambda_2^2} s(\lambda_1, \lambda_2) |_{\lambda_1=\lambda, \lambda_2=1}.$$

Evaluating the right-hand sides from differentiating the closed-form expression (41), we get

$$\begin{aligned} \lim_{d, n \rightarrow \infty, d/n \rightarrow \gamma} \mathbf{I}_{n,d} &= \frac{1}{2\gamma} \cdot \frac{\lambda + 1 + \gamma}{\sqrt{(\lambda + 1 + \gamma)^2 - 4\gamma}} - \frac{1}{2\gamma}, \\ \lim_{d, n \rightarrow \infty, d/n \rightarrow \gamma} \Pi_{n,d} &= \frac{(\gamma - 1)^2 + (\gamma + 1)\lambda}{((\lambda + 1 + \gamma)^2 - 4\gamma)^{5/2}}. \end{aligned}$$

Substituting back to (40) yields that

$$\begin{aligned} \lim_{d, n \rightarrow \infty, d/n \rightarrow \gamma} \text{AsymMSE}(\widehat{\mathbf{w}}_{0,T}^{\text{tr-tr}}(n; \lambda)) &= \lim_{d, n \rightarrow \infty, d/n \rightarrow \gamma} R^2 \cdot \mathbf{I}_{n,d} / \Pi_{n,d}^2 \\ &= R^2 \cdot \frac{4\gamma^2 [(\gamma - 1)^2 + (\gamma + 1)\lambda]}{((\lambda + 1 + \gamma)^2 - 4\gamma)^{5/2} \cdot \left(\frac{\lambda + 1 + \gamma}{\sqrt{(\lambda + 1 + \gamma)^2 - 4\gamma}} - 1 \right)^2} \\ &= R^2 \cdot \frac{4\gamma^2 [(\gamma - 1)^2 + (\gamma + 1)\lambda]}{((\lambda + 1 + \gamma)^2 - 4\gamma)^{3/2} \cdot \left(\lambda + 1 + \gamma - \sqrt{(\lambda + 1 + \gamma)^2 - 4\gamma} \right)^2}. \end{aligned}$$

This proves the desired result. \square

C.6 Rate of train-train in the proportional limit

In order to compare the two methods, we need to further simplify the asymptotic rate of the train-train method in Theorem 5 and optimize it over λ . We perform such calculation under the *proportional limit*, where we let $d, n \rightarrow \infty$ and $d/n \rightarrow \gamma$, where $\gamma \in (0, \infty)$ is a fixed constant¹. The proportional limit allows a better understanding of the limiting spectrum of $\frac{1}{n} \mathbf{X}_t^\top \mathbf{X}_t$ (Bai & Silverstein, 2010; Anderson et al., 2010), and is also an accurate approximation of the finite (n, d) setting when (n, d) are large and roughly proportional.

We begin by calculating the asymptotic rate of the train-train method in the proportional limit.

¹Recall that the definition of AsymMSE already takes the $T \rightarrow \infty$ limit first, and thus our proportional limit can be thought as d, n being large, but still much smaller than T .

Theorem 10 (Exact rates of the train-train method in the proportional limit). *In the high-dimensional limiting regime $d, n \rightarrow \infty$, $d/n \rightarrow \gamma$ where $\gamma \in (0, \infty)$ is a fixed shape parameter, for any $\lambda > 0$*

$$\lim_{d, n \rightarrow \infty, d/n = \gamma} \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-tr}}(n; \lambda)) = \rho_{\lambda, \gamma} R^2.$$

where $\rho_{\lambda, \gamma} = 4\gamma^2 [(\gamma - 1)^2 + (\gamma + 1)\lambda] / (\lambda + 1 + \gamma - \sqrt{(\lambda + \gamma + 1)^2 - 4\gamma})^2 / ((\lambda + \gamma + 1)^2 - 4\gamma)^{3/2}$.

The proof builds on the Steiltjes transform and its “derivative trick”, and is deferred to Appendix C.5. With this exact rate at hand, we are now able to optimize the rate for the train-train method over λ , and compare with the train-val method.

C.7 Proof of Theorem 6

Throughout this proof we assume that $R^2 = 1$ without loss of generality (as all the rates are constant multiples of R^2).

Part I: Optimal rate for $L^{\text{tr-tr}}$ By Theorem 10, we have

$$\begin{aligned} & \inf_{\lambda > 0} \lim_{d, n \rightarrow \infty, d/n = \gamma} \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-tr}}(n; \lambda)) \\ &= \inf_{\lambda > 0} \frac{4\gamma^2 [(\gamma - 1)^2 + (\gamma + 1)\lambda]}{\underbrace{(\lambda + 1 + \gamma - \sqrt{(\lambda + \gamma + 1)^2 - 4\gamma})^2 \cdot ((\lambda + \gamma + 1)^2 - 4\gamma)^{3/2}}_{:=f(\lambda, \gamma)}}. \end{aligned}$$

In order to bound $\inf_{\lambda > 0} f(\lambda, \gamma)$, picking any $\lambda = \lambda(\gamma)$ gives $f(\lambda(\gamma), \gamma)$ as a valid upper bound, and our goal is to choose λ that yields a bound as tight as possible. Here we consider the choice

$$\lambda = \lambda(\gamma) = \max \{1 - \gamma/2, \gamma - 1/2\} = (1 - \gamma/2)\mathbf{1}\{\gamma \leq 1\} + (\gamma - 1/2)\mathbf{1}\{\gamma > 1\}$$

which we now show yields the claimed upper bound.

Case 1: $\gamma \leq 1$ Substituting $\lambda = 1 - \gamma/2$ into $f(\lambda, \gamma)$ and simplifying, we get

$$f(1 - \gamma/2, \gamma) = \frac{2(\gamma^2 - 3\gamma + 4)}{(2 - \gamma/2)^3} =: g_1(\gamma).$$

Clearly, $g_1(0) = 1$ and $g_1(1) = 32/27$. Further differentiating g_1 twice gives

$$g_1''(\gamma) = \frac{\gamma^2 + 7\gamma + 4}{(2 - \gamma/2)^5} > 0 \quad \text{for all } \gamma \in [0, 1].$$

Thus g_1 is convex on $[0, 1]$, from which we conclude that

$$g_1(\gamma) \leq (1 - \gamma) \cdot g_1(0) + \gamma \cdot g_1(1) = 1 + \frac{5}{27}\gamma.$$

Case 2: $\gamma > 1$ Substituting $\lambda = \gamma - 1/2$ into $f(\lambda, \gamma)$ and simplifying, we get

$$f(\gamma - 1/2, \gamma) = \frac{2\gamma^2(4\gamma^2 - 3\gamma + 1)}{(2\gamma - 1/2)^3} =: g_2(\gamma).$$

We have $g_2(1) = g_1(1) = 32/27$. Further differentiating g_2 gives

$$g_2'(\gamma) = -\frac{1}{(4\gamma - 1)^2} - \frac{6}{(4\gamma - 1)^3} - \frac{6}{(4\gamma - 1)^4} + 1 < 1 \quad \text{for all } \gamma > 1.$$

Therefore we have for all $\gamma > 1$ that

$$g_2(\gamma) = g_2(1) + \int_1^\gamma g_2'(t) dt \leq g_2(1) + \gamma - 1 = \gamma + \frac{5}{27}.$$

624 Combining Case 1 and 2, we get

$$\begin{aligned} & \inf_{\lambda > 0} f(\lambda, \gamma) \leq g_1(\gamma) \\ & \leq \mathbf{1}\{\gamma \leq 1\} + g_2(\gamma) \mathbf{1}\{\gamma > 1\} \leq \left(1 + \frac{5}{27}\gamma\right) \mathbf{1}\{\gamma \leq 1\} + \left(\frac{5}{27} + \gamma\right) \mathbf{1}\{\gamma > 1\} \\ & = \max \left\{1 + \frac{5}{27}\gamma, \frac{5}{27} + \gamma\right\}. \end{aligned}$$

625 This is the desired upper bound for $L^{\text{tr-tr}}$.

626 **Equality at $\gamma = 1$** We finally show that the above upper bound becomes an equality when $\gamma = 1$.
 627 At $\gamma = 1$, we have

$$f(\lambda, 1) = \frac{8\lambda}{(\lambda + 2 - \sqrt{\lambda^2 + 4\lambda})^2(\lambda^2 + 4\lambda)^{3/2}} = \frac{8\lambda^{-4}}{(1 + 2/\lambda - \sqrt{1 + 4/\lambda})^2(1 + 4/\lambda)^{3/2}}.$$

628 Make the change of variable $t = \sqrt{1 + 4/\lambda}$ so that $\lambda^{-1} = (t^2 - 1)/4$, minimizing the above
 629 expression is equivalent to minimizing

$$\frac{(t^2 - 1)^4/32}{(t^2/2 - t + 1/2)^2 t^3} = \frac{(t + 1)^4}{8t^3}$$

630 over $t > 1$. It is straightforward to check (by computing the first and second derivatives) that the
 631 above quantity is minimized at $t = 3$ with value $32/27$. In other words, we have shown

$$\inf_{\lambda > 0} f(\lambda, 1) = \frac{32}{27} = \max \left\{1 + \frac{5}{27}\gamma, \frac{5}{27} + \gamma\right\} \Big|_{\gamma=1},$$

632 that is, the equality holds at $\gamma = 1$.

633 **Part II: Optimal rate for $L^{\text{tr-val}}$** We now prove the result on $L^{\text{tr-val}}$, that is,

$$\begin{aligned} & \inf_{\lambda > 0, s \in (0,1)} \lim_{d, n \rightarrow \infty, d/n = \gamma} \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}(ns, n(1-s); \lambda)) \\ & \stackrel{(i)}{=} \lim_{d, n \rightarrow \infty, d/n = \gamma} \underbrace{\inf_{\lambda > 0, n_1 + n_2 = n} \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}(n_1, n_2; \lambda))}_{\frac{d+n+1}{n}} \stackrel{(ii)}{=} 1 + \gamma. \end{aligned}$$

634 First, equality (ii) follows from Corollary 9 and the fact that $(d + n + 1)/n \rightarrow 1 + \gamma$. Second, the
 635 “ \geq ” direction of equality (i) is trivial (since we always have “ $\inf \lim \geq \lim \inf$ ”). Therefore we get
 636 the “ \geq ” direction of the overall equality, and it remains to prove the “ \leq ” direction.

637 For the “ \leq ” direction, we fix any $\lambda > 0$, and bound $\text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}(n_1, n_2; \lambda))$ (and conse-
 638 quently its limit as $d, n \rightarrow \infty$.) We have from Theorem 5 that

$$\begin{aligned} & \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}(n_1, n_2; \lambda)) \\ & = \frac{d}{n_2} \cdot \frac{\mathbb{E} \left[\left(\sum_{i=1}^d \lambda^2 / (\sigma_i^{(n_1)} + \lambda)^2 \right)^2 + (n_2 + 1) \sum_{i=1}^d \lambda^4 / (\sigma_i^{(n_1)} + \lambda)^4 \right]}{\left(\mathbb{E} \left[\sum_{i=1}^d \lambda^2 / (\sigma_i^{(n_1)} + \lambda)^2 \right] \right)^2} \\ & \leq \frac{d}{n_2} \cdot \frac{d^2 + (n_2 + 1)d}{\left(\mathbb{E} \left[\sum_{i=1}^d \lambda^2 / (\sigma_i^{(n_1)} + \lambda)^2 \right] \right)^2} \\ & = \frac{d + n_2 + 1}{n_2} \cdot \frac{1}{\left(\mathbb{E} \left[\frac{1}{d} \sum_{i=1}^d \lambda^2 / (\sigma_i^{(n_1)} + \lambda)^2 \right] \right)^2} \end{aligned}$$

639 Observe that

$$\begin{aligned} & \mathbb{E} \left[\frac{1}{d} \sum_{i=1}^d \frac{\lambda^2}{(\sigma_i^{(n_1)} + \lambda)^2} \right] \stackrel{(i)}{\geq} \mathbb{E} \left[\frac{\lambda^2}{\left(\sum_{i=1}^d \sigma_i^{(n_1)} / d + \lambda \right)^2} \right] \\ & \stackrel{(ii)}{\geq} \frac{\lambda^2}{\left(\mathbb{E} \left[\sum_{i=1}^d \sigma_i^{(n_1)} / d \right] + \lambda \right)^2} \stackrel{(iii)}{=} \frac{\lambda^2}{(1 + \lambda)^2}, \end{aligned}$$

640 where (i) follows from the convexity of $t \mapsto \lambda^2/(t + \lambda)^2$ on $t \geq 0$; (ii) follows from the
 641 same convexity and Jensen's inequality, and (iii) is since $\mathbb{E} \left[\sum_{i=1}^d \sigma_i^{(n_1)} \right] = \mathbb{E} \left[\text{tr} \left(\frac{1}{n_1} \mathbf{X}_t^\top \mathbf{X}_t \right) \right] =$
 642 $\mathbb{E} \left[\|\mathbf{X}_t\|_{F_r}^2 / n_1 \right] = d$. Applying this in the preceeding bound yields

$$\text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}(n_1, n_2; \lambda)) \leq \frac{d + n_2 + 1}{n_2} \cdot \frac{(1 + \lambda)^2}{\lambda^2}.$$

643 Further plugging in $n_1 = ns$ and $n_2 = n(1 - s)$ for any $s \in (0, 1)$ yields

$$\lim_{d, n \rightarrow \infty, d/n \rightarrow \gamma} \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}(ns, n(1 - s); \lambda)) \leq \frac{\gamma + 1 - s}{1 - s} \cdot \frac{(1 + \lambda)^2}{\lambda^2}.$$

644 Finally, the right-hand side is minimized at $\lambda \rightarrow \infty$ and $s = 0$, from which we conclude that

$$\inf_{\lambda > 0, s \in (0, 1)} \lim_{d, n \rightarrow \infty, d/n \rightarrow \gamma} \text{AsymMSE}(\hat{\mathbf{w}}_{0,T}^{\text{tr-val}}(n_1, n_2; \lambda)) \leq 1 + \gamma,$$

645 which is the desired “ \leq ” direction. □