
Revisiting the Bias-Variance Tradeoff of Multi-Task Learning in High Dimensions

Anonymous Author(s)

Affiliation

Address

email

Abstract

Multi-task learning is a powerful approach in many applications such as image and text classification. Yet, there is little rigorous understanding of when multi-task learning outperforms single-task learning. In this work, we provide a rigorous study of this question in the setting of high-dimensional linear regression. We show that the bias-variance tradeoff of multi-task learning determines the effect of transfer and develop new concentration bounds to analyze the tradeoff. The high-dimensional linear regression setting allows us to define three properties that measure the difference between task data, including *task similarity*, *sample size*, and *covariate shift*. We relate each property to the bias and variance of multi-task learning to explain three negative effects as a result of decreased task similarity, increased source sample size, and covariate shift under increased source sample size. We validate the three effects on text classification tasks. Inspired by our theory, we show two practical connections of interest. First, single-task performance can help understand multi-task performance. Second, incrementally adding training data can mitigate negative transfer and improve multi-task training efficiency.

1 Introduction

Multi-task learning is a powerful approach for improving performance on prediction tasks in computer vision [1, 2], natural language processing [3, 4], and many other areas [5]. In many settings, multiple source tasks are available for predicting a particular target task. The performance of multi-task learning depends on the relationship between the source and target tasks [6]. When the sources are relatively different from the target, multi-task learning has often been observed to perform worse than single-task learning [7, 8], which is referred to as *negative transfer* [9]. While many empirical approaches have been proposed to mitigate negative transfer [5], a precise understanding of when negative transfer occurs remains elusive in the literature [10].

Understanding negative transfer requires developing generalization bounds that scale tightly with properties of each task data such as its sample size. This presents a technical challenge in the multi-task setting because of the difference between task features, even for two tasks. For Rademacher complexity or VC-based techniques, the generalization error scales down as the sample sizes of all tasks increase, when applied to the multi-task setting [11, 12, 13, 14, 15]. Without a tight lower bound for multi-task learning, comparing its performance to single-task learning results in vacuous bounds. From a practical standpoint, developing a better understanding of multi-task learning in terms of properties of task data can provide guidance for downstream applications [16].

In this work, we study the bias and variance of multi-task learning in the high-dimensional linear regression setting [17, 18]. Our key observation is that three properties of task data, including *task similarity*, *sample size*, and *covariate shift* can affect whether multi-task learning outperforms single-task learning (which we refer to as *positive transfer*). As an example, we vary each property in Figure 1 for two linear regression tasks and measure the loss of single-task learning minus multi-task learning for predicting the second task. We observe that the effect of transfer can be positively or negatively affected as we vary each property. Moreover, these phenomena cannot be explained using previous techniques [15]. The high-dimensional linear regression setting allows us to measure the

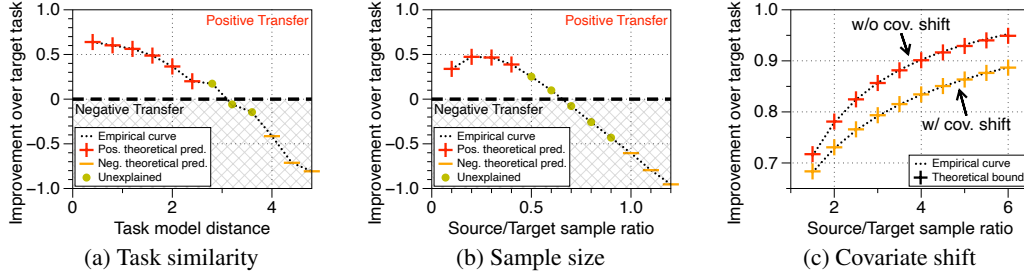


Figure 1: We observe a transition from positive to negative transfer as (a) *task model distance* increases and (b) *source/target sample ratio* increases. For the special case of having the same task model, we observe in (c) that as *source/target sample ratio* increases, having *covariate shift* worsens the performance of MTL. The *y*-axis measures the loss of STL minus MTL.

three properties precisely. We define each property for two tasks and our definition applies to general settings. We refer to the first task as the **source** and the second as the **target**.

- **Task similarity:** Assume that both tasks follow a linear model with parameters $\beta_1, \beta_2 \in \mathbb{R}^p$, respectively. We measure the distance between them by $\|\beta_1 - \beta_2\|$.
- **Sample size:** Let $n_1 = \rho_1 \cdot p, n_2 = \rho_2 \cdot p$ be the sample size of each task, where $\rho_1, \rho_2 > 1$ are both fixed values that do not grow with p . We measure the source/target sample ratio by ρ_1/ρ_2 .
- **Covariate shift:** Assume that the task features are random vectors with positive semidefinite covariance matrix $\Sigma_1, \Sigma_2 \in \mathbb{R}^{p \times p}$, respectively. We define covariate shift as the matrix $\Sigma_1^{1/2} \Sigma_2^{-1/2}$.

We consider a multi-task estimator obtained using a shared linear layer for all tasks and a separate output layer for each task [15]. This two-layer model is inspired by a commonly used idea of hard parameter sharing in multi-task learning [10, 22]. We consider the bias and variance of the multi-task estimator for predicting a **particular** task and compare its performance to single-task learning.

Main results. First, we develop tight bounds for the bias and variance of the multi-task estimator for two tasks by applying recent development from the random matrix theory literature [19, 20, 21]. We observe that the variance of the multi-task estimator is *always smaller* than single-task learning, because of additional **source task samples**. On the other hand, the bias of the multi-task estimator is *always larger* than single-task learning, due to model differences. **Hence**, the tradeoff between bias and variance determines whether the transfer is positive or negative. We provide a sharp analysis of the variance that scales only with **sample size and covariate shift** in Lemma 3.1 and extend it to the bias. Combining both, we analyze the bias-variance tradeoff for two tasks in Theorem 3.2 and extend the analysis to many tasks with the same features in Theorem 3.6.

Second, we explain the phenomena in Figure 1 in isotropic and covariate shifted settings.

- We provide conditions to predict the effect of transfer as a parameter of **model distance** (Section 3.2). As task model distance increases, the bias becomes larger, resulting in negative transfer.
- We provide conditions to predict transfer as a parameter of sample ratio n_1/n_2 (Section 3.3). Adding source task samples helps initially by reducing variance, but hurts eventually due to bias.
- We study a special case of having the same models and show that MTL performs best without covariate shift (Section 3.5). The bias is zero but the variance reduces less due to covariate shift.

Along the way, we analyze the benefit of MTL for reducing labeled data to achieve comparable performance to STL, which has been empirically observed in Taskonomy by Zamir et al. [2].

Our study also leads to several algorithmic consequences with practical interest. First, we show that single-task learning results can help predict positive or negative transfer for multi-task learning. We validate this observation on ChestX-ray14 [1] and sentiment analysis datasets [23]. **Second, we propose to train multi-task models by incrementally adding source task data for predicting a target task.** This is inspired by our observation in Figure 1b where adding more source task data helps initially, but hurts eventually. Using our incremental training schedule, we reduce the computational cost by 65% compared to baseline multi-task training over six sentiment analysis **dataset** while keeping the accuracy the same. Third, we provide a fine-grained insight on a covariance alignment procedure proposed in [15]. We show that the alignment procedure provides more significant improvement when the source/target sample ratio is large. Finally, we validate our three theoretical findings on sentiment analysis tasks.

2 Problem Formulation for Multi-Task Learning

We begin by defining our problem setup including the multi-task estimator we study. Then, we describe the bias-variance tradeoff of the multi-task estimator and connect the bias and variance of the estimator to *task similarity*, *sample size*, and *covariate shift*.

Problem setup. Suppose we have t datasets, where t is a fixed value that does not grow with the feature dimension p . In the high-dimensional linear regression setting (e.g. [17, 18]), the features of the k -th task, denoted by $X_k \in \mathbb{R}^{n_k \times p}$, consists of n_k feature vectors given by x_1, x_2, \dots, x_{n_k} . And each feature $x_i = \Sigma^{1/2} z_i$, where $z_i \in \mathbb{R}^p$ consists of i.i.d. entries with mean zero and unit variance. The sample size n_k equals $\rho_k \cdot p$ for a fixed value ρ_k . The labels $Y_k = X_k \beta_k + \varepsilon_k$, where β_k denotes the linear model parameters and ε_k denotes i.i.d. noise with mean zero and variance σ^2 .

We focus on the commonly used hard parameter sharing model for multi-task learning [10]. When specialized to the linear regression setting, the model consists of a linear layer $B \in \mathbb{R}^{p \times r}$ that is shared by all tasks and t output layers W_1, \dots, W_t that are in \mathbb{R}^r . The width of B , denoted by r , plays an important regularization effect. As observed in Proposition 1 of [15], if $r \geq t$, there is no regularization effect. Hence, we assume that $r < t$ in our study. For example, when there are only two tasks, $r = 1$ and B reduces a vector whereas W_1, W_2 become scalars. We study the following procedure inspired by how hard parameter sharing models are trained in practice (e.g. [22]).

- Separate each dataset (X_i, Y_i) randomly into a training set (X_i^{tr}, Y_i^{tr}) and a validation set (X_i^{val}, Y_i^{val}) . The size of each set is described below.
- Learning the shared layer B : minimize the training loss over B and W_1, \dots, W_t , leading to a closed form equation for \hat{B} that depends on W_1, \dots, W_t .

$$f(B; W_1, \dots, W_t) = \sum_{k=1}^t \|X_k^{tr} B W_k - Y_k^{tr}\|^2. \quad (2.1)$$

- Tuning the output layers W_i : set $B = \hat{B}$ and minimize the validation loss over W_1, \dots, W_t .

$$g(W_1, \dots, W_t) = \sum_{k=1}^t \|X_k^{val} \hat{B} W_k - Y_k^{val}\|^2. \quad (2.2)$$

We make several remarks. In general, the objective $f(\cdot)$ is non-convex in B and the W_k 's. Therefore, we first minimize B in equation (2.1) and then minimize W_k given B in equation (2.2). For our purpose, a validation set of size $\rho_i \cdot p^{0.99}$ that is much larger than the number of output layer parameters $r \cdot t$ suffices. The size of the training set is then $\rho_i(p - p^{0.99})$. The advantage of tuning the output layers on the validation set is to reduce the effect of noise from \hat{B} .

Problem statement. We focus on predicting a particular task, say the t -th task without loss of generality. Let $\hat{\beta}_t^{\text{MTL}}$ denote the multi-task estimator obtained from the procedure above. Our goal is to compare the prediction loss of $\hat{\beta}_t^{\text{MTL}}$, defined by

$$L(\hat{\beta}_t^{\text{MTL}}) = \mathbb{E}_{\{\varepsilon_i\}_{i=1}^t} \mathbb{E}_{x \sim \Sigma^{1/2} z} \left[(x^\top \hat{\beta} - x^\top \beta_t)^2 \right] = \mathbb{E}_{\{\varepsilon_i\}_i^t} \left\| \Sigma_2^{1/2} (\hat{\beta}_t^{\text{MTL}} - \beta_t) \right\|^2,$$

to the prediction loss $L(\hat{\beta}_t^{\text{STL}})$ of the single-task estimator $\hat{\beta}_t^{\text{STL}} = (X_t^\top X_t)^{-1} X_t^\top Y_t$. We say there is negative transfer if $L(\hat{\beta}_t^{\text{MTL}}) > L(\hat{\beta}_t^{\text{STL}})$, or positive transfer otherwise.

As an example, for the setting of two tasks, we can decompose $L(\hat{\beta}_t^{\text{MTL}}) - L(\hat{\beta}_t^{\text{STL}})$ into a bias term and a variance term as follows (derived in Appendix A).

$$L(\hat{\beta}_t^{\text{MTL}}) - L(\hat{\beta}_t^{\text{STL}}) = \hat{v}^2 \left\| \Sigma_2^{1/2} (\hat{v}^2 X_1^\top X_1 + X_2^\top X_2)^{-1} X_1^\top X_1 (\beta_1 - \hat{v} \beta_2) \right\|^2 \quad (2.3)$$

$$+ \sigma^2 \left(\text{Tr} \left[(\hat{v}^2 X_1^\top X_1 + X_2^\top X_2)^{-1} \Sigma_2 \right] - \text{Tr} \left[(X_2^\top X_2)^{-1} \Sigma_2 \right] \right). \quad (2.4)$$

In the above, $\hat{v} = W_1/W_2$ where W_1, W_2 are obtained from solving equation (2.2) (recalling that W_1, W_2 are scalars for two tasks). The role of \hat{v} is to scale the shared subspace B to fit each task.

Equation (2.3) corresponds to the bias of $\hat{\beta}_t^{\text{MTL}}$. Hence, the bias term introduces a negative effect that depends on the *similarity* between β_1 and β_2 . Equation (2.4) corresponds to the variance of $\hat{\beta}_t^{\text{MTL}}$ minus the variance of $\hat{\beta}_t^{\text{STL}}$, which is always negative. Intuitively, the more *samples* we have, the smaller the variance is. Meanwhile, *covariate shift* also affects how small the variance can be.

3 Comparing Multi-Task Learning to Single-Task Learning

We provide tight bounds on the bias and variance of the multi-task estimator for two tasks. We show theoretical implications for understanding the performance of multi-task learning. (a) *Task similarity*: we explain the phenomenon of negative transfer precisely as task models become different. (b) *Sample size*: we further explain a curious phenomenon where increasing the source sample size helps initially, but hurts eventually. (c) *Covariate shift*: as the source sample size increases, we show that the covariate shift worsens the performance of the multi-task estimator. Finally, we extend our results from two tasks to many tasks with the same features.

3.1 Analyzing the Bias-Variance Tradeoff using Random Matrix Theory

A well-known result in the high-dimensional linear regression setting states that $\text{Tr}[(X_2^\top X_2)^{-1} \Sigma_2]$ is concentrated around $1/(\rho_2 - 1)$ (e.g. Chapter 6 of [24]), which scales with the sample size of the target task. Our main technical contribution is to extend this result to two tasks. We show how the variance of the multi-task estimator scales with sample size and covariate shift in the following result.

Lemma 3.1 (Variance bound). *In the setting of two tasks, let $n_1 = \rho_1 \cdot p$ and $n_2 = \rho_2 \cdot p$ be the sample size of the two tasks. Let $\lambda_1, \dots, \lambda_p$ be the singular values of the covariate shift matrix $\Sigma_1^{1/2} \Sigma_2^{-1/2}$ in decreasing order. With high probability, the variance of the multi-task estimator $\hat{\beta}_t^{\text{MTL}}$ equals*

$$\frac{\sigma^2}{n_1 + n_2} \cdot \text{Tr} \left[(\hat{v}^2 a_1 \Sigma_2^{-1/2} \Sigma_1 \Sigma_2^{-1/2} + a_2 \text{Id})^{-1} \right] + O \left(p^{-1/2+o(1)} \right),$$

where a_1, a_2 are solutions of the following equations:

$$a_1 + a_2 = 1 - \frac{1}{\rho_1 + \rho_2}, \quad a_1 + \frac{1}{\rho_1 + \rho_2} \cdot \frac{1}{p} \sum_{i=1}^p \frac{\hat{v} \lambda_i^2 a_1}{\hat{v} \lambda_i^2 a_1 + a_2} = \frac{\rho_1}{\rho_1 + \rho_2}.$$

Lemma 3.1 allows us to get a tight bound on equation (2.4), that only depends on *sample size*, *covariate shift* and the scalar \hat{v} . As a remark, the concentration error $O(p^{-1/2+o(1)})$ of our result is nearly optimal. For the bias term of equation (2.3), a similar result that scales with task model distance in addition to sample size and covariate shift holds (cf. Lemma B.3 in Appendix B). Combining the two lemmas, we provide a sharp analysis of the bias-variance tradeoff of the multi-task estimator. For a matrix X , let $\lambda_{\min}(X)$ denote its smallest singular value and $\|X\|$ denote its spectral norm.

Theorem 3.2 (Two tasks). *For the setting of two tasks, let $\delta > 0$ be a fixed error margin, $\rho_2 > 1$ and $\rho_1 \gtrsim \delta^{-2} \cdot \lambda_{\min}(\Sigma_1^{1/2} \Sigma_2^{-1/2})^{-4} \|\Sigma_1\| \max(\|\beta_1\|^2, \|\beta_2\|^2)$, and . There exists two deterministic functions Δ_{bias} and Δ_{var} that only depend on $\{\hat{v}, \Sigma_1, \Sigma_2, \rho_1, \rho_2, \beta_1, \beta_2\}$ such that*

- If $\Delta_{\text{bias}} - \Delta_{\text{var}} < -\delta$, then w.h.p. over the randomness of X_1, X_2 , we have $L(\hat{\beta}_t^{\text{MTL}}) < L(\hat{\beta}_t^{\text{STL}})$.
- If $\Delta_{\text{bias}} - \Delta_{\text{var}} > \delta$, then w.h.p. over the randomness of X_1, X_2 , we have $L(\hat{\beta}_t^{\text{MTL}}) > L(\hat{\beta}_t^{\text{STL}})$.

Theorem 3.2 applies to settings where large amounts of source task data is available but the target sample size is small. For such settings, we obtain a sharp transition from positive transfer to negative transfer determined by $\Delta_{\text{bias}} - \Delta_{\text{var}}$. While the general form of these functions can be complex (as are previous generalization bounds for MTL), they admit interpretable forms for simplified settings. The proof of Theorem 3.2 is presented in Appendix B and the proof of Lemma 3.1 is in Appendix E.

3.2 Task Similarity

It is well-known since the seminal work of Caruana [6] that how well multi-task learning performs depends on task relatedness. We formalize this connection in the following simplified setting, where we can perform explicit calculations. We show that as we increase the distance between β_1 and β_2 , there is a transition from positive transfer to negative transfer in MTL.

The isotropic model. Consider two tasks with isotropic covariances $\Sigma_1 = \Sigma_2 = \text{Id}$. Each task has sample size $n_1 = \rho_1 \cdot p$ and $n_2 \rho_2 \cdot p$. Assume that for the target task, β_2 has i.i.d. entries with mean zero and variance κ^2 . For the source task, β_1 equals β_2 plus i.i.d. entries with mean 0 and variance d^2 . The labels are $Y_i = X_i \beta_i + \varepsilon_i$, where ε_i consists of i.i.d. entries with mean zero and variance σ^2 . For our purpose, it is enough to think of the order of d being $1/\sqrt{p}$ and $p d^2 / \sigma^2$ being constant.

We introduce the following notations.

$$\Psi(\beta_1, \beta_2) = \mathbb{E} [\|\beta_1 - \beta_2\|^2] / \sigma^2, \quad \Phi(\rho_1, \rho_2) = \frac{(\rho_1 + \rho_2 - 1)^2}{\rho_1(\rho_1 + \rho_2)(\rho_2 - 1)}.$$

168

169 **Proposition 3.3** (Task model distance). *In the isotropic model, suppose that ρ_1 and $\rho_2 > 1$. Then*

- 170 • *If $\Psi(\beta_1, \beta_2) < \frac{1}{\nu} \cdot \Phi(\rho_1, \rho_2)$, then w.h.p. over the randomness of X_1, X_2 , $L(\hat{\beta}_t^{MTL}) < L(\hat{\beta}_t^{STL})$.*
- 171 • *If $\Psi(\beta_1, \beta_2) > \nu \cdot \Phi(\rho_1, \rho_2)$, then w.h.p. over the randomness of X_1, X_2 , $L(\hat{\beta}_t^{MTL}) > L(\hat{\beta}_t^{STL})$.*

172 *Here $\nu = (1 - o(1)) \min((1 - 1/\sqrt{\rho_1})^{-4}, (1 + 1/\sqrt{\rho_1})^4)$. Concretely, if $\rho_1 > 40$, then $\nu \in (1, 2)$.*

173 Proposition 3.3 simplifies Theorem 3.2 in the isotropic model, allowing for a more explicit statement
 174 of the bias-variance tradeoff. Concretely, $\Psi(\beta_1, \beta)$ and $\Phi(\rho_1, \rho_2)$ corresponds to Δ_{bias} and Δ_{var} ,
 175 respectively. Roughly speaking, the transition threshold scales as $\frac{pd^2}{\sigma^2} - \frac{1}{\rho_1} - \frac{1}{\rho_2}$. We apply Proposition
 176 3.3 to the parameter setting of Figure 1a (the details are left to Appendix F.1). We can see that our
 177 result is able to predict positive or negative transfer accurately that matches the empirical curve.
 178 There are several unexplained observations near the transition threshold 0, which are caused by the
 179 concentration error ν . The proof of Proposition 3.3 can be found in Appendix C.1. A key part of the
 180 analysis shows that $\hat{\nu} \approx 1$ in the isotropic model, thus simplifying the result of Theorem 3.2.

181 **Algorithmic consequence.** We can in fact extend the result to the cases where the noise variances
 182 are different. In this case, we will see that MTL is particularly effective. Concretely, suppose the
 183 noise variance σ_1^2 of task 1 differs from the noise variance σ_2^2 of task 2. If σ_1^2 is too large, the source
 184 task provides a negative transfer to the target. If σ_1^2 is small, the source task is more helpful. We leave
 185 the result to Proposition C.2 in Appendix C.1. Inspired by the observation, we propose a single-task
 186 based metric to help understand MTL results using STL results.

- 187 • For each task, we train a single-task model. Let z_s and z_t be the prediction accuracy of each task,
 188 respectively. Let $\tau \in (0, 1)$ be a fixed threshold.
- 189 • If $z_s - z_t > \tau$, then we predict that there will be positive transfer when combining the two tasks
 190 using MTL. If $z_s - z_t < -\tau$, then we predict negative transfer.

191

3.3 Sample Size

192 In classical Rademacher or VC based theory of multi-task learning, the generalization bounds are
 193 usually presented for settings where the sample sizes are equal for all tasks [11, 13, 14]. On the other
 194 hand, uneven sample sizes between different tasks (or even dominating tasks) have been empirically
 195 observed as a cause of negative transfer [25]. For such settings, we have also observed that adding
 196 more labeled data from one task does not always help. In the isotropic model, we consider what
 197 happens if we vary the source task sample size. Our theory accurately predicts a curious phenomenon,
 198 where increasing the sample size of the source task results in negative transfer!

199 **Proposition 3.4** (Source/target sample ratio). *In the isotropic model, suppose that $\rho_1 > 40$ and
 200 $\rho_2 > 110$ are fixed constants, and $\Psi(\beta_1, \beta_2) > 2/(\rho_2 - 1)$. Then we have that*

- 201 • *If $\frac{n_1}{n_2} = \frac{\rho_1}{\rho_2} < \frac{1}{\nu} \cdot \frac{1-2\rho_2^{-1}}{\Psi(\beta_1, \beta_2)(\rho_2-1)-\nu^{-1}}$, then w.h.p. $L(\hat{\beta}_t^{MTL}) < L(\hat{\beta}_t^{STL})$.*
- 202 • *If $\frac{n_1}{n_2} = \frac{\rho_1}{\rho_2} > \nu \cdot \frac{1-2\rho_2^{-1}}{\Psi(\beta_1, \beta_2)(\rho_2-1.5)-\nu}$, then w.h.p. $L(\hat{\beta}_t^{MTL}) > L(\hat{\beta}_t^{STL})$.*

203 Proposition 3.4 describes the bias-variance tradeoff in terms of the sample ratio n_1/n_2 . We apply the
 204 result to the parameter setting of Figure 1b (described in Appendix F.1). There are several unexplained
 205 observations near $y = 0$ caused by ν . The proof of Proposition 3.4 can be found in Appendix C.2.

206 **Connection to Taskonomy.** We use our tools to explain a key result of Taskonomy by Zamir et
 207 al. (2018) [2], which shows that MTL can reduce the amount of labeled data needed to achieve
 208 comparable performance to STL. For $i = 1, 2$, let $\hat{\beta}_i^{MTL}(x)$ denote the estimator trained using $x \cdot n_i$
 209 datapoints from every task. The data efficiency ratio is defined as

$$\arg \min_{x \in (0,1)} L_1(\hat{\beta}_1^{MTL}(x)) + L_2(\hat{\beta}_2^{MTL}(x)) \leq L_1(\hat{\beta}_1^{STL}) + L_2(\hat{\beta}_2^{STL}).$$

210 For example, the data efficiency ratio is 1 if there is negative transfer. Using our tools, we show that
 211 in the isotropic model, the data efficiency ratio is roughly

$$\frac{1}{\rho_1 + \rho_2} + \frac{2}{(\rho_1 + \rho_2)(\rho_1^{-1} + \rho_2^{-1} - \Theta(\Psi(\beta_1, \beta_2)))}.$$

212 Compared with Proposition 3.3, we see that when $\Psi(\beta_1, \beta_2)$ is smaller than $\rho_1^{-1} + \rho_2^{-1}$ (up to a
 213 constant multiple), the transfer is positive. Moreover, the data efficiency ratio quantifies how effective
 214 the positive transfer is using MTL. The result can be found in Proposition C.3 in Appendix C.2.

215 **Algorithmic consequence.** An interesting consequence of Proposition 3.4 is that $L(\hat{\beta}_t^{\text{MTL}})$ is not
 216 monotone in ρ_1 . In particular, Figure 1b (and our analysis) shows that $L(\hat{\beta}_t^{\text{MTL}})$ behaves as a
 217 quadratic function over ρ_1 . More generally, depending on how large $\Psi(\beta_1, \beta_2)$ is, $L(\hat{\beta}_t^{\text{MTL}})$ may
 218 also be monotonically increasing or decreasing. Based on this insight, we propose an incremental
 219 optimization schedule to improve MTL training efficiency.

- 220 • We divide the source task data into S batches. For S rounds, we incrementally add the source
 221 task data by adding one batch at a time.
- 222 • After training T epochs, if the validation accuracy becomes worse than the previous round's
 223 result, we terminate. Algorithm 1 in Appendix F describes the procedure in detail.

224 3.4 Covariate Shift

225 So far we have considered the isotropic model where $\Sigma_1 = \Sigma_2$. This setting is relevant for settings
 226 where different tasks share the same input features such as multi-class image classification. In general,
 227 the covariance matrices of the two tasks may be different such as in text classification. In this part,
 228 we consider what happens when $\Sigma_1 \neq \Sigma_2$. We show that when n_1/n_2 is large, MTL with covariate
 229 shift can be suboptimal compared to MTL without covariate shift.

230 *Example.* We measure covariate shift by $M = \Sigma_1^{1/2} \Sigma_2^{-1/2}$. Assume that $\Psi(\beta_1, \beta_2) = 0$ for simplicity.
 231 We compare two cases: (i) when $M = \text{Id}$; (ii) when M has $p/2$ singular values that are equal to λ
 232 and $p/2$ singular values that are equal to $1/\lambda$. Hence, λ measures the severity of the covariate shift.
 233 Figure 1c shows a simulation of this setting by varying λ . We observe that as source/target sample
 234 ratio increases, the performance gap between the two cases increases.

235 We compare different choices of M that belong to the following bounded set. Let λ_i be the i -th
 236 singular value of M . Let $\mu_{\min} < \mu < \mu_{\max}$ be fixed values that do not grow with p .

$$\mathcal{S}_\mu := \left\{ M \left| \prod_{i=1}^p \lambda_i \leq \mu^p, \mu_{\min} \leq \lambda_i \leq \mu_{\max}, \text{ for all } 1 \leq i \leq p \right. \right\},$$

237 **Proposition 3.5** (Covariate shift). *Assume that $\Psi(\beta_1, \beta_2) = 0$ and $\rho_1, \rho_2 > 1$. Let $g(M)$ denote the*
 238 *prediction loss of $\hat{\beta}_t^{\text{MTL}}$ when $M = \Sigma_1^{1/2} \Sigma_2^{-1/2} \in \mathcal{S}_\mu$. We have that*

$$g(\mu \text{Id}) \leq (1 + O(\rho_2/\rho_1)) \min_{M \in \mathcal{S}_\mu} g(M).$$

239 This proposition shows that when source/target sample ratio is large, then having no covariate shift is
 240 optimal. The proof of Proposition 3.5 is left to Appendix C.3.

241 **Algorithmic consequence.** Our observation highlights the need to correct covariate shift when
 242 n_1/n_2 is large. Hence for such settings, we expect procedures that aim at correcting covariate shift to
 243 provide more significant gains. We consider a covariance alignment procedure proposed in Wu et
 244 al. (2020) [15], which is designed for the purpose of correcting covariate shift. The idea is to add
 245 an alignment module between the input and the shared module B . This new module is then trained
 246 together with B and the output layers. We validate our insight on this procedure in the experiments.

247 3.5 Extensions

248 Next, we describe our result for more than two tasks with same features, i.e. $X_i = X$ for any i . This
 249 setting is prevalent in applications of multi-task learning to image classification, where there are
 250 multiple prediction labels/tasks for every image [1, 26].

251 **Theorem 3.6** (Many tasks). *For the setting of t tasks where $X_i = X$, for all $1 \leq i \leq t$. Let*
 252 *$B^* := [\beta_1, \beta_2, \dots, \beta_t]$ and $U_r \in \mathbb{R}^{t \times r}$ denote the linear model parameters. Let $U_r U_r^\top$ denote the*
 253 *best rank- r subspace approximation of $(B^*)^\top \Sigma B^*$. Assume that $\lambda_{\min}(B^{*\top} \Sigma B^*) \gtrsim \sigma^2$. Let v_i*
 254 *denote the i -th row vector of U_r . There exists a value $\delta = o(\|B^*\|^2 + \sigma^2)$ such that*

- 255 • If $(1 - \|v_t\|^2) \frac{\sigma^2}{\rho-1} - \|\Sigma(B^* U_r v_t - \beta_t)\|^2 > \delta$, then w.h.p $L(\hat{\beta}_t^{\text{MTL}}) < L(\hat{\beta}_t^{\text{STL}})$.
- 256 • If $(1 - \|v_t\|^2) \frac{\sigma^2}{\rho-1} - \|\Sigma(B^* U_r v_t - \beta_t)\|^2 < -\delta$, then w.h.p. $L(\hat{\beta}_t^{\text{MTL}}) > L(\hat{\beta}_t^{\text{STL}})$.

257 Theorem 3.6 provides a sharp analysis of the bias-variance tradeoff beyond two tasks. Specially,
 258 $(1 - \|v_t\|^2) \sigma^2 / (\rho - 1)$ shows the amount of reduced variance and $\|\Sigma(B^* U_r v_t - \beta_t)\|^2$ shows the bias
 259 of the multi-task estimator. The proof of 3.6 can be found in Appendix D.

4 Experiments

We validate our theory and algorithmic insights. First, we validate the single-task based metric on sentiment analysis and ChestX-ray14 datasets. We show that single-task learning results can help predict positive or negative transfer for both datasets. Second, our proposed incremental training schedule improves the training efficiency of standard multi-task training on sentiment analysis tasks. Third, when the sample ratio is large, performing the alignment procedure of [15] provides more improvement for MTL. Finally, we validate our theoretical results on text classification tasks.

4.1 Experimental Setup

We consider a text classification task and an image classification task as follows.

Sentiment Analysis. We consider six tasks: movie review sentiment (MR), sentence subjectivity (SUBJ), customer reviews polarity (CR), question type (TREC), opinion polarity (MPQA), and the Stanford sentiment treebank (SST) tasks. The question is to predict positive or negative sentiment expressed in the text. We use an embedding layer with GloVe embeddings followed by an LSTM, MLP or CNN layer proposed by [27].

ChestX-ray14. This dataset contains 112,120 frontal-view X-ray images. There are 14 diseases (tasks) for every image that we would like to predict. We use densenet121 as the shared module.

For all models, we share the main module across all tasks and assign a separate regression or classification layer on top of the shared module for each tasks. The baseline training schedule for MTL is the round-robin training schedule. We measure the test accuracy of predicting a target task.

4.2 Experimental Results

Predicting transfer effect via STL results. We show that the single-task based metric proposed in Section 3.2 can predict positive or negative transfer in MTL. A common challenge in the study of multi-task learning is that the results can be hard to understand. It is difficult to predict when MTL performs well without running extensive trials. Our insight is that we can use STL results to help understand MTL results. Table 2 shows the result on both the sentiment analysis and the ChestX-ray14 tasks. We find that using a threshold of $\tau = 0.1$, the STL results correctly predict positive or negative transfer with 75.6% accuracy and 38.8% recall among 30 times 5 (random seeds) task pairs! We observe similar results for 91 task pairs from the ChestX-ray14 dataset.

Mitigating negative transfer via incremental training. First, we show that our proposed incremental training schedule (Algorithm 1) can help mitigate negative transfer for predicting a particular target task. Our insight is that since adding more samples from the source task does not always help, we can improve efficiency by adding source samples *incrementally* during training. Over six randomly selected pairs from the sentiment analysis tasks, we find that Algorithm 1 requires only 45% of the computational cost to achieve similar performance on the target task, compared to the MTL baseline. Our next result shows the incremental training schedule applies to multiple tasks as well. In Table 1, we find that over all six sentiment analysis tasks, incremental training requires less than 35% of the computational cost compared to baseline MTL training, while achieving the same accuracy averaged over all six tasks. As a further validation, excluding TREC, we observe similar comparative results.

4.3 Validating the Theoretical Results

We validate our three theoretical results in Section 3 on the sentiment analysis tasks. In Figure 2a, we compare the performance training with a semantically similar task versus a dissimilar task with a target task. We select each task pair based on our domain knowledge. We observe that adding a similar task helps the target task whereas adding a dissimilar task hurts. In Figure 2b, we validate

Models	Sentiment analysis	
	all tasks	w/o TREC
MLP	31%	29%
LSTM	35%	34%
CNN	30%	28%

Table 1: Efficiency of incremental training compared to baseline MTL.

Threshold	Sentiment analysis		ChestX-ray14	
	Precision	Recall	Precision	Recall
0.0	0.596	1.000	0.593	1.000
0.1	0.756	0.388	0.738	0.462
0.2	0.919	0.065	0.875	0.044

Table 2: Single-task learning results can help predict positive or negative transfer in multi-task learning.

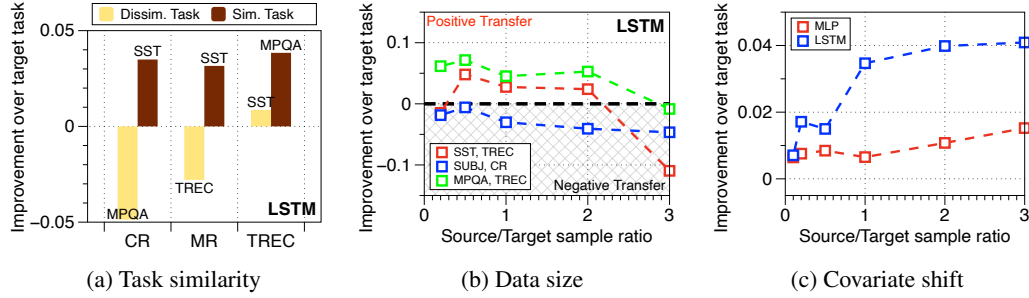


Figure 2: Validating the three results of Section 3 on sentiment analysis tasks. (a) Adding a semantically similar source task in MTL performs better than adding a dissimilar task. (b) As source/target sample ratio increases, we observe a transition from positive to negative transfer. (c) As source/target sample ratio increases, aligning task covariances [15] improves more over the baseline.

the phenomenon that adding more source data samples does not always improve performance on the target task. In Figure 2c, we measure the performance gains from performing the alignment procedure proposed in [15] minus baseline MTL performance. We average the results over all 15 task pairs. The result shows that as the source samples increases, the alignment procedure shows a bigger improvement over MTL. The rest of experimental procedures are left to Appendix F.2.

5 Related Work

We refer the interested readers to several excellent surveys on multi-task learning for a comprehensive survey [9, 10, 5, 28]. Below, we describe several lines of work that are most related to this work.

Multi-task learning theory. Some of the earliest works on multi-task learning are Baxter [11], Ben-David and Schuller [29]. Maurer [13] studies generalization bounds for linear separation settings of MTL. Ben-David et al. [30] provides uniform convergence bounds that combines source and target errors in an optimal way. The benefit of learning multi-task representations is studied for learning certain half-spaces [14] and sparse regression [31, 32]. Our work is closely related to Wu et al. [15]. While Wu et al. provide generalization bounds to show that adding more labeled helps learn the target task more accurately, their techniques do not explain the phenomena of negative transfer.

Multi-task learning methodology. Ando and Zhang [12] introduces an alternating minimization framework for learning multiple tasks. Argyriou et al. [33] present a convex algorithm which learns common sparse representations across a pool of related tasks. Evgeniou et al. [34] develop a framework for multi-task learning in the context of kernel methods. The multi-task learning model that we have focused on is based on the idea of hard parameter sharing [35, 36, 10]. We believe that the technical tools we have developed can also be applied to many other multi-task learning models.

Random matrix theory. The random matrix theory tool and related proof of our work fall into a paradigm of the so-called local law of random matrices [19]. For a sample covariance matrix $X^\top X$ with $\Sigma = \text{Id}$, such a local law was proved in [20]. It was later extended to sample covariance matrices with non-identity Σ [21], and separable covariance matrices [37]. On the other hand, one may derive the asymptotic result in Theorem 3.1 with error $o(1)$ using the free addition of two independent random matrices in free probability theory [38]. To the best of my knowledge, we do not find an *explicit result* for the sum of two sample covariance matrices with general covariates in the literature.

6 Conclusions and Open Problems

In this work, we analyzed the bias and variance of multi-task learning versus single-task learning. We provided tight concentration bounds for the bias and the variance. Based on these bounds, we analyzed the impact of three properties, including task similarity, sample size, and covariate shift on the bias and variance, to derive conditions for transfer. We validated our theoretical results. Based on the theory, we proposed to train multi-task models by incrementally adding labeled data and showed encouraging results inspired by our theory. We describe several open questions for future work. First, our bound on the bias term (cf. Lemma B.3) involves an error term that scales down with ρ_1 . Tightening this error bound can potentially cover the unexplained observations in Figure 1. Second, it would be interesting to extend our results to non-linear settings. We remark that this likely requires addressing significant technical challenges to deal with non-linearity.

Broader Impacts

In this work, we provide a theoretical study to help understand when multi-task learning performs well. We approach this question by studying the bias-variance tradeoff of multi-task learning. We provide new technical tools to bound the bias and variance. We relate the bounds to three properties of task data. Overall, our theoretical study provides a framework to help understand multi-task learning performance. We further provide guidance for detecting and mitigating negative transfer on image and text classification tasks.

Our theoretical framework has the potential to impact many other neighboring areas in the ML community. Multi-task learning connects to a wide range of areas [28]. To name a few, transfer learning, meta learning, multimodal learning, semi-supervised learning, and representation learning are all closely related areas to multi-task learning. Any learning scenario such as reinforcement learning [25] that combines multiple datasets to supervise a model is using multi-task learning. While the theoretical results that we have provided are not directly applicable to these different settings, we believe that the tools we have developed and the framework we have provided can inspire followup works in different settings. For one specific example, we have developed new concentration bounds that may apply to many settings such as soft parameter sharing [10], kernel methods [34], and convex formulation of multi-task learning [39]. For another example, our results also allow extensions to transfer learning and domain adaptation [40]. The insights we have developed on positive and negative transfer can potentially find applications in multimodal learning, where the data sources are usually heterogeneous. Our fine-grained study on sample sizes have the potential to provide new insight in meta learning, where scarce labeled samples presents a significant challenge.

Our algorithmic consequences that have from our theory have the potential to impact downstream applications of multi-task learning. For example, many medical applications use multi-task learning to train large-scale image classification models by combining multiple datasets [1, 26]. Unlike the applications of multi-task learning in text classification where large amounts of labeled data are collected [3], in medical applications it is typically difficult to acquire large amounts of labeled data. For such settings, training multi-task models can be very challenging. Our insight on using single-task learning results to help understand multi-task learning can be valuable for helping practitioners understand their results.

References

- [1] Pranav Rajpurkar, Jeremy Irvin, Kaylie Zhu, Brandon Yang, Hershel Mehta, Tony Duan, Daisy Ding, Aarti Bagul, Curtis Langlotz, Katie Shpanskaya, et al. Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning. *arXiv preprint arXiv:1711.05225*, 2017.
- [2] Amir R Zamir, Alexander Sax, William Shen, Leonidas J Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Disentangling task transfer learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3712–3722, 2018.
- [3] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, 2018.
- [4] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. Superglue: A stickier benchmark for general-purpose language understanding systems. In *Advances in Neural Information Processing Systems*, pages 3261–3275, 2019.
- [5] Yu Zhang and Qiang Yang. A survey on multi-task learning. *arXiv preprint arXiv:1707.08114*, 2017.
- [6] Rich Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997.
- [7] Héctor Martínez Alonso and Barbara Plank. When is multitask learning effective? semantic sequence prediction under varying data conditions. *arXiv preprint arXiv:1612.02251*, 2016.
- [8] Joachim Bingel and Anders Søgaard. Identifying beneficial task relations for multi-task learning in deep neural networks. *arXiv preprint arXiv:1702.08303*, 2017.

- [9] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2009.
- [10] Sebastian Ruder. An overview of multi-task learning in deep neural networks. *arXiv preprint arXiv:1706.05098*, 2017.
- [11] Jonathan Baxter. A model of inductive bias learning. *Journal of artificial intelligence research*, 12:149–198, 2000.
- [12] Rie Kubota Ando and Tong Zhang. A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6(Nov):1817–1853, 2005.
- [13] Andreas Maurer. Bounds for linear multi-task learning. *Journal of Machine Learning Research*, 7(Jan):117–139, 2006.
- [14] Andreas Maurer, Massimiliano Pontil, and Bernardino Romera-Paredes. The benefit of multitask representation learning. *The Journal of Machine Learning Research*, 17(1):2853–2884, 2016.
- [15] Sen Wu, Hongyang R. Zhang, and Christopher Ré. Understanding and improving information transfer in multi-task learning. In *International Conference on Learning Representations*, 2020.
- [16] Alexander Ratner, Braden Hancock, Jared Dunnmon, Frederic Sala, Shreyash Pandey, and Christopher Ré. Training complex models with multi-task weak supervision. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4763–4771, 2019.
- [17] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. *arXiv preprint arXiv:1903.08560*, 2019.
- [18] Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 2020.
- [19] László Erdos and Horng-Tzer Yau. A dynamical approach to random matrix theory. *Courant Lecture Notes in Mathematics*, 28, 2017.
- [20] A. Bloemendal, L. Erdős, A. Knowles, H.-T. Yau, and J. Yin. Isotropic local laws for sample covariance and generalized Wigner matrices. *Electron. J. Probab.*, 19(33):1–53, 2014.
- [21] Antti Knowles and Jun Yin. Anisotropic local laws for random matrices. *Probability Theory and Related Fields*, pages 1–96, 2016.
- [22] Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. Multi-task deep neural networks for natural language understanding. *arXiv preprint arXiv:1901.11504*, 2019.
- [23] Tao Lei, Yu Zhang, Sida I Wang, Hui Dai, and Yoav Artzi. Simple recurrent units for highly parallelizable recurrence. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4470–4481, 2018.
- [24] Vadim Ivanovich Serdobolskii. *Multiparametric statistics*. Elsevier, 2007.
- [25] Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. Gradient surgery for multi-task learning. *arXiv preprint arXiv:2001.06782*, 2020.
- [26] Sabri Eyuboglu, Geoffrey Angus, Bhavik N. Patel, Anuj Pareek, Guido Davidzon, JaredDunnmon, and Matthew P. Lungren. Multi-task weak supervision enables automated abnormality localization in whole-body fdg-pet/ct. 2020.
- [27] Tao Lei, Yu Zhang, Sida I Wang, Hui Dai, and Yoav Artzi. Simple recurrent units for highly parallelizable recurrence. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4470–4481, 2018.
- [28] Simon Vandenhende, Stamatios Georgoulis, Marc Proesmans, Dengxin Dai, and Luc Van Gool. Revisiting multi-task learning in the deep learning era. *arXiv preprint arXiv:2004.13379*, 2020.
- [29] Shai Ben-David and Reba Schuller. Exploiting task relatedness for multiple task learning. In *Learning Theory and Kernel Machines*, pages 567–580. Springer, 2003.
- [30] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79(1-2):151–175, 2010.
- [31] Karim Lounici, Massimiliano Pontil, Alexandre B Tsybakov, and Sara Van De Geer. Taking advantage of sparsity in multi-task learning. *arXiv preprint arXiv:0903.1468*, 2009.

- [32] Karim Lounici, Massimiliano Pontil, Sara Van De Geer, Alexandre B Tsybakov, et al. Oracle inequalities and optimal inference under group sparsity. *The annals of statistics*, 39(4):2164–2204, 2011.
- [33] Andreas Argyriou, Theodoros Evgeniou, and Massimiliano Pontil. Convex multi-task feature learning. *Machine learning*, 73(3):243–272, 2008.
- [34] Theodoros Evgeniou, Charles A Micchelli, and Massimiliano Pontil. Learning multiple tasks with kernel methods. *Journal of machine learning research*, 6(Apr):615–637, 2005.
- [35] Rich Caruana. Multitask learning: A knowledge-based source of inductive bias. *Proceedings of the Tenth International Conference on Machine Learning*.
- [36] Abhishek Kumar and Hal Daumé III. Learning task grouping and overlap in multi-task learning. In *Proceedings of the 29th International Conference on Machine Learning*, 2012.
- [37] Fan Yang. Edge universality of separable covariance matrices. *Electron. J. Probab.*, 24:57 pp., 2019.
- [38] Alexandru Nica and Roland Speicher. *Lectures on the combinatorics of free probability*, volume 13. Cambridge University Press, 2006.
- [39] Yu Zhang and Dit-Yan Yeung. A regularization approach to learning task relationships in multitask learning. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 8(3):1–31, 2014.
- [40] Wouter M Kouw. An introduction to domain adaptation and transfer learning. *arXiv preprint arXiv:1812.11806*, 2018.
- [41] Theodore W Anderson. *An Introduction to Multivariate Statistical Analysis*. Wiley New York, 2003.
- [42] Zhidong Bai and Jack W Silverstein. *Spectral analysis of large dimensional random matrices*. Springer Series in Statistics. Springer, New York, 2nd edition, 2010.
- [43] Terence Tao. *Topics in Random Matrix Theory*, volume 132. American Mathematical Soc., 2012.
- [44] V. A. Marčenko and L. A. Pastur. Distribution of eigenvalues for some sets of random matrices. *Mathematics of the USSR-Sbornik*, 1:457, 1967.
- [45] Z. D. Bai and Jack W. Silverstein. No eigenvalues outside the support of the limiting spectral distribution of large-dimensional sample covariance matrices. *Ann. Probab.*, 26(1):316–345, 1998.
- [46] Shurong Zheng, Zhidong Bai, and Jianfeng Yao. Clt for eigenvalue statistics of large-dimensional general fisher matrices with applications. *Bernoulli*, 23(2):1130–1178, 2017.
- [47] L. Erdős, A. Knowles, and H.-T. Yau. Averaging fluctuations in resolvents of random band matrices. *Ann. Henri Poincaré*, 14:1837–1926, 2013.
- [48] Johannes Alt, L. Erdős, and Torben Krüger. Local law for random Gram matrices. *Electron. J. Probab.*, 22:41 pp., 2017.
- [49] Haokai Xi, Fan Yang, and Jun Yin. Local circular law for the product of a deterministic matrix with a random matrix. *Electron. J. Probab.*, 22:77 pp., 2017.
- [50] Natesh S. Pillai and Jun Yin. Universality of covariance matrices. *Ann. Appl. Probab.*, 24:935–1001, 2014.
- [51] Xiucui Ding and Fan Yang. A necessary and sufficient condition for edge universality at the largest singular values of covariance matrices. *Ann. Appl. Probab.*, 28(3):1679–1738, 2018.
- [52] L. Erdős, A. Knowles, H.-T. Yau, and J. Yin. Spectral statistics of Erdős-Rényi graphs I: Local semicircle law. *Ann. Probab.*, 41(3B):2279–2375, 2013.

A Missing Details of Problem Formulation

Assumptions on task data generation. First, we give the basic assumption for our main objects—the random matrices $X_i, i = 1, 2$.

Assumption A.1. We will consider $n \times p$ random matrices of the form $X = Z\Sigma^{1/2}$, where Σ is a $p \times p$ deterministic positive definite symmetric matrices, and $Z = (z_{ij})$ is an $n \times p$ random matrix with real i.i.d. entries with mean zero and variance one. Note that the rows of X are i.i.d. centered random vectors with covariance matrix Σ . For simplicity, we assume that all the moments of z_{ij} exists, that is, for any fixed $k \in \mathbb{N}$, there exists a constant $C_k > 0$ such that

$$\mathbb{E}|z_{ij}|^k \leq C_k, \quad 1 \leq i \leq n, \quad 1 \leq j \leq p. \quad (\text{A.1})$$

We assume that $n = \rho p$ for some fixed constant $\rho > 1$. Without loss of generality, after a rescaling we can assume that the norm of Σ is bounded by a constant $C > 0$. Moreover, we assume that Σ is well-conditioned: $\kappa(\Sigma) \leq C$, where $\kappa(\cdot)$ denotes the condition number.

Here we have assumed (A.1) solely for simplicity of representation. If the entries of Z only have finite a -th moment for some $a > 4$, then all the results below still hold except that we need to replace $O(p^{-\frac{1}{2}+\epsilon})$ with $O(p^{-\frac{1}{2}+\frac{2}{a}+\epsilon})$ in some error bounds. We will not get deeper into this issue in this section, but refer the reader to Corollary E.8 below.

Then we make the following assumptions on the data models.

Assumption A.2. For for some fixed $t \in \mathbb{N}$, let $Y_i = X_i\beta_i + \varepsilon_i, 1 \leq i \leq t$, be independent data models, where X_i, β_i and ε_i are also independent of each other. Suppose that $X_i = Z_i\Sigma_i^{1/2} \in \mathbb{R}^{n_i \times p}$ satisfy Assumption A.1 with $\rho_i := n_i/p > 1$ being fixed constants, and $\varepsilon_i \in \mathbb{R}^{n_i}$ are random vectors with i.i.d. entries with mean zero, variance σ_i^2 and all moments as in (A.1).

Throughout the appendix, we shall say an event Ξ holds with high probability (whp) if for any fixed $D > 0, \mathbb{P}(\Xi) \geq 1 - p^{-D}$ for large enough p . Moreover, we shall use $o(1)$ to mean a small positive number that converges to 0 as $p \rightarrow \infty$.

The multi-task learning estimator. From [15], we know that we need to explicitly restrict the capacity r of B so that there is transfer between the two tasks. for the rest of the section, we shall consider the case of two tasks with $r = 1$. Then equation (2.1) simplifies to

$$f(B; w_1, w_2) = \|X_1 B w_1 - Y_1\|^2 + \|X_2 B w_2 - Y_2\|^2, \quad (\text{A.2})$$

where $B \in \mathbb{R}^p$ and w_1, w_2 are both real numbers. To solve the above problem, suppose that w_1, w_2 are fixed, by local optimality, we find the optimal B as

$$\begin{aligned} \hat{B}(w_1, w_2) &= (w_1^2 X_1^\top X_1 + w_2^2 X_2^\top X_2)^{-1} (w_1 X_1^\top Y_1 + w_2 X_2^\top Y_2) \\ &= \frac{1}{w_2} \left(\frac{w_1^2}{w_2^2} X_1^\top X_1 + X_2^\top X_2 \right)^{-1} \left(\frac{w_1}{w_2} X_1^\top Y_1 + X_2^\top Y_2 \right) \\ &= \frac{1}{w_2} \left[\beta_t + \left(\frac{w_1^2}{w_2^2} X_1^\top X_1 + X_2^\top X_2 \right)^{-1} \left(X_1^\top X_1 \left(\frac{w_1}{w_2} \beta_1 - \frac{w_1^2}{w_2^2} \beta_2 \right) + \left(\frac{w_1}{w_2} X_1^\top \varepsilon_1 + X_2^\top \varepsilon_2 \right) \right) \right]. \end{aligned} \quad (\text{A.3})$$

As a remark, when $w_1 = w_2 = 1$, we recover the linear regression estimator. The advantage of using $f(B; w_1, w_2)$ is that if β_1 is a scaling of β_2 , then this case can be solved optimally using equation (A.2) [36].

Next we consider N_i independent samples of the training set $\{(\tilde{x}_k^{(i)}, \tilde{y}_k^{(i)}) : 1 \leq k \leq N_i\}$ from task- $i, i = 1, 2$. With these sample, we form the random matrices $\tilde{X}_i \in \mathbb{R}^{N_i \times p}$ and $\tilde{Y}_i \in \mathbb{R}^{N_i \times p}, i = 1, 2$, whose row vectors are given by $\tilde{x}_k^{(i)}$ and $\tilde{y}_k^{(i)}$. Here we assume that N_1 and N_2 satisfies $N_1/N_2 = n_1/n_2$ and $N_i \geq n_i^{1-\varepsilon_0}$ for some constant $\varepsilon_0 > 0$. Then we define the validation loss as

$$\tilde{f}(\hat{B}; w_1, w_2) = \|\tilde{X}_1 \hat{B} w_1 - \tilde{Y}_1\|^2 + \|\tilde{X}_2 \hat{B} w_2 - \tilde{Y}_2\|^2. \quad (\text{A.4})$$

Inserting (A.3) into (A.4), one can see that \tilde{f} only depends on the ratio $v := w_1/w_2$. Hence we will also write $\tilde{f}(\hat{B}; v)$ in the following discussion.

Let $\hat{v} = \hat{w}_1/\hat{w}_2$ be the global minimizer of $\tilde{f}(\hat{B}; v)$. We will define the multi-task learning estimator for the target task as

$$\hat{\beta}_t^{\text{MTL}} = \hat{w}_2 \hat{B}(\hat{w}_1, \hat{w}_2),$$

where $t = 2$ since we are considering the two task case, and it also stands for the “target task”. The intuition for deriving $\hat{\beta}_t^{\text{MTL}}$ is akin to performing multi-task training in practice. Then the test loss of using $\hat{\beta}_t^{\text{MTL}}$ for the target task is

$$\begin{aligned} L(\hat{\beta}_t^{\text{MTL}}) = & \hat{v}^2 \left\| \Sigma_2^{1/2} (\hat{v}^2 X_1^\top X_1 + X_2^\top X_2)^{-1} X_1^\top X_1 (\beta_1 - \hat{v} \beta_2) \right\|^2 \\ & + \text{Tr} \left[\Sigma_2 (\hat{v}^2 X_1^\top X_1 + X_2^\top X_2)^{-2} (\sigma_1^2 \cdot \hat{v}^2 X_1^\top X_1 + \sigma_2^2 \cdot X_2^\top X_2) \right], \end{aligned} \quad (\text{A.5})$$

which is well-defined since it only depends on \hat{v} , but otherwise does not depend on \hat{w}_1 or \hat{w}_2 separately. Our goal is to study under model and covariate shifts, whether multi-task learning helps to learn the target task better than single-task learning. The baseline where we solve the target task with its own data is

$$L(\hat{\beta}_t^{\text{STL}}) = \sigma_2^2 \cdot \text{Tr} [\Sigma_2 (X_2^\top X_2)^{-1}], \quad \text{where } \hat{\beta}_t^{\text{STL}} = (X_2^\top X_2)^{-1} X_2^\top Y_2.$$

One may observe that we can reduce \tilde{f} to an expression that is easier to handle using concentration of random vectors with i.i.d. entries. Before doing that, we first need to fix the setting for the following discussions, because we want to keep track of the error rate carefully instead of obtaining an asymptotic result only.

Now suppose $Y_i = X_i \beta_i + \varepsilon_i$ and $\tilde{Y}_i = \tilde{X}_i \tilde{\beta}_i + \tilde{\varepsilon}_i$, $i = 1, 2$, all satisfy Assumption A.2. Then we rewrite (A.4) as

$$\tilde{f}(\hat{B}; v) = \sum_{i=1}^2 \left\| \tilde{X}_i \tilde{\beta}_i - \tilde{\varepsilon}_i \right\|^2, \quad \tilde{\beta} := \hat{B} w_i - \beta_i.$$

Since $\tilde{X}_i \tilde{\beta}_i$ and $\tilde{\varepsilon}_i$ are independent random vectors with i.i.d. centered entries, we can use the concentration estimate, Lemma E.14, to get that for any constant $\varepsilon > 0$,

$$\begin{aligned} \left| \left\| \tilde{X}_i \tilde{\beta}_i - \tilde{\varepsilon}_i \right\|^2 - \mathbb{E}_{\tilde{X}_i, \tilde{\varepsilon}_i} \left[\left\| \tilde{X}_i \tilde{\beta}_i - \tilde{\varepsilon}_i \right\|^2 \right] \right| &= \left| \left\| \tilde{X}_i \tilde{\beta}_i - \tilde{\varepsilon}_i \right\|^2 - N_i (\tilde{\beta}_i^\top \Sigma_i \tilde{\beta}_i + \sigma_i^2) \right| \\ &\leq N_i^{1/2+\varepsilon} (\tilde{\beta}_i^\top \Sigma_i \tilde{\beta}_i + \sigma_i^2), \end{aligned}$$

with high probability. Thus we obtain that

$$\tilde{f}(\hat{B}; v) = \left[\sum_{i=1}^2 N_i \left\| \Sigma_i^{1/2} (\hat{B} w_i - \beta_i) \right\|^2 + (N_1 \sigma_1^2 + N_2 \sigma_2^2) \right] \cdot \left(1 + O(p^{-(1-\varepsilon_0)/2+\varepsilon}) \right),$$

where we also used $N_i \geq p^{-1+\varepsilon_0}$. Inserting (A.3) into the above expression and using again the concentration result, Lemma E.14, we obtain that

$$\sum_{i=1}^2 N_i \left\| \Sigma_i^{1/2} (\hat{B} w_i - \beta_i) \right\|^2 = \text{val}(\hat{B}; v) \cdot \left(1 + O(p^{-1/2+\varepsilon}) \right)$$

with high probability, where

$$\begin{aligned} \text{val}(\hat{B}; v) &:= \mathbb{E}_{\varepsilon_1, \varepsilon_2} \left[\sum_{i=1}^2 \left\| \Sigma_i^{1/2} (\hat{B} w_i - \beta_i) \right\|^2 \right] \\ &= N_1 \cdot \left\| \Sigma_1^{1/2} (v^2 X_1^\top X_1 + X_2^\top X_2)^{-1} X_2^\top X_2 (\beta_1 - v \beta_2) \right\|^2 \\ &\quad + N_2 \cdot v^2 \left\| \Sigma_2^{1/2} (v^2 X_1^\top X_1 + X_2^\top X_2)^{-1} X_1^\top X_1 (\beta_1 - v \beta_2) \right\|^2 \\ &\quad + N_1 \cdot v^2 \text{Tr} \left[\Sigma_1 (v^2 X_1^\top X_1 + X_2^\top X_2)^{-2} (\sigma_1^2 \cdot v^2 X_1^\top X_1 + \sigma_2^2 \cdot X_2^\top X_2) \right] \\ &\quad + N_2 \cdot \text{Tr} \left[\Sigma_2 (v^2 X_1^\top X_1 + X_2^\top X_2)^{-2} (\sigma_1^2 \cdot v^2 X_1^\top X_1 + \sigma_2^2 \cdot X_2^\top X_2) \right]. \end{aligned}$$

In sum, we have obtained that

$$\tilde{f}(\hat{B}; v) = \left[\text{val}(\hat{B}; v) + (N_1 \sigma_1^2 + N_2 \sigma_2^2) \right] \cdot \left(1 + O(p^{-(1-\varepsilon_0)/2+\varepsilon}) \right). \quad (\text{A.6})$$

Hence to minimize \tilde{f} , it suffices to minimize $\text{val}(\hat{B}; v)$ over v .

B Proof of Theorem 3.2

We now state several helper lemmas to get estimates on $L(\hat{\beta}_t^{\text{STL}})$ and $L(\hat{\beta}_t^{\text{MTL}})$. The first lemma, which is a folklore result in random matrix theory, helps to determine the asymptotic limit of $L(\hat{\beta}_t^{\text{STL}})$, as $p \rightarrow \infty$. When the entries of X are multivariate Gaussian, this lemma recovers the classical result for the mean of inverse Wishart distribution [41]. For general non-Gaussian random matrices, it can be obtained from Stieltjes transform method; see e.g., Lemma 3.11 of [42]. Here we shall state a result obtained from Theorem 2.4 in [20], which gives an almost sharp error bound.

Lemma B.1. *Suppose X satisfies assumption A.1. Let A be any $p \times p$ matrix that is independent of X . We have that for any constant $\varepsilon > 0$,*

$$\text{Tr}[(X^\top X)^{-1}A] = \frac{1}{\rho - 1} \frac{1}{p} \text{Tr}(\Sigma^{-1}A) + O\left(\|A\|p^{-1/2+\varepsilon}\right) \quad (\text{B.1})$$

with high probability.

We shall refer to random matrices of the form $X^\top X$ as sample covariance matrices following the standard notations in high-dimensional statistics. The second lemma extends Lemma B.1 for a single sample covariance matrix to the sum of two independent sample covariance matrices. It is the main random matrix theoretical input of this paper.

Lemma B.2. *Suppose $X_1 = Z_1 \Sigma_1^{1/2} \in \mathbb{R}^{n_1 \times p}$ and $X_2 = Z_2 \Sigma_2^{1/2} \in \mathbb{R}^{n_2 \times p}$ satisfy Assumption A.1 with $\rho_1 := n_1/p > 1$ and $\rho_2 := n_2/p > 1$ being fixed constants. Denote by $M = \Sigma_1^{1/2} \Sigma_2^{-1/2}$ and let $\lambda_1, \lambda_2, \dots, \lambda_p$ be the singular values of $M^\top M$ in descending order. Let A be any $p \times p$ matrix that is independent of X_1 and X_2 . We have that for any constant $\varepsilon > 0$,*

$$\text{Tr}[(X_1^\top X_1 + X_2^\top X_2)^{-1}A] = \frac{1}{\rho_1 + \rho_2} \frac{1}{p} \text{Tr}[(a_1 \Sigma_1 + a_2 \Sigma_2)^{-1}A] + O\left(\|A\|p^{-1/2+\varepsilon}\right) \quad (\text{B.2})$$

with high probability, where (a_1, a_2) is the solution to the following deterministic equations:

$$a_1 + a_2 = 1 - \frac{1}{\rho_1 + \rho_2}, \quad a_1 + \frac{1}{\rho_1 + \rho_2} \cdot \frac{1}{p} \sum_{i=1}^p \frac{\lambda_i^2 a_1}{\lambda_i^2 a_1 + a_2} = \frac{\rho_1}{\rho_1 + \rho_2}. \quad (\text{B.3})$$

Proof Overview. We first describe the proof of Theorem 3.1. We use the Stieltjes transform method (or the resolvent method) in random matrix theory [42, 43, 19]. Roughly speaking, we study the resolvent $R(z) := [\Sigma_2^{-1/2}(X_1^\top X_1 + X_2^\top X_2)\Sigma_2^{-1/2} - z]^{-1}$ for $z \in \mathbb{C}$ around $z = 0$. Using the methods in [21, 37], we find the asymptotic limit, say $R_\infty(z)$, of $R(z)$ for any z as $p \rightarrow \infty$ with an almost optimal convergence rate. In particular, when $z = 0$, $\text{Tr}[R_\infty(0)]$ gives the expression in Theorem 3.1. The details can be found in Appendix E and E.3.

Finally, the last lemma describes the asymptotic limit of $(X_1^\top X_1 + X_2^\top X_2)^{-1} \Sigma_2 (X_1^\top X_1 + X_2^\top X_2)^{-1}$, which will be needed when we estimate the first term on the right-hand side of (A.5).

Lemma B.3. *In the setting of Lemma B.2, let $\beta \in \mathbb{R}^p$ be any vector that is independent of X_1 and X_2 . We have that for any constant $\varepsilon > 0$,*

$$\begin{aligned} & (n_1 + n_2)^2 \left\| \Sigma_2^{1/2} (X_1^\top X_1 + X_2^\top X_2)^{-1} \beta \right\|^2 \\ &= \beta^\top \Sigma_2^{-1/2} \frac{(1 + a_3) \text{Id} + a_4 M^\top M}{(a_1 M^\top M + a_2)^2} \Sigma_2^{-1/2} \beta + O(p^{-1/2+\varepsilon} \|\beta\|^2), \end{aligned} \quad (\text{B.4})$$

with high probability, where a_3 and a_4 satisfy the following system of linear equations:

$$(\rho_2 a_2^{-2} - b_0) \cdot a_3 - b_1 \cdot a_4 = b_0, \quad (\rho_1 a_1^{-2} - b_2) \cdot a_4 - b_1 \cdot a_3 = b_1. \quad (\text{B.5})$$

Here b_0, b_1 and b_2 are defined as

$$b_k := \frac{1}{p} \sum_{i=1}^p \frac{\lambda_i^{2k}}{(a_2 + \lambda_i^2 a_1)^2}, \quad k = 0, 1, 2.$$

The proof of Lemma B.2 and Lemma B.3 is a main focus of Section E. We remark that one can probably derive the same asymptotic result using free probability theory (see e.g. [38]), but our results (B.2) and (B.4) also give an almost sharp error bound $O(p^{-1/2+\varepsilon})$.

578 In this section, we state and prove the formal version of Theorem 3.2, which covers the two tasks
 579 case with $t = 2$. In this section, we consider the case where the entries of ε_1 and ε_2 have the same
 580 variance $\sigma_1^2 = \sigma_2^2 = \sigma^2$.

581 First, we introduce several quantities that will be used in our statement, and they are also related
 582 to the quantities in Lemma B.2 and Lemma B.3. Given the optimal ratio \hat{v} , let $\hat{M} = \hat{v}\Sigma_1^{1/2}\Sigma_2^{-1/2}$
 583 denote the weighted covariate shift matrix, and $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p$ be the eigenvalues of $\hat{M}^\top \hat{M}$.
 584 Define (\hat{a}_1, \hat{a}_2) as the solution to the following system of deterministic equations,

$$\hat{a}_1 + \hat{a}_2 = 1 - \frac{1}{\rho_1 + \rho_2}, \quad \hat{a}_1 + \frac{1}{\rho_1 + \rho_2} \cdot \frac{1}{p} \sum_{i=1}^p \frac{\hat{\lambda}_i^2 \hat{a}_1}{\hat{\lambda}_i^2 \hat{a}_1 + \hat{a}_2} = \frac{\rho_1}{\rho_1 + \rho_2}. \quad (\text{B.6})$$

585 After obtaining (\hat{a}_1, \hat{a}_2) , we can solve the following linear equations to get (\hat{a}_3, \hat{a}_4) :

$$(\rho_2 \hat{a}_2^{-2} - \hat{b}_0) \cdot \hat{a}_3 - \hat{b}_1 \cdot \hat{a}_4 = \hat{b}_0, \quad (\rho_1 \hat{a}_1^{-2} - \hat{b}_2) \cdot \hat{a}_4 - \hat{b}_1 \cdot \hat{a}_3 = \hat{b}_1. \quad (\text{B.7})$$

where we denoted

$$\hat{b}_k := \frac{1}{p} \sum_{i=1}^p \frac{\hat{\lambda}_i^{2k}}{(\hat{a}_2 + \hat{\lambda}_i^2 \hat{a}_1)^2}, \quad k = 0, 1, 2.$$

586 Then we introduce the following matrix

$$\Pi = \frac{\rho_1^2}{(\rho_1 + \rho_2)^2} \cdot \hat{M} \frac{(1 + \hat{a}_3) \text{Id} + \hat{a}_4 \hat{M}^\top \hat{M}}{(\hat{a}_1 \hat{M}^\top \hat{M} + \hat{a}_2)^2} \hat{M}^\top. \quad (\text{B.8})$$

We introduce two factors that will appear often in our statements and discussions:

$$\alpha_-(\rho_1) := \left(1 - \rho_1^{-1/2}\right)^2, \quad \alpha_+(\rho_1) := \left(1 + \rho_1^{-1/2}\right)^2.$$

587 In fact, $\alpha_-(\rho_1)$ and $\alpha_+(\rho_1)$ correspond to the largest and smallest singular values of $Z_1/\sqrt{n_1}$,
 588 respectively, as given by the famous Marčenko-Pastur law [44]. In particular, as ρ_1 increases, both
 589 α_- and α_+ will converge to 1 and $Z_1/\sqrt{n_1}$ will be more close to an isometry. Finally, we introduce
 590 the error term

$$\delta \equiv \delta(\hat{v}) := \frac{\alpha_+(\rho_1) - 1}{\alpha_-^2(\rho_1) \lambda_{\min}^2(\hat{M})} \cdot \|\Sigma_1^{1/2}(\beta_1 - \hat{v}\beta_2)\|^2, \quad (\text{B.9})$$

591 where $\lambda_{\min}(\hat{M})$ is the smallest singular value of \hat{M} . Note that this factor converges to 0 as ρ_1
 592 increases.

593 Now we are ready to state our main result for two tasks with both covariate and model shift. It shows
 594 that the information transfer is determined by two deterministic quantities Δ_{bias} and Δ_{var} , which give
 595 the change of model shift bias and the change of variance, respectively.

596 **Theorem B.4.** Consider two data models $Y_i = X_i \beta_i + \varepsilon_i$, $i = 1, 2$, that satisfy Assumption A.2.
 597 With high probability, we have

$$L(\hat{\beta}_t^{\text{MTL}}) \leq L(\hat{\beta}_t^{\text{STL}}) \quad \text{when: } \Delta_{\text{var}} - \Delta_{\text{bias}} \geq \delta \quad (\text{B.10})$$

$$L(\hat{\beta}_t^{\text{MTL}}) \geq L(\hat{\beta}_t^{\text{STL}}) \quad \text{when: } \Delta_{\text{var}} - \Delta_{\text{bias}} \leq -\delta, \quad (\text{B.11})$$

598 where

$$\Delta_{\text{var}} := \sigma^2 \left(\frac{1}{\rho_2 - 1} - \frac{1}{\rho_1 + \rho_2} \cdot \frac{1}{p} \text{Tr} \left[(\hat{a}_1 \hat{M}^\top \hat{M} + \hat{a}_2 \text{Id})^{-1} \right] \right) \quad (\text{B.12})$$

$$\Delta_{\text{bias}} := (\beta_1 - \hat{v}\beta_2)^\top \Sigma_1^{1/2} \Pi \Sigma_1^{1/2} (\beta_1 - \hat{v}\beta_2). \quad (\text{B.13})$$

599 For the isotropic model in Section 3, we actually have an easier and sharper bound than Theorem B.4
 600 as follows.

601 **Lemma B.5.** In the setting of Theorem B.4, assume that $\Sigma_1 = \text{Id}$, β_2 is a random vector with
 602 i.i.d. entries with mean 0, variance κ^2 and all moments, and β_1 is a random vector such that
 603 $(\beta_1 - \beta_2)$ is a random vector with i.i.d. entries with mean 0, variance d^2 and all moments. Denote
 604 $\Delta_{\text{bias}}^* := ((1 - \hat{v})^2 \kappa^2 + d^2) \text{Tr}[\Pi]$. Then we have

$$L(\hat{\beta}_t^{\text{MTL}}) \leq L(\hat{\beta}_t^{\text{STL}}) \quad \text{when: } \Delta_{\text{var}} \geq (\alpha_+(\rho_1) + o(1)) \cdot \Delta_{\text{bias}}^*,$$

$$L(\hat{\beta}_t^{\text{MTL}}) \geq L(\hat{\beta}_t^{\text{STL}}) \quad \text{when: } \Delta_{\text{var}} \leq (\alpha_-(\rho_1) - o(1)) \cdot \Delta_{\text{bias}}^*.$$

Now we give the proof of Theorem B.4 based on Lemma B.2 and Lemma B.3.

Proof of Theorem B.4. Note that

$$\begin{aligned} L(\hat{\beta}_t^{\text{STL}}) - L(\hat{\beta}_t^{\text{MTL}}) &= \sigma^2 \left(\text{Tr} [(X_2^\top X_2)^{-1} \Sigma_2] - \text{Tr} [(\hat{v}^2 X_1^\top X_1 + X_2^\top X_2)^{-1} \Sigma_2] \right) \\ &\quad - \hat{v}^2 \left\| \Sigma_2^{1/2} (\hat{v}^2 X_1^\top X_1 + X_2^\top X_2)^{-1} X_1^\top X_1 (\beta_1 - \hat{v} \beta_2) \right\|^2 =: \delta_{\text{var}}(\hat{v}) - \delta_{\text{bias}}(\hat{v}). \end{aligned}$$

The proof is divided into the following four steps.

- (i) We first consider $\hat{M} \equiv \hat{M}(v) = v \Sigma_1^{1/2} \Sigma_2^{-1/2}$ for a fixed $v \in \mathbb{R}$. Then we use Lemma B.1 and Lemma B.2 to calculate the variance reduction $\delta_{\text{var}}(v)$, which will lead to the Δ_{var} term.
- (ii) Using the approximate isometry property of X_1 (see (B.16) below), we will bound the bias term $\delta_{\text{bias}}(v)$ through

$$\tilde{\delta}_{\text{bias}}(v) := v^2 n_1^2 \left\| \Sigma_2^{1/2} (v^2 X_1^\top X_1 + X_2^\top X_2)^{-1} \Sigma_1 (\beta_1 - v \beta_2) \right\|^2. \quad (\text{B.14})$$

- (iii) We use Lemma B.3 to calculate (B.14), which will lead to the Δ_{bias} term.

- (iv) Finally we use a standard ε -net argument to extend the above results to $\hat{M} = \hat{v} \Sigma_1^{1/2} \Sigma_2^{-1/2}$ for a possibly random \hat{v} which depends on Y_1 and Y_2 .

Step I: Variance reduction. Let $\hat{M} = v \Sigma_1^{1/2} \Sigma_2^{-1/2}$ for any fixed constant $v \in \mathbb{R}$. Using Lemma B.2, we can obtain that for any constant $\varepsilon > 0$,

$$\sigma^2 \cdot \text{Tr} [(X_2^\top X_2)^{-1} \Sigma_2] = \frac{\sigma^2}{\rho_2 - 1} \left(1 + O(p^{-1/2+\varepsilon}) \right),$$

and

$$\sigma^2 \cdot \text{Tr} [(v^2 X_1^\top X_1 + X_2^\top X_2)^{-1} \Sigma_2] = \frac{\sigma^2}{\rho_1 + \rho_2} \cdot \frac{1}{p} \text{Tr} [(\hat{a}_1 \hat{M}^\top \hat{M} + \hat{a}_2 \text{Id})^{-1}] \left(1 + O(p^{-1/2+\varepsilon}) \right),$$

with high probability, where \hat{a}_1 and \hat{a}_2 satisfy (B.6). Combining them, we get

$$\delta_{\text{var}}(v) = \Delta_{\text{var}}(v) + O(\sigma^2 p^{-1/2+\varepsilon}) \quad \text{whp}, \quad (\text{B.15})$$

where $\Delta_{\text{var}}(v)$ is defined as in (B.12) but with \hat{v} replaced by v .

Step II: Bounding the bias term. In this step, we shall use the following the following bounds on the singular values of Z_1 : for any fixed $\varepsilon > 0$, we have

$$\alpha_-(\rho_1) - O(p^{-1/2+\varepsilon}) \preceq \frac{Z_1^T Z_1}{n_1} \preceq \alpha_+(\rho_1) + O(p^{-1/2+\varepsilon}) \quad (\text{B.16})$$

with high probability. In fact, $Z_1^T Z_1$ is a standard sample covariance matrix, and it is well-known that its nonzero eigenvalues are all inside the support of the Marchenko-Pastur law $[\alpha_-(\rho_1) - o(1), \alpha_+(\rho_1) + o(1)]$ with probability $1 - o(1)$ [45]. For the estimate (B.16) we used [20, Theorem 2.10] to get a stronger probability bound.

Next we shall use (B.16) to approximate $\delta_{\text{bias}}(v)$ with $\tilde{\delta}_{\text{bias}}(v)$ in (B.14).

Lemma B.6. In the setting of Theorem B.4, we denote by $K = (v^2 X_1^\top X_1 + X_2^\top X_2)^{-1}$, and

$$\delta_\varepsilon(v) := n_1^2 v^2 \left\| \Sigma_1^{1/2} K \Sigma_2 K \Sigma_1^{1/2} \right\| \cdot \left\| \Sigma_1^{1/2} (\beta_1 - v \beta_2) \right\|^2.$$

Then we have whp,

$$\left| \delta_{\text{bias}}(v) - \tilde{\delta}_{\text{bias}}(v) \right| \leq \left(\alpha_+^2(\rho_1) - 1 + O(p^{-1/2+\varepsilon}) \right) \delta_\varepsilon.$$

626 *Proof.* Denote by $\mathcal{E} = Z_1^\top Z_1 - n_1 \text{Id}$. Then we can write

$$\begin{aligned} \delta_{\text{bias}}(v) - \tilde{\delta}_{\text{bias}}(v) &= 2v^2 n_1 (\beta_1 - v\beta_2)^\top \Sigma_1^{1/2} \mathcal{E} \left(\Sigma_1^{1/2} K \Sigma_2 K \Sigma_1^{1/2} \right) \Sigma_1^{1/2} (\beta_1 - v\beta_2) \\ &\quad + v^2 \left\| \Sigma_2^{1/2} K \Sigma_1^{1/2} \mathcal{E} \Sigma_1^{1/2} (\beta_1 - v\beta_2) \right\|^2. \end{aligned} \quad (\text{B.17})$$

Using (B.16), we can bound

$$\|\mathcal{E}\| \leq \left(\alpha_+(\rho_1) - 1 + O(p^{-1/2+\varepsilon}) \right) n_1, \quad \text{whp.}$$

627 Thus we can estimate that

$$\begin{aligned} |\delta_{\text{bias}}(v) - \tilde{\delta}_{\text{bias}}(v)| &\leq v^2 (2n_1 \|\mathcal{E}\| + \|\mathcal{E}\|^2) \left\| \Sigma_1^{1/2} K \Sigma_2 K \Sigma_1^{1/2} \right\| \left\| \Sigma_1^{1/2} (\beta_1 - v\beta_2) \right\|^2 \\ &= v^2 \left[(n_1 + \|\mathcal{E}\|)^2 - n_1^2 \right] \left\| \Sigma_1^{1/2} K \Sigma_2 K \Sigma_1^{1/2} \right\| \left\| \Sigma_1^{1/2} (\beta_1 - v\beta_2) \right\|^2 \\ &\leq v^2 n_1^2 \left[\alpha_+^2(\rho_1) + O(p^{-1/2+\varepsilon}) - 1 \right] \left\| \Sigma_1^{1/2} K \Sigma_2 K \Sigma_1^{1/2} \right\| \left\| \Sigma_1^{1/2} (\beta_1 - v\beta_2) \right\|^2, \end{aligned}$$

628 which concludes the proof by the definition of δ_ε . \square

629 Note by (B.16), we have with high probability,

$$\begin{aligned} v^2 n_1^2 \Sigma_1^{1/2} K \Sigma_2 K \Sigma_1^{1/2} &= \hat{M} \frac{1}{(\hat{M}^\top Z_1^\top Z_1 \hat{M} + Z_2^\top Z_2)^2} \hat{M}^\top \\ &\preceq n_1^2 \hat{M} \frac{1}{\left[n_1 \alpha_-(\rho_1) \hat{M}^\top \hat{M} + n_2 \alpha_-(\rho_2) + O(p^{1/2+\varepsilon}) \right]^2} \hat{M}^\top \\ &\preceq \left[\alpha_-^2(\rho_1) \hat{M} \hat{M}^\top + 2 \frac{\rho_2}{\rho_1} \alpha_-(\rho_1) \alpha_-(\rho_2) + 2 \left(\frac{\rho_2}{\rho_1} \right)^2 \alpha_-^2(\rho_2) (\hat{M} \hat{M}^\top)^{-1} \right]^{-1} + O(p^{-1/2+\varepsilon}) \\ &\prec [\alpha_-^2(\rho_1) \lambda_{\min}^2(\hat{M})]^{-1} \cdot (1 - c) \end{aligned}$$

630 for some small enough constant $c > 0$. Together with Lemma B.6, we get with high probability,

$$\left| \delta_{\text{bias}}(v) - \tilde{\delta}_{\text{bias}}(v) \right| \leq (1 - c) \delta(v) \quad (\text{B.18})$$

631 for some small constant $c > 0$, where recall $\delta(v)$ defined in (B.9).

632 **Step III: The limit of $\tilde{\delta}_{\text{bias}}(v)$.** Using Lemma B.3 with Σ_1 and M replaced by $v^2 \Sigma_1$ and \hat{M} , we
633 obtain that

$$\begin{aligned} \tilde{\delta}_{\text{bias}}(v) &= \frac{\rho_1^2}{(\rho_1 + \rho_2)^2} \cdot v^2 (\beta_1 - v\beta_2)^\top \Sigma_1 \Sigma_2^{-1/2} \frac{(1 + \hat{a}_3) \text{Id} + \hat{a}_4 \hat{M}^\top \hat{M}}{(a_1 \hat{M}^\top \hat{M} + a_2)^2} \Sigma_2^{-1/2} \Sigma_1 (\beta_1 - v\beta_2) + O(p^{-1/2+\varepsilon}) \\ &= (\beta_1 - v\beta_2)^\top \Sigma_1^{1/2} \Pi \Sigma_1^{1/2} (\beta_1 - v\beta_2) + O(p^{-1/2+\varepsilon}) =: \Delta_{\text{bias}}(v) + O(p^{-1/2+\varepsilon}), \end{aligned}$$

634 with high probability. Together with and (B.15) and (B.18), we obtain that whp,

$$\begin{cases} \delta_{\text{var}}(v) > \delta_{\text{bias}}(v), & \text{if } \Delta_{\text{var}}(v) - \Delta_{\text{bias}}(v) \geq \delta(v), \\ \delta_{\text{var}}(v) < \delta_{\text{bias}}(v), & \text{if } \Delta_{\text{var}}(v) - \Delta_{\text{bias}}(v) \leq -\delta(v). \end{cases} \quad (\text{B.19})$$

Step IV: An ε -net argument. Finally, it remains to extend the above result to $v = \hat{v}$, which is random and depends on X_1 and X_2 . We first show that for any fixed constant $C_0 > 0$, there exists a high probability event Ξ on which (B.19) holds uniformly for all $v \in [-C_0, C_0]$. In fact, for a large constant $C_1 > 0$, we consider v belonging to a discrete set

$$V := \{v_k = kp^{-1} : -(C_0 p + 1) \leq k \leq C_0 p + 1\}.$$

Then using the arguments for the first three steps and a simple union bound, we get that (B.19) holds simultaneously for all $v \in V$ with high probability. On the other hand, by (B.16) the event

$$\Xi_1 := \left\{ \alpha_-(\rho_1)/2 \preceq \frac{Z_1^\top Z_1}{n_1} \preceq 2\alpha_+(\rho_1), \alpha_-(\rho_2)/2 \preceq \frac{Z_2^\top Z_2}{n_2} \preceq 2\alpha_+(\rho_2) \right\}$$

holds with high probability. Now it is easy to check that on Ξ_1 , for all $v_k \leq v \leq v_{k+1}$ we have the following estimates:

$$|\delta_{\text{var}}(v) - \delta_{\text{var}}(v_k)| \lesssim p^{-1} \delta_{\text{var}}(v_k), \quad |\delta_{\text{bias}}(v) - \delta_{\text{bias}}(v_k)| \lesssim p^{-1} \delta_{\text{bias}}(v_k), \quad |\delta(v) - \delta(v_k)| \lesssim p^{-1} \delta(v_k), \\ |\Delta_{\text{bias}}(v) - \Delta_{\text{bias}}(v_k)| \lesssim p^{-1} \Delta_{\text{bias}}(v_k), \quad |\Delta_{\text{var}}(v) - \Delta_{\text{var}}(v_k)| \lesssim p^{-1} \Delta_{\text{var}}(v_k).$$

Then a simple application of triangle inequality gives that the event

$$\Xi_2 = \{(\text{B.19}) \text{ holds simultaneously for all } -C_0 \leq v \leq C_0\}$$

holds with high probability. On the other hand, on Ξ_1 one can see that for any small constant $\varepsilon > 0$,

$$|\delta_{\text{var}}(v) - \delta_{\text{var}}(C_0)| \leq \varepsilon \delta_{\text{var}}(C_0), \quad |\delta_{\text{bias}}(v) - \delta_{\text{bias}}(C_0)| \leq \varepsilon \delta_{\text{bias}}(C_0), \quad |\delta(v) - \delta(C_0)| \leq \varepsilon \delta(C_0), \\ |\Delta_{\text{bias}}(v) - \Delta_{\text{bias}}(C_0)| \leq \varepsilon \Delta_{\text{bias}}(C_0), \quad |\Delta_{\text{var}}(v) - \Delta_{\text{var}}(C_0)| \leq \varepsilon \Delta_{\text{var}}(C_0),$$

for all $v \geq C_0$ as long as C_0 is chosen large enough depending on ε . Similar estimates hold for $v \leq -C_0$ if we replace C_0 with $-C_0$ in the above estimates. Together with the estimate at $\pm C_0$, we get that (B.19) holds simultaneously for all $v \in \mathbb{R}$ on the high probability event $\Xi_1 \cap \Xi_2$. This concludes the proof since v must be one of the real values. \square

Remark B.7. One can see from the above proof that the main error, δ , of Theorem B.4 comes from approximating δ_{bias} by $\tilde{\delta}_{\text{bias}}$ in (B.18). In order to improve this estimate and obtain an exact asymptotic result as for the δ_{var} term, one needs to study the singular value distribution of the following random matrix:

$$(X_1^\top X_1)^{-1} X_2^\top X_2 + v^2.$$

In fact, the eigenvalues of $\mathcal{X} := (X_1^\top X_1)^{-1} X_2^\top X_2$ have been studied in the name of Fisher matrices; see e.g. [46]. However, since \mathcal{X} is not symmetric, it is known that the singular values of \mathcal{X} are different from its eigenvalues. To the best of our knowledge, the asymptotic singular value behavior of \mathcal{X} is still unknown in random matrix theory literature, and the study of the singular values of $\mathcal{X} + v^2$ will be even harder. We leave this problem to future study.

By replacing (B.18) with a tighter bound in Step II of the above proof, we can conclude the proof of Lemma B.5.

Proof of Lemma B.5. For any fixed $v \in \mathbb{R}$, $\beta_1 - v\beta_2$ is a random vector with i.i.d. entries with mean 0 and variance $(1-v)^2\kappa^2 + d^2$. Then using the concentration result, Lemma E.14, we get that for any constant $\varepsilon > 0$,

$$|\delta_{\text{bias}}(v) - [(1-v)^2\kappa^2 + d^2] \text{Tr}(\mathcal{K}^\top \mathcal{K})| \\ = |(\beta_1 - v\beta_2)^\top \mathcal{K}^\top \mathcal{K} (\beta_1 - v\beta_2) - [(1-v)^2\kappa^2 + d^2] \text{Tr}(\mathcal{K}^\top \mathcal{K})| \\ \leq p^\varepsilon [(1-v)^2\kappa^2 + d^2] \{\text{Tr}[(\mathcal{K}^\top \mathcal{K})^2]\}^{1/2} \lesssim p^{1/2+\varepsilon} [(1-v)^2\kappa^2 + d^2], \quad (\text{B.20})$$

where we denoted $\mathcal{K} := v\Sigma_2^{1/2}(v^2 X_1^\top X_1 + X_2^\top X_2)^{-1} X_1^\top X_1$, and in the last step we used $\|\mathcal{K}\| = O(1)$ by (B.16). Now for $\text{Tr}(\mathcal{K}^\top \mathcal{K})$, we rewrite it as

$$v^2[(1-v)^2\kappa^2 + d^2] \text{Tr}[(v^2 X_1^\top X_1 + X_2^\top X_2)^{-1} \Sigma_2 (v^2 X_1^\top X_1 + X_2^\top X_2)^{-1} (X_1^\top X_1)^2].$$

Recalling that $\Sigma_1 = \text{Id}$ and bounding $(X_1^\top X_1)^2 = (Z_1^\top Z_1)^2$ using (B.16) again, we obtain that

$$\delta_{\text{bias}}^*(v) \cdot (\alpha_-^2(\rho_1) - O(p^{-1/2+\varepsilon})) \leq [(1-v)^2\kappa^2 + d^2] \text{Tr}(\mathcal{K}^\top \mathcal{K}) \leq \delta_{\text{bias}}^*(v) \cdot (\alpha_+^2(\rho_1) + O(p^{-1/2+\varepsilon})), \quad (\text{B.21})$$

where

$$\delta_{\text{bias}}^*(v) := n_1^2 v^2 [(1-v)^2\kappa^2 + d^2] \text{Tr}[(v^2 X_1^\top X_1 + X_2^\top X_2)^{-1} \Sigma_2 (v^2 X_1^\top X_1 + X_2^\top X_2)^{-1}].$$

Note that $\delta_{\text{bias}}^*(v) \sim 1$, hence combining (B.20) and (B.21) we get

$$\delta_{\text{bias}}^*(v) \cdot (\alpha_-^2(\rho_1) - O(p^{-1/2+\varepsilon})) \leq \delta_{\text{bias}}(v) \leq \delta_{\text{bias}}^*(v) \cdot (\alpha_+^2(\rho_1) + O(p^{-1/2+\varepsilon})). \quad (\text{B.22})$$

Now we can replace the estimate (B.18) with this stronger estimate, and repeat all the other parts of the proof of Theorem B.4 to conclude Lemma B.5. In particular, one can calculate $\delta_{\text{bias}}^*(v)$ using Lemma B.3 and get the $\Delta_{\text{bias}}^*(v)$ term. We omit the details. \square

C Proofs for Isotropic and Covariate Shifted Settings

C.1 Missing Proofs of Section 3.2

We define the function

$$\begin{aligned} \text{val}(v) &= \frac{\rho_1}{\rho_2} \left[d^2 + (v-1)^2 \kappa^2 \right] \cdot \text{Tr} \left[(v^2 X_1^\top X_1 + X_2^\top X_2)^{-2} (X_2^\top X_2)^2 \right] \\ &\quad + v^2 \left[d^2 + (v-1)^2 \kappa^2 \right] \cdot \text{Tr} \left[(v^2 X_1^\top X_1 + X_2^\top X_2)^{-2} (X_1^\top X_1)^2 \right] \\ &\quad + \left(\frac{\rho_1}{\rho_2} v^2 + 1 \right) \sigma^2 \cdot \text{Tr} \left[(v^2 X_1^\top X_1 + X_2^\top X_2)^{-1} \right]. \end{aligned}$$

For the isotropic model with $\sigma_1^2 = \sigma_2^2 = \sigma^2$, using concentration for random vectors with i.i.d. entries, Lemma E.14, we can obtain that $\text{val}(\hat{B}; w_1, w_2) = \text{val}(v) \cdot (1 + O(p^{-1/2+\varepsilon}))$. Hence the validation loss in (A.6) reduces to

$$\tilde{f}(\hat{B}; v) = [N_2 \cdot \text{val}(v) + (N_1 \sigma_1^2 + N_2 \sigma_2^2)] \cdot \left(1 + O(p^{-(1-\varepsilon_0)/2+\varepsilon}) \right) \quad (\text{C.1})$$

with high probability for any constant $\varepsilon > 0$. Thus for the following discussions, it suffices to focus on the behavior of $\text{val}(v)$. Let \hat{w} the minimizer of $\text{val}(v)$. The proof will consist of two main steps.

- First, we show that \hat{w} is close to 1, and then (C.1) implies that \hat{v} is also close to 1.
- Second, we plug \hat{v} back into $L(\hat{\beta}_2^{\text{MTL}})$ and use Lemma B.5 to show the result.

For the first step, we will prove the following result.

Lemma C.1. *For the isotropic model, the minimizer for $\text{val}(v)$ satisfies*

$$|\hat{w} - 1| \leq C \left(\frac{d^2}{\kappa^2} + \frac{\sigma^2}{p\kappa^2} \right) \quad \text{whp} \quad (\text{C.2})$$

for some constant $C > 0$.

Proof. To be consistent with the notation \hat{w} , we shall change the name of the argument to w in the proof. First it is easy to observe that $\text{val}(w) < \text{val}(-w)$ for $w > 0$. Hence it suffices to assume that $w \geq 0$.

We first consider the case $w \geq 1$. We write

$$\begin{aligned} \text{val}(w) &= \frac{\rho_1}{\rho_2} \left[\frac{d^2}{w^4} + \frac{(w-1)^2}{w^4} \kappa^2 \right] \cdot \text{Tr} \left[(X_1^\top X_1 + w^{-2} X_2^\top X_2)^{-2} (X_2^\top X_2)^2 \right] \\ &\quad + \left[\frac{d^2}{w^2} + \frac{(w-1)^2}{w^2} \kappa^2 \right] \cdot \text{Tr} \left[(X_1^\top X_1 + w^{-2} X_2^\top X_2)^{-2} (X_1^\top X_1)^2 \right] \\ &\quad + \frac{\rho_1}{\rho_2} \sigma^2 \cdot \text{Tr} \left[(X_1^\top X_1 + w^{-2} X_2^\top X_2)^{-1} \right] + \sigma^2 \cdot \text{Tr} \left[(w^2 X_1^\top X_1 + X_2^\top X_2)^{-1} \right]. \end{aligned}$$

Notice that

$$\text{Tr} \left[(X_1^\top X_1 + w^{-2} X_2^\top X_2)^{-2} (X_i^\top X_i)^2 \right], \quad i = 1, 2, \quad \text{and} \quad \text{Tr} \left[(X_1^\top X_1 + w^{-2} X_2^\top X_2)^{-1} \right]$$

are increasing functions in w . Hence taking derivative of $\text{val}(w)$ with respect to w , we obtain that

$$\begin{aligned} \text{val}'(w) &\geq \frac{\rho_1}{\rho_2} \left[\frac{2(w-1)(2-w)}{w^5} \kappa^2 - \frac{4d^2}{w^5} \right] \text{Tr} \left[(X_1^\top X_1 + w^{-2} X_2^\top X_2)^{-2} (X_2^\top X_2)^2 \right] \\ &\quad + \left[\frac{2(w-1)}{w^3} \kappa^2 - \frac{2d^2}{w^3} \right] \cdot \text{Tr} \left[(X_1^\top X_1 + w^{-2} X_2^\top X_2)^{-2} (X_1^\top X_1)^2 \right] \\ &\quad - 2 \frac{\sigma^2}{w^3} \cdot \text{Tr} \left[(X_1^\top X_1 + w^{-2} X_2^\top X_2)^{-2} X_1^\top X_1 \right] = \text{Tr} \left[(X_1^\top X_1 + w^{-2} X_2^\top X_2)^{-2} \mathcal{A} \right], \end{aligned}$$

where the matrix \mathcal{A} is

$$\mathcal{A} := \frac{\rho_1}{\rho_2} \left[\frac{2(w-1)(2-w)}{w^5} \kappa^2 - \frac{4d^2}{w^5} \right] (X_2^\top X_2)^2 + \left[\frac{2(w-1)}{w^3} \kappa^2 - \frac{2d^2}{w^3} \right] (X_1^\top X_1)^2 - 2 \frac{\sigma^2}{w^3} X_1^\top X_1.$$

678 Using the estimate (B.16), we get that \mathcal{A} is lower bounded as

$$\begin{aligned}\mathcal{A} \succeq & -\frac{4d^2}{w^5}n_1n_2(\alpha_+(\rho_2) + o(1))^2 + \left[\frac{2(w-1)}{w^3}\kappa^2 - \frac{2d^2}{w^3}\right]n_1^2(\alpha_-(\rho_1) - o(1))^2 \\ & - 2\frac{\sigma^2}{w^3}n_1(\alpha_+(\rho_1) + o(1)) \succ 0,\end{aligned}$$

as long as

$$w > w_1 := 1 + \frac{d^2}{\kappa^2} + \frac{\sigma^2}{n_1\kappa^2} \frac{\alpha_+(\rho_1) + o(1)}{\alpha_-^2(\rho_1)} + \frac{2d^2}{\kappa^2} \frac{\rho_2(\alpha_+^2(\rho_2) + o(1))}{\rho_1\alpha_-^2(\rho_1)}.$$

679 Hence $val'(w) > 0$ on (w_1, ∞) , i.e. $val(w)$ is strictly increasing for $w > w_1$. Hence we must have
680 $\hat{w} \leq w_1$.

681 Then we consider the case $w \leq 1$, and the proof is similar as above. Taking derivative of $val(w)$, we
682 obtain that

$$\begin{aligned}val'(w) & \leq \frac{\rho_1}{\rho_2} [2(w-1)\kappa^2] \cdot \text{Tr}[(w^2X_1^\top X_1 + X_2^\top X_2)^{-2}(X_2^\top X_2)^2] \\ & + [2wd^2 + 2w(w-1)(2w-1)\kappa^2] \cdot \text{Tr}[(w^2X_1^\top X_1 + X_2^\top X_2)^{-2}(X_1^\top X_1)^2] \\ & + \frac{\rho_1}{\rho_2}(2w\sigma^2) \cdot \text{Tr}[(w^2X_1^\top X_1 + X_2^\top X_2)^{-2}X_2^\top X_2] \\ & = \frac{\rho_1}{\rho_2} \text{Tr}[(w^2X_1^\top X_1 + X_2^\top X_2)^{-1}\mathcal{B}],\end{aligned}\tag{C.3}$$

where the matrix \mathcal{B} is

$$\mathcal{B} = 2(w-1)\kappa^2(X_2^\top X_2)^2 + \frac{\rho_2}{\rho_1} [2wd^2 + 2w(w-1)(2w-1)\kappa^2] (X_1^\top X_1)^2 + 2w\sigma^2 X_2^\top X_2.$$

683 Using the estimate (B.16), we get that \mathcal{B} is upper bounded as

$$\mathcal{B} \preceq -2(1-w)\kappa^2n_2^2(\alpha_-(\rho_2) - o(1))^2 + 2wd^2n_1n_2(\alpha_+(\rho_1) + o(1))^2 + 2w\sigma^2n_2(\alpha_+(\rho_2) + o(1)) \prec 0,$$

as long as

$$w < w_2 := 1 - \frac{d^2}{\kappa^2} \frac{\rho_1(\alpha_+(\rho_1) + o(1))^2}{\rho_2\alpha_-^2(\rho_2)} - \frac{\sigma^2}{n_2\kappa^2} \frac{\alpha_+(\rho_2) + o(1)}{\alpha_-^2(\rho_2)}.$$

684 Hence $val'(w) < 0$ on $[0, w_2)$, i.e. $val(w)$ is strictly decreasing for $w < w_2$. Hence we must have
685 $\hat{w} \geq w_2$.

In sum, we obtain that $w_2 \leq w \leq w_1$. Note that under our assumptions, we have

$$\max(|w_1 - 1|, |w_2 - 1|) = O\left(\frac{d^2}{\kappa^2} + \frac{\sigma^2}{p\kappa^2}\right),$$

686 which concludes the proof. \square

687 For the rest of this section, we choose the parameters that satisfy the following relations:

$$pd^2 \sim \sigma^2 \sim 1, \quad p^{-1+c_0}\sigma^2 \leq \kappa^2 \leq p^{-\varepsilon_0-c_0}\sigma^2,\tag{C.4}$$

688 for some small constant $c_0 > 0$. We will explain below why we make this choice. Before that, we
689 first show the following estimate on the optimizer \hat{v} : with high probability,

$$|\hat{v} - 1| = O(\mathcal{E}), \quad \mathcal{E} := \frac{d^2}{\kappa^2} + \frac{\sigma^2}{p\kappa^2} + p^{-1/2+\varepsilon_0/2+2\varepsilon}.\tag{C.5}$$

In fact, from the proof of Lemma C.1 above, one can check that if $C\mathcal{E} \leq |w - \hat{w}| \leq 2C\mathcal{E}$ for a large enough constant $C > 1$, then $|val'(w)| \gtrsim pd^2$. Moreover, under the choice (C.4) we have

$$val(w) = O(pd^2), \quad \text{for } |w - \hat{w}| \leq 2C\mathcal{E}.$$

Thus we obtain that for $|w - \hat{w}| \geq 2C\mathcal{E}$,

$$|val(w) - val(\hat{w})| \geq |val(w) - \min(val(w_1), val(w_2))| \gtrsim pd^2\mathcal{E} \gtrsim \mathcal{E} \cdot val(\hat{w}),$$

690 which leads to $\tilde{f}(\hat{B}; w) > \tilde{f}(\hat{B}; \hat{w})$ whp by (C.1). Thus w cannot be a minimizer of $\tilde{f}(\hat{B}; v)$, and we
 691 must have $|\hat{v} - \hat{w}| \leq 2C\mathcal{E}$. Together with (C.2), we conclude (C.5).
 692 Inserting (C.5) into (A.5) and applying Lemma E.14 to $(\beta_1 - \hat{v}\beta_s)$ again, we get whp,

$$\begin{aligned} L(\hat{\beta}_t^{\text{MTL}}) &= (1 + O(\mathcal{E})) \cdot [d^2 + O(\mathcal{E}^2\kappa^2)] \text{Tr} [(X_1^\top X_1 + X_2^\top X_2)^{-2} (X_1^\top X_1)^2] \\ &\quad + (1 + O(\mathcal{E})) \cdot \sigma^2 \text{Tr} [(X_1^\top X_1 + X_2^\top X_2)^{-1}]. \end{aligned} \quad (\text{C.6})$$

In order to study the phenomenon of bias-variance trade-off, we need the bias term with d^2 and the variance term with σ^2 to be of the same order. With estimate (B.16), we see that

$$\text{Tr} [(X_1^\top X_1 + X_2^\top X_2)^{-2} (X_1^\top X_1)^2] \sim p, \quad \text{Tr} [(X_1^\top X_1 + X_2^\top X_2)^{-1}] \sim \frac{p}{n_1 + n_2}.$$

693 Hence we need to choose that $p \cdot d^2 \sim \sigma^2$. On the other hand, we want the error term $\mathcal{E}^2\kappa^2$ to be much
 694 smaller than d^2 , which leads to the condition $p^{-1+\varepsilon_0+4\varepsilon}\kappa^2 \ll d^2 \ll \kappa^2$. The above considerations
 695 lead to the choices of parameters in (C.4). Moreover, under (C.4) we can simplify (C.6) to

$$\begin{aligned} L(\hat{\beta}_2^{\text{MTL}}) &= (1 + O(n^{-\varepsilon})) \cdot d^2 \text{Tr} [(X_1^\top X_1 + X_2^\top X_2)^{-2} (X_1^\top X_1)^2] \\ &\quad + (1 + O(n^{-\varepsilon})) \cdot \sigma^2 \text{Tr} [(X_1^\top X_1 + X_2^\top X_2)^{-1}] \end{aligned} \quad (\text{C.7})$$

696 whp for some constant $\varepsilon > 0$.

697 With (C.7) and Lemma B.5, we can prove Proposition 3.3, which gives a transition threshold with
 698 respect to the ratio between the model bias and the noise level. With slight abuse of notations, we
 699 shall write \hat{a}_i, \hat{b}_k and \hat{M} as a_i, b_k and M throughout the rest of this section.

700 *Proof of Proposition 3.3.* In the setting of Proposition 3.3, we have $M = \Sigma_1^{1/2} \Sigma_2^{-1/2} = \text{Id}$. Then
 701 solving equations (B.6) and (B.7) with $\hat{\lambda}_i = 1$, we get that

$$a_1 = \frac{\rho_1(\rho_1 + \rho_2 - 1)}{(\rho_1 + \rho_2)^2}, \quad a_2 = \frac{\rho_2(\rho_1 + \rho_2 - 1)}{(\rho_1 + \rho_2)^2}, \quad (\text{C.8})$$

$$a_3 = \frac{\rho_2}{(\rho_1 + \rho_2)(\rho_1 + \rho_2 - 1)}, \quad a_4 = \frac{\rho_1}{(\rho_1 + \rho_2)(\rho_1 + \rho_2 - 1)}. \quad (\text{C.9})$$

702 Using Lemma B.1 and Lemma B.2, we can track the reduction of variance from $\hat{\beta}_2^{\text{MTL}}$ to $\hat{\beta}_2^{\text{STL}}$ as

$$\begin{aligned} \delta_{\text{var}} &:= \sigma^2 \text{Tr} [(X_2^\top X_2)^{-1}] - (1 + O(n^{-\varepsilon})) \cdot \sigma^2 \text{Tr} [(X_1^\top X_1 + X_2^\top X_2)^{-1}] \\ &= \Delta_{\text{var}} \cdot (1 + O(n^{-\varepsilon})) \end{aligned} \quad (\text{C.10})$$

703 with high probability, where

$$\Delta_{\text{var}} := \sigma^2 \left(\frac{1}{\rho_2 - 1} - \frac{1}{\rho_1 + \rho_2} \cdot \frac{1}{a_1 + a_2} \right) = \sigma^2 \cdot \frac{\rho_1}{(\rho_2 - 1)(\rho_1 + \rho_2 - 1)}.$$

Next for the model shift bias

$$\delta_{\text{bias}} := (1 + O(n^{-\varepsilon})) \cdot d^2 \text{Tr} [(X_1^\top X_1 + X_2^\top X_2)^{-2} (X_1^\top X_1)^2],$$

704 we can get from Lemma B.5 (or rather the proof of it) that

$$\alpha_-^2(\rho_1) - o(1) \leq \frac{\delta_{\text{bias}}}{\Delta_{\text{bias}}} \leq \alpha_+^2(\rho_1) + o(1), \quad (\text{C.11})$$

where

$$\Delta_{\text{bias}} := pd^2 \cdot \frac{\rho_1^2}{(\rho_1 + \rho_2)^2} \frac{1 + a_3 + a_4}{(a_1 + a_2)^2} = pd^2 \cdot \frac{\rho_1^2(\rho_1 + \rho_2)}{(\rho_1 + \rho_2 - 1)^3}.$$

705 Note that

$$L(\hat{\beta}_2^{\text{STL}}) - L(\hat{\beta}_2^{\text{MTL}}) = \delta_{\text{var}} - \delta_{\text{bias}}. \quad (\text{C.12})$$

706 Then we can track its sign using (C.10) and (C.11).

707 **Positive transfer.** With (C.10) and (C.11), we conclude that if

$$pd^2 \cdot \frac{\rho_1^2(\rho_1 + \rho_2)}{(\rho_1 + \rho_2 - 1)^3} \cdot (\alpha_+^2(\rho_1) + o(1)) < \sigma^2 \cdot \frac{\rho_1}{(\rho_2 - 1)(\rho_1 + \rho_2 - 1)}, \quad (\text{C.13})$$

708 we have that $\delta_{\text{var}} > \delta_{\text{bias}}$, which implies $L(\hat{\beta}_2^{\text{MTL}}) < L(\hat{\beta}_2^{\text{STL}})$. We can simplify (C.13) to

$$\frac{pd^2}{\sigma^2} < \Phi(\rho_1, \rho_2) \cdot (\alpha_+^2(\rho_1) + o(1))^{-1}, \quad (\text{C.14})$$

709 Since $\Psi(\beta_1, \beta_2) = pd^2/\sigma^2$, it gives the first statement of Proposition 3.3.

710 **Negative transfer.** On the other hand, if

$$pd^2 \cdot \frac{\rho_1^2(\rho_1 + \rho_2)}{(\rho_1 + \rho_2 - 1)^3} \cdot (\alpha_-^2(\rho_1) - o(1)) > \sigma^2 \cdot \frac{\rho_1}{(\rho_2 - 1)(\rho_1 + \rho_2 - 1)}, \quad (\text{C.15})$$

711 we have that $\delta_{\text{var}} < \delta_{\text{bias}}$, which implies $L(\hat{\beta}_2^{\text{MTL}}) > L(\hat{\beta}_2^{\text{STL}})$. We can simplify (C.15) to

$$\Psi(\beta_1, \beta_2) = \frac{pd^2}{\sigma^2} > \Phi(\rho_1, \rho_2) \cdot (\alpha_-^2(\rho_1) - o(1))^{-1}, \quad (\text{C.16})$$

712 which gives the second statement of Proposition 3.3. \square

713 Next we consider the case where the two tasks have different noise variances $\sigma_1^2 \neq \sigma_2^2$. In particular,
714 we show Proposition C.2, which gives a transition threshold with respect to the difference between
715 the noise levels of the two tasks.

716 **Proposition C.2.** *In the isotropic model, assume that $\rho_1 > 40$ and $\mathbb{E}[\|\beta_1 - \beta_2\|^2] < \frac{1}{2}\sigma_2^2 \cdot \Phi(\rho_1, \rho_2)$.
717 Then we have the following transition with respect to σ_1^2 :*

- 718 • If $\sigma_1^2 < -\gamma_+^{1/2} \rho_1 \cdot pd^2 + \left(1 + \gamma_+^{-1/2} \rho_1 \Phi(\rho_1, \rho_2)\right) \cdot \sigma_2^2$, then whp $L(\hat{\beta}_2^{\text{MTL}}) < L(\hat{\beta}_2^{\text{STL}})$.
- 719 • If $\sigma_1^2 > -\gamma_-^{1/2} \rho_1 \cdot pd^2 + \left(1 + \gamma_-^{-1/2} \rho_1 \Phi(\rho_1, \rho_2)\right) \cdot \sigma_2^2$, then whp $L(\hat{\beta}_2^{\text{MTL}}) > L(\hat{\beta}_2^{\text{STL}})$.

720 As a corollary, if $\sigma_1^2 \leq \sigma_2^2$, then we always get positive transfer.

721 *Proof of Proposition C.2.* In the setting of Proposition C.2, the test loss is given by (A.5). In the
722 isotropic model, using again the concentration of random vector with i.i.d. entries, Lemma E.14, we
723 can rewrite $L(\hat{\beta}_2^{\text{MTL}})$ as

$$\begin{aligned} L(\hat{\beta}_2^{\text{MTL}}) = & \hat{v}^2 \left[d^2 + (\hat{v} - 1)^2 \kappa^2 \right] \text{Tr} \left[(\hat{v}^2 X_1^\top X_1 + X_2^\top X_2)^{-2} (X_1^\top X_1)^2 \right] \cdot \left(1 + O(p^{-1/2+\epsilon}) \right) \\ & + \sigma_2^2 \cdot \text{Tr} \left[(\hat{v}^2 X_1^\top X_1 + X_2^\top X_2)^{-1} \right] + (\sigma_1^2 - \sigma_2^2) \cdot \text{Tr} \left[(\hat{v}^2 X_1^\top X_1 + X_2^\top X_2)^{-2} \hat{v}^2 X_1^\top X_1 \right] \end{aligned}$$

724 with high probability for any constant $\epsilon > 0$.

725 In the current setting, we can also show that (C.5) holds for \hat{v} . Since the proof is almost the same as
726 the one for Lemma C.1, we omit the details. Thus under the choice parameters in (C.4), $L(\hat{\beta}_2^{\text{MTL}})$
727 can be simplified as in (C.7):

$$\begin{aligned} L(\hat{\beta}_2^{\text{MTL}}) = & (1 + O(n^{-\epsilon})) \cdot d^2 \text{Tr} \left[(X_1^\top X_1 + X_2^\top X_2)^{-2} (X_1^\top X_1)^2 \right] \\ & + (1 + O(n^{-\epsilon})) \cdot \sigma_2^2 \text{Tr} \left[(X_1^\top X_1 + X_2^\top X_2)^{-1} \right] \\ & + (1 + O(n^{-\epsilon})) \cdot (\sigma_1^2 - \sigma_2^2) \text{Tr} \left[(X_1^\top X_1 + X_2^\top X_2)^{-2} X_1^\top X_1 \right]. \end{aligned} \quad (\text{C.17})$$

Then we write

$$L(\hat{\beta}_2^{\text{STL}}) - L(\hat{\beta}_2^{\text{MTL}}) = \delta_{\text{var}} - \delta_{\text{bias}} - \delta_{\text{var}}^{(2)},$$

where

$$\delta_{\text{var}} := \sigma_2^2 \text{Tr} \left[(X_2^\top X_2)^{-1} \right] - (1 + O(n^{-\epsilon})) \cdot \sigma_2^2 \text{Tr} \left[(X_1^\top X_1 + X_2^\top X_2)^{-1} \right]$$

satisfies (C.10) but with σ^2 replaced with σ_2^2 ,

$$\delta_{\text{bias}} := (1 + O(n^{-\epsilon})) \cdot d^2 \text{Tr} \left[(X_1^\top X_1 + X_2^\top X_2)^{-2} (X_1^\top X_1)^2 \right]$$

728 satisfies (C.11), and

$$\delta_{\text{var}}^{(2)} := (1 + O(n^{-\epsilon})) \cdot (\sigma_1^2 - \sigma_2^2) \text{Tr} \left[(X_1^\top X_1 + X_2^\top X_2)^{-2} X_1^\top X_1 \right].$$

729 To estimate this new term $\delta_{\text{var}}^{(2)}$, we use the same arguments as in the proof of Lemma B.5: we first
 730 replace $X_1^\top X_1$ with $n_1 \text{Id}$ up to some error using (B.16), and then apply Lemma B.3 to calculate
 731 $\text{Tr}[(X_1^\top X_1 + X_2^\top X_2)^{-2}]$. This process leads to the following estimates on $\delta_{\text{var}}^{(2)}$:

$$\alpha_-(\rho_1) - o(1) \leq \frac{\delta_{\text{var}}^{(2)}}{\Delta_{\text{var}}^{(2)}} \leq \alpha_+(\rho_1) + o(1), \quad (\text{C.18})$$

where

$$\Delta_{\text{var}}^{(2)} := (\sigma_1^2 - \sigma_2^2) \frac{\rho_1(\rho_1 + \rho_2)}{(\rho_1 + \rho_2 - 1)^3}.$$

732 Next we compare δ_{var} with $\delta_{\text{bias}} + \delta_{\text{var}}^{(2)}$. Our main focus is to see how the extra $\delta_{\text{var}}^{(2)}$ affects the
 733 information transfer in this case.

734 Note that the condition $\mathbb{E}[\|\beta_1 - \beta_2\|^2] < \frac{1}{2}\sigma_2^2 \cdot \Phi(\rho_1, \rho_2)$ for $\rho_1 > 40$ means that we have $\delta_{\text{var}} > \delta_{\text{bias}}$
 735 by Proposition 3.3. Hence if $\sigma_1^2 \leq \sigma_2^2$, then $\delta_{\text{var}}^{(2)} < 0$ and we always have $\delta_{\text{var}} > \delta_{\text{bias}} + \delta_{\text{var}}^{(2)}$, which
 736 gives $L(\hat{\beta}_2^{\text{MTL}}) < L(\hat{\beta}_2^{\text{STL}})$. It remains to consider the case $\sigma_1^2 \geq \sigma_2^2$.

737 **Positive transfer.** By (C.10), (C.11) and (C.18), if the following inequality holds,

$$\begin{aligned} & \sigma_2^2 \cdot \frac{\rho_1}{(\rho_2 - 1)(\rho_1 + \rho_2 - 1)} \cdot (1 - o(1)) \\ & > pd^2 \cdot \frac{\rho_1^2(\rho_1 + \rho_2)}{(\rho_1 + \rho_2 - 1)^3} \alpha_+^2(\rho_1) + (\sigma_1^2 - \sigma_2^2) \cdot \frac{\rho_1(\rho_1 + \rho_2)}{(\rho_1 + \rho_2 - 1)^3} \alpha_+(\rho_1), \end{aligned} \quad (\text{C.19})$$

738 then we have $\delta_{\text{var}} > \delta_{\text{bias}} + \delta_{\text{var}}^{(2)}$ whp, which gives $L(\hat{\beta}_t^{\text{MTL}}) < L(\hat{\beta}_t^{\text{STL}})$. We can solve (C.19) to get

$$\sigma_1^2 < -pd^2 \cdot \rho_1 \alpha_+(\rho_1) + \sigma_2^2 [1 + \rho_1 \Phi(\rho_1, \rho_2) \alpha_+^{-1}(\rho_1)] \cdot (1 - o(1)).$$

739 This proves the first claim of Proposition C.2 for positive transfer.

740 **Negative transfer.** On the other hand, if the following inequality holds,

$$\begin{aligned} & \sigma_2^2 \cdot \frac{\rho_1}{(\rho_2 - 1)(\rho_1 + \rho_2 - 1)} \cdot (1 + o(1)) \\ & < pd^2 \cdot \frac{\rho_1^2(\rho_1 + \rho_2)}{(\rho_1 + \rho_2 - 1)^3} \alpha_-^2(\rho_1) + (\sigma_1^2 - \sigma_2^2) \cdot \frac{\rho_1(\rho_1 + \rho_2)}{(\rho_1 + \rho_2 - 1)^3} \alpha_-(\rho_1), \end{aligned} \quad (\text{C.20})$$

741 then we have $\delta_{\text{var}} < \delta_{\text{bias}} + \delta_{\text{var}}^{(2)}$ whp, which gives $L(\hat{\beta}_t^{\text{MTL}}) > L(\hat{\beta}_t^{\text{STL}})$. We can solve (C.20) to get

$$\sigma_1^2 > -pd^2 \cdot \rho_1 \alpha_-(\rho_1) + \sigma_2^2 [1 + \rho_1 \Phi(\rho_1, \rho_2) \alpha_-^{-1}(\rho_1)] \cdot (1 + o(1)).$$

742 This proves the second claim of Proposition C.2 for negative transfer. \square

743 C.2 Missing Proofs of Section 3.3

744 We first prove Proposition 3.4, which describes the effect of source/task data ratio on the information
 745 transfer.

746 *Proof of Proposition 3.4.* Following the above proof of Proposition 3.3, we see that $L(\hat{\beta}_2^{\text{MTL}}) <$
 747 $L(\hat{\beta}_2^{\text{STL}})$ whp if (C.14) holds, while $L(\hat{\beta}_2^{\text{MTL}}) > L(\hat{\beta}_2^{\text{STL}})$ whp if (C.16) holds.

748 We first explain the meaning of the condition

$$\Psi(\beta_1, \beta_2) > 2/(\rho_2 - 1). \quad (\text{C.21})$$

Notice that the function

$$\Phi(\rho_1, \rho_2) = \frac{(\rho_1 + \rho_2 - 1)^2}{\rho_1(\rho_1 + \rho_2)(\rho_2 - 1)} = \frac{1}{\rho_2 - 1} \left(1 + \frac{\rho_2 - 2}{\rho_1} + \frac{1}{\rho_1(\rho_1 + \rho_2)} \right)$$

749 is strictly decreasing with respect to ρ_1 as long as $\rho_2 > 2$, and $\Phi(\rho_1, \rho_2)$ converges to $(\rho_2 - 1)^{-1}$
 750 as $\rho_1 \rightarrow \infty$. Moreover, we notice that $(\alpha_-^2(\rho_1) - o(1))^{-1} < 2$ for $\rho_1 > 40$. Hence (C.21) implies
 751 that (C.16) holds for all large enough ρ_1 . The transition from positive transfer when ρ_1 is small to
 752 negative transfer when ρ_1 is large is described by the two bounds in Proposition 3.4.

753 The two bounds follows directly from (C.14) and (C.16). We will use the following trivial inequalities
754

$$\frac{(\rho_2 - 1)\rho_1}{\rho_1 + \rho_2 - 2} \cdot \left(1 - \frac{1}{(\rho_1 + \rho_2 - 2)^2}\right) \leq \Phi(\rho_1, \rho_2) \leq \frac{(\rho_2 - 1)\rho_1}{\rho_1 + \rho_2 - 2}. \quad (\text{C.22})$$

755 **Positive transfer.** With (C.22), we see that (C.14) is implied by the following inequality:

$$\Psi(\beta_1, \beta_2) \cdot \frac{(\rho_2 - 1)\rho_1}{\rho_1 + \rho_2 - 2} < \gamma_+^{-1}. \quad (\text{C.23})$$

756 then we can solve (C.23) to get

$$\rho_1 < \frac{\rho_2 - 2}{\Psi(\beta_1, \beta_2) \cdot \gamma_+(\rho_2 - 1) - 1}. \quad (\text{C.24})$$

757 This gives the first statement of Proposition 3.4.

Note that if we require the RHS of (C.24) to be larger than 40, that is, (C.24) is not a null condition. Then together with (C.21), we get

$$\rho_2 - 2 > (2\gamma_+ - 1)\rho_1.$$

758 Plugging into $\rho_1 > 40$, we get $\rho_2 \geq 106$. This gives a constraint on ρ_2 .

759 **Negative transfer.** With (C.22), we see that (C.16) is implied by the following inequality:

$$\Psi(\beta_1, \beta_2) \cdot \frac{(\rho_2 - 1)\rho_1}{\rho_1 + \rho_2 - 2} \left(1 - \frac{1}{(\rho_1 + \rho_2 - 2)^2}\right) > \Psi(\beta_1, \beta_2) \cdot \frac{(\rho_2 - 1.5)\rho_1}{\rho_1 + \rho_2 - 2} > \gamma_-^{-1}. \quad (\text{C.25})$$

760 where we used $(1 - (\rho_1 + \rho_2 - 2)^{-2})(\rho_2 - 1) > \rho_2 - 1.5$ for $\rho_1 > 40$ and $\rho_2 > 110$. Then we can
761 solve (C.25) to get

$$\rho_1 > \frac{(\rho_2 - 2)\sigma^2}{\Psi(\beta_1, \beta_2) \cdot \gamma_-(\rho_2 - 1.5) - 1}, \quad (\text{C.26})$$

762 which gives the second statement of Proposition 3.4. We remark that condition (C.21) implies

763 $\Psi(\beta_1, \beta_2) \cdot \gamma_-(\rho_2 - 1.5) > 1$, so (C.26) does not give a trivial bound. \square

764 Next we state Proposition C.3, which gives precise upper and lower bounds on the data efficiency
765 ratio for taskonomy.

766 **Proposition C.3** (Labeled data efficiency). *In the isotropic model, assume that $\rho_1, \rho_2 \geq 9$ and*
767 *$\Psi(\beta_1, \beta_2) < (5(\rho_1 - 1))^{-1} + (5(\rho_2 - 1))^{-1}$. Then the data efficiency ratio x^* satisfies*

$$x_l \leq x^* \leq \frac{1}{\rho_1 + \rho_2} \left(\frac{2}{(\rho_1 - 1)^{-1} + (\rho_2 - 1)^{-1} - 5\Psi(\beta_1, \beta_2)} + 1 \right), \quad (\text{C.27})$$

where we denoted

$$x_l := \frac{1}{\rho_1 + \rho_2} \left(\frac{2}{(\rho_1 - 1)^{-1} + (\rho_2 - 1)^{-1}} + 1 \right).$$

768 *Proof of Proposition C.3.* Suppose we have reduced number of datapoints— xn_1 for task 1 and xn_2
769 for task 2 with $n_1 = \rho_1 p$ and $n_2 = \rho_2 p$. Then all the results in the proof of Proposition 3.3 still hold,
770 except that we need to replace (ρ_1, ρ_2) with $(x\rho_1, x\rho_2)$. More precisely, we have

$$\begin{aligned} a_1 &= \frac{\rho_1(x\rho_1 + x\rho_2 - 1)}{x(\rho_1 + \rho_2)^2}, & a_2 &= \frac{\rho_2(x\rho_1 + \alpha\rho_2 - 1)}{x(\rho_1 + \rho_2)^2}, \\ a_3 &= \frac{\rho_2}{(\rho_1 + \rho_2)(x\rho_1 + x\rho_2 - 1)}, & a_4 &= \frac{\rho_1}{(\rho_1 + \rho_2)(x\rho_1 + x\rho_2 - 1)}. \end{aligned}$$

771 Moreover, with high probability,

$$L_i(\hat{\beta}_i^{\text{MTL}}(x)) = \frac{\sigma^2}{x(\rho_1 + \rho_2) - 1} (1 + o(1)) + \delta_{\text{bias}}^{(i)}, \quad i = 1, 2. \quad (\text{C.28})$$

772 Here the model shift biases $\delta_{\text{bias}}^{(i)}$ satisfy that

$$\alpha_-^2(\alpha\rho_i) - o(1) \leq \delta_{\text{bias}}^{(i)} / \Delta_{\text{bias}}^{(i)} \leq \alpha_+^2(\alpha\rho_i) + o(1), \quad i = 1, 2,$$

773 where $\Delta_{\text{bias}}^{(i)}$ are defined as

$$\Delta_{\text{bias}}^{(i)} := pd^2 \frac{(x\rho_i)^2 \cdot x(\rho_1 + \rho_2)}{[x(\rho_1 + \rho_2) - 1]^3}, \quad i = 1, 2, .$$

774 On the other hand, using Lemma B.1 we have whp,

$$L_i(\hat{\beta}_i^{\text{STL}}) = \frac{\sigma^2}{\rho_i - 1} (1 + o(1)), \quad i = 1, 2. \quad (\text{C.29})$$

Comparing (C.28) and (C.29), we immediately obtain the lower bound $x^* \geq x_l$. In fact, one can see that if $x < x_l$, then we have

$$\frac{2\sigma^2}{x(\rho_1 + \rho_2) - 1} > \frac{\sigma^2}{\rho_1 - 1} + \frac{\sigma^2}{\rho_2 - 1},$$

775 that is, $L_1(\hat{\beta}(\alpha)) + L_2(\hat{\beta}(\alpha))$ is larger than $L_1(\hat{\beta}_l^{\text{STL}}) + L_2(\hat{\beta}_l^{\text{STL}})$ even if we do not take into account
776 the model shift bias terms $\delta_{\text{bias}}^{(i)}$.

777 Then we try to obtain an upper bound on x^* . In the following discussions, we only consider x such
778 that $x > x_l$. In particular, we have $x\rho > x_l\rho \geq \min(\rho_1, \rho_2)$, where we abbreviated $\rho := \rho_1 + \rho_2$.

779 **The upper bound.** From (C.28) and (C.29), we see that $x^* \leq x$ if x satisfies

$$(1 + o(1)) \cdot \sum_{i=1}^2 pd^2 \frac{(x\rho_i)^2 \cdot x\rho}{(x\rho - 1)^3} \left(1 + \sqrt{\frac{1}{x\rho_i}}\right)^4 \leq \frac{\sigma^2}{\rho_1 - 1} + \frac{\sigma^2}{\rho_2 - 1} - \frac{2\sigma^2}{x\rho - 1}.$$

780 We rewrite the inequality as

$$(1 + o(1)) \cdot \frac{\Psi(\beta_1, \beta_2)}{[1 - (x\rho)^{-1}]^3} \sum_{i=1}^2 \left(\sqrt{\frac{\rho_i}{\rho}} + \sqrt{\frac{1}{x\rho}}\right)^4 \leq \frac{1}{\rho_1 - 1} + \frac{1}{\rho_2 - 1} - \frac{2}{x\rho - 1}. \quad (\text{C.30})$$

With $x\rho \geq \min(\rho_1, \rho_2) > 9$, we can get the simple bound

$$\frac{1 + o(1)}{[1 - (x\rho)^{-1}]^3} \sum_{i=1}^2 \left(\sqrt{\frac{\rho_i}{\rho}} + \sqrt{\frac{1}{x\rho}}\right)^4 < 5.$$

781 Inserting it into (C.30), we can solve for the upper bound in (C.27).

We can get better bounds if the values of ρ_1 and ρ_2 increase. For example, if we consider the case $\min(\rho_1, \rho_2) \geq 100$, then with some basic calculations, one can show that in this case

$$\frac{1}{[1 - (x\rho)^{-1}]^3} \sum_{i=1}^2 \left(\sqrt{\frac{\rho_i}{\rho}} + \sqrt{\frac{1}{x\rho}}\right)^4 < \frac{\rho_1^2 + \rho_2^2}{\rho^2} + 0.52.$$

782 Thus the following inequality implies (C.30):

$$\left(\frac{\rho_1^2 + \rho_2^2}{\rho^2} + 0.52\right) \Psi(\beta_1, \beta_2) < \frac{1}{\rho_1 - 1} + \frac{1}{\rho_2 - 1} - \frac{2}{x\rho - 1},$$

783 from which we can solve for the following upper bound on α^* :

$$\alpha^* < \frac{1}{\rho} \frac{2}{\frac{1}{\rho_1 - 1} + \frac{1}{\rho_2 - 1} - \left(\frac{\rho_1^2 + \rho_2^2}{\rho^2} + 0.52\right) \Psi(\beta_1, \beta_2)} + \frac{1}{\rho}.$$

784 Similarly, we can get a better lower bound. From (C.28) and (C.29), we see that $x^* \geq x$ if x satisfies

$$(1 - o(1)) \cdot \frac{\Psi(\beta_1, \beta_2)}{[1 - (x\rho)^{-1}]^3} \sum_{i=1}^2 \left(\sqrt{\frac{\rho_i}{\rho}} - \sqrt{\frac{1}{x\rho}}\right)^4 \geq \frac{1}{\rho_1 - 1} + \frac{1}{\rho_2 - 1} - \frac{2}{x\rho - 1}. \quad (\text{C.31})$$

Then in the case $\min(\rho_1, \rho_2) \geq 100$, with some basic calculations, one can show that the sum on the left-hand side of (C.31) satisfies

$$\frac{1}{[1 - (x\rho)^{-1}]^3} \sum_{i=1}^2 \left(\sqrt{\frac{\rho_i}{\rho}} - \sqrt{\frac{1}{x\rho}}\right)^4 > \frac{\rho_1^2 + \rho_2^2}{\rho^2} - 0.33.$$

Thus the following inequality implies (C.31):

$$\left(\frac{\rho_1^2 + \rho_2^2}{\rho^2} - 0.33\right) p d^2 > \frac{\sigma^2}{\rho_1 - 1} + \frac{\sigma^2}{\rho_2 - 1} - \frac{2\sigma^2}{x\rho - 1}, \quad (\text{C.32})$$

from which we can solve for the following lower bound on x^* :

$$x^* > \frac{1}{\rho} \frac{2}{\frac{1}{\rho_1 - 1} + \frac{1}{\rho_2 - 1} - \left(\frac{\rho_1^2 + \rho_2^2}{\rho^2} - 0.33\right) \Psi(\beta_1, \beta_2)} + \frac{1}{\rho}.$$

This gives a lower bound above x_l . \square

C.3 Missing Proofs of Section 3.4

We now prove Proposition 3.5, which shows that $L(\hat{\beta}^{\text{MTL}})$ is minimized approximately when M is a scalar matrix where there is enough source data.

Proof of Proposition 3.5. Let

$$M_0 := \arg \min_{M \in \mathcal{S}_\mu} g(M).$$

We now calculate $g(M_0)$. With the same arguments as in Lemma C.1 we can show that (C.5) holds. Moreover, if the parameters are chosen such that $p^{-1+c_0}\sigma^2 \leq \kappa^2 \leq p^{-\varepsilon_0-c_0}\sigma^2$ as in (C.4), we can simplify

$$g(M_0) = (1 + O(p^{-\varepsilon})) \cdot \sigma^2 \text{Tr} [\Sigma_2(X_1^\top X_1 + X_2^\top X_2)^{-1}],$$

with high probability for some constant $\varepsilon > 0$. In fact, Lemma C.1 was proved assuming that $M = \text{Id}$, but its proof can be easily extended to the case with general $M \in \mathcal{S}_\mu$ by using that $\mu_{\min} \leq \lambda_p(M) \leq \lambda_1(M) \leq \mu_{\max}$. We omit the details here.

Now using Lemma B.2, we obtain that with high probability,

$$g(M_0) = \frac{\sigma^2}{\rho_1 + \rho_2} \cdot \frac{1}{p} \text{Tr} \left(\frac{1}{a_1(M_0) \cdot M_0^\top M_0 + a_2(M_0)} \right) \cdot (1 + O(p^{-\varepsilon})). \quad (\text{C.33})$$

From equation (B.3), it is easy to obtain the following estimates on $a_1(M)$ and $a_2(M)$ for any $M \in \mathcal{S}_\mu$:

$$\frac{\rho_1 - 1}{\rho_1 + \rho_2} < a_1(M) < \frac{\rho_1 + \rho_2 - 1}{\rho_1 + \rho_2}, \quad a_2(M) < \frac{\rho_2}{\rho_1 + \rho_2}. \quad (\text{C.34})$$

Inserting (C.34) into (C.33) and using $\lambda(M_0^\top M_0) \geq \mu_{\min}^2$, we obtain that with high probability,

$$\left(1 + \frac{\rho_2}{(\rho_1 - 1)\mu_{\min}^2}\right)^{-1} h(M_0) \cdot (1 - O(p^{-\varepsilon})) \leq g(M_0) \leq h(M_0) \cdot (1 + O(p^{-\varepsilon})), \quad (\text{C.35})$$

where

$$h(M_0) := \frac{\sigma^2}{(\rho_1 + \rho_2)a_1(M_0)} \cdot \frac{1}{p} \text{Tr} \left(\frac{1}{M_0^\top M_0} \right).$$

By AM-GM inequality, we observe that

$$\text{Tr} \left(\frac{1}{M^\top M} \right) = \sum_{i=1}^p \frac{1}{\lambda_i^2}$$

is minimized when $\lambda_1 = \dots = \lambda_p = \mu$ under the restriction $\prod_{i=1}^p \lambda_i \leq \mu^p$. Hence we get that

$$h(M_0) \leq \frac{\sigma^2}{\mu^2(\rho_1 + \rho_2)a_1(M_0)}. \quad (\text{C.36})$$

On the other hand, when $M = \mu \text{Id}$, applying Lemma B.2 we obtain that with high probability,

$$\begin{aligned} g(\mu \text{Id}) &= \frac{\sigma^2}{\rho_1 + \rho_2} \cdot \frac{1}{p} \text{Tr} \left(\frac{1}{\mu^2 a_1(\mu \text{Id}) + a_2(\mu \text{Id})} \right) \cdot (1 + O(p^{-\varepsilon})) \\ &\leq \frac{\sigma^2}{\mu^2(\rho_1 + \rho_2)a_1(\mu \text{Id})}. \end{aligned} \quad (\text{C.37})$$

Combining (C.34), (C.35), (C.36) and (C.37), we conclude the proof. \square

804 **D Proof of Theorem 3.6**

805 *Proof of Theorem 3.6.* In this setting, we need to study the following loss function:

$$f(B; W_1, \dots, W_t) = \sum_{i=1}^t \|XBW_i - Y_i\|^2. \quad (\text{D.1})$$

806 For any fixed $W_1, W_2, \dots, W_t \in \mathbb{R}^r$, we can derive a closed form solution for B as

$$\begin{aligned} \hat{B}(W_1, \dots, W_t) &= (X^\top X)^{-1} X^\top \left(\sum_{i=1}^t Y_i W_i^\top \right) (\mathcal{W} \mathcal{W}^\top)^{-1} \\ &= (B^* \mathcal{W}^\top) (\mathcal{W} \mathcal{W}^\top)^{-1} + (X^\top X)^{-1} X^\top \left(\sum_{i=1}^t \varepsilon_i W_i^\top \right) (\mathcal{W} \mathcal{W}^\top)^{-1}, \end{aligned}$$

807 where we denote $\mathcal{W} \in \mathbb{R}^{r \times t}$ as $\mathcal{W} = [W_1, W_2, \dots, W_t]$. Then as in (A.6), we pick N independent
808 samples of the training set for each task with $N \geq n^{1-\varepsilon_0}$, and use concentration to get the validation
809 loss as

$$\tilde{f}(\hat{B}; \mathcal{W}) = N [\text{val}(\mathcal{W}) + t\sigma^2] \cdot \left(1 + O(p^{-(1-\varepsilon_0)/2+\varepsilon}) \right). \quad (\text{D.2})$$

Here $\text{val}(\mathcal{W})$ is defined as

$$\text{val}(\mathcal{W}) := \mathbb{E}_{\varepsilon_j, \forall 1 \leq j \leq t} \left[\sum_{i=1}^t \left\| \Sigma^{1/2} (\hat{B} W_i - \beta_i) \right\|^2 \right] = \delta_{\text{bias}}(\mathcal{W}) + \delta_{\text{var}}(\mathcal{W}),$$

810 where the model shift bias term $\delta_{\text{bias}}(\mathcal{W})$ is given by

$$\delta_{\text{bias}}(\mathcal{W}) := \sum_{i=1}^t \left\| \Sigma^{1/2} ((B^* \mathcal{W}^\top) (\mathcal{W} \mathcal{W}^\top)^{-1} W_i - \beta_i) \right\|^2,$$

811 and the variance term $\delta_{\text{var}}(\mathcal{W})$ can be calculated as

$$\delta_{\text{var}}(\mathcal{W}) := \sigma^2 \cdot \text{Tr} [\Sigma (X^\top X)^{-1}].$$

812 It suffices to minimize $\delta_{\text{bias}}(\mathcal{W})$ over \mathcal{W} , since both $tN\sigma^2$ and $\delta_{\text{var}}(\mathcal{W})$ do not depend on the weights.

813 We denote $Q := \mathcal{W}^\top (\mathcal{W} \mathcal{W}^\top)^{-1} \mathcal{W} \in \mathbb{R}^{k \times k}$, whose (i, j) -th entry is equal to $W_i^\top (\mathcal{W} \mathcal{W}^\top)^{-1} W_j$.

814 Now we can write $\delta_{\text{bias}}(\mathcal{W})$ succinctly as

$$\delta_{\text{bias}}(\mathcal{W}) = \left\| \Sigma^{1/2} B^* (Q - \text{Id}) \right\|_F^2.$$

815 From this equation we can solve the minimizer optimally as $Q_0 = U_r U_r^\top$. On the other hand, let $\hat{\mathcal{W}}$

816 be the minimizer of \tilde{f} , and denote $\hat{Q} := \hat{\mathcal{W}}^\top (\hat{\mathcal{W}} \hat{\mathcal{W}}^\top)^{-1} \hat{\mathcal{W}}$. We claim that \hat{Q} satisfies

$$\|Q_0^{-1} \hat{Q} - \text{Id}\| = o(1) \quad \text{whp.} \quad (\text{D.3})$$

In fact, if (D.3) does not hold, then using the condition $\lambda_{\min}((B^*)^\top \Sigma B^*) \gtrsim \sigma^2$ and that $\delta_{\text{var}}(\mathcal{W}) = O(\sigma^2)$ by (B.16), we obtain that

$$\text{val}(\hat{\mathcal{W}}) + t\sigma^2 > (\text{val}(\mathcal{W}_0) + t\sigma^2) \cdot (1 + o(1)) \Rightarrow \tilde{f}(\hat{B}; \hat{\mathcal{W}}) > \tilde{f}(\hat{B}; \mathcal{W}_0),$$

817 that is, $\hat{\mathcal{W}}$ is not a minimizer. This leads to a contradiction.

818 In sum, we have solved that $\hat{\beta}_i^{\text{MTL}} = B^* (U_r U_r(i) + o(1))$. Inserting it into the definition of the test
819 loss, we get that

$$\begin{aligned} L(\hat{\beta}_t^{\text{MTL}}) &= \left\| \Sigma^{1/2} \left((B^* \hat{\mathcal{W}}^\top) (\hat{\mathcal{W}} \hat{\mathcal{W}}^\top)^{-1} \hat{W}_t - \beta_2 \right) \right\|^2 + \sigma^2 \hat{W}_t^\top (\hat{\mathcal{W}} \hat{\mathcal{W}}^\top)^{-1} \hat{W}_t \cdot \text{Tr} [\Sigma (X^\top X)^{-1}] \\ &= \left\| \Sigma^{1/2} (B^* U_r U_r(t) - \beta_2) \right\|^2 + o(\|B^*\|^2) + \sigma^2 \|U_r(t)\|^2 \text{Tr} [\Sigma (X^\top X)^{-1}] \cdot (1 + o(1)) \\ &= \left\| \Sigma^{1/2} (B^* U_r U_r(t) - \beta_2) \right\|^2 + \frac{\sigma^2}{\rho - 1} \|U_r(t)\|^2 + o(\|B^*\|^2 + \sigma^2), \end{aligned}$$

with high probability, where we used Lemma B.1 in the last step. Similar, by Lemma B.1 we have

$$L(\hat{\beta}_t^{\text{MTL}}) = \frac{\sigma^2}{\rho - 1} \cdot (1 + o(1)).$$

820 Combining the above two estimates, we conclude the proof. \square

E Proof of Lemma B.2 and Lemma B.3

In random matrix theory, it is more convenient to rescale the matrices Z_1 and Z_2 such that their entries have variance n^{-1} , where $n := n_1 + n_2$. The advantage of this scaling is that the singular eigenvalues of Z_1 and Z_2 all lie in a bounded support that does grow with n .

E.1 Notations and basic tools

We denote the two sample covariance matrices by $\mathcal{Q}_1 := X_1^\top X_1$ and $\mathcal{Q}_2 := X_2^\top X_2$. We assume that $Z_1 = (z_{ij}^{(1)})$ and $Z_2 = (z_{ij}^{(2)})$ are $n_1 \times p$ and $n_2 \times p$ random matrices with i.i.d. entries satisfying

$$\mathbb{E} z_{ij}^{(\alpha)} = 0, \quad \mathbb{E} |z_{ij}^{(\alpha)}|^2 = n^{-1}. \quad (\text{E.1})$$

Moreover, we assume that the fourth moments exist:

$$\mathbb{E} |\sqrt{n} z_{ij}^{(\alpha)}|^4 \leq C \quad (\text{E.2})$$

for some constant $C > 0$. Let $0 < \tau < 1$ be a small constant. We assume that the aspect ratios $d_1 := p/n_1$ and $d_2 := p/n_2$ satisfy that

$$0 \leq d_1 \leq \tau^{-1}, \quad 1 + \tau \leq d_2 \leq \tau^{-1}. \quad (\text{E.3})$$

Here the lower bound $1 + \tau \leq d_2$ is to ensure that the sample covariance matrix \mathcal{Q}_2 is non-singular with high probability; see Lemma E.10 below.

We assume that Σ_1 and Σ_2 have eigendecompositions

$$\Sigma_1 = O_1 \Lambda_1 O_1^\top, \quad \Sigma_2 = O_2 \Lambda_2 O_2^\top, \quad \Lambda_1 = \text{diag}(\sigma_1^{(1)}, \dots, \sigma_n^{(1)}), \quad \Lambda_2 = \text{diag}(\sigma_1^{(2)}, \dots, \sigma_N^{(2)}), \quad (\text{E.4})$$

where the eigenvalues satisfy that

$$\tau^{-1} \geq \sigma_1^{(1)} \geq \sigma_2^{(1)} \geq \dots \geq \sigma_p^{(1)} \geq 0, \quad \tau^{-1} \geq \sigma_1^{(2)} \geq \sigma_2^{(2)} \geq \dots \geq \sigma_p^{(2)} \geq \tau. \quad (\text{E.5})$$

We assume that $M = \Sigma_1^{1/2} \Sigma_2^{-1/2}$ has singular value decomposition

$$M = U \Lambda V^\top, \quad \Lambda = \text{diag}(\sigma_1, \dots, \sigma_p), \quad (\text{E.6})$$

where the singular values satisfy that

$$\tau \leq \sigma_p \leq \sigma_1 \leq \tau^{-1}. \quad (\text{E.7})$$

We summarize our basic assumptions here for future reference. Note that this assumption is in accordance with Assumption A.1, except that we rescale the entries of Z_1 and Z_2 here.

Assumption E.1. We assume that Z_1 and Z_2 are independent $n_1 \times p$ and $n_2 \times p$ random matrices with real i.i.d. entries satisfying (E.1) and (E.2), Σ_1 and Σ_2 are deterministic non-negative definite symmetric matrices satisfying (E.4)-(E.7), and $d_{1,2}$ satisfy (E.3).

We will use the following notion of stochastic domination, which was first introduced in [47] and subsequently used in many works on random matrix theory. It simplifies the presentation of the results and their proofs by systematizing statements of the form “ ξ is bounded by ζ with high probability up to a small power of n ”.

Definition E.2 (Stochastic domination). (i) Let

$$\xi = \left(\xi^{(n)}(u) : n \in \mathbb{N}, u \in U^{(n)} \right), \quad \zeta = \left(\zeta^{(n)}(u) : n \in \mathbb{N}, u \in U^{(n)} \right)$$

be two families of nonnegative random variables, where $U^{(n)}$ is a possibly n -dependent parameter set. We say ξ is stochastically dominated by ζ , uniformly in u , if for any fixed (small) $\varepsilon > 0$ and (large) $D > 0$,

$$\sup_{u \in U^{(n)}} \mathbb{P} \left[\xi^{(n)}(u) > N^\varepsilon \zeta^{(n)}(u) \right] \leq N^{-D}$$

for large enough $n \geq n_0(\varepsilon, D)$, and we shall use the notation $\xi \prec \zeta$. If for some complex family ξ we have $|\xi| \prec \zeta$, then we will also write $\xi \prec \zeta$ or $\xi = O_{\prec}(\zeta)$.

(ii) We say an event Ξ holds with high probability if for any constant $D > 0$, $\mathbb{P}(\Xi) \geq 1 - n^{-D}$ for large enough n . We say Ξ holds with high probability on an event Ω if for any constant $D > 0$, $\mathbb{P}(\Omega \setminus \Xi) \leq n^{-D}$ for large enough n .

855 The following lemma collects basic properties of stochastic domination \prec , which will be used tacitly
 856 in the proof.

857 **Lemma E.3** (Lemma 3.2 in [20]). *Let ξ and ζ be families of nonnegative random variables.*

858 (i) *Suppose that $\xi(u, v) \prec \zeta(u, v)$ uniformly in $u \in U$ and $v \in V$. If $|V| \leq n^C$ for some constant C ,
 859 then $\sum_{v \in V} \xi(u, v) \prec \sum_{v \in V} \zeta(u, v)$ uniformly in u .*

860 (ii) *If $\xi_1(u) \prec \zeta_1(u)$ and $\xi_2(u) \prec \zeta_2(u)$ uniformly in $u \in U$, then $\xi_1(u)\xi_2(u) \prec \zeta_1(u)\zeta_2(u)$
 861 uniformly in u .*

862 (iii) *Suppose that $\Psi(u) \geq n^{-C}$ is deterministic and $\xi(u)$ satisfies $\mathbb{E}\xi(u)^2 \leq n^C$. If $\xi(u) \prec \Psi(u)$
 863 uniformly in u , then we also have $\mathbb{E}\xi(u) \prec \Psi(u)$ uniformly in u .*

864 **Definition E.4** (Bounded support condition). *We say a random matrix Z satisfies the bounded
 865 support condition with q , if*

$$\max_{i,j} |x_{ij}| \prec q. \quad (\text{E.8})$$

866 Here $q \equiv q(N)$ is a deterministic parameter and usually satisfies $n^{-1/2} \leq q \leq n^{-\phi}$ for some (small)
 867 constant $\phi > 0$. Whenever (E.8) holds, we say that X has support q .

868 Our main goal is to study the matrix inverse $(\mathcal{Q}_1 + \mathcal{Q}_2)^{-1}$. Using (E.6), we can rewrite it as

$$(\mathcal{Q}_1 + \mathcal{Q}_2)^{-1} = \Sigma_2^{-1/2} V (\Lambda U^\top Z_1^\top Z_1 U \Lambda + V^\top Z_2^\top Z_2 V)^{-1} V^\top \Sigma_2^{-1/2}. \quad (\text{E.9})$$

869 For this purpose, we shall study the following matrix for $z \in \mathbb{C}_+$,

$$\mathcal{G}(z) := (\Lambda U^\top Z_1^\top Z_1 U \Lambda + V^\top Z_2^\top Z_2 V - z)^{-1}, \quad z \in \mathbb{C}_+, \quad (\text{E.10})$$

870 which we shall refer to as resolvent (or Green's function).

871 Next we introduce a convenient self-adjoint linearization trick. It has been proved to be useful in
 872 studying the local laws of random matrices of the Gram type [21, 48, 49]. We define the following
 873 $(p+n) \times (p+n)$ self-adjoint block matrix, which is a linear function of Z_1 and Z_2 :

$$H \equiv H(Z_1, Z_2) := \begin{pmatrix} 0 & \Lambda U^\top Z_1^\top & V^\top Z_2^\top \\ Z_1 U \Lambda & 0 & 0 \\ Z_2 V & 0 & 0 \end{pmatrix}. \quad (\text{E.11})$$

874 Then we define its resolvent (Green's function) as

$$G \equiv G(Z_1, Z_2, z) := \left[H(Z_1, Z_2) - \begin{pmatrix} z I_{p \times p} & 0 & 0 \\ 0 & I_{n_1 \times n_1} & 0 \\ 0 & 0 & I_{n_2 \times n_2} \end{pmatrix} \right]^{-1}, \quad z \in \mathbb{C}_+. \quad (\text{E.12})$$

For simplicity of notations, we define the index sets

$$\mathcal{I}_1 := \llbracket 1, p \rrbracket, \quad \mathcal{I}_2 := \llbracket p+1, p+n_1 \rrbracket, \quad \mathcal{I}_3 := \llbracket p+n_1+1, p+n_1+n_2 \rrbracket, \quad \mathcal{I} := \mathcal{I}_1 \cup \mathcal{I}_2 \cup \mathcal{I}_3.$$

We will consistently use the latin letters $i, j \in \mathcal{I}_1$, greek letters $\mu, \nu \in \mathcal{I}_2 \cup \mathcal{I}_3$, and $\mathbf{a}, \mathbf{b} \in \mathcal{I}$. We label the indices of the matrices according to

$$Z_1 = (z_{\mu i} : i \in \mathcal{I}_1, \mu \in \mathcal{I}_2), \quad Z_2 = (z_{\nu i} : i \in \mathcal{I}_1, \nu \in \mathcal{I}_3).$$

875 Then we denote the $\mathcal{I}_1 \times \mathcal{I}_1$ block of $G(z)$ by $\mathcal{G}_L(z)$, the $\mathcal{I}_1 \times (\mathcal{I}_2 \cup \mathcal{I}_3)$ block by \mathcal{G}_{LR} , the $(\mathcal{I}_2 \cup \mathcal{I}_3) \times \mathcal{I}_1$
 876 block by \mathcal{G}_{RL} , and the $(\mathcal{I}_2 \cup \mathcal{I}_3) \times (\mathcal{I}_2 \cup \mathcal{I}_3)$ block by \mathcal{G}_R . For simplicity, we abbreviate $Y_1 := Z_1 U \Lambda$,
 877 $Y_2 := Z_2 V$ and $W := (Y_1^\top, Y_2^\top)$. By Schur complement formula, one can find that

$$\mathcal{G}_L = (W W^\top - z)^{-1} = \mathcal{G}, \quad \mathcal{G}_{LR} = \mathcal{G}_{RL}^\top = \mathcal{G} W, \quad \mathcal{G}_R = z (W^\top W - z)^{-1}. \quad (\text{E.13})$$

878 Thus a control of G yields directly a control of the resolvent \mathcal{G} . We also introduce the following
 879 random quantities (some partial traces and weighted partial traces):

$$\begin{aligned} m(z) &:= \frac{1}{p} \sum_{i \in \mathcal{I}_1} G_{ii}(z), & m_1(z) &:= \frac{1}{p} \sum_{i \in \mathcal{I}_1} \sigma_i^2 G_{ii}(z), \\ m_2(z) &:= \frac{1}{n_1} \sum_{\mu \in \mathcal{I}_2} G_{\mu\mu}(z), & m_3(z) &:= \frac{1}{n_2} \sum_{\mu \in \mathcal{I}_3} G_{\mu\mu}(z). \end{aligned} \quad (\text{E.14})$$

Our proof will use the spectral decomposition of G . Let $W = \sum_{k=1}^p \sqrt{\lambda_k} \xi_k \zeta_k^\top$ be a singular value decomposition of W , where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0 = \lambda_{p+1} = \dots = \lambda_n$ are the eigenvalues, $\{\xi_k\}_{k=1}^p$ are the left-singular vectors, and $\{\zeta_k\}_{k=1}^n$ are the right-singular vectors. Then using (E.13), we get that for $i, j \in \mathcal{I}_1$ and $\mu, \nu \in \mathcal{I}_2 \cup \mathcal{I}_3$,

$$G_{ij} = \sum_{k=1}^p \frac{\xi_k(i) \xi_k^\top(j)}{\lambda_k - z}, \quad G_{\mu\nu} = z \sum_{k=1}^n \frac{\zeta_k(\mu) \zeta_k^\top(\nu)}{\lambda_k - z}, \quad G_{i\mu} = G_{\mu i} = \sum_{k=1}^p \frac{\sqrt{\lambda_k} \xi_k(i) \zeta_k^\top(\mu)}{\lambda_k - z}. \quad (\text{E.15})$$

E.2 Local laws

We now describe the asymptotic limit of $\mathcal{G}(z)$. First define the deterministic limits of $(m_2(z), m_3(z))$, denoted by $(m_{2c}(z), m_{3c}(z))$, as the (unique) solution to the following system of equations

$$\frac{1}{m_{2c}} = \frac{\gamma_n}{p} \sum_{i=1}^p \frac{\sigma_i^2}{z + \sigma_i^2 r_1 m_{2c} + r_2 m_{3c}} - 1, \quad \frac{1}{m_{3c}} = \frac{\gamma_n}{p} \sum_{i=1}^p \frac{1}{z + \sigma_i^2 r_1 m_{2c} + r_2 m_{3c}} - 1, \quad (\text{E.16})$$

such that $(m_{2c}(z), m_{3c}(z)) \in \mathbb{C}_+^2$ for $z \in \mathbb{C}_+$, where, for simplicity, we introduce the parameters

$$\gamma_n := \frac{p}{n}, \quad r_1 \equiv r_1(n) := \frac{n_1}{n}, \quad r_2 \equiv r_2(n) := \frac{n_2}{n}. \quad (\text{E.17})$$

We then define the matrix limit of $G(z)$ as

$$\Pi(z) := \begin{pmatrix} -(z + r_1 m_{2c} \Lambda^2 + r_2 m_{3c})^{-1} & 0 & 0 \\ 0 & m_{2c}(z) I_{n_1} & 0 \\ 0 & 0 & m_{3c}(z) I_{n_2} \end{pmatrix}. \quad (\text{E.18})$$

In particular, the matrix limit of $\mathcal{G}(z)$ is given by $-(z + r_1 m_{2c} \Lambda^2 + r_2 m_{3c})^{-1}$.

If $z = 0$, then the equations (E.16) are reduced to

$$r_1 b_2 + r_2 b_3 = 1 - \gamma_n, \quad b_2 + \frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^2 b_2}{\sigma_i^2 r_1 b_2 + (1 - \gamma_n - r_1 b_2)} = 1. \quad (\text{E.19})$$

where $b_2 := -m_{2c}(0)$ and $b_3 := -m_{3c}(0)$. Note that the function

$$f(b_2) := b_2 + \frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^2 b_2}{\sigma_i^2 r_1 b_2 + (1 - \gamma_n - r_1 b_2)}$$

is a strictly increasing function on $[0, r_1^{-1}(1 - \gamma_n)]$. Moreover, we have $f(0) = 0 < 1$ and $f(r_1^{-1}(1 - \gamma_n)) = 1 + \gamma_n > 1$. Hence by mean value theorem, there exists a unique solution $b_2 \in (0, r_1^{-1}(1 - \gamma_n))$. Moreover, it is easy to check that $f'(a) = O(1)$ for $a \in [0, r_1^{-1}(1 - \gamma_n)]$, and $f(1) > 1$ if $1 \leq r_1^{-1}(1 - \gamma_n)$. Hence there exists a constant $\tau > 0$, such that

$$r_1 \tau \leq r_1 b_2 < \min\{(1 - \gamma_n) - r_1 \tau, r_1(1 - \tau)\}, \quad \tau < r_2 b_3 \leq 1 - \gamma_n - r_1 \tau. \quad (\text{E.20})$$

For general z around $z = 0$, the existence and uniqueness of the solution $(m_{2c}(z), m_{3c}(z))$ is given by the following lemma. Moreover, we will also include some basic estimates on it.

Lemma E.5. *There exist constants $c_0, C_0 > 0$ depending only on τ in (E.3), (E.5), (E.7) and (E.20) such that the following statements hold. There exists a unique solution to (E.16) under the conditions*

$$|z| \leq c_0, \quad |m_{2c}(z) - m_{2c}(0)| + |m_{3c}(z) - m_{3c}(0)| \leq c_0. \quad (\text{E.21})$$

Moreover, the solution satisfies

$$\max_{\alpha=2}^3 |m_{\alpha c}(z) - m_{\alpha c}(0)| \leq C_0 |z|. \quad (\text{E.22})$$

The proof is a standard application of the contraction principle. For reader's convenience, we will include its proof in Appendix E.3.4. As a byproduct of the contraction mapping argument there, we also obtain the following stability result that will be useful for our proof of Theorem E.7 below.

903 **Lemma E.6.** *There exist constants $c_0, C_0 > 0$ depending only on τ in (E.3), (E.5), (E.7) and (E.20)*
 904 *such that the self-consistent equations in (E.16) are stable in the following sense. Suppose $|z| \leq c_0$*
 905 *and $m_\alpha : \mathbb{C}_+ \mapsto \mathbb{C}_+$, $\alpha = 2, 3$, are analytic functions of z such that*

$$|m_2(z) - m_{2c}(0)| + |m_3(z) - m_{3c}(0)| \leq c_0.$$

906 *Suppose they satisfy the system of equations*

$$\frac{1}{m_2} + 1 - \frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^2}{z + \sigma_i^2 r_1 m_2 + r_2 m_3} = \mathcal{E}_2, \quad \frac{1}{m_3} + 1 - \frac{1}{n} \sum_{i=1}^p \frac{1}{z + \sigma_i^2 r_1 m_2 + r_2 m_3} = \mathcal{E}_3, \quad (\text{E.23})$$

907 *for some (random) errors satisfying $\max_{\alpha=2}^3 |\mathcal{E}_\alpha| \leq \delta(z)$, where $\delta(z)$ is a deterministic z -dependent*
 908 *function with $\delta(z) \leq (\log n)^{-1}$. Then we have*

$$\max_{\alpha=2}^3 |m_\alpha(z) - m_{\alpha c}(z)| \leq C_0 \delta(z). \quad (\text{E.24})$$

909 In the following proof, we choose a sufficiently small constants $c_0 > 0$ such that Lemma E.5 and
 910 Lemma E.6 hold. Then we define a domain of the spectral parameter z as

$$\mathbf{D} := \{z = E + i\eta \in \mathbb{C}_+ : |z| \leq (\log n)^{-1}\}. \quad (\text{E.25})$$

911 The following theorem gives almost optimal estimates on the resolvent G , which are conventionally
 912 called local laws.

913 **Theorem E.7.** *Suppose Assumption E.1 holds, and Z_1, Z_2 satisfy the bounded support condition*
 914 *(E.8) for a deterministic parameter $q \equiv q(n)$ satisfying $n^{-1/2} \leq q \leq n^{-\phi}$ for some (small) constant*
 915 *$\phi > 0$. Then there exists a sufficiently small constant $c_0 > 0$ such that the following **anisotropic***
 916 ***local law** holds uniformly for all $z \in \mathbf{D}$. For any deterministic unit vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^{\mathcal{I}}$, we have*

$$|\mathbf{u}^\top (G(z) - \Pi(z)) \mathbf{v}| \prec q. \quad (\text{E.26})$$

917 The proof of this theorem will be given in Section E.3. We now use it to complete the proof of Lemma
 918 B.2 and Lemma B.3.

Proof of Lemma B.2. In the setting of Lemma B.2, we write

$$\mathcal{R} := (X_1^\top X_1 + X_2^\top X_2)^{-1} = n^{-1} \left(\Sigma_1^{1/2} Z_1^\top Z_1 \Sigma_1^{1/2} + \Sigma_2^{1/2} Z_2^\top Z_2 \Sigma_2^{1/2} \right)^{-1},$$

where the extra n^{-1} is due to the choice of the scaling—in the setting of Lemma B.2 the variances of the entries of $Z_{1,2}$ are equal to 1, while here they are taken to be n^{-1} . Then as in (E.9), we can write

$$\mathcal{R} = n^{-1} \Sigma_2^{-1/2} V \mathcal{G}(0) V^\top \Sigma_2^{-1/2}, \quad \mathcal{G}(0) = (\Lambda U^\top Z_1^\top Z_1 U \Lambda + V^\top Z_2^\top Z_2 V)^{-1}.$$

919 If the entries of $\sqrt{n} Z_1$ and $\sqrt{n} Z_2$ have arbitrarily high moments as in (A.1), then Z_1 and Z_2 have
 920 bounded support $q = n^{-1/2}$. Using Theorem E.7, we obtain that for any small constant $\varepsilon > 0$,

$$\max_{1 \leq i \leq p} |(A\mathcal{R} - n^{-1} A \Sigma_2^{-1/2} V \Pi(0) V^\top \Sigma_2^{-1/2})_{ii}| \prec n^{-3/2} \|A\|, \quad (\text{E.27})$$

where by (E.18), we have

$$\Pi(0) = -(r_1 m_{2c}(0) \Lambda^2 + r_2 m_{3c}(0))^{-1} = (r_1 b_2 V^\top M^\top M V + r_2 b_3)^{-1},$$

with (b_2, b_3) satisfying (E.19). Thus from (E.27) we get that

$$\text{Tr}(A\mathcal{R}) = n^{-1} \text{Tr}(r_1 b_2 M^\top M + r_2 b_3)^{-1} + O_\prec(n^{-1/2} \|A\|).$$

921 This concludes (B.2) if we rename $r_1 b_2 \rightarrow a_1$ and $r_2 b_3 \rightarrow a_2$.

922 Note that if we set $n_1 = 0$ and $n_2 = n$, then $a_1 = 0$ and $a_2 = (n_2 - p)/n_2$ is the solution to (B.3).
 923 This gives (B.1) using (B.2). \square

924 *Proof of Lemma B.3.* In the setting of Lemma B.3, we can write

$$\Delta := n^2 \left\| \Sigma_2^{-1/2} (X_1^\top X_1 + X_2^\top X_2)^{-1} \beta \right\|^2 = \beta^\top \Sigma_2^{-1/2} (M^\top Z_1^\top Z_1 M + Z_2^\top Z_2)^{-2} \Sigma_2^{-1/2} \beta.$$

925 Here again the n^2 factor disappears due to the choice of scaling. With (E.6), we can write the above
 926 expression as $\Delta := \mathbf{v}^\top (\mathcal{G}^2)(0) \mathbf{v}$ where $\mathbf{v} := V^\top \Sigma_2^{-1/2} \beta$. Note that $\mathcal{G}^2(0) = \partial_z \mathcal{G}|_{z=0}$. Now using
 927 Cauchy's integral formula and Theorem E.7, we get that

$$\begin{aligned} \mathbf{v}^\top \mathcal{G}^2(0) \mathbf{v} &= \frac{1}{2\pi i} \oint_{\mathcal{C}} \frac{\mathbf{v}^\top \mathcal{G}(z) \mathbf{v}}{z^2} dz = \frac{1}{2\pi i} \oint_{\mathcal{C}} \frac{\mathbf{v}^\top \Pi(z) \mathbf{v}}{z^2} dz + O_{\prec}(n^{-1/2} \|\beta\|^2) \\ &= \mathbf{v}^\top \Pi'(0) \mathbf{v} + O_{\prec}(n^{-1/2} \|\beta\|^2), \end{aligned} \quad (\text{E.28})$$

928 where \mathcal{C} is the contour $\{z \in \mathbb{C} : |z| \leq (\log n)^{-1}\}$. Hence it remains to study the derivatives

$$\mathbf{v}^\top \Pi'(0) \mathbf{v} = \mathbf{v}^\top \frac{1 + r_1 m'_{2c}(0) \Lambda^2 + r_2 m'_{3c}(0)}{(r_1 m_{2c}(0) \Lambda^2 + r_2 m_{3c}(0))^2} \mathbf{v}, \quad (\text{E.29})$$

929 where we need to calculate the derivatives $m'_{2c}(0)$ and $m'_{3c}(0)$.

930 By the implicit differentiation of (E.16), we obtain that

$$\begin{aligned} \frac{1}{m_{2c}^2(0)} m'_{2c}(0) &= \frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^2 (1 + \sigma_i^2 r_1 m'_{2c}(0) + r_2 m'_{3c}(0))}{(\sigma_i^2 r_1 m_{2c}(0) + r_2 m_{3c}(0))^2}, \\ \frac{1}{m_{3c}^2(0)} m'_{3c}(0) &= \frac{1}{n} \sum_{i=1}^p \frac{1 + \sigma_i^2 r_1 m'_{2c}(0) + r_2 m'_{3c}(0)}{(\sigma_i^2 r_1 m_{2c}(0) + r_2 m_{3c}(0))^2}. \end{aligned}$$

931 If we rename $-r_1 m_{2c}(0) \rightarrow a_1$, $-r_2 m_{3c}(0) \rightarrow a_2$, $r_2 m'_{3c}(0) \rightarrow a_3$ and $r_1 m'_{2c}(0) \rightarrow a_4$, then this
 932 equation becomes

$$\begin{aligned} \left(\frac{r_2}{a_2^2} - \frac{1}{n} \sum_{i=1}^p \frac{1}{(\sigma_i^2 a_1 + a_2)^2} \right) a_3 - \left(\frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^2}{(\sigma_i^2 a_1 + a_2)^2} \right) a_4 &= \frac{1}{n} \sum_{i=1}^p \frac{1}{(\sigma_i^2 a_1 + a_2)^2}, \\ \left(\frac{r_1}{a_1^2} - \frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^4}{(\sigma_i^2 a_1 + a_2)^2} \right) a_4 - \left(\frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^2}{(\sigma_i^2 a_1 + a_2)^2} \right) a_3 &= \frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^2}{(\sigma_i^2 a_1 + a_2)^2}. \end{aligned} \quad (\text{E.30})$$

933 Then by (E.28) and (E.29), we get

$$\Delta = \beta^\top \Sigma_2^{-1/2} V \frac{1 + a_3 + a_4 \Lambda^2}{(a_1 \Lambda^2 + a_2)} V^\top \Sigma_2^{-1/2} \beta = \beta^\top \Sigma_2^{-1/2} \frac{1 + a_3 + a_4 M^\top M}{(a_1 M^\top M + a_2)} \Sigma_2^{-1/2} \beta,$$

934 where we used $M^\top M = V \Lambda^2 V^\top$ in the second step. This concludes Lemma B.3. \square

935 Using a simple cutoff argument, it is easy to obtain from Theorem E.7 the following corollary under
 936 weaker moment assumptions.

937 **Corollary E.8.** *Suppose Assumption E.1 holds. Moreover, assume that the entries of Z_1 and Z_2 are*
 938 *i.i.d. random variables satisfying (E.1) and*

$$\max_{i,j} \mathbb{E} |\sqrt{n} z_{ij}^{(\alpha)}|^a = O(1), \quad \alpha = 1, 2, \quad (\text{E.31})$$

939 for some fixed $a > 4$. Then (E.26) holds for $q = n^{2/a-1/2}$ on an event with probability $1 - o(1)$.

Proof of Corollary E.8. Fix any sufficiently small constant $\varepsilon > 0$. We choose $q = n^{-c_a + \varepsilon}$ with
 $c_a = 1/2 - 2/a$. Then we introduce the truncated matrices \tilde{Z}_1 and \tilde{Z}_2 , with entries

$$\tilde{z}_{ij}^{(\alpha)} := \mathbf{1} \left\{ |\tilde{z}_{ij}^{(\alpha)}| \leq q \right\} \cdot z_{ij}^{(\alpha)}, \quad \alpha = 1, 2.$$

940 By the moment conditions (E.31) and a simple union bound, we have

$$\mathbb{P}(\tilde{Z}_1 = Z_1, \tilde{Z}_2 = Z_2) = 1 - O(n^{-a\varepsilon}). \quad (\text{E.32})$$

Using (E.31) and integration by parts, it is easy to verify that

$$|\mathbb{E}\tilde{z}_{ij}^{(\alpha)}| = O(n^{-2-\varepsilon}), \quad \mathbb{E}|\tilde{z}_{ij}^{(\alpha)}|^2 = n^{-1} + O(n^{-2-\varepsilon}), \quad \alpha = 1, 2, \quad (\text{E.33})$$

Then we can centralize and rescale \tilde{Z}_1 and \tilde{Z}_2 as $\hat{Z}_\alpha := (\tilde{Z}_\alpha - \mathbb{E}\tilde{Z}_\alpha)/(\mathbb{E}|\tilde{z}_{11}^{(\alpha)}|^2)^{1/2}$, $\alpha = 1, 2$. Now \hat{Z}_1 and \hat{Z}_2 satisfy the assumptions in Theorem E.7 with $q = n^{-c_a+\varepsilon}$, and (E.26) gives that

$$\left| \mathbf{u}^\top (G(\hat{Z}_1, \hat{Z}_2, z) - \Pi(z)) \mathbf{v} \right| \prec q.$$

Then using (E.33) and (E.37) below, we obtain that

$$\left| \mathbf{u}^\top (G(\hat{Z}_1, \hat{Z}_2, z) - G(\tilde{Z}_1, \tilde{Z}_2, z)) \mathbf{v} \right| \prec n^{-1-\varepsilon},$$

where we also used the bound $\|\mathbb{E}\tilde{Z}_\alpha\| = O(n^{-1-\varepsilon})$ by (E.33). This shows that (E.26) also holds for $G(\tilde{Z}_1, \tilde{Z}_2, z)$ with $q = n^{-c_a+\varepsilon}$, and hence concludes the proof by (E.32). \square

With this corollary, we can easily extend Lemma B.2 and Lemma B.3 to the case with weaker moment assumptions. Due to length constraints, we will not go into further details here.

E.3 Proof of Theorem E.7

The main difficulty for the proof of Theorem E.7 is due to the fact that the entries of $Y_1 = Z_1 U \Lambda$ and $Y_2 = Z_2 V$ are not independent. However, notice that if the entries of $Z_1 \equiv Z_1^{Gauss}$ and $Z_2 \equiv Z_2^{Gauss}$ are i.i.d. Gaussian, then by the rotational invariance of the multivariate Gaussian distribution, we have

$$Z_1^{Gauss} U \Lambda \stackrel{d}{=} Z_1^{Gauss} \Lambda, \quad Z_2^{Gauss} V \stackrel{d}{=} Z_2^{Gauss}.$$

In this case, the problem is reduced to proving the anisotropic local law for G with $U = \text{Id}$ and $V = \text{Id}$, such that the entries of Y_1 and Y_2 are independent. This can be handled using the standard resolvent methods as in e.g. [20, 37, 50]. To go from the Gaussian case to the general X case, we will adopt a continuous self-consistent comparison argument developed in [21].

For the case $U = \text{Id}$ and $V = \text{Id}$, we need to deal with the following resolvent:

$$G_0(z) := \begin{pmatrix} -zI_{p \times p} & \Lambda Z_1^\top & Z_2^\top \\ Z_1 \Lambda & -I_{n_1 \times n_1} & 0 \\ Z_2 & 0 & -I_{n_2 \times n_2} \end{pmatrix}^{-1}, \quad z \in \mathbb{C}_+, \quad (\text{E.34})$$

and prove the following result.

Proposition E.9. *Suppose Assumption E.1 holds, and Z_1, Z_2 satisfy the bounded support condition (E.8) with $q = n^{-1/2}$. Suppose U and V are identity. Then the estimate (E.26) holds for $G_0(z)$.*

In Section E.3.1, we collect some a priori estimates and resolvent identities that will be used in the proof of Theorem E.7 and Proposition E.9. Then in Section E.3.2 we give the proof of Proposition E.9, which concludes Theorem E.7 for i.i.d. Gaussian Z_1 and Z_2 . Finally, in Section E.3.3, we describe how to extend the result in Theorem E.7 from the Gaussian case to the case with generally distributed entries of Z_1 and Z_2 . In the proof, we always denote the spectral parameter by $z = E + i\eta$.

E.3.1 Basic estimates

The estimates in this section work for general G , that is, we do not require U and V to be identity.

First, note that $Z_1^\top Z_1$ (resp. $Z_2^\top Z_2$) is a standard sample covariance matrix, and it is well-known that its nonzero eigenvalues are all inside the support of the Marchenko-Pastur law $[(1 - \sqrt{d_1})^2, (1 + \sqrt{d_1})^2]$ (resp. $[(1 - \sqrt{d_2})^2, (1 + \sqrt{d_2})^2]$) with probability $1 - o(1)$ [45]. In our proof, we shall need a slightly stronger probability bound, which is given by the following lemma. Denote the eigenvalues of $Z_1^\top Z_1$ and $Z_2^\top Z_2$ by $\lambda_1(Z_1^\top Z_1) \geq \dots \geq \lambda_p(Z_1^\top Z_1)$ and $\lambda_1(Z_2^\top Z_2) \geq \dots \geq \lambda_p(Z_2^\top Z_2)$.

Lemma E.10. *Suppose Assumption E.1 holds, and Z_1, Z_2 satisfy the bounded support condition (E.8) for some deterministic parameter $q \equiv q(n)$ satisfying $n^{-1/2} \leq q \leq n^{-\phi}$ for some (small) constant $\phi > 0$. Then for any constant $\varepsilon > 0$, we have with high probability,*

$$\lambda_1(Z_1^\top Z_1) \leq (1 + \sqrt{d_1})^2 + \varepsilon, \quad (\text{E.35})$$

and

$$(1 - \sqrt{d_2})^2 - \varepsilon \leq \lambda_p(Z_2^\top Z_2) \leq \lambda_1(Z_2^\top Z_2) \leq (1 + \sqrt{d_2})^2 + \varepsilon. \quad (\text{E.36})$$

971 *Proof.* This lemma essentially follows from [20, Theorem 2.10], although the authors considered
 972 the case with $q \prec n^{-1/2}$ only. The results for larger q follows from [51, Lemma 3.12], but only
 973 the bounds for the largest eigenvalues are given there in order to avoid the issue with the smallest
 974 eigenvalue when d_2 is close to 1. However, under the assumption (E.3), the lower bound for the
 975 smallest eigenvalue follows from the same arguments as in [51]. Hence we omit the details. \square

976 With this lemma, we can obtain the following a priori estimate on the resolvent $G(z)$ for $z \in \mathbf{D}$.

977 **Lemma E.11.** *Suppose the assumptions of Lemma E.10 holds. Then there exists a constant $C > 0$
 978 such that the following estimates hold uniformly in $z, z' \in \mathbf{D}$ with high probability:*

$$\|G(z)\| \leq C, \quad (\text{E.37})$$

979 and for any deterministic unit vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^{\mathcal{I}}$,

$$|\mathbf{u}^\top [G(z) - G(z')] \mathbf{v}| \leq C|z - z'|. \quad (\text{E.38})$$

980 *Proof.* As in (E.15), we let $\{\lambda_k\}_{1 \leq k \leq p}$ be the eigenvalues of WW^\top . By Lemma E.10 and the
 981 assumption (E.3), we obtain that $\lambda_p \geq \lambda_p(Z_2^\top Z_2) \gtrsim 1$, which further implies the estimate that
 982 $\inf_{z \in \mathbf{D}} \min_{1 \leq k \leq p} |\lambda_k - z| \gtrsim 1$. Together with (E.15), we obtain (E.37) and (E.38). \square

983 Now we introduce the concept of minors, which are defined by removing certain rows and columns
 984 of the matrix H .

985 **Definition E.12 (Minors).** *For any $(p+n) \times (p+n)$ matrix \mathcal{A} and $\mathbb{T} \subseteq \mathcal{I}$, we define the minor
 986 $\mathcal{A}^{(\mathbb{T})} := (\mathcal{A}_{\mathbf{ab}} : \mathbf{a}, \mathbf{b} \in \mathcal{I} \setminus \mathbb{T})$ as the $(p+n-|\mathbb{T}|) \times (p+n-|\mathbb{T}|)$ matrix obtained by removing
 987 all rows and columns indexed by \mathbb{T} . Note that we keep the names of indices when defining $\mathcal{A}^{(\mathbb{T})}$, i.e.
 988 $(\mathcal{A}^{(\mathbb{T})})_{ab} = \mathcal{A}_{ab}$ for $a, b \notin \mathbb{T}$. Correspondingly, we define the resolvent minor as (recall (E.13))*

$$G^{(\mathbb{T})} := \left[\left(H - \begin{pmatrix} zI_p & 0 \\ 0 & I_n \end{pmatrix} \right)^{(\mathbb{T})} \right]^{-1} = \begin{pmatrix} \mathcal{G}^{(\mathbb{T})} & \mathcal{G}^{(\mathbb{T})} W^{(\mathbb{T})} \\ (W^{(\mathbb{T})})^\top \mathcal{G}^{(\mathbb{T})} & \mathcal{G}_R^{(\mathbb{T})} \end{pmatrix},$$

989 and the partial traces $m^{(\mathbb{T})}$, $m_1^{(\mathbb{T})}$, $m_2^{(\mathbb{T})}$ and $m_3^{(\mathbb{T})}$ by replacing G with $G^{(\mathbb{T})}$ in (E.14). For convenience,
 990 we will adopt the convention that for any minor $\mathcal{A}^{(\mathbb{T})}$ defined as above, $\mathcal{A}_{ab}^{(\mathbb{T})} = 0$ if $a \in \mathbb{T}$ or
 991 $b \in \mathbb{T}$. Moreover, we will abbreviate $(\{a\}) \equiv (a)$ and $(\{a, b\}) \equiv (ab)$.

992 **Lemma E.13.** *We have the following resolvent identities.*

993 (i) For $i \in \mathcal{I}_1$ and $\mu \in \mathcal{I}_2 \cup \mathcal{I}_3$, we have

$$\frac{1}{G_{ii}} = -z - \left(W G^{(i)} W^\top \right)_{ii}, \quad \frac{1}{G_{\mu\mu}} = -1 - \left(W^\top G^{(\mu)} W \right)_{\mu\mu}. \quad (\text{E.39})$$

994 (ii) For $i \in \mathcal{I}_1$, $\mu \in \mathcal{I}_2 \cup \mathcal{I}_3$, $\mathbf{a} \in \mathcal{I} \setminus \{i\}$ and $\mathbf{b} \in \mathcal{I} \setminus \{\mu\}$, we have

$$G_{i\mathbf{a}} = -G_{ii} \left(W G^{(i)} \right)_{i\mathbf{a}}, \quad G_{\mu\mathbf{b}} = -G_{\mu\mu} \left(W^\top G^{(\mu)} \right)_{\mu\mathbf{b}}. \quad (\text{E.40})$$

995 (iii) For $\mathbf{a} \in \mathcal{I}$ and $\mathbf{b}, \mathbf{c} \in \mathcal{I} \setminus \{\mathbf{a}\}$,

$$G_{\mathbf{bc}}^{(\mathbf{a})} = G_{\mathbf{bc}} - \frac{G_{\mathbf{ba}} G_{\mathbf{ac}}}{G_{\mathbf{aa}}}, \quad \frac{1}{G_{\mathbf{bb}}} = \frac{1}{G_{\mathbf{bb}}^{(\mathbf{a})}} - \frac{G_{\mathbf{ba}} G_{\mathbf{ab}}}{G_{\mathbf{bb}}^{(\mathbf{a})} G_{\mathbf{aa}}}. \quad (\text{E.41})$$

996 *Proof.* All these identities can be proved directly using Schur's complement formula. The reader can
 997 also refer to, for example, [21, Lemma 4.4]. \square

998 The following lemma gives large deviation bounds for bounded supported random variables.

999 **Lemma E.14** (Lemma 3.8 of [52]). *Let (x_i) , (y_j) be independent families of centered and inde-
 1000 pendent random variables, and (A_i) , (B_{ij}) be families of deterministic complex numbers. Suppose*

1001 the entries x_i, y_j have variances at most n^{-1} and satisfy the bounded support condition (E.8) with
 1002 $q \leq n^{-\phi}$ for some constant $\phi > 0$. Then we have the following bound:

$$\left| \sum_i A_i x_i \right| \prec q \max_i |A_i| + \frac{1}{\sqrt{n}} \left(\sum_i |A_i|^2 \right)^{1/2}, \quad \left| \sum_{i,j} x_i B_{ij} y_j \right| \prec q^2 B_d + q B_o + \frac{1}{n} \left(\sum_{i \neq j} |B_{ij}|^2 \right)^{1/2},$$

$$\left| \sum_i \bar{x}_i B_{ii} x_i - \sum_i (\mathbb{E}|x_i|^2) B_{ii} \right| \prec q B_d, \quad \left| \sum_{i \neq j} \bar{x}_i B_{ij} x_j \right| \prec q B_o + \frac{1}{n} \left(\sum_{i \neq j} |B_{ij}|^2 \right)^{1/2},$$

1003 where $B_d := \max_i |B_{ii}|$ and $B_o := \max_{i \neq j} |B_{ij}|$. Moreover, if all the moments of $\sqrt{n}x_i$ and $\sqrt{n}y_j$
 1004 exist in the sense of (A.1), then we have stronger bounds

$$\left| \sum_i A_i x_i \right| \prec \frac{1}{\sqrt{n}} \left(\sum_i |A_i|^2 \right)^{1/2}, \quad \left| \sum_{i,j} x_i B_{ij} y_j \right| \prec \frac{1}{n} \left(\sum_{i \neq j} |B_{ij}|^2 \right)^{1/2},$$

$$\left| \sum_i \bar{x}_i B_{ii} x_i - \sum_i (\mathbb{E}|x_i|^2) B_{ii} \right| \prec \frac{1}{n} \left(\sum_i |B_{ii}|^2 \right)^{1/2}, \quad \left| \sum_{i \neq j} \bar{x}_i B_{ij} x_j \right| \prec \frac{1}{n} \left(\sum_{i \neq j} |B_{ij}|^2 \right)^{1/2}.$$

1005 E.3.2 Entrywise local law

1006 The main goal of this subsection is to prove the following entrywise local law. The anisotropic local
 1007 law (E.26) then follows from the entrywise local law combined with a polynomialization method as
 1008 we will explain later. Recall that in the setting of Proposition E.9, we have $q = n^{-1/2}$ and

$$W = (\Lambda Z_1^\top, Z_2^\top). \quad (\text{E.42})$$

1009 **Lemma E.15.** Suppose the assumptions in Proposition E.9 hold. Then the following estimate holds
 1010 uniformly for $z \in \mathbf{D}$:

$$\max_{\mathbf{a}, \mathbf{b} \in \mathcal{I}} |(G_0)_{\mathbf{ab}}(z) - \Pi_{\mathbf{ab}}(z)| \prec n^{-1/2}. \quad (\text{E.43})$$

1011 *Proof.* The proof of Lemma E.15 is divided into three steps. For simplicity, we will still denote
 1012 $G \equiv G_0$ in the following proof, while keeping in mind that W takes the form in (E.42).

1013 **Step 1: Large deviations estimates.** In this step, we prove some (almost) optimal large deviation
 1014 estimates on the off-diagonal entries of G , and on the following Z variables. In analogy to [52,
 1015 Section 3] and [21, Section 5], we introduce the Z variables

$$Z_{\mathbf{a}}^{(\mathbb{T})} := (1 - \mathbb{E}_{\mathbf{a}})(G_{\mathbf{aa}}^{(\mathbb{T})})^{-1}, \quad \mathbf{a} \notin \mathbb{T},$$

1016 where $\mathbb{E}_{\mathbf{a}}[\cdot] := \mathbb{E}[\cdot \mid H^{(\mathbf{a})}]$, i.e. it is the partial expectation over the randomness of the \mathbf{a} -th row and
 1017 column of H . Using (E.39), we get that for $i \in \mathcal{I}_1$, $\mu \in \mathcal{I}_2$ and $\nu \in \mathcal{I}_3$,

$$Z_i = \sigma_i^2 \sum_{\mu, \nu \in \mathcal{I}_2} G_{\mu\nu}^{(i)} \left(\frac{1}{n} \delta_{\mu\nu} - z_{\mu i} z_{\nu i} \right) + \sum_{\mu, \nu \in \mathcal{I}_3} G_{\mu\nu}^{(i)} \left(\frac{1}{n} \delta_{\mu\nu} - z_{\mu i} z_{\nu i} \right), \quad (\text{E.44})$$

$$Z_\mu = \sum_{i, j \in \mathcal{I}_1} \sigma_i \sigma_j G_{ij}^{(\mu)} \left(\frac{1}{n} \delta_{ij} - z_{\mu i} z_{\mu j} \right), \quad Z_\nu = \sum_{i, j \in \mathcal{I}_1} G_{ij}^{(\nu)} \left(\frac{1}{n} \delta_{ij} - z_{\nu i} z_{\nu j} \right). \quad (\text{E.45})$$

1018 For simplicity, we introduce the random error $\Lambda_o := \max_{\mathbf{a} \neq \mathbf{b}} |G_{\mathbf{aa}}^{-1} G_{\mathbf{ab}}|$. The following lemma
 1019 gives the desired large deviations estimates on Λ_o and the Z variables.

1020 **Lemma E.16.** Suppose the assumptions in Proposition E.9 hold. Then the following estimates hold
 1021 uniformly for all $z \in \mathbf{D}$:

$$\Lambda_o + \max_{\mathbf{a} \in \mathcal{I}} |Z_{\mathbf{a}}| \prec n^{-1/2}. \quad (\text{E.46})$$

1022 *Proof.* Note that for any $\mathbf{a} \in \mathcal{I}$, $H^{(\mathbf{a})}$ and $G^{(\mathbf{a})}$ also satisfies the assumptions for Lemma E.11. Hence
 1023 (E.37) and (E.38) also hold for $G^{(\mathbf{a})}$. Now applying Lemma E.14 to (E.44) and (E.45), and using the
 1024 a priori bound (E.37), we get that for any $i \in \mathcal{I}_1$,

$$|Z_i| \lesssim \sum_{\alpha=2}^3 \left| \sum_{\mu, \nu \in \mathcal{I}_\alpha} G_{\mu\nu}^{(i)} \left(\frac{1}{n} \delta_{\mu\nu} - z_{\mu i} z_{\nu i} \right) \right| \prec n^{-1/2} + \frac{1}{n} \left(\sum_{\mu, \nu \in \mathcal{I}_2 \cup \mathcal{I}_3} |G_{\mu\nu}^{(i)}|^2 \right)^{1/2} \prec n^{-1/2},$$

1025 where in the last step we used (E.37) to get that for any μ ,

$$\sum_{\nu \in \mathcal{I}_2 \cup \mathcal{I}_3} \left| G_{\mu\nu}^{(i)} \right|^2 \leq \sum_{\mathbf{a} \in \mathcal{I}} \left| G_{\mu\mathbf{a}}^{(i)} \right|^2 = \left[G^{(i)} (G^{(i)})^* \right]_{\mu\mu} = O(1). \quad (\text{E.47})$$

1026 Similarly, applying Lemma E.14 to Z_μ and Z_ν in (E.45) and using (E.37), we obtain the same bound.
1027 we have

$$G_{i\mathbf{a}} = -G_{ii} \left(W G^{(i)} \right)_{i\mathbf{a}}, \quad G_{\mu\mathbf{b}} = -G_{\mu\mu} \left(W^\top G^{(\mu)} \right)_{\mu\mathbf{b}}. \quad (\text{E.48})$$

1028 Then we prove the off-diagonal estimate on Λ_o . For $i \in \mathcal{I}_1$ and $\mathbf{a} \in \mathcal{I} \setminus \{i\}$, using (E.40), Lemma
1029 E.14 and (E.37), we obtain that

$$\left| G_{ii}^{-1} G_{i\mathbf{a}} \right| \prec n^{-1/2} + \frac{1}{\sqrt{n}} \left(\sum_{\mu \in \mathcal{I}_2 \cup \mathcal{I}_3} \left| G_{\mu\mathbf{a}}^{(i)} \right|^2 \right)^{1/2} \prec n^{-1/2}.$$

1030 We have a similar estimate for $\left| G_{\mu\mu}^{-1} G_{\mu\mathbf{b}} \right|$ with $\mu \in \mathcal{I}_2 \cup \mathcal{I}_3$ and $\mathbf{b} \in \mathcal{I} \setminus \{\mu\}$. Thus we obtain that
1031 $\Lambda_o \prec n^{-1/2}$, which concludes (E.46). \square

1032 Note that combining (E.37) and (E.46), we immediately conclude (E.43) for $\mathbf{a} \neq \mathbf{b}$.

1033 **Step 2: Self-consistent equations.** This is the key step of the proof for Proposition E.15, which
1034 derives approximate self-consistent equations satisfied by $m_2(z)$ and $m_3(z)$. More precisely, we
1035 will show that $(m_2(z), m_3(z))$ satisfies (E.23) for some small error $|\mathcal{E}_{2,3}| \prec n^{-1/2}$. Then in Step 3
1036 we will apply Lemma E.6 to show that $(m_2(z), m_3(z))$ is close to $(m_{2c}(z), m_{3c}(z))$.

1037 We define the following z -dependent event

$$\Xi(z) := \left\{ |m_2(z) - m_{2c}(z)| + |m_3(z) - m_{3c}(z)| \leq (\log n)^{-1/2} \right\}. \quad (\text{E.49})$$

1038 Note that by (E.22), we have $|m_{2c} + b_2| \lesssim (\log n)^{-1}$ and $|m_{3c} + b_3| \lesssim (\log n)^{-1}$. Together with
1039 (E.16), (E.20) and (E.7), we obtain the following basic estimates

$$|m_{2c}| \sim |m_{3c}| \sim 1, \quad |z + \sigma_i^2 r_1 m_{2c} + r_2 m_{3c}| \sim 1, \quad |1 + \gamma_n m_c| \sim |1 + \gamma_n m_{1c}| \sim 1, \quad (\text{E.50})$$

uniformly in $z \in \mathbf{D}$, where we abbreviated

$$m_c(z) := -\frac{1}{p} \sum_{i \in \mathcal{I}_1} \frac{1}{z + \sigma_i^2 r_1 m_{2c} + r_2 m_{3c}}, \quad m_{1c}(z) := -\frac{1}{p} \sum_{i \in \mathcal{I}_1} \frac{\sigma_i^2}{z + \sigma_i^2 r_1 m_{2c} + r_2 m_{3c}}.$$

1040 Plugging (E.50) into (E.18), we get

$$|\Pi_{\mathbf{a}\mathbf{a}}(z)| \sim 1 \quad \text{uniformly in } z \in \mathbf{D}, \mathbf{a} \in \mathcal{I}. \quad (\text{E.51})$$

1041 Then we claim the following result.

1042 **Lemma E.17.** *Suppose the assumptions in Proposition E.9 hold. Then the following estimates hold*
1043 *uniformly in $z \in \mathbf{D}$:*

$$\begin{aligned} \mathbf{1}(\Xi) \left| \frac{1}{m_2} + 1 - \frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^2}{z + \sigma_i^2 r_1 m_2 + r_2 m_3} \right| &\prec n^{-1/2}, \\ \mathbf{1}(\Xi) \left| \frac{1}{m_3} + 1 - \frac{1}{n} \sum_{i=1}^p \frac{1}{z + \sigma_i^2 r_1 m_2 + r_2 m_3} \right| &\prec n^{-1/2}. \end{aligned} \quad (\text{E.52})$$

1044 *Proof.* By (E.39), (E.44) and (E.45), we obtain that for $i \in \mathcal{I}_1$, $\mu \in \mathcal{I}_2$ and $\nu \in \mathcal{I}_3$,

$$\frac{1}{G_{ii}} = -z - \frac{\sigma_i^2}{n} \sum_{\mu \in \mathcal{I}_2} G_{\mu\mu}^{(i)} - \frac{1}{n} \sum_{\mu \in \mathcal{I}_3} G_{\mu\mu}^{(i)} + Z_i = -z - \sigma_i^2 r_1 m_2 - r_2 m_3 + \varepsilon_i, \quad (\text{E.53})$$

$$\frac{1}{G_{\mu\mu}} = -1 - \frac{1}{n} \sum_{i \in \mathcal{I}_1} \sigma_i^2 G_{ii}^{(\mu)} + Z_\mu = -1 - \gamma_n m_1 + \varepsilon_\mu, \quad (\text{E.54})$$

$$\frac{1}{G_{\nu\nu}} = -1 - \frac{1}{n} \sum_{i \in \mathcal{I}_1} G_{ii}^{(\nu)} + Z_\nu = -1 - \gamma_n m + \varepsilon_\nu, \quad (\text{E.55})$$

where we recall Definition E.12, and

$$\varepsilon_i := Z_i + \sigma_i r_1 \left(m_2 - m_2^{(i)} \right) + r_2 \left(m_3 - m_3^{(i)} \right), \quad \varepsilon_\mu := \begin{cases} Z_\mu + \gamma_n (m_1 - m_1^{(\mu)}), & \text{if } \mu \in \mathcal{I}_2 \\ Z_\mu + \gamma_n (m - m^{(\mu)}), & \text{if } \mu \in \mathcal{I}_3 \end{cases}.$$

1045 By (E.41) we can bound that

$$|m_2 - m_2^{(i)}| \leq \frac{1}{n_1} \sum_{\mu \in \mathcal{I}_2} \left| \frac{G_{\mu i} G_{i \mu}}{G_{ii}} \right| \prec n^{-1},$$

1046 where we used (E.46) in the second step. Similarly, we can get that

$$|m - m^{(\mu)}| + |m_1 - m_1^{(\mu)}| + |m_2 - m_2^{(i)}| + |m_3 - m_3^{(i)}| \prec n^{-1} \quad (\text{E.56})$$

1047 for any $i \in \mathcal{I}_1$ and $\mu \in \mathcal{I}_2 \cup \mathcal{I}_3$. Together with (E.46), we obtain that for all i and μ ,

$$|\varepsilon_i| + |\varepsilon_\mu| \prec n^{-1/2}. \quad (\text{E.57})$$

1048 With (E.50) and the definition of Ξ , we get that $\mathbf{1}(\Xi) |z + \sigma_i^2 r_1 m_2 + r_2 m_3| \sim 1$. Hence using (E.53),
1049 (E.57) and (E.46), we obtain that

$$\mathbf{1}(\Xi) G_{ii} = \mathbf{1}(\Xi) \left[-\frac{1}{z + \sigma_i^2 r_1 m_2 + r_2 m_3} + O_{\prec} \left(n^{-1/2} \right) \right]. \quad (\text{E.58})$$

1050 Plugging it into the definitions of m and m_1 in (E.14), we get

$$\mathbf{1}(\Xi) m = \mathbf{1}(\Xi) \left[-\frac{1}{p} \sum_{i \in \mathcal{I}_1} \frac{1}{z + \sigma_i^2 r_1 m_2 + r_2 m_3} + O_{\prec} \left(n^{-1/2} \right) \right], \quad (\text{E.59})$$

$$\mathbf{1}(\Xi) m_1 = \mathbf{1}(\Xi) \left[-\frac{1}{p} \sum_{i \in \mathcal{I}_1} \frac{\sigma_i^2}{z + \sigma_i^2 r_1 m_2 + r_2 m_3} + O_{\prec} \left(n^{-1/2} \right) \right]. \quad (\text{E.60})$$

1051 As a byproduct, we obtain from these two estimates that

$$\mathbf{1}(\Xi) (|m - m_c| + |m_1 - m_{1c}|) \lesssim (\log n)^{-1/2}, \quad \text{with high probability.} \quad (\text{E.61})$$

1052 Together with (E.50), we get that

$$|1 + \gamma_n m_1| \sim 1, \quad |1 + \gamma_n m| \sim 1, \quad \text{with high probability on } \Xi. \quad (\text{E.62})$$

1053 Now using (E.54), (E.55), (E.57), (E.46) and (E.62), we obtain that for $\mu \in \mathcal{I}_2$ and $\nu \in \mathcal{I}_3$,

$$\mathbf{1}(\Xi) \left(G_{\mu\mu} + \frac{1}{1 + \gamma_n m_1} \right) = O_{\prec} (n^{-1/2}), \quad \mathbf{1}(\Xi) \left(G_{\nu\nu} + \frac{1}{1 + \gamma_n m} \right) = O_{\prec} (n^{-1/2}). \quad (\text{E.63})$$

1054 Taking average over $\mu \in \mathcal{I}_2$ and $\nu \in \mathcal{I}_3$, we get that with high probability,

$$\mathbf{1}(\Xi) \left(m_2 + \frac{1}{1 + \gamma_n m_1} \right) = O_{\prec} (n^{-1/2}), \quad \mathbf{1}(\Xi) \left(m_3 + \frac{1}{1 + \gamma_n m} \right) = O_{\prec} (n^{-1/2}). \quad (\text{E.64})$$

1055 Finally, plugging (E.59) and (E.60) into (E.64), we conclude (E.52). \square

1056 **Step 3: Ξ holds with high probability.** In this step, we show that the event $\Xi(z)$ in fact holds with
1057 high probability for all $z \in \mathbf{D}$. Once we have proved this fact, then applying Lemma E.6 to (E.52)
1058 immediately shows that $(m_2(z), m_3(z))$ is equal to $(m_{2c}(z), m_{3c}(z))$ up to an error of order $n^{-1/2}$.
1059 We claim that it suffices to show

$$|m_2(0) - m_{2c}(0)| + |m_3(0) - m_{3c}(0)| \prec n^{-1/2}. \quad (\text{E.65})$$

1060 Once we know (E.65), then by (E.22) and (E.38), we get $\max_{\alpha=2}^3 |m_{\alpha c}(z) - m_{\alpha c}(0)| =$
1061 $O((\log n)^{-1})$ and $\max_{\alpha=2}^3 |m_{\alpha}(z) - m_{\alpha}(0)| = O((\log n)^{-1})$ with high probability for all $z \in \mathbf{D}$.
1062 Together with (E.65), we obtain that

$$\sup_{z \in \mathbf{D}} (|m_2(z) - m_{2c}(z)| + |m_3(z) - m_{3c}(z)|) \lesssim (\log n)^{-1} \quad \text{with high probability,} \quad (\text{E.66})$$

1063 and

$$\sup_{z \in \mathbf{D}} (|m_2(z) - m_{2c}(0)| + |m_3(z) - m_{3c}(0)|) \lesssim (\log n)^{-1} \quad \text{with high probability.} \quad (\text{E.67})$$

1064 The condition (E.66) shows that Ξ holds with high probability, and the condition (E.67) verifies the
1065 condition (E.21) of Lemma E.6. Then applying Lemma E.6 to (E.52), we obtain that

$$|m_2(z) - m_{2c}(z)| + |m_3(z) - m_{3c}(z)| \prec n^{-1/2} \quad (\text{E.68})$$

1066 for all $z \in \mathbf{D}$. Plugging (E.68) into (E.53)-(E.55), we get the diagonal estimate

$$\max_{a \in \mathcal{I}} |G_{aa}(z) - \Pi_{aa}(z)| \prec n^{-1/2}. \quad (\text{E.69})$$

1067 Together with the off-diagonal estimate in (E.46), we conclude (E.43). \square

1068 Now we give the proof of (E.65).

Proof of (E.65). By (E.15), we get

$$m(0) = \frac{1}{p} \sum_{i \in \mathcal{I}_1} G_{ii}(0) = \frac{1}{p} \sum_{k=1}^p \frac{|\xi_k(i)|^2}{\lambda_k} \geq \lambda_1^{-1} \gtrsim 1.$$

Similarly, we can also get that $m_1(0)$ is positive and has size $m_1(0) \sim 1$. Hence we have

$$1 + \gamma_n m_1(0) \sim 1, \quad 1 + \gamma_n m_1(0) \sim 1.$$

Together with (E.54), (E.55) and (E.57), we obtain that (E.64) holds at $z = 0$ even without the indicator function $\mathbf{1}(\Xi)$. Furthermore, it gives that

$$|\sigma_i^2 r_1 m_2(0) + r_2 m_3(0)| = \left| \frac{\sigma_i^2 r_1}{1 + \gamma_n m_1(0)} + \frac{r_2}{1 + \gamma_n m(0)} + O_{\prec}(n^{-1/2}) \right| \sim 1$$

1069 with high probability. Then using (E.53) and (E.57), we obtain that (E.59) and (E.60) hold at
1070 $z = 0$ even without the indicator function $\mathbf{1}(\Xi)$. Finally, plugging (E.59) and (E.60) into (E.64), we
1071 conclude (E.52) holds at $z = 0$, that is,

$$\begin{aligned} \left| \frac{1}{m_2(0)} + 1 - \frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^2}{\sigma_i^2 r_1 m_2(0) + r_2 m_3(0)} \right| &\prec n^{-1/2}, \\ \left| \frac{1}{m_3(0)} + 1 - \frac{1}{n} \sum_{i=1}^p \frac{1}{\sigma_i^2 r_2 m_2(0) + r_2 m_3(0)} \right| &\prec n^{-1/2}. \end{aligned} \quad (\text{E.70})$$

Denoting $\omega_2 = -m_{2c}(0)$ and $\omega_3 = -m_{3c}(0)$. By (E.64), we have

$$\omega_2 = \frac{1}{1 + \gamma_n m_1(0)} + O_{\prec}(n^{-1/2}), \quad \omega_3 = \frac{1}{1 + \gamma_n m(0)} + O_{\prec}(n^{-1/2}).$$

1072 Hence there exists a sufficiently small constant $c > 0$ such that

$$c \leq \omega_2 \leq 1, \quad c \leq \omega_3 \leq 1, \quad \text{with high probability.} \quad (\text{E.71})$$

1073 Also one can verify from (E.70) that (ω_2, ω_3) satisfy approximately the same equations as (E.19):

$$r_1 \omega_2 + r_2 \omega_3 = 1 - \gamma_n + O_{\prec}(n^{-1/2}), \quad f(\omega_2) = 1 + O_{\prec}(n^{-1/2}). \quad (\text{E.72})$$

1074 The first equation and (E.71) together implies that $\omega_2 \in [0, r_1^{-1}(1 - \gamma_n)]$ with high probability. Since
1075 f is strictly increasing and has bounded derivatives on $[0, r_1^{-1}(1 - \gamma_n)]$, by basic calculus the second
1076 equation in (E.72) gives that $|\omega_2 - b_2| \prec n^{-1/2}$. Together with the first equation in (E.72), we get
1077 $|\omega_3 - b_3| \prec n^{-1/2}$. This concludes (E.65). \square

1078 With Lemma E.15, we can complete the proof of Proposition E.9.

1079 *Proof of Proposition E.9.* With (E.43), one can use the polynomialization method in [20, Section 5]
1080 to get the anisotropic local law (E.26) for G_0 with $q = n^{-1/2}$. The proof is exactly the same, except
1081 for some minor differences in notations, so we omit the details. \square

1082 E.3.3 Anisotropic local law

1083 In this subsection, we finish the proof of Theorem E.7 for a general X satisfying the bounded support
 1084 condition (E.8) with $q \leq n^{-\phi}$ for some constant $\phi > 0$. Proposition E.9 implies that (E.26) holds for
 1085 Gaussian Z_1^{Gauss} and Z_2^{Gauss} as discussed before. Thus the basic idea is to prove that for Z_1 and Z_2
 1086 satisfying the assumptions in Theorem E.7,

$$\mathbf{u}^\top (G(Z, z) - G(Z^{Gauss}, z)) \mathbf{v} \prec q$$

1087 for any deterministic unit vectors $\mathbf{u}, \mathbf{v} \in \mathbb{C}^{\mathcal{I}}$ and $z \in \mathbf{D}$. Here we abbreviated $Z := \begin{pmatrix} Z_1 \\ Z_2 \end{pmatrix}$ and
 1088 $Z^{Gauss} := \begin{pmatrix} Z_1^{Gauss} \\ Z_2^{Gauss} \end{pmatrix}$. We prove the above statement using a continuous comparison argument
 1089 introduced in [21]. The proof is similar to the ones in Sections 7-8 of [21], so we only give a rough
 1090 description of the basic idea, without writing down all the details.

1091 **Definition E.18** (Interpolation). We denote $Z^0 := Z^{Gauss}$ and $Z^1 := Z$. Let $\rho_{\mu i}^0$ and $\rho_{\mu i}^1$ be the laws
 1092 of $Z_{\mu i}^0$ and $Z_{\mu i}^1$, respectively. For $\theta \in [0, 1]$, we define the interpolated law $\rho_{\mu i}^\theta := (1 - \theta)\rho_{\mu i}^0 + \theta\rho_{\mu i}^1$.
 1093 We shall work on the probability space consisting of triples (Z^0, Z^θ, Z^1) of independent $n \times p$
 1094 random matrices, where the matrix $Z^\theta = (Z_{\mu i}^\theta)$ has law

$$\prod_{i \in \mathcal{I}_1} \prod_{\mu \in \mathcal{I}_2 \cup \mathcal{I}_3} \rho_{\mu i}^\theta(dZ_{\mu i}^\theta). \quad (\text{E.73})$$

1095 For $\lambda \in \mathbb{R}$, $i \in \mathcal{I}_1$ and $\mu \in \mathcal{I}_2 \cup \mathcal{I}_3$, we define the matrix $Z_{(\mu i)}^{\theta, \lambda}$ through

$$\left(Z_{(\mu i)}^{\theta, \lambda} \right)_{\nu j} := \begin{cases} Z_{\mu i}^\theta, & \text{if } (j, \nu) \neq (i, \mu) \\ \lambda, & \text{if } (j, \nu) = (i, \mu) \end{cases}.$$

1096 We also introduce the matrices $G^\theta(z) := G(Z^\theta, z)$, $G_{(\mu i)}^{\theta, \lambda}(z) := G(Z_{(\mu i)}^{\theta, \lambda}, z)$.

1097 We shall prove (E.26) through interpolation matrices Z^θ between Z^0 and Z^1 . We have seen that (E.26)
 1098 holds for Z^0 by Proposition E.9. Using (E.73) and fundamental calculus, we get the following basic
 1099 interpolation formula: for $F : \mathbb{R}^{n \times p} \rightarrow \mathbb{C}$,

$$\frac{d}{d\theta} \mathbb{E} F(Z^\theta) = \sum_{i \in \mathcal{I}_1} \sum_{\mu \in \mathcal{I}_2 \cup \mathcal{I}_3} \left[\mathbb{E} F \left(Z_{(\mu i)}^{\theta, Z_{\mu i}^1} \right) - \mathbb{E} F \left(Z_{(\mu i)}^{\theta, Z_{\mu i}^0} \right) \right] \quad (\text{E.74})$$

1100 provided all the expectations exist. We shall apply (E.74) to $F(Z) := F_{\mathbf{u}\mathbf{v}}^s(Z, z)$ for (large) $s \in 2\mathbb{N}$
 1101 and $F_{\mathbf{u}\mathbf{v}}(Z, z)$ defined as

$$F_{\mathbf{u}\mathbf{v}}(Z, z) := |\mathbf{u}^\top (G(Z, z) - \Pi(z)) \mathbf{v}|.$$

1102 The main part of the proof is to show the following self-consistent estimate for the right-hand side of
 1103 (E.74) for any fixed $s \in 2\mathbb{N}$ and constant $\varepsilon > 0$:

$$\sum_{i \in \mathcal{I}_1} \sum_{\mu \in \mathcal{I}_2 \cup \mathcal{I}_3} \left[\mathbb{E} F_{\mathbf{u}\mathbf{v}}^s \left(Z_{(\mu i)}^{\theta, Z_{\mu i}^1}, z \right) - \mathbb{E} F_{\mathbf{u}\mathbf{v}}^s \left(Z_{(\mu i)}^{\theta, Z_{\mu i}^0}, z \right) \right] = O \left((n^\varepsilon q)^s + \mathbb{E} F_{\mathbf{u}\mathbf{v}}^s(Z^\theta, z) \right) \quad (\text{E.75})$$

for all $\theta \in [0, 1]$. If (E.75) holds, then combining (E.74) with a Grönwall's argument we obtain that
 for any fixed $s \in 2\mathbb{N}$ and constant $\varepsilon > 0$:

$$\mathbb{E} |G_{\mathbf{u}\mathbf{v}}(Z^1, z) - \Pi_{\mathbf{u}\mathbf{v}}(z)|^p \leq (n^\varepsilon q)^p.$$

1104 Together with Markov's inequality, we conclude (E.26). Underlying the proof of (E.75) is an
 1105 expansion approach, which is very similar to the ones for Lemma 7.10 of [21] and Lemma 6.11 of
 1106 [37]. So we omit the details.

1107 E.3.4 Proofs of Lemma E.5 and Lemma E.6

1108 Finally, we give the proof of Lemma E.5 and Lemma E.6 using the contraction principle.

1109 *Proof of Lemma E.5.* One can check that the equations in (E.16) are equivalent to the following ones:
 1110

$$r_1 m_{2c} = -(1 - \gamma_n) - r_2 m_{3c} - z (m_{3c}^{-1} + 1), \quad g_z(m_{3c}(z)) = 1, \quad (\text{E.76})$$

where

$$g_z(m_{3c}) := -m_{3c} + \frac{1}{n} \sum_{i=1}^p \frac{m_{3c}}{z - \sigma_i^2(1 - \gamma_n) + (1 - \sigma_i^2)r_2 m_{3c} - \sigma_i^2 z (m_{3c}^{-1} + 1)}.$$

1111 We first show that there exists a unique solution $m_{3c}(z)$ to the equation $g_z(m_{3c}(z)) = 1$ under the
 1112 conditions in (E.21), and the solution satisfies (E.22). Now we abbreviate $\varepsilon(z) := m_{3c}(z) - m_{3c}(0)$,
 1113 and from (E.76) we obtain that

$$0 = [g_z(m_{3c}(z)) - g_0(m_{3c}(0)) - g'_z(m_{3c}(0))\varepsilon(z)] + g'_z(m_{3c}(0))\varepsilon(z),$$

1114 which implies

$$\varepsilon(z) = -\frac{g_z(m_{3c}(0)) - g_0(m_{3c}(0))}{g'_z(m_{3c}(0))} - \frac{g_z(m_{3c}(0) + \varepsilon(z)) - g_z(m_{3c}(0)) - g'_z(m_{3c}(0))\varepsilon(z)}{g'_z(m_{3c}(0))}.$$

1115 Inspired by this equation, we define iteratively a sequence $\varepsilon^{(k)} \in \mathbb{C}$ such that $\varepsilon^{(0)} = 0$, and

$$\varepsilon^{(k+1)} = -\frac{g_z(m_{3c}(0)) - g_0(m_{3c}(0))}{g'_z(m_{3c}(0))} - \frac{g_z(m_{3c}(0) + \varepsilon^{(k)}) - g_z(m_{3c}(0)) - g'_z(m_{3c}(0))\varepsilon^{(k)}}{g'_z(m_{3c}(0))}. \quad (\text{E.77})$$

1116 Then (E.77) defines a mapping $h : \mathbb{C} \rightarrow \mathbb{C}$, which maps $\varepsilon^{(k)}$ to $\varepsilon^{(k+1)} = h(\varepsilon^{(k)})$.

With direct calculation, one can get the derivative

$$g'_z(m_{3c}(0)) = -1 - \frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^2(1 - \gamma_n) - z [1 - \sigma_i^2 (2m_{3c}^{-1}(0) + 1)]}{[z - \sigma_i^2(1 - \gamma_n) + (1 - \sigma_i^2)r_2 m_{3c}(0) - \sigma_i^2 z (m_{3c}^{-1}(0) + 1)]^2}.$$

1117 Then it is easy to check that there exist constants $\tilde{c}, \tilde{C} > 0$ depending only on τ in (E.7) and (E.20)
 1118 such that

$$|[g'_z(m_{3c}(0))]^{-1}| \leq \tilde{C}, \quad \left| \frac{g_z(m_{3c}(0)) - g_0(m_{3c}(0))}{g'_z(m_{3c}(0))} \right| \leq \tilde{C}|z|, \quad (\text{E.78})$$

1119 and

$$\left| \frac{g_z(m_{3c}(0) + \varepsilon_1) - g_z(m_{3c}(0) + \varepsilon_2) - g'_z(m_{3c}(0))(\varepsilon_1 - \varepsilon_2)}{g'_z(m_{3c}(0))} \right| \leq \tilde{C}|\varepsilon_1 - \varepsilon_2|^2, \quad (\text{E.79})$$

for all $|z| \leq \tilde{c}$ and $|\varepsilon_1| \leq \tilde{c}, |\varepsilon_2| \leq \tilde{c}$. Then with (E.78) and (E.79), it is easy to see that there exists a
 sufficiently small constant $\delta > 0$ depending only on \tilde{C} , such that h is a self-mapping

$$h : B_r \rightarrow B_r, \quad B_r := \{\varepsilon \in \mathbb{C} : |\varepsilon| \leq r\},$$

1120 as long as $r \leq \delta$ and $|z| \leq c_\delta$ for some constant $c_\delta > 0$ depending only on \tilde{C} and δ . Now it suffices
 1121 to prove that h restricted to B_r is a contraction, which then implies that $\varepsilon := \lim_{k \rightarrow \infty} \varepsilon^{(k)}$ exists and
 1122 $m_{3c}(0) + \varepsilon$ is a unique solution to the second equation of (E.76) subject to the condition $\|\varepsilon\|_\infty \leq r$.
 1123 From the iteration relation (E.77), using (E.78) one can readily check that

$$\varepsilon^{(k+1)} - \varepsilon^{(k)} = h(\varepsilon^{(k)}) - h(\varepsilon^{(k-1)}) \leq \tilde{C}|\varepsilon^{(k)} - \varepsilon^{(k-1)}|^2. \quad (\text{E.80})$$

1124 Hence as long as r is chosen to be sufficiently small such that $2r\tilde{C} \leq 1/2$, then h is indeed
 1125 a contraction mapping on B_r , which proves both the existence and uniqueness of the solution
 1126 $m_{3c}(z) = m_{3c}(0) + \varepsilon$, if we choose c_0 in (E.21) as $c_0 = \min\{c_\delta, r\}$. After obtaining $m_{3c}(z)$, we
 1127 can then find $m_{2c}(z)$ using the first equation in (E.76).

1128 Note that with (E.79) and $\varepsilon^{(0)} = 0$, we get from (E.77) that $|\varepsilon^{(1)}| \leq \tilde{C}|z|$. With the contraction
 1129 mapping, we have the bound

$$|\varepsilon| \leq \sum_{k=0}^{\infty} |\varepsilon^{(k+1)} - \varepsilon^{(k)}| \leq 2\tilde{C}|z|. \quad (\text{E.81})$$

Algorithm 1 An incremental training schedule for efficient multi-task learning with two tasks

Input: Two tasks (X_1, Y_1) and (X_2, Y_2) .

Parameter: A shared module B , output layers W_1, W_2 as in the hard parameter sharing architecture.

Require: # batches S , epochs T , task 2's validation accuracy $\hat{g}(B; W_2)$, a threshold $\tau \in (0, 1)$.

Output: The trained modules B, W_2 optimized for task 2.

```
1: Divide  $(X_1, Y_1)$  randomly into  $S$  batches:  $(x^{(1)}, y^{(1)}), \dots, (x^{(S)}, y^{(S)})$ .
2: for  $i = 1, \dots, S$  do
3:   for  $j = 1, \dots, T$  do
4:     Update  $B, W_1, W_2$  using the training data  $\{x^{(k)}, x^{(k)}\}_{k=1}^i$  and  $(X_2, Y_2)$ .
5:   end for
6:   Let  $a_i = \hat{g}(B; W_2)$  be the validation accuracy.
7:   if  $a_i < a_{i-1}$  or  $a_i > \tau$  then
8:     break
9:   end if
10: end for
```

This gives the bound (E.22) for $m_{3c}(z)$. Using the first equation in (E.76), we immediately obtain the bound $r_1|m_{2c}(z) - m_{2c}(0)| \leq C|z|$. This gives (E.22) for $m_{2c}(z)$ as long as if $r_1 \gtrsim 1$. To deal with the small r_1 case, we go back to the first equation in (E.16) and treat $m_{2c}(z)$ as the solution to the following equation:

$$\tilde{g}_z(m_{2c}(z)) = 1, \quad \tilde{g}_z(x) := -x + \frac{\gamma_n}{p} \sum_{i=1}^p \frac{\sigma_i^2 x}{z + \sigma_i^2 r_1 x + r_2 m_{3c}(z)}.$$

1130 Then with similar arguments as above between (E.76) and (E.81), we can conclude (E.22) for $m_{2c}(z)$.
1131 This concludes the proof of Lemma E.5. \square

1132 *Proof of Lemma E.6.* Under (E.21), we can obtain equation (E.76) approximately up to some small
1133 error

$$r_1 m_{2c} = -(1 - \gamma_n) - r_2 m_{3c} - z(m_{3c}^{-1} + 1) + \mathcal{E}'_2(z), \quad g_z(m_{3c}(z)) = 1 + \mathcal{E}'_3(z), \quad (\text{E.82})$$

1134 with $|\mathcal{E}'_2(z)| + |\mathcal{E}'_3(z)| = O(\delta(z))$. Then we subtract the equations (E.76) from (E.82), and consider
1135 the contraction principle for the functions $\varepsilon(z) := m_3(z) - m_{3c}(z)$. The rest of the proof is exactly
1136 the same as the one for Lemma E.5, so we omit the details. \square

1137 F Missing Details from the Experiments

1138 F.1 Synthetic Settings

1139 In Figure 1 (c), we plot the test error of the target task for $n_2 = 4p$ and n_1 ranging from p to $20p$.

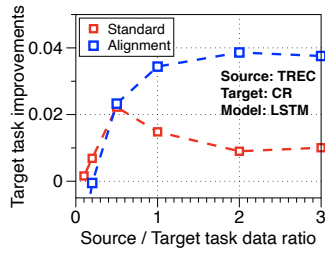
1140 F.2 Image and Text Classification Settings

1141 Note: For text classification tasks, the source task training data size ranges from 500 to 1,500 and
1142 target task training data size is 1000; For ChestX-ray14, the training data size is 10,000.

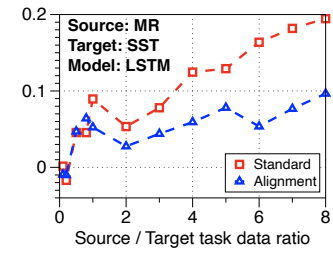
1143 **Task similarity.** We validate that MTL performs better when the source task is more similar to the
1144 target task. We show the result on the sentiment analysis tasks. For a target task, we manually select
1145 a similar task and a dissimilar task based on prior knowledge. Figure 2a confirms the result. Recall
1146 that Section 3.3 shows that increasing the data size of the source task does not always improve the
1147 performance of MTL for the target task. In Figure 2b, we show that for source task MR and target
1148 task SST, there is a transition from positive to negative transfer as we increase the data size of the
1149 source task. When the source task data size is particularly large compared to the target task, we show
1150 that applying the covariance alignment algorithm results in more significant gains. In Figure 2c, we
1151 observe that the benefit from aligning task covariances becomes more significant for LSTM and MLP
1152 as we increase the number of datapoints of the source task.

1153 ¹

¹<http://nlp.stanford.edu/data/wordvecs/glove.6B.zip>



(a) Task pair TREC and CR



(b) Task pair MR and SST

Figure 3: The performance of aligning task covariances depends on data size. As the ratio between source task data size and target task data size increases, the performance improvement from aligning task covariances increases.