
Generalization Effects of Transferring in High-Dimensions using Multi-Task Learning

Anonymous Author(s)

Affiliation

Address

email

1 Introduction

Multi-task learning is an inductive learning mechanism to improve generalization performance using related task data. Many state-of-the-art results in computer vision and natural language processing are obtained using multi-task learning. On the other hand, unless the structures across task data are well-understood, applying multi-task learning on several different datasets often result in suboptimal models (or negative transfer in more technical terms).

[**Todo: Motivate high-dimensional data.**] The technical challenge is to deal with high-dimensional data, in particular when the size of the training set is only a small constant times the feature dimension. Prior theory is unable to explain a phase transition that occurs when comparing multi-task learning to single-task learning. In particular in Figure 1, we observe a shift from positive transfer to negative transfer as a parameter of task relatedness. The theory we develop will provide a precise explanation to such a phenomenon (and more).

To gain insight into the working of multi-task learning, we consider a simplified setting for learning multiple high-dimensional linear regression tasks. A typical process to do multi-task learning involves two steps: (i) Jointly learn a shared representation for all the tasks; (ii) Fine-tune the learnt model on a specific target task. We focus on a hard parameter sharing model proposed in [23] and identify conditions on when multi-task and transfer learning works, and when it doesn't. The high-dimensional linear regression setting where the target task data size is limited captures the intuition that the target task only contains limited labeled data.

Concretely, our input consists of k tasks $(X_1, Y_1), (X_2, Y_2), \dots, (X_k, Y_k)$. We shall assume that each task data follows a linear model, i.e. $y_i = X_i \beta_i + \varepsilon_i$, $1 \leq i \leq k$. Here $\beta_i \in \mathbb{R}^p$ is the model parameter for the i -th task. Each row of $X_i \in \mathbb{R}^{n_i \times p}$ is assumed to be drawn i.i.d. from a fixed distribution with covariance matrix Σ_i . We use a shared body $B \in \mathbb{R}^{p \times r}$ for all tasks and a separate prediction head $\{W_i \in \mathbb{R}^r\}_{i=1}^k$ for each task. This corresponds to minimizing the following optimization objective.

$$f(B; W_1, \dots, W_k) = \sum_{i=1}^k \|X_i B W_i - Y_i\|^2. \quad (1.1)$$

Note that we consider the natural parameterization without reweighting the tasks above. The shared body B plays an important role because it allows information transfer between different task data. This is known as the hard parameter sharing architecture in the literature, where we control the capacity r of B , e.g. [18, 23].

We focus on comparing the test performance on a target task using estimators from doing multi-task training and transfer learning. We compare the test performance of these estimators to the single-task baseline. The details are described in Algorithm 1.

By using random matrix theory, we can explain several phenomena that are not explained by the techniques of [23]. [**Todo: list those here**] We refine the task covariance part of [23] into three factors:



Figure 1: Phase transition as a parameter of model distance.

Algorithm 1 Multi-task learning using a hard-parameter sharing architecture

Input: Two regression tasks $(X_1, Y_2), (X_2, Y_2)$.

Parameter: Shared body B , task-specific prediction heads W_1, W_2 .

- 1: Training the shared body B .
 - 2: Optimizing the task heads on the validation set
 - Jointly optimizing both tasks: $\hat{\beta}_t^{\text{MTL}}$
 - Optimizing the target task: $\hat{\beta}_t^{\text{TL-FT}}$
 - Single-task training baseline: $\hat{\beta}_t^{\text{STL}}$
 - 3: Problem statement: how can we compare the test error of the three estimators on the target task?
-

model distance, covariate shift matrix and data ratio [21, 17]. These are achieved through tight generalization bounds established in the high-dimensional regression setting. Our main results are described as follows.

Variance reduction from information transfer. We show that the benefit from doing multi-task or transfer learning stems from reducing the variance of the estimator for the target task through newly added source task data. We derive this result for the setting of two tasks with general inputs and k tasks with the same covariates for any $k \geq 2$. The latter setting is prevalent in applications of multi-task learning to image classification, where there are multiple prediction labels/tasks for every image [12]. On the other hand, the difference between task models causes a negative effect that we call the *model shift bias*. We show bounds on the trade-off between the amount of variance reduced and the amount of model shift bias incurred, which become tighter and tighter as the number of source task data points increases.

Insight 1: $\hat{\beta}_t^{\text{MTL}}$ vs $\hat{\beta}_t^{\text{STL}}$. With multi-task training, since the training objective balances the losses from both the source and target tasks, the trained model can have worse performance for the target task. In particular, if model shift is too large, we get negative transfer from multi-task training.

Insight 2: $\hat{\beta}_t^{\text{TL-FT}}$ vs $\hat{\beta}_t^{\text{MTL}}$. Finally, we show that the transfer learning estimator always improves over multi-task training. The amount of improvement becomes more significant as the model distance becomes larger.

Covariate shift and data ratio. For the case of two tasks with general inputs, we further study the factors that determine the transfer rate. We identify two factors, the covariate shift matrix and the ratio between number of task data.

Insight 3: $\hat{\beta}_t^{\text{TL-FT}}$ vs $\hat{\beta}_t^{\text{STL}}$. Our result has implications on the following question. Is it better for two tasks to have the same covariance matrix or complementary covariance matrices. For our setting, we show that when the data ratio is large, having the same covariance matrix provably yields the lowest test performance on the target task. On the other hand, when data ratio is small, we find that there are cases when having complementary covariance matrices is better. The result provides insight into why the covariance alignment algorithm can help improve performance in [23].

Experimental results.

2 Preliminaries

We assume that for every row x of X_i , we have $\mathbb{E}[xx^\top] = \Sigma_i$. We also write $x = \Sigma_i^{1/2} z_i$, where z_i is a random vector with mean 0 and variance 1. We will designate the k -th task as the target. Our goal is to come up with an estimator $\hat{\beta}$ to provide accurate predictions for the target task, provided with the other auxiliary task data. Concretely, we focus on the test error for the target task:

$$\begin{aligned} te_k(\hat{\beta}) &:= \mathbb{E}_{x \sim \Sigma_k} \left[\mathbb{E}_{\varepsilon_i, \forall 1 \leq i \leq k} \left[(x^\top \hat{\beta} - x^\top \beta_t)^2 \right] \right] \\ &= \mathbb{E}_{\varepsilon_i, \forall 1 \leq i \leq k} \left[(\hat{\beta} - \beta_t)^\top \Sigma_k (\hat{\beta} - \beta_t) \right]. \end{aligned}$$

[Todo: show that $te_k(\hat{\beta}_t^{\text{TL-FT}})$ is less than both $te_k(\hat{\beta}_t^{\text{MTL}})$ and $te_k(\hat{\beta}_t^{\text{STL}})$.]

Hypothesis on Heterogeneous Task Data Our hypothesis is that the heterogeneity among the multiple tasks can be categorized into two classes, *covariate shift* and *model shift*. We consider two natural questions within each category.

- **Model shift.** In general the single-task models can also be different across different tasks. We shall argue that in addition to the bias and variance terms of generalization error, model shift introduces a third term which is the bias caused by model shift. Under model shift, when do we get positive vs. negative transfer? How does the type of transfer depend on the number of data points, the distance of the task models etc?
- **Covariate shift.** The covariance matrices Σ_i may be different across tasks, i.e. having different spectrum or singular vectors. This is also known as covariate shift in the literature. How does covariate shift affect the rate of information transfer? For example, is it better to have the same covariance matrix or not?

The High-Dimensional Setting. We would like to get insight on how covariate and model shifts affect the rate of transfer. We will consider the high-dimensional setting where for the target task, its number of data points is a small constant times p . This setting captures a wide range of applications of multi-task learning where we would like to use auxiliary task data to help train tasks with limited labeled data. Furthermore, this setting is particularly suited to our study since there is need for adding more data to help learn the target task.

For the case of two tasks, we can get precise rates using random matrix theory. For the sake of clarity, we call task 1 the source task and task 2 the target task, i.e. $\beta_1 = \beta_s$ and $\beta_2 = \beta_t$. We introduce the following notations for the high-dimensional setting

$$c_{n_1} := \frac{n_1}{p} \rightarrow c_1, \quad c_{n_2} := \frac{n_2}{p} \rightarrow c_2, \quad \text{as } n_1, n_2 \rightarrow \infty,$$

for some constants $c_1, c_2 \in (1, \infty)$. A crucial quantity is what we call the *covariate shift* matrix $M = \Sigma_1^{1/2} \Sigma_2^{-1/2}$. Let $\lambda_1, \lambda_2, \dots, \lambda_p$ denote the singular values of M .

3 Illustrative Examples and Their Insights on Transfer

We illustrate our main results (to be presented in Section 4) by considering a few special cases, namely special settings of the task models $\{\beta_i\}_{i=1}^k$, covariance matrices $\{\Sigma_i\}_{i=1}^k$, and number of data points $\{n_i\}_{i=1}^k$. We show that our results explain several phenomenon that cannot be explained before. [Todo: list those here]

3.1 Task Model Similarity versus Noise

We compare the test error of $\hat{\beta}_t^{\text{MTL}}$ to that of $\hat{\beta}_t^{\text{STL}}$. For a simple example, we show that whether $\hat{\beta}_t^{\text{MTL}}$ performs better than $\hat{\beta}_t^{\text{STL}}$ is determined by the distance of the task models. We derive a sharp threshold when positive transfer transitions to negative transfer, as a ratio between the model distance and the noise level.

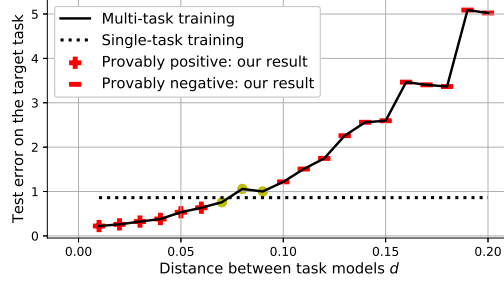


Figure 2: Positive vs negative transfer as a parameter of the task model distances.

Example 3.1. Consider a setting where $\Sigma_1 = \Sigma_2 = \text{Id}$, in other words there is no covariate shift between the two tasks. For the task models, suppose that β_2 has i.i.d. entries with mean zero and variance κ^2 and $\beta_1 - \beta_2$ also has i.i.d. entries with mean 0 and variance d^2 . We have $n_i = c_i \cdot p$ data points from each task, for $i = 1, 2$.

We illustrate the example in a synthetic setting. We demonstrate our result with a simulation. ([**Todo:** uses the tighter bound Proposition A.3?]) We consider a setting where $p = 200$, $n_1 = 90p$, $n_2 = 30p$. [**Todo:** Fill in other params.] We fix the target task and vary the source task, by varying the task model distance parameter d . We show that Theorem 4.1 predicts whether we can get positive or negative transfer. Figure 2 shows the result.

Specifically, the transition threshold is derived in the following proposition.

Proposition 3.2. *In the setting of Example 3.1 with $\sigma_1 = \sigma_2 = \sigma$, assume that $c > 100$ is a fixed constant. Whether $te(\hat{\beta}_t^{MTL})$ is lower than $te(\hat{\beta}_t^{STL})$ is determined by the ratio between the model distance and the noise level:*

- If $d^2 < \frac{2\sigma^2}{3p} \frac{(c_1+c_2-1)^2}{c_1(c_1+c_2)(c_2-1)} - o(1)$, then whp we have that $te(\hat{\beta}_t^{MTL}) < te(\hat{\beta}_t^{STL})$.
- If $d^2 \geq \frac{3\sigma^2}{2p} \frac{(c_1+c_2-1)^2}{c_1(c_1+c_2)(c_2-1)} + o(1)$, then whp we have that $te(\hat{\beta}_t^{MTL}) \geq te(\hat{\beta}_t^{STL})$.

The proof of Proposition 3.2 involves two parts. First, adding the source task has a positive effect of reducing the variance of the estimator, which scales with $n_1 = c_1p$, the number of source task data points. Second, the difference between task models β_1 and β_2 introduces an additional bias term, which scales with pd^2 , the distance between β_1 and β_2 . Hence, the type of transfer is determined by the trade-off between model shift bias and the reduction of variance. The proof can be found in Appendix A.1, which is based on our main result described later in Theorem 4.1.

Consider a more general setting where the noise level σ_1 of task 1 differs from the noise level σ_2 of task 2. We derive a sharp transition similar to Proposition 3.2.

Proposition 3.3. *In the setting of Example 3.1 with d being fixed but σ_1 varies, we derive the following transition as a parameter of σ_1 :*

- If $\sigma_1^2 \leq pd^2 \cdot c_1 + \frac{(c_1+c_2-1)^2}{(c_1+c_2)(c_2-1)} \cdot \sigma_2^2$, then whp $te(\hat{\beta}_t^{MTL}) < te(\hat{\beta}_t^{STL})$.
- If $\sigma_1^2 > pd^2 \cdot c_1 + \frac{(c_1+c_2-1)^2}{(c_1+c_2)(c_2-1)} \cdot \sigma_2^2$, then whp $te(\hat{\beta}_t^{MTL}) > te(\hat{\beta}_t^{STL})$.

Implications. [**Todo:** add connection to taskonomy]

3.2 Data Size and Efficiency

In classical Rademacher or VC based theory of multi-task learning, adding more labeled data improves the generalization performance of a model. On the other hand, we have observed that adding more labeled data does not always improve performance in multi-task learning. Using Example 3.1, we analyze the effect of varying the source task data size.

Proposition 3.4. *In the setting of Example 3.1, assume that $c_1 \geq 100$ and $c_2 \geq 3$ are fixed constants. We have the following conditions to determine whether $te(\hat{\beta}_t^{MTL})$ is lower than $te(\hat{\beta}_t^{STL})$:*

- If $d^2 < \frac{2\sigma^2}{3(c_2-1)p}$, then whp $te(\hat{\beta}_t^{MTL}) < te(\hat{\beta}_t^{STL})$, for any c_1 larger than ??.
- If $d^2 > \frac{3\sigma^2}{2(c_2-1)p}$, then we have the following transition depending on c_1 :
 - If $100 \leq c_1 < \frac{(c_2-2)\sigma^2}{3(c_2-1)pd^2/2-\sigma^2}$, then whp $te(\hat{\beta}_t^{MTL}) < te(\hat{\beta}_t^{STL})$.
 - If $c_1 > \frac{c_2\sigma^2}{2(c_2-1)pd^2/3-\sigma^2}$ and $c_1 \geq 110$, then whp $te(\hat{\beta}_t^{MTL}) > te(\hat{\beta}_t^{STL})$.

Furthermore, as a function of c_1 over the range $[??, \infty]$, the maximum of $te(\hat{\beta}_t^{STL}) - te(\hat{\beta}_t^{MTL})$ is attained when $c_1 = c_2\sigma^2 / \max(2(c_2-1)pd^2 - \sigma^2, 0)$. (cannot get this because we only have some bounds. If we let $c_1 \rightarrow \infty$, then the curve of $te(\hat{\beta}_t^{STL}) - te(\hat{\beta}_t^{MTL})$ already becomes flat, and it is meaningless to discuss the minimum of this function at this point?)

Implications.

3.3 Covariance Alignment

Now we apply Theorem A.5 to show how covariate shift affects the rate of transfer.

Example 3.5 (When $\Sigma_1 = \Sigma_2/\lambda$). In this case, equations (4.2) become

$$a_1 + a_2 = 1 - p/(n_1 + n_2), \quad a_1 + \frac{p}{n_1 + n_2} \cdot \frac{a_1}{a_1 + \lambda^2 a_2} = \frac{n_1}{n_1 + n_2}.$$

By solving these, we can get the test errors (the estimation error behaves similarly). Figure ?? shows how they grow as we increase the number of source task data points. Here $n_2 = 4p$ and n_1 ranges from p to $20p$. We can see that the smaller λ is, the lower the test errors will be.

Example 3.6 (When Σ_1 and Σ_2 are complementary). We now consider another case when Σ_1 and Σ_2 have complementary eigenspaces. Suppose Σ_1 and Σ_2 have the eigendecomposition

$$\Sigma_1^{1/2} = 1 + U\Lambda U^\top, \quad \Sigma_2^{1/2} = 1 + V\Lambda V,$$

where

$$\Lambda = \text{Diag}(\tilde{\lambda}_1, \dots, \tilde{\lambda}_{p/2}), \quad U = (u_1, \dots, u_{p/2}), \quad V = (v_1, \dots, v_{p/2}).$$

If $V = U_\perp$, i.e. the vectors $v_1, \dots, v_{p/2}$ are perpendicular to the vectors $u_1, \dots, u_{p/2}$, then

$$M = \Sigma_1^{1/2} \Sigma_2^{-1/2} = (1 + \Lambda)U U^\top + (1 + \Lambda)^{-1}V V^\top.$$

As a concrete example we consider the case where $\tilde{\lambda}_1 = \dots = \tilde{\lambda}_{p/2}$ and we denote $\lambda := 1 + \tilde{\lambda}_1$. Thus for M , the first $p/2$ singular values are equal to λ and the rest ones are equal to λ^{-1} . In this case, equations in (4.2) become

$$a_1 + a_2 = 1 - \frac{p}{n_1 + n_2}, \quad a_1 + \frac{p}{2(n_1 + n_2)} \cdot \left(\frac{a_1}{a_1 + \lambda^2 a_2} + \frac{a_1}{a_1 + \frac{a_2}{\lambda^2}} \right) = \frac{n_1}{n_1 + n_2}. \quad (3.1)$$

It's not hard to verify that there is only one valid solution (a_1, a_2) to (3.1). After solving these, we get the test error for the target task as follows.

$$te(\lambda) = \frac{p}{2(n_1 + n_2)} \cdot \left(\frac{1}{\frac{a_1}{\lambda^2} + a_2} + \frac{1}{a_1 \lambda^2 + a_2} \right). \quad (3.2)$$

In Figure 3, we plot the test error of the target task for $n_2 = 4p$ and n_1 ranging from p to $20p$. First we notice that the curves in Figure 3 all cross at the point $n_1 = n_2$. In fact, if $n_1 = n_2$, then it is easy to observe that $a_1 = a_2 = (1 - \gamma)/2$ is the solution to equation (3.1), where we denote $\gamma = p/(n_1 + n_2)$. Then for any λ , the test error in (3.2) takes the value

$$te(\lambda) = \frac{\gamma}{2} \frac{1}{(1 - \gamma)/2} = \frac{p}{n_1 + n_2 - p}.$$

Second, we observe the following two phases as we increase n_1/p .

- When $n_1 \leq n_2$, having complementary covariance matrices leads to lower test error compared to the case when $\Sigma_1 = \Sigma_2$.
- When $n_1 > n_2$, having complementary covariance matrices leads to higher test error compared to the case when $\Sigma_1 = \Sigma_2$.

This phenomenon can be also explained using our theory. With (3.1), we can write

$$te(\lambda) = \frac{\gamma}{2} \cdot \left(\frac{1}{\frac{a_1}{\lambda^2} + (1 - \gamma - a_1)} + \frac{1}{a_1 \lambda^2 + (1 - \gamma - a_1)} \right).$$

We can compute that

$$\begin{aligned} te(\lambda) - te(1) &= \frac{\gamma}{2(1 - \gamma)} (\lambda^2 - 1) a_1 \cdot \left(\frac{1}{-a_1(\lambda^2 - 1) + (1 - \gamma)\lambda^2} - \frac{1}{a_1(\lambda^2 - 1) + (1 - \gamma)} \right) \\ &= \frac{\gamma}{2(1 - \gamma)} (\lambda^2 - 1)^2 a_1 \cdot \frac{2a_1 - (1 - \gamma)}{[-a_1(\lambda^2 - 1) + (1 - \gamma)\lambda^2][a_1(\lambda^2 - 1) + (1 - \gamma)]}. \end{aligned}$$

If $n_1 > n_2$, we have $a_1 > (1 - \gamma)/2$ (because $a_1 > a_2$ as observed from the equation (3.1)), and hence $te(\lambda) > te(1)$. Otherwise if $n_1 < n_2$, we have $a_1 < (1 - \gamma)/2$, and hence $te(\lambda) < te(1)$.

We can extend the above observation for Example 3.6 to more general cases with arbitrary covariate shift. To compare different choices of M , we need to fix a scaling for them, because otherwise aM always achieves a smaller error than M for $a > 1$. For this purpose, we introduce the following condition

$$\det(M^\top M) = a^p \Leftrightarrow \prod_{i=1}^p \lambda_i = a^p \quad (3.3)$$

for some constant $a > 0$, and compare different choices of M under this constraint. The next proposition shows that, roughly speaking, as long as there are sufficiently many source task datas, then $M = a \text{Id}$ always gives (approximately) the smallest test error.

Proposition 3.7. *We have that the test error satisfies*

$$te(M) \left(1 - \frac{n_2}{n_1 - p} \frac{1}{\lambda_p^2 + \frac{n_2}{n_1 - p}} \right) \leq \frac{\sigma^2}{n_1 + n_2} \text{Tr} \left(\frac{1}{a_1 M^\top M + a_2} \right) \leq te(M), \quad (3.4)$$

where λ_p is the smallest singular value of p and

$$te(M) := \frac{\sigma^2}{a_1(n_1 + n_2)} \text{Tr} \left(\frac{1}{M^\top M} \right).$$

Moreover, for all M satisfying (3.3), the minimum of $te(M)$ is attained when $M = a \text{Id}$.

Proof. From equation (4.2), we get

$$a_1 \geq \frac{n_1 - p}{n_1 + n_2}, \quad a_2 \leq \frac{n_2}{n_1 + n_2}.$$

With these two bounds, we can easily conclude (3.4).

For the second statement, we find that we need to minimize

$$\text{Tr} \left(\frac{1}{M^\top M} \right) = \sum_{i=1}^p \frac{1}{\lambda_i}, \quad \text{under } \prod_{i=1}^p \lambda_i = a^p.$$

Then using AM-GM inequality, we conclude that the sum $\sum_{i=1}^p \lambda_i^{-1}$ is smallest when $\lambda_1 = \dots = \lambda_p = a$. \square

From the above proposition, we see that when $n_1 \gg n_2$, we have that the test error is approximately given by $te(M)$ (as long as we impose a proper lower bound on λ_p). Moreover, $te(M)$ is minimized when Σ_1 and Σ_2 are proportional to each other, i.e. there is no covariate shift between the source task data and target task data. This provides a theoretical evidence that in general covariate shift is unfavored in transfer learning if we have enough source task data, although Example 3.6 show that this may not be true when the number of source task data is small.

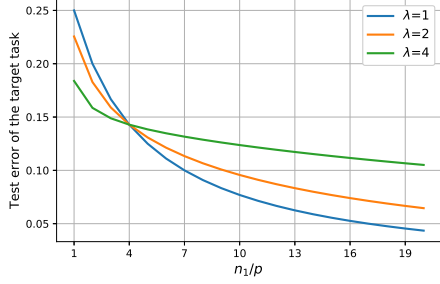


Figure 3: When Σ_1 and Σ_2 are complementary. The number of target task data points is $n_2 = 4p$.

4 Main Results: A Model Shift Bias versus Variance Trade-off

In this part, we consider the case of two tasks to show establish the intuition that adding more data helps in multi-task learning by reducing the variance of the estimator. We achieve this through tight generalization bounds obtained from random matrix theory. For the case of two tasks, we identify three factors that determine the type of transfer between tasks: model distance, covariate shift matrix, and data ratio.

Recall that the test error of $\hat{\beta}_t^{\text{MTL}}$ consists of two parts

$$te(\hat{\beta}_t^{\text{MTL}}) = \hat{w}^2 \left\| \Sigma_2^{1/2} (\hat{w}^2 X_1^\top X_1 + X_2^\top X_2)^{-1} X_1^\top X_1 (\beta_s - \hat{w} \beta_t) \right\|^2 + \sigma^2 \cdot \text{Tr} [(\hat{w}^2 X_1^\top X_1 + X_2^\top X_2)^{-1} \Sigma_2]. \quad (4.1)$$

It is not hard to show that the variance of $\hat{\beta}_t^{\text{MTL}}$ is reduced compared to $\hat{\beta}_t^{\text{STL}}$ (following the argument of Proposition ??), i.e.

$$\sigma^2 \cdot \text{Tr} [(\hat{w}^2 X_1^\top X_1 + X_2^\top X_2)^{-1} \Sigma_2] \leq \sigma^2 \cdot \text{Tr} [(X_2^\top X_2)^{-1} \Sigma_2].$$

Because of model shift however, i.e. $\beta_s \neq \beta_t$. We can no longer guarantee that $te(\hat{\beta}_t^{\text{MTL}}) \leq te(\hat{\beta}_t^{\text{STL}})$. The main result of this part show deterministic conditions under which we get positive or negative transfer. And the conditions depend only on the covariate shift matrix M , the difference of the task models, and the number of per-task data points. In order to characterize $te(\hat{\beta}_t^{\text{MTL}})$ and $te(\hat{\beta}_t^{\text{STL}})$, the technical crux of our approach relies on deriving the limit of the trace of matrix inverse in the high-dimensional setting. To illustrate the idea, we observe that by using Lemma A.1, we have that

$$te(\hat{\beta}_t^{\text{STL}}) = \frac{\sigma^2}{n_2 - p} \text{Tr} [\Sigma_2^{-1}].$$

We shall also derive the limit of $te(\hat{\beta}_t^{\text{MTL}})$. [Todo: write a brief technical overview]

Let $M = \hat{w} \Sigma_1^{1/2} \Sigma_2^{-1/2}$ denote the weighted covariate shift matrix. Denote by $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$ the singular values of $M^\top M$. Let (a_1, a_2) be the solutions to the following system of equations

$$a_1 + a_2 = 1 - \frac{p}{n_1 + n_2}, \quad a_1 + \sum_{i=1}^p \frac{a_1}{(n_1 + n_2)(a_1 + a_2/\lambda_i^2)} = \frac{n_1}{n_1 + n_2}. \quad (4.2)$$

After obtaining (a_1, a_2) , we can solve the following linear equations to get (a_3, a_4) :

$$\left(\frac{n_2}{a_2^2} - \sum_{i=1}^p \frac{1}{(a_2 + \lambda_i^2 a_1)^2} \right) a_3 - \left(\sum_{i=1}^p \frac{\lambda_i^2}{(a_2 + \lambda_i^2 a_1)^2} \right) a_4 = \sum_{i=1}^p \frac{1}{(a_2 + \lambda_i^2 a_1)^2}, \quad (4.3)$$

$$\left(\frac{n_1}{a_1^2} - \sum_{i=1}^p \frac{\lambda_i^4}{(a_2 + \lambda_i^2 a_1)^2} \right) a_4 - \left(\sum_{i=1}^p \frac{\lambda_i^2}{(a_2 + \lambda_i^2 a_1)^2} \right) a_3 = \sum_{i=1}^p \frac{\lambda_i^2}{(a_2 + \lambda_i^2 a_1)^2}. \quad (4.4)$$

Then we introduce the following matrix

$$Z = \frac{n_1^2}{(n_1 + n_2)^2} \cdot M \frac{(1 + a_3) \text{Id} + a_4 M^\top M}{(a_2 + a_1 M^\top M)^2} M^\top,$$

which can be regarded as the asymptotic limit of $\hat{w}\Sigma_2^{1/2}(\hat{w}^2 X_1^\top X_1 + X_2^\top X_2)^{-1}\Sigma_1^{1/2}$. Finally we introduce

$$\delta := \left[\frac{n_1 \lambda_1}{(\sqrt{n_1} - \sqrt{p})^2 \lambda_p + (\sqrt{n_2} - \sqrt{p})^2} \right]^2 \cdot \|\Sigma_1^{1/2}(\beta_s - \hat{w}\beta_t)\|^2.$$

may be able to get a better bound, but the statement will be long

We now state our main result for two tasks with both covariate and model shift in the following theorem.

Theorem 4.1. *Let n_1, n_2 be the number of data points for the source, target task, respectively. Let \hat{w} denote the optimal solution for the ratio w_1/w_2 in equation (A.2). The information transfer is solely determined by two deterministic quantities Δ_β and Δ_{var} , which show the change of model shift bias and variance, respectively. With high probability we have*

$$te(\hat{\beta}_t^{MTL}) \leq te(\hat{\beta}_t^{STL}) \text{ when: } \Delta_{var} - \Delta_\beta \geq \left(\left(1 + \sqrt{\frac{p}{n_1}} \right)^4 - 1 \right) \delta \quad (4.5)$$

$$te(\hat{\beta}_t^{MTL}) \geq te(\hat{\beta}_t^{STL}) \text{ when: } \Delta_{var} - \Delta_\beta \leq -2 \left(2\sqrt{\frac{p}{n_1}} + \frac{p}{n_1} \right) \delta, \quad (4.6)$$

where

$$\Delta_{var} := \sigma^2 \left(\frac{p}{n_2 - p} - \frac{1}{n_1 + n_2} \text{Tr}[(a_1 M^\top M + a_2 \text{Id})^{-1}] \right)$$

$$\Delta_\beta := (\beta_s - \hat{w}\beta_t)^\top \Sigma_1^{1/2} Z \Sigma_1^{1/2} (\beta_s - \hat{w}\beta_t).$$

Theorem 4.1 shows upper and lower bounds that guarantee positive transfer, which is determined by the change of variance Δ_{var} and a certain model shift bias parameter Δ_β determined by the covariate shift matrix and the model shift. The bounds get tighter and tighter as n_1/p increases.

4.1 Extension to Many Tasks of the Same Covariates

In this section we consider the setting with k many that have the same covariates. Since every task has the same number of data points as well as the same covariance, the only differences between different tasks are their models $\{\beta_i\}_{i=1}^k$. For this setting, we derive solutions for the multi-task training and the transfer learning setting that match our insights qualitatively from Section ??.

5 Related Work

Adding a regularization over B , e.g. [19, 20]. Moreover, [18] observed that controlling the capacity can outperform the implicit capacity control of adding regularization over B .

References

- [1] Johannes Alt. Singularities of the density of states of random Gram matrices. *Electron. Commun. Probab.*, 22:13 pp., 2017.
- [2] Johannes Alt, L. Erdős, and Torben Krüger. Local law for random Gram matrices. *Electron. J. Probab.*, 22:41 pp., 2017.
- [3] 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.
- [4] 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.
- [5] A. Bloemendal, A. Knowles, H.-T. Yau, and J. Yin. On the principal components of sample covariance matrices. *Prob. Theor. Rel. Fields*, 164(1):459–552, 2016.
- [6] P. Bourgade, H.-T. Yau, and J. Yin. Local circular law for random matrices. *Probab. Theory Relat. Fields*, 159:545–595, 2014.

- [7] Xiucan 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.
- [8] L. Erdős, A. Knowles, and H.-T. Yau. Averaging fluctuations in resolvents of random band matrices. *Ann. Henri Poincaré*, 14:1837–1926, 2013.
- [9] L. Erdős, A. Knowles, H.-T. Yau, and J. Yin. Delocalization and diffusion profile for random band matrices. *Commun. Math. Phys.*, 323:367–416, 2013.
- [10] L. Erdős, A. Knowles, H.-T. Yau, and J. Yin. The local semicircle law for a general class of random matrices. *Electron. J. Probab.*, 18:1–58, 2013.
- [11] 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.
- [12] Sabri Eyuboglu, Geoffrey Angus, Bhavik N. Patel, Anuj Pareek, Guido Davidzon, Jared Dunmon, and Matthew P. Lungren. Multi-task weak supervision enables automated abnormality localization in whole-body fdg-pet/ct. 2020.
- [13] Viacheslav Leonidovich Girko. *Theory of random determinants*, volume 45. Springer Science & Business Media, 2012.
- [14] VL Girko. Random matrices. *Handbook of Algebra*, ed. Hazewinkel, 1:27–78, 1975.
- [15] Vyacheslav L Girko. Spectral theory of random matrices. *Russian Mathematical Surveys*, 40(1):77, 1985.
- [16] Antti Knowles and Jun Yin. Anisotropic local laws for random matrices. *Probability Theory and Related Fields*, pages 1–96, 2016.
- [17] Wouter M Kouw. An introduction to domain adaptation and transfer learning. *arXiv preprint arXiv:1812.11806*, 2018.
- [18] Abhishek Kumar and Hal Daume III. Learning task grouping and overlap in multi-task learning. *arXiv preprint arXiv:1206.6417*, 2012.
- [19] 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.
- [20] 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.
- [21] Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2009.
- [22] Natesh S. Pillai and Jun Yin. Universality of covariance matrices. *Ann. Appl. Probab.*, 24:935–1001, 2014.
- [23] Sen Wu, Hongyang Zhang, and Christopher Ré. Understanding and improving information transfer in multi-task learning. In *International Conference on Learning Representations*, 2020.
- [24] 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.
- [25] Fan Yang. Edge universality of separable covariance matrices. *Electron. J. Probab.*, 24:57 pp., 2019.

A Supplementary Materials for Section 3

From [23], 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 when $r = 1$ we are considering the case of two tasks. Here, equation (1.1) simplifies to the following

$$f(B; w_1, w_2) = \|X_1 B w_1 - Y_1\|^2 + \|X_2 B w_2 - Y_2\|^2, \quad (\text{A.1})$$

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

$$\begin{aligned} \hat{B}(w_1, w_2) &= (w_1^2 X_1^\top X_1 + w_2^2 X_2^\top X_2)^{-1} (w_1 X_1^\top Y_1 + w_2 X_2^\top Y_2) \\ &= \frac{1}{w_2} \left(\left(\frac{w_1}{w_2} \right)^2 X_1^\top X_1 + X_2^\top X_2 \right)^{-1} \left(\frac{w_1}{w_2} X_1^\top Y_1 + X_2^\top Y_2 \right) \\ &= \frac{1}{w_2} \left(\beta_t + \left(\left(\frac{w_1}{w_2} \right)^2 X_1^\top X_1 + X_2^\top X_2 \right)^{-1} \left(X_1^\top X_1 \left(\frac{w_1}{w_2} \beta_s - w^2 \beta_t \right) + \left(\frac{w_1}{w_2} X_1^\top \varepsilon_1 + X_2^\top \varepsilon_2 \right) \right) \right). \end{aligned}$$

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 θ_1 is a scaling of θ_2 , then this case can be solved optimally using equation (A.1) [18].

Defining the multi-task learning estimator. Using a validation set that is sub-sampled from the original training dataset, we get a validation loss as follows

$$\begin{aligned} val(\hat{B}; w_1, w_2) &= n_1 \cdot \left\| \Sigma_1^{1/2} \left(\left(\frac{w_1}{w_2} \right)^2 X_1^\top X_1 + X_2^\top X_2 \right)^{-1} X_2^\top X_2 \left(\beta_s - \frac{w_1}{w_2} \beta_t \right) \right\|^2 \\ &\quad + n_1 \sigma^2 \cdot \text{Tr} \left[\left(\left(\frac{w_1}{w_2} \right)^2 X_1^\top X_1 + X_2^\top X_2 \right)^{-1} \Sigma_1 \right] \\ &\quad + n_2 \cdot w^2 \left\| \Sigma_2^{1/2} \left(\left(\frac{w_1}{w_2} \right)^2 X_1^\top X_1 + X_2^\top X_2 \right)^{-1} X_1^\top X_1 \left(\beta_s - \frac{w_1}{w_2} \beta_t \right) \right\|^2 \\ &\quad + n_2 \cdot \sigma^2 \cdot \text{Tr} \left[\left(\left(\frac{w_1}{w_2} \right)^2 X_1^\top X_1 + X_2^\top X_2 \right)^{-1} \Sigma_2 \right]. \end{aligned} \quad (\text{A.2})$$

Let \hat{w}_1/\hat{w}_2 be the global minimizer of $val(\hat{B}; w_1, w_2)$. We will define the multi-task learning estimator for the target task as

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

The intuition for deriving $\hat{\beta}_t^{\text{MTL}}$ is akin to performing multi-task training in practice. Let $\hat{v} = \hat{w}_1/\hat{w}_2$ for the simplicity of notation. The test loss of using $\hat{\beta}_t^{\text{MTL}}$ for the target task is

$$\begin{aligned} te(\hat{\beta}_t^{\text{MTL}}) &= \hat{v}^2 \left\| \Sigma_2^{1/2} (\hat{v}^2 X_1^\top X_1 + X_2^\top X_2)^{-1} X_1^\top X_1 (\beta_s - \hat{v} \beta_t) \right\|^2 \\ &\quad + \sigma^2 \cdot \text{Tr} [(\hat{v}^2 X_1^\top X_1 + X_2^\top X_2)^{-1} \Sigma_2]. \end{aligned} \quad (\text{A.3})$$

Our goal is to study under model and covariate shifts, whether multi-task learning helps learn the target task better than single-task learning. The baseline where we solve the target task with its own data is

$$te(\hat{\beta}_t^{\text{STL}}) = \sigma^2 \cdot \text{Tr} [(X_1^\top X_1 + X_2^\top X_2)^{-1}], \text{ where } \hat{\beta}_t^{\text{STL}} = (X_2^\top X_2)^{-1} X_2^\top Y_2.$$

We state several helper lemmas to get a bound on the variance of $\hat{\beta}_t^{\text{STL}}$ and $\hat{\beta}_t^{\text{MTL}}$. The first lemma, which is a folklore result in random matrix theory, helps determine the asymptotic limit of $te(\hat{\beta}_t^{\text{STL}})$, as p goes to infinity.

Lemma A.1. *[[TODO: ref?](#)] Let $X \in \mathbb{R}^{n \times p}$ be a random matrix that contains i.i.d. row vectors with mean 0 and covariance $\Sigma \in \mathbb{R}^{p \times p}$. In the setting when $n = cp$ we have that as p goes to infinity,*

$$\text{Tr} [(X^\top X)^{-1} \Sigma] = \frac{1}{c-1}.$$

The second lemma, which deals the inverse of the sum of two random matrices, can be viewed as a special case of Theorem 4.1.

Lemma A.2. *Let $X_i \in \mathbb{R}^{n_i \times p}$ be a random matrix that contains i.i.d. row vectors with mean 0 and variance $\Sigma_i \in \mathbb{R}^{p \times p}$, for $i = 1, 2$. 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 decreasing order. When $n_1 = c_1 p$ and $n_2 = c_2 p$, we have that with high probability over the randomness of X_1 and X_2 , the following equation holds*

$$\text{Tr}[(X_1^\top X_1 + X_2^\top X_2)^{-1} \Sigma_2] = \frac{1}{c_1 + c_2} \cdot \text{Tr}\left[\frac{1}{p} \cdot (a_1 M^\top M + a_2)^{-1}\right] + O(n^{-1/2+\varepsilon}), \quad (\text{A.4})$$

for any constant $\varepsilon > 0$, where a_1, a_2 are solutions to the following deterministic equations:

$$a_1 + a_2 = 1 - \frac{1}{c_1 + c_2}, \quad a_1 + \frac{1}{p} \cdot \sum_{i=1}^p \frac{a_1}{(c_1 + c_2)(a_1 + a_2/\lambda_i^2)} = \frac{c_1}{c_1 + c_2}. \quad (\text{A.5})$$

We will give the proof of Lemma A.2 in Section C.

Finally, we can get a bound tighter than Theorem 4.1 as follows. [Todo: restate this]

Lemma A.3. *In the setting of Theorem 4.1, assume that $\Sigma_1 = \text{Id}$, β_t is i.i.d. with mean 0 and variance κ^2 and $\beta_s - \beta_t$ is i.i.d. with mean 0 and variance d^2 . We set $\Delta_\beta = ((1 - \hat{w})^2 \kappa^2 + d^2) \text{Tr}[Z]$ and we have*

$$\begin{aligned} te(\hat{\beta}_t^{MTL}) &\leq te(\hat{\beta}_t^{STL}) \text{ when: } \Delta_{var} \geq \left(1 + \sqrt{\frac{p}{n_1}}\right)^4 \Delta_\beta, \\ te(\hat{\beta}_t^{MTL}) &\geq te(\hat{\beta}_t^{STL}) \text{ when: } \Delta_{var} \leq \left(1 - \sqrt{\frac{p}{n_1}}\right)^4 \Delta_\beta. \end{aligned}$$

A.1 Proof of Proposition 3.2 and Proposition 3.3

The proof will consist of two main steps.

- First, we show that \hat{v} is close to 1.
- Second, we plug \hat{v} back into $te(\hat{\beta}_t^{MTL})$ to show the result.

We denote

$$\begin{aligned} val(w) &= n_1 \left[d^2 + (w - 1)^2 \kappa^2 \right] \cdot \text{Tr}[(w^2 X_1^\top X_1 + X_2^\top X_2)^{-2} (X_2^\top X_2)^2] \\ &\quad + n_2 w^2 \left[d^2 + (w - 1)^2 \kappa^2 \right] \cdot \text{Tr}[(w^2 X_1^\top X_1 + X_2^\top X_2)^{-2} (X_1^\top X_1)^2] \\ &\quad + (n_1 + n_2) \sigma^2 \cdot \text{Tr}[(w^2 X_1^\top X_1 + X_2^\top X_2)^{-1}]. \end{aligned}$$

Under the setting of Proposition A.3, using concentration of random vector with i.i.d. Gaussian entries, Lemma D.6, we have that

$$val(\hat{B}; w_1, w_2) = val(w) \left(1 + O(p^{-1/2})\right) \quad \text{with probability } 1 - o(1).$$

Thus it suffices to study the behavior of $val(w)$. For the minimizer \hat{w} of $val(w)$, we have a similar result as in Proposition A.5.

Lemma A.4. *Suppose the assumptions of Proposition A.3 hold. Assume that $\kappa^2 \sim p d^2 \sim \sigma^2$ are of the same order. Then we have that the optimal ratio for $val(w)$ satisfies*

$$|\hat{w} - 1| = O(p^{-1}).$$

Proof. The proof is also similar to the one for Proposition A.5. First it is easy to observe that $val(w) \leq val(-w)$ for $w \geq 0$. Hence it suffices to consider the $w \geq 0$ case.

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

$$\begin{aligned} \text{val}(w) &= n_1 \left[\frac{d^2}{w^4} + \frac{(w-1)^2}{w^4} \kappa^2 \right] \cdot \text{Tr} [(X_1^\top X_1 + w^{-2} X_2^\top X_2)^{-2} (X_2^\top X_2)^2] \\ &\quad + n_2 \left[\frac{d^2}{w^2} + \frac{(w-1)^2}{w^2} \kappa^2 \right] \cdot \text{Tr} [(X_1^\top X_1 + w^{-2} X_2^\top X_2)^{-2} (X_1^\top X_1)^2] \\ &\quad + \sigma^2 (n_1 + n_2) \cdot \text{Tr} [(w^2 X_1^\top X_1 + X_2^\top X_2)^{-1}]. \end{aligned}$$

Taking derivative of $\text{val}(w)$ with respect to w , we obtain that

$$\begin{aligned} \text{val}'(w) &\geq n_1 \left[\frac{2(w-1)(2-w)}{w^5} \kappa^2 - \frac{4d^2}{w^5} \right] \text{Tr} [(X_1^\top X_1 + w^{-2} X_2^\top X_2)^{-2} (X_2^\top X_2)^2] \\ &\quad + n_2 \left[\frac{2(w-1)}{w^3} \kappa^2 - \frac{2d^2}{w^3} \right] \cdot \text{Tr} [(X_1^\top X_1 + w^{-2} X_2^\top X_2)^{-2} (X_1^\top X_1)^2] \\ &\quad - 2(n_1 + n_2) \frac{\sigma^2}{w^3} \cdot \text{Tr} [(X_1^\top X_1 + w^{-2} X_2^\top X_2)^{-2} X_1^\top X_1] = n_1 \text{Tr} [(X_1^\top X_1 + w^{-2} X_2^\top X_2)^{-2} \mathcal{A}], \end{aligned}$$

where the matrix \mathcal{A} is

$$\mathcal{A} := \left[\frac{2(w-1)(2-w)}{w^5} \kappa^2 - \frac{4d^2}{w^5} \right] (X_2^\top X_2)^2 + \frac{n_2}{n_1} \left[\frac{2(w-1)}{w^3} \kappa^2 - \frac{2d^2}{w^3} \right] (X_1^\top X_1)^2 - 2 \frac{n_1 + n_2}{n_1} \frac{\sigma^2}{w^3} X_1^\top X_1$$

Using the estimate (B.3), we get that \mathcal{A} is lower bounded as

$$\mathcal{A} \succeq -\frac{4d^2}{w^5} (\sqrt{n_2} + \sqrt{p})^4 + \frac{n_2}{n_1} \left[\frac{2(w-1)}{w^3} \kappa^2 - \frac{2d^2}{w^3} \right] (\sqrt{n_1} - \sqrt{p})^4 - 2 \frac{n_1 + n_2}{n_1} \frac{\sigma^2}{w^3} (\sqrt{n_1} + \sqrt{p})^2 \succ 0,$$

as long as

$$w > w_1 := 1 + \frac{d^2}{\kappa^2} + \frac{\sigma^2 (n_1 + n_2) (\sqrt{n_1} + \sqrt{p})^2}{\kappa^2 n_2 (\sqrt{n_1} - \sqrt{p})^4} + \frac{2d^2 n_1 (\sqrt{n_2} + \sqrt{p})^4}{\kappa^2 n_2 (\sqrt{n_1} - \sqrt{p})^4}.$$

Hence $\text{val}'(w) > 0$ on (w_1, ∞) , i.e. $\text{val}(w)$ is strictly increasing for $w > w_1$. Hence we must have $\hat{w} \leq w_1$. Note that under our assumptions, we have $w_1 = 1 + O(p^{-1})$.

Then we consider the case $w \leq 1$. Taking derivative of $\text{val}(w)$ with respect to w , we obtain that

$$\begin{aligned} \text{val}'(w) &\leq n_1 [2(w-1)\kappa^2] \cdot \text{Tr} [(w^2 X_1^\top X_1 + X_2^\top X_2)^{-2} (X_2^\top X_2)^2] \\ &\quad + n_2 [2d^2 w + 2w(w-1)(2w-1)\kappa^2] \cdot \text{Tr} [(w^2 X_1^\top X_1 + X_2^\top X_2)^{-2} (X_1^\top X_1)^2] \\ &= n_1 \text{Tr} [(w^2 X_1^\top X_1 + X_2^\top X_2)^{-1} \mathcal{B}], \end{aligned}$$

where the matrix \mathcal{B} is

$$\mathcal{B} = 2(w-1)\kappa^2 (X_2^\top X_2)^2 + \frac{n_2}{n_1} [2d^2 w + 2w(w-1)(2w-1)\kappa^2] (X_1^\top X_1)^2.$$

Using the estimate (B.3), we get that \mathcal{B} is upper bounded as

$$\mathcal{B} \preceq -2(1-w)\kappa^2 (\sqrt{n_2} - \sqrt{p})^4 + 2d^2 w \frac{n_2}{n_1} (\sqrt{n_1} + \sqrt{p})^4 \prec 0,$$

as long as

$$w < w_2 := 1 - \frac{d^2 n_2 (\sqrt{n_1} + \sqrt{p})^4}{\kappa^2 n_1 (\sqrt{n_2} - \sqrt{p})^4}.$$

Hence $\text{val}'(w) < 0$ on $[0, w_2)$, i.e. $\text{val}(w)$ is strictly increasing for $w < w_2$. Hence we must have $\hat{w} \leq w_2$. Note that under our assumptions, we have $w_2 = 1 - O(p^{-1})$. \square

Based on Lemma A.4, we can prove Proposition 3.2.

Proof of Proposition 3.2. Since $\Sigma_1 = \Sigma_2$, we know that $M = \Sigma_1^{1/2} \Sigma_2^{-1/2} = \hat{v} \text{Id}$. Using Lemma A.1 and A.2, we can track the reduction of variance from $\hat{\beta}_t^{\text{MTL}}$ to $\hat{\beta}_t^{\text{STL}}$ as whp

$$\begin{aligned} \Delta_{\text{var}} &:= \sigma^2 \left(\frac{1}{c_2 - 1} - \frac{1}{c_1 + c_2} \cdot \frac{1}{a_1 \hat{v}^2 + a_2} + O\left(p^{-1/2+\varepsilon}\right) \right) \\ &= \sigma^2 \left(\frac{c_1}{(c_2 - 1)(c_1 + c_2 - 1)} + O\left(p^{-1/2+\varepsilon}\right) \right), \end{aligned} \quad (\text{A.6})$$

where we use equation (A.5) and Lemma A.4. Next we consider the model shift bias

$$\begin{aligned} \Delta_{\beta} &:= \hat{v}^2 \left\| (\hat{v}^2 X_1^\top X_1 + X_2^\top X_2)^{-1} X_1^\top X_1 (\beta_s - \hat{v} \beta_t) \right\|^2 \\ &= d^2 \cdot \left\| (\hat{v}^2 X_1^\top X_1 + X_2^\top X_2)^{-1} X_1^\top X_1 \right\|_F^2 + O\left(p^{-1/2+\varepsilon} d^2\right) \end{aligned}$$

Using Lemma A.3, we get an upper and lower bound on Δ_{β} as

$$\left(1 - \sqrt{\frac{1}{c_1}}\right)^4 \leq \Delta_{\beta} / \left(p \cdot d^2 \cdot \frac{c_1^2}{(c_1 + c_2)^2} \cdot \frac{1 + a_3 + a_4}{(a_1 + a_2)^2} + O\left(p^{1/2+\varepsilon} d^2\right)\right) \leq \left(1 + \sqrt{\frac{1}{c_1}}\right)^4. \quad (\text{A.7})$$

By solving equations (A.5), (4.3) and (4.4), we get

$$\begin{aligned} a_1 &= \frac{c_1(c_1 + c_2 - 1)}{(c_1 + c_2)^2} + O(p^{-1/2+\varepsilon}), a_2 = \frac{c_2(c_1 + c_2 - 1)}{(c_1 + c_2)^2} + O(p^{-1/2+\varepsilon}), \\ a_3 &= \frac{c_2}{(c_1 + c_2)(c_1 + c_2 - 1)} + O(p^{-1/2+\varepsilon}), a_4 = \frac{c_1}{(c_1 + c_2)(c_1 + c_2 - 1)} + O(p^{-1/2+\varepsilon}). \end{aligned}$$

Hence we obtain that

$$\frac{1 + a_3 + a_4}{(a_1 + a_2)^2} = \frac{(c_1 + c_2)^3}{(c_1 + c_2 - 1)^3} + O\left(p^{-1/2+\varepsilon}\right).$$

To sum up, we have shown that when

$$\frac{p \cdot d^2}{\sigma^2} < \frac{(c_1 + c_2 - 1)^2}{c_1(c_1 + c_2)(c_2 - 1)} \cdot \left(1 + \sqrt{\frac{1}{c_1}}\right)^{-4},$$

we have that $\Delta_{\text{var}} > \Delta_{\beta}$, which implies that $te(\hat{\beta}_t^{\text{MTL}})$ is lower than $te(\hat{\beta}_t^{\text{STL}})$. When

$$\frac{p \cdot d^2}{\sigma^2} > \frac{(c_1 + c_2 - 1)^2}{c_1(c_1 + c_2)(c_2 - 1)} \cdot \left(1 - \sqrt{\frac{1}{c_1}}\right)^{-4},$$

we have that $\Delta_{\text{var}} < \Delta_{\beta}$, which implies that $te(\hat{\beta}_t^{\text{MTL}})$ is higher than $te(\hat{\beta}_t^{\text{STL}})$. The result follows by taking $c_1 = c_2 = c$ for the above equation. \square

A.2 Proof of Proposition 3.4

Recall that under the setting of Example 3.1, (A.6) and (A.7) hold. Then we get that $te(\hat{\beta}_t^{\text{MTL}}) < te(\hat{\beta}_t^{\text{STL}})$ whp if

$$\left(p d^2 \cdot \frac{c_1(c_1 + c_2)(c_2 - 1)}{(c_1 + c_2 - 1)^2} + O\left(p^{1/2+\varepsilon} d^2\right)\right) \cdot \left(1 + \sqrt{\frac{1}{c_1}}\right)^4 < \sigma^2 \left(1 + O\left(p^{-1/2+\varepsilon}\right)\right); \quad (\text{A.8})$$

otherwise if

$$\left(p d^2 \cdot \frac{c_1(c_1 + c_2)(c_2 - 1)}{(c_1 + c_2 - 1)^2} + O\left(p^{1/2+\varepsilon} d^2\right)\right) \cdot \left(1 - \sqrt{\frac{1}{c_1}}\right)^4 > \sigma^2 \left(1 + O\left(p^{-1/2+\varepsilon}\right)\right), \quad (\text{A.9})$$

then we have $te(\hat{\beta}_t^{\text{MTL}}) > te(\hat{\beta}_t^{\text{STL}})$ whp.

Now we prove the first statement of Proposition 3.4. Notice that the function

$$\frac{c_1(c_1 + c_2)(c_2 - 1)}{(c_1 + c_2 - 1)^2} = (c_2 - 1) \left(1 + \frac{c_2 - 2}{c_1} + \frac{1}{c_1(c_1 + c_2)}\right)^{-1}$$

is strictly increasing with respect to c_1 as long as $c_2 \geq 3$. In particular, for $c_1 \geq 100$, we get the bound by taking $c_1 \rightarrow \infty$:

$$\frac{c_1(c_1 + c_2)(c_2 - 1)}{(c_1 + c_2 - 1)^2} \cdot \left(1 + \sqrt{\frac{1}{c_1}}\right)^4 < 1.1^4(c_2 - 1).$$

Hence we conclude that for $c_1 \geq 100$, we always have $te(\hat{\beta}_t^{\text{MTL}}) < te(\hat{\beta}_t^{\text{STL}})$ whp as long as $pd^2 \leq \frac{2\sigma^2}{3(c_2-1)}$. This gives the first statement.

The second statement can be proved in a similar way. Suppose $c_1 \geq 110$ and $pd^2 > \frac{3\sigma^2}{2(c_2-1)}$. If $c_1 > \frac{(c_2-2)\sigma^2}{2(c_2-1)pd^2/3-\sigma^2}$, we have

$$pd^2 \cdot \frac{c_1(c_1 + c_2)(c_2 - 1)}{(c_1 + c_2 - 1)^2} \cdot \left(1 - \sqrt{\frac{1}{c_1}}\right)^4 > \frac{2}{3}pd^2(c_2 - 1) \cdot \frac{c_1}{c_1 + (c_2 - 2)} > \sigma^2,$$

where in the first step we used that

$$\frac{(c_1 + c_2)(c_1 + c_2 - 2)}{(c_1 + c_2 - 1)^2} \cdot \left(1 - \sqrt{\frac{1}{c_1}}\right)^4 > \frac{2}{3}$$

for $c_1 \geq 110$. This shows that (A.9) holds as long as p is large enough, and hence $te(\hat{\beta}_t^{\text{MTL}}) > te(\hat{\beta}_t^{\text{STL}})$ holds. On the other hand, if $100 \leq c_1 \leq \frac{(c_2-2)\sigma^2}{3(c_2-1)pd^2/2-\sigma^2}$, then we have

$$pd^2 \cdot \frac{c_1(c_1 + c_2)(c_2 - 1)}{(c_1 + c_2 - 1)^2} \cdot \left(1 + \sqrt{\frac{1}{c_1}}\right)^4 < \frac{3}{2}pd^2(c_2 - 1) \cdot \frac{c_1}{c_1 + (c_2 - 2)} < \sigma^2,$$

where in the first step we used that

$$\frac{(c_1 + c_2)(c_1 + c_2 - 2)}{(c_1 + c_2 - 1)^2} \cdot \left(1 + \sqrt{\frac{1}{c_1}}\right)^4 < \frac{3}{2}$$

for $c_1 \geq 100$. This shows that (A.8) holds as long as p is large enough, and hence $te(\hat{\beta}_t^{\text{MTL}}) < te(\hat{\beta}_t^{\text{STL}})$ holds.

A.3 Proof of ???

Note that for the case of k tasks with the same covariates, since there is no covariate shift and the data ratio is always equal to one, the main factor is model distance.

A precise bound when there is no model shift. As Proposition ?? shows, if β_s and β_t are equal, then adding the source task dataset always helps learn the target task. The goal of this section is to understand how covariate shift affects the rate of transfer. [Todo: add conceptual msg]

A simple observation here is that when $\beta_s = \beta_t$, the optimal \hat{w} for minimizing equation (??) is equal to 1. Based on this observation, we can get a more precise result than Theorem 4.1 on the improvement of adding the source task data that only depends on the covariance matrices Σ_1, Σ_2 and the number of data points n_1, n_2 .

Proposition A.5 (Transfer rate without model shift). *Suppose $\beta_s = \beta_t$ and $\|\beta_t\|_2^2 \sim p\sigma^2$ (i.e. the l^2 -norm of the vector β_t is of the same order as that of the error vector). Assume that the condition numbers of Σ_1, Σ_2 and $M := \Sigma_1^{1/2}\Sigma_2^{-1/2}$ are all bounded by a constant $C > 0$. Then we have that the optimal ratio for W_1/W_2 in equation (??) satisfies*

$$1 \leq \hat{w} \leq 1 + O(p^{-1}).$$

Moreover, we have

$$te(\hat{\beta}_t^{\text{TL-FT}}) = \sigma^2 \cdot \text{Tr} \left[((n_1 + n_2)a_1 M^\top M + (n_1 + n_2)a_2 \text{Id})^{-1} \right] \cdot (1 + O(p^{-1})), \quad (\text{A.10})$$

where a_1, a_2 are the solutions to equations (4.2).

Proof. We abbreviate $\text{val}(w_2 \hat{B}(w)) := \text{val}(w)$. Note that $\text{val}(w) \leq \text{val}(-w)$ for $w \geq 0$. Hence we have $\hat{w} \geq 0$. Moreover, we notice that $\text{val}(w) < \text{val}(1)$ for all $0 \leq w < 1$. Thus we have $\hat{w} \geq 1$. It suffices to consider the case with $w > 1$. Under the assumption on β_s and β_t , we can write

$$\begin{aligned} \text{val}(w) &= \left(1 - \frac{1}{w}\right)^2 \left\| (M^\top Z_1^\top Z_1 M + w^{-2} Z_2^\top Z_2)^{-1} M^\top Z_1^\top Z_1 \Sigma_1^{1/2} \beta_t \right\|^2 \\ &\quad + \frac{\sigma^2}{w^2} \cdot \text{Tr} \left[(M^\top Z_1^\top Z_1 M + w^{-2} Z_2^\top Z_2)^{-1} \right]. \end{aligned}$$

Since

$$\begin{aligned} &\left\| (M^\top Z_1^\top Z_1 M + w^{-2} Z_2^\top Z_2)^{-1} M^\top Z_1^\top Z_1 \Sigma_1^{1/2} \beta_t \right\|^2 \\ &= \text{Tr} \left[(M^\top Z_1^\top Z_1 M + w^{-2} Z_2^\top Z_2)^{-2} M^\top Z_1^\top Z_1 \Sigma_1^{1/2} \beta_t \beta_t^\top \Sigma_1^{1/2} Z_1^\top Z_1 M \right] \end{aligned}$$

is increasing with respect to w , then the derivative of $\text{val}(w)$ can be bounded from below as

$$\begin{aligned} \text{val}'(w) &\geq 2 \frac{w-1}{w^3} \left\| (M^\top Z_1^\top Z_1 M + w^{-2} Z_2^\top Z_2)^{-1} M^\top Z_1^\top Z_1 \Sigma_1^{1/2} \beta_t \right\|^2 \\ &\quad - 2 \frac{\sigma^2}{w^3} \cdot \text{Tr} \left[(M^\top Z_1^\top Z_1 M + w^{-2} Z_2^\top Z_2)^{-1} M^\top Z_1^\top Z_1 M (M^\top Z_1^\top Z_1 M + w^{-2} Z_2^\top Z_2)^{-1} \right] \\ &\geq 2 \frac{w-1}{w^3} \left\| (M^\top Z_1^\top Z_1 M + Z_2^\top Z_2)^{-1} M^\top Z_1^\top Z_1 \Sigma_1^{1/2} \beta_t \right\|^2 - 2 \frac{\sigma^2}{w^3} \cdot \text{Tr} \left[(M^\top Z_1^\top Z_1 M)^{-1} \right]. \end{aligned}$$

Hence $\text{val}'(w) \geq 0$ if $w > 1 + \varepsilon_0$, where

$$\varepsilon_0 := \frac{\sigma^2 \text{Tr} \left[(M^\top Z_1^\top Z_1 M)^{-1} \right]}{\left\| (M^\top Z_1^\top Z_1 M + Z_2^\top Z_2)^{-1} M^\top Z_1^\top Z_1 \Sigma_1^{1/2} \beta_t \right\|^2}.$$

In other words, $\text{val}(w)$ is strictly increasing function on $[1 + \varepsilon_0, \infty]$. Thus we get that \hat{w} satisfies

$$1 \leq w \leq 1 + \varepsilon_0. \quad (\text{A.11})$$

Using (B.3), we get that

$$\varepsilon_0 = O(\sigma^2 / \|\beta_t\|_2^2) = O(p^{-1}).$$

Finally, plugging (A.11) into the expression $te(\hat{\beta}_t^{\text{TL-FT}})$, we obtain (A.10). \square

As a remark, we see that Proposition ?? follows from Theorem A.5. The amount of reduction on test error for the target task is given as

$$te(\hat{\beta}_t) - te(\hat{\beta}_{s,t}) = \sigma^2 \cdot \left(\frac{p}{n_2 - p} - \text{Tr} \left[\left((n_1 + n_2) a_1 \Sigma_2^{-1/2} \Sigma_1 \Sigma_2^{-1/2} + (n_1 + n_2) a_2 \text{Id} \right)^{-1} \right] \right).$$

Because

$$\begin{aligned} te(\hat{\beta}_{s,t}) \leq te(\hat{\beta}_t) &\Leftrightarrow (n_2 - p) \Sigma_2 \preceq (n_1 + n_2) a_1 \Sigma_1 + (n_1 + n_2) a_2 \Sigma_2 \\ &\Leftrightarrow \mathbf{0} \preceq (n_1 + n_2) a_1 \Sigma_1 + (n_1 - (n_1 + n_2) \cdot a_1) \Sigma_2, \end{aligned}$$

which is true since $a_1 \leq n_1 / (n_1 + n_2)$ by equation (4.2). The proof for $te(\hat{\beta}_{s,t}) \leq te(\hat{\beta}_t)$ follows by multiplying $\Sigma_2^{-1/2}$ on both sides of the inequalities above.

B Supplementary Materials for Section 4

B.1 Proof of Theorem 4.1

[**Todo: A proof outline; including the following key lemma.**] To prove Theorem A.5, we study the spectrum of the random matrix model:

$$Q = \Sigma_1^{1/2} Z_1^\top Z_1 \Sigma_1^{1/2} + \Sigma_2^{1/2} Z_2^\top Z_2 \Sigma_2^{1/2},$$

where $\Sigma_{1,2}$ are $p \times p$ deterministic covariance matrices, and $X_1 = (x_{ij})_{1 \leq i \leq n_1, 1 \leq j \leq p}$ and $X_2 = (x_{ij})_{n_1+1 \leq i \leq n_1+n_2, 1 \leq j \leq p}$ are $n_1 \times p$ and $n_2 \times p$ random matrices, respectively, where the entries x_{ij} , $1 \leq i \leq n_1 + n_2 \equiv n$, $1 \leq j \leq p$, are real independent random variables satisfying

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

The proof of Theorem 4.1 involves two parts.

Part I: Bounding the bias from model shift. We relate the first term in equation (4.1) to Δ_β .

Proposition B.1. *In the setting of Theorem 4.1, denote by $K = (\hat{w}^2 X_1^\top X_1 + X_2^\top X_1)^{-1}$, and*

$$\begin{aligned}\delta_1 &= \hat{w}^2 \left\| \Sigma_2^{1/2} K X_1^\top X_1 (\beta_s - \hat{w} \beta_t) \right\|^2, \\ \delta_2 &= n_1^2 \cdot \hat{w}^2 \left\| \Sigma_2^{1/2} K \Sigma_1 (\beta_s - \hat{w} \beta_t) \right\|^2, \\ \delta_3 &= n_1^2 \cdot \hat{w}^2 \left\| \Sigma_1^{1/2} K \Sigma_2 K \Sigma_1^{1/2} \right\| \cdot \left\| \Sigma_1^{1/2} (\beta_s - \hat{w} \beta_t) \right\|^2.\end{aligned}$$

We have that

$$-2n_1^2 \left(2\sqrt{\frac{p}{n_1}} + \frac{p}{n_1} \right) \delta_3 \leq \delta_1 - \delta_2 \leq n_1^2 \left(2\sqrt{\frac{p}{n_1}} + \frac{p}{n_1} \right) \left(2 + 2\sqrt{\frac{p}{n_1}} + \frac{p}{n_1} \right) \delta_3.$$

For the special case when $\Sigma_1 = \text{Id}$ and $\beta_s - \beta_t$ is i.i.d. with mean 0 and variance d^2 , we further have

$$\left(1 - \sqrt{\frac{p}{n_1}} \right)^4 \Delta_\beta \leq \left\| \Sigma_2^{1/2} (X_1^\top X_1 + X_2^\top X_2)^{-1} X_1^\top X_1 (\beta_s - \beta_t) \right\|^2.$$

Proof. The proof follows by applying equation (B.3). Recall that $X_1^\top X_1 = \Sigma_1^{1/2} Z_1^\top Z_1 \Sigma_1^{1/2}$. Denote by $\mathcal{E} = Z_1^\top Z_1 - n_1 \text{Id}$. Let We have

$$\delta_1 = \delta_2 + 2\hat{w}^2 n_1 (\beta_s - \hat{w} \beta_t)^\top \Sigma_1^{1/2} \mathcal{E} \Sigma_1^{1/2} K \Sigma_2 K \Sigma_1 (\beta_s - \hat{w} \beta_t) + \hat{w}^2 \left\| \Sigma_2^{1/2} K \Sigma_1^{1/2} \mathcal{E} \Sigma_1^{1/2} (\beta_s - \hat{w} \beta_t) \right\|^2 \quad (\text{B.2})$$

Here we use the following on the second term in equation (B.2)

$$\begin{aligned}& \left| (\beta_s - \hat{w} \beta_t)^\top \Sigma_1^{1/2} \mathcal{E} \Sigma_1^{1/2} K \Sigma_2 K \Sigma_1 (\beta_s - \hat{w} \beta_t) \right| \\&= \left| \text{Tr} \left[\mathcal{E} \Sigma_1^{1/2} K \Sigma_2 K \Sigma_1 (\beta_s - \hat{w} \beta_t) (\beta_s - \hat{w} \beta_t)^\top \Sigma_1^{1/2} \right] \right| \\&\leq \|\mathcal{E}\| \cdot \left\| \Sigma_1^{1/2} K \Sigma_2 K \Sigma_1 (\beta_s - \hat{w} \beta_t) (\beta_s - \hat{w} \beta_t)^\top \Sigma_1^{1/2} \right\|_* \\&\leq n_1 \left(2\sqrt{\frac{p}{n_1}} + \frac{p}{n_1} \right) \cdot \left\| \Sigma_1^{1/2} K \Sigma_2 K \Sigma_1 (\beta_s - \hat{w} \beta_t) (\beta_s - \hat{w} \beta_t)^\top \Sigma_1^{1/2} \right\|_* \quad (\text{by equation (B.3)}) \\&\leq n_1 \left(2\sqrt{\frac{p}{n_1}} + \frac{p}{n_1} \right) \left\| \Sigma_1^{1/2} K \Sigma_2 K \Sigma_1^{1/2} \right\| \cdot \left\| \Sigma_1^{1/2} (\beta_s - \hat{w} \beta_t) \right\|^2 \\&\quad (\text{since the matrix inside is rank 1})\end{aligned}$$

The third term in equation (B.2) can be bounded with

$$\left\| \Sigma_2^{1/2} K \Sigma_1^{1/2} \mathcal{E} \Sigma_1^{1/2} (\beta_s - \hat{w} \beta_t) \right\|^2 \leq n_1^2 \left(2\sqrt{\frac{p}{n_1}} + \frac{p}{n_1} \right)^2 \left\| \Sigma_1^{1/2} K \Sigma_2 K \Sigma_1^{1/2} \right\| \cdot \left\| \Sigma_1^{1/2} (\beta_s - \hat{w} \beta_t) \right\|^2.$$

Combined together we have shown the right direction for $\delta_1 - \delta_2$. For the left direction, we simply note that the third term in equation (B.2) is positive. And the second term is bigger than $-2n_1^2 (2\sqrt{\frac{p}{n_1}} + \frac{p}{n_1}) \alpha$ using equation (B.3). \square

Part II: The limit of $\left\| \Sigma_2^{1/2} (\hat{w}^2 X_1^\top X_1 + X_2^\top X_2)^{-1} \Sigma_1 (\beta_s - \hat{w} \beta_t) \right\|^2$ using random matrix theory. We consider the same setting as in previous subsection:

$$X_1^\top X_1 := \Sigma_1^{1/2} Z_1^\top Z_1 \Sigma_1^{1/2}, \quad X_2^\top X_2 = \Sigma_2^{1/2} Z_2^\top Z_2 \Sigma_2^{1/2},$$

where z_{ij} , $1 \leq i \leq n_1 + n_2 \equiv n$, $1 \leq j \leq p$, are real independent random variables satisfying (B.1). For now, we assume that the random variables z_{ij} are i.i.d. Gaussian, but we know that universality holds for generally distributed entries. Assume that p/n_1 is a small number such that $Z_1^\top Z_1$ is roughly an isometry, that is, under (B.1), *If we assume the variances of the entries of Z_1 are 1, then we have*

$$-n_1 \left(2\sqrt{\frac{p}{n_1}} - \frac{p}{n_1} \right) \leq Z_1^\top Z_1 - n_1 \text{Id} \leq n_1 \left(2\sqrt{\frac{p}{n_1}} + \frac{p}{n_1} \right). \quad (\text{B.3})$$

Lemma B.2. *In the setting of Theorem 4.1, we have with high probability $1 - o(1)$,*

$$\begin{aligned} & \hat{w}^2 (n_1 + n_2)^2 \left\| \Sigma_2^{1/2} (\hat{w} X_1^\top X_1 + X_2^\top X_2)^{-1} \Sigma_1 (\beta_s - w \beta_t) \right\|^2 \\ &= (\beta_s - w \beta_t)^\top \Sigma_1^{1/2} M \frac{(1 + a_3) \text{Id} + a_4 M^\top M}{(a_2 + a_1 M^\top M)^2} M^\top \Sigma_1^{1/2} (\beta_s - \hat{w} \beta_t) + O(n^{-1/2+\varepsilon}), \end{aligned} \quad (\text{B.4})$$

for any constant $\varepsilon > 0$.

We will give the proof of this lemma in Section C.

Proof of Proposition A.3. the proof for tighter bound □

add some arguments with ε -net.

B.2 Proof of Theorem ???

Extension to Many Tasks of the Same Covariates In this section we consider the setting with k many that have the same covariates. Since every task has the same number of data points as well as the same covariance, the only differences between different tasks are their models $\{\beta_i\}_{i=1}^k$. For this setting, we derive solutions for the multi-task training and the transfer learning setting that match our insights qualitatively from Section ??.

Concretely we will consider the following problem.

$$f(B; W_1, \dots, W_k) = \sum_{i=1}^k \|X B W_i - Y_i\|^2. \quad (\text{B.5})$$

By fixing W_1, W_2, \dots, W_k , we can derive a closed form solution for B as

$$\begin{aligned} \hat{B}(W_1, \dots, W_k) &= (X^\top X)^{-1} X^\top \left(\sum_{i=1}^k Y_i W_i^\top \right) (Z Z^\top)^{-1} \\ &= \sum_{i=1}^k (\beta_i W_i^\top) (Z Z^\top)^{-1} + (X^\top X)^{-1} X^\top \left(\sum_{i=1}^k \varepsilon_i W_i^\top \right) (Z Z^\top)^{-1} \end{aligned}$$

where we denote $Z \in \mathbb{R}^{r \times k}$ as the k vectors W_1, W_2, \dots, W_k stacked together. Similar to Section 2, we consider minimizing the validation loss over W_1, W_2, \dots, W_k provided with \hat{B} .

Jointly minimizing over all tasks. Denote by $\varepsilon(W) = \sum_{i=1}^k \varepsilon_i W_i^\top$. We shall decompose the validation loss $\text{val}(\hat{B}; W_1, \dots, W_k)$ into two parts. The first part is the model shift bias, which is equal to

$$\sum_{j=1}^k \left(\left\| \Sigma^{1/2} \left(\sum_{i=1}^k (\beta_i W_i^\top) (Z Z^\top)^{-1} W_j - \beta_j \right) \right\|^2 \right)$$

The second part is the variance, which is equal to

$$\begin{aligned} & \sum_{j=1}^k \mathbb{E}_{\varepsilon_i, \forall i} \left[\left(\left(\sum_{i=1}^k \varepsilon_i W_i^\top \right) (Z Z^\top)^{-1} W_j \right)^2 \right] \cdot \text{Tr} [\Sigma (X^\top X)^{-1}] \\ &= \sigma^2 \cdot \text{Tr} [\Sigma (X^\top X)^{-1}]. \end{aligned}$$

Therefore we shall focus on the minimizer for the model shift bias since the variance part does not depend the weights. Let us denote $Q \in \mathbb{R}^{k \times k}$ where the (i, j) -th entry is equal to $W_i^\top (Z Z^\top)^{-1} W_j$, for any $1 \leq i, j \leq k$. Let $B^* = [\beta_1, \beta_2, \dots, \beta_k] \in \mathbb{R}^{p \times k}$ denote the true model parameters. We can now write the validation loss succinctly as follows.

$$\text{val}(\hat{B}; W_1, \dots, W_k) = \sum_{j=1}^k \left\| \Sigma^{1/2} (B^* Q - \beta_j) \right\|^2 + \sigma^2 \cdot \text{Tr} [\Sigma (X^\top X)^{-1}]$$

From the above we can solve for Q optimally as [Todo: this].

Minimizing over the target task alone. If we only minimize over the the validation loss for the target task, we shall get the following.

$$\text{val}_j(\hat{B}W_j) = \left\| \Sigma^{1/2} \left(\sum_{i=1}^k W_i^\top (ZZ^\top)^{-1} W_j \beta_i - \beta_j \right) \right\|^2 + \sigma^2 \cdot W_j^\top (ZZ^\top)^{-1} W_j \cdot \text{Tr} [\Sigma (X^\top X)^{-1}].$$

From the above we can obtain three conceptual insights that are consistent with Section ?? and 3.

- The de-noising effect of multi-task learning.
- Multi-task training vs single-task training can be either positive or negative.
- Transfer learning is better than the other two. And the improvement over multi-task training increases as the model distances become larger.

C Proof of Lemma A.2 and Lemma B.2

We consider two $p \times p$ random sample covariance matrices $Q_1 := \Sigma_1^{1/2} Z_1^\top Z_1 \Sigma_1^{1/2}$ and $Q_2 := \Sigma_2^{1/2} Z_2^\top Z_2 \Sigma_2^{1/2}$, where Σ_1 and Σ_2 are $p \times p$ deterministic non-negative definite (real) symmetric matrices. 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 matrix with (real) i.i.d. entries satisfying

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

where we denote $n := n_1 + n_2$. Here we have chosen the scaling that is more standard in the random matrix theory literature—under this $n^{-1/2}$ scaling, the eigenvalues of Q_1 and Q_2 are all of order 1. Moreover, we assume that the fourth moment exists:

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

for some constant $C > 0$. We assume that the aspect ratios $d_1 := p/n_1$ and $d_2 := p/n_2$ satisfy that

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

for some small constant $0 < \tau < 1$. Here the lower bound $1 + \tau \leq d_2$ is to ensure that the covariance matrix Q_2 for the target task is non-singular with high probability; see Lemma D.2 below.

We assume that Σ_1 and Σ_2 have eigendecompositions

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

where the eigenvalues satisfy that

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

for some small constant $0 < \tau < 1$. We assume that $M := \Sigma_1^{1/2} \Sigma_2^{-1/2}$ has singular value decomposition

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

where the singular values satisfy that

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

for some small constant $0 < \tau < 1$.

We summarize our basic assumptions here for future reference.

Assumption C.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 (C.1) and (C.2), Σ_1 and Σ_2 are deterministic non-negative definite symmetric matrices satisfying (C.4)-(C.7), and $d_{1,2}$ satisfy (C.3).

Before giving the main proof, we first introduce some notations and tools.

C.1 Notations

We will use the following notion of stochastic domination, which was first introduced in [8] and subsequently used in many works on random matrix theory, such as [4, 5, 6, 9, 10, 16]. 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 C.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 $\xi \prec \zeta$. Throughout this paper, the stochastic domination will always be uniform in all parameters that are not explicitly fixed (such as matrix indices, and z that takes values in some compact set). 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 .

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

Lemma C.3 (Lemma 3.2 in [4]). Let ξ and ζ be families of nonnegative random variables.

(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 , then $\sum_{v \in V} \xi(u, v) \prec \sum_{v \in V} \zeta(u, v)$ uniformly in u .

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

(iii) Suppose that $\Psi(u) \geq n^{-C}$ is deterministic and $\xi(u)$ satisfies $\mathbb{E}\xi(u)^2 \leq n^C$ for all u . Then if $\xi(u) \prec \Psi(u)$ uniformly in u , we have $\mathbb{E}\xi(u) \prec \Psi(u)$ uniformly in u .

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

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

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

Our main goal is to study the following matrix inverse

$$(\mathcal{Q}_1 + \mathcal{Q}_2)^{-1} = \left(\Sigma_1^{1/2} Z_1^\top Z_1 \Sigma_1^{1/2} + \Sigma_2^{1/2} Z_2^\top Z_2 \Sigma_2^{1/2} \right)^{-1}.$$

Using (C.6), we can rewrite it as

$$\Sigma_2^{-1/2} V \left(\Lambda U^\top Z_1^\top Z_1 U \Lambda + V^\top Z_2^\top Z_2 V \right)^{-1} V^\top \Sigma_2^{-1/2}. \quad (\text{C.9})$$

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

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

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

Now we introduce a convenient self-adjoint linearization trick. This idea dates back at least to Girko, see e.g., the works [13, 14, 15] and references therein. It has been proved to be useful in studying the local laws of random matrices of the Gram type [1, 2, 16, 24]. We define the following $(p+n) \times (p+n)$ self-adjoint block matrix, which is a linear function of X :

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

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

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

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 $i, j \in \mathcal{I}_1$ and greek letters $\mu, \nu \in \mathcal{I}_2 \cup \mathcal{I}_3$. Moreover, we shall use the notations $\mathbf{a}, \mathbf{b} \in \mathcal{I} := \cup_{i=1}^3 \mathcal{I}_i$. We label the indices of the matrices according to

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

Then we denote the $\mathcal{I}_1 \times \mathcal{I}_1$ block of $G(z)$ by $\mathcal{G}_L(z)$, the $\mathcal{I}_1 \times (\mathcal{I}_2 \cup \mathcal{I}_3)$ by \mathcal{G}_{LR} , the $(\mathcal{I}_2 \cup \mathcal{I}_3) \times \mathcal{I}_1$ 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 $Y_1 := Z_1 U \Lambda$, $Y_2 := Z_2 V$ and $W := (Y_1^\top, Y_2^\top)$. By Schur complement formula, one can find that (recall (C.10))

$$\mathcal{G}_{11} = (WW^\top - z)^{-1} = \mathcal{G}, \quad \mathcal{G}_{LR} = \mathcal{G}_{RL}^\top = \mathcal{G}W, \quad \mathcal{G}_R := \begin{pmatrix} \mathcal{G}_{22} & \mathcal{G}_{23} \\ \mathcal{G}_{32} & \mathcal{G}_{33} \end{pmatrix} = z(W^\top W - z)^{-1}. \quad (\text{C.13})$$

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

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

Next we introduce 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

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0 = \lambda_{p+1} = \dots = \lambda_n,$$

$\{\xi_k\}_{k=1}^p$ are the left-singular vectors, and $\{\zeta_k\}_{k=1}^n$ are the right-singular vectors. Then using (C.13), we can get that for $i, j \in \mathcal{I}_1$ and $\mu, \nu \in \mathcal{I}_2$,

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

We now define the deterministic limit of $\mathcal{G}(z)$. We first define the deterministic limits of $(m_2(z), m_3(z))$, that is $(m_{2c}(z), m_{3c}(z))$, as the (unique) solution to the following system of self-consistent equations

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

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

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

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

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

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 (C.16) is reduced to

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

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

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

is a strictly increasing function on $[0, r_1^{-1}(1 - \gamma_n)]$, and $f(0) = 0 < 1$, $f(r_1^{-1}(1 - \gamma_n)) = 1 + \gamma_n > 1$. Hence there exists a unique solution (b_2, b_3) to (C.19). 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

$$r_1 \tau \leq r_1 b_2 < \min\{(1 - \gamma_n) - r_1 \tau, r_1(1 - \tau)\}, \quad \tau < r_3 b_3 \leq 1 - \gamma_n - r_1 \tau. \quad (\text{C.20})$$

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. (say something about the previous work)

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

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

Moreover, the solution satisfies

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

The proof is a standard application of the contraction principle. For reader's convenience, we will gives its proof in Appendix D.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 C.7 below.

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

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

Suppose they satisfy the system of equations

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

for some (random) errors satisfying

$$\max_{\alpha=2}^3 |\mathcal{E}_\alpha| \leq \delta(z),$$

where $\delta(z)$ is any deterministic z -dependent function $\delta(z) \leq (\log n)^{-1}$. Then we have

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

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

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

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

Theorem C.7. *Suppose Assumption C.1 holds, and Z_1, Z_2 satisfy the bounded support condition (C.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 there exists a sufficiently small constant $c_0 > 0$ such that the following **anisotropic 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*

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

The proof of this theorem will be given in Section D. Using a simple cutoff argument, it is easy to obtain the following corollary under certain moment assumptions.

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

$$\max_{i,j} \mathbb{E} |\sqrt{n} z_{ij}^{(\alpha)}|^a = O(1), \quad \alpha = 1, 2, \quad (\text{C.27})$$

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

Proof of Corollary C.8. Fix any sufficiently small constant $\varepsilon > 0$. We then 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

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

By the moment conditions (C.27) and a simple union bound, we have

$$\mathbb{P}(\tilde{Z}_1 = Z_1, \tilde{Z}_2 = Z_2) = 1 - O(n^{-a\varepsilon}). \quad (\text{C.28})$$

Using (C.27) and integration by parts, it is easy to verify that

$$\mathbb{E} |z_{ij}^{(\alpha)}| \mathbf{1}_{|z_{ij}^{(\alpha)}| > q} = O(n^{-2-\varepsilon}), \quad \mathbb{E} |z_{ij}^{(\alpha)}|^2 \mathbf{1}_{|z_{ij}^{(\alpha)}| > q} = O(n^{-2-\varepsilon}), \quad \alpha = 1, 2,$$

which imply that

$$|\mathbb{E} \tilde{z}_{ij}^{(\alpha)}| = O(n^{-2-\varepsilon}), \quad \mathbb{E} |\tilde{z}_{ij}^{(\alpha)}|^2 = n^{-1} + O(n^{-2-\varepsilon}), \quad \alpha = 1, 2. \quad (\text{C.29})$$

Moreover, we trivially have

$$\mathbb{E} |\tilde{z}_{ij}^{(\alpha)}|^4 \leq \mathbb{E} |z_{ij}^{(\alpha)}|^4 = O(n^{-2}), \quad \alpha = 1, 2.$$

Then we centralize and rescale \tilde{Z}_1 and \tilde{Z}_2 as

$$\hat{Z}_\alpha := \frac{\tilde{Z}_\alpha - \mathbb{E} \tilde{Z}_\alpha}{(\mathbb{E} |\tilde{z}_{11}^{(\alpha)}|^2)^{1/2}}, \quad \alpha = 1, 2.$$

Now \hat{Z}_1 and \hat{Z}_2 satisfy the assumptions in Theorem C.7 with $q = n^{-c_a+\varepsilon}$, and (C.26) gives that

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

Then using (C.29) and (D.4) below, we can easily get that

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

where we also used the bound $\|\mathbb{E} \tilde{Z}_\alpha\| = O(n^{-1-\varepsilon})$. This shows that (C.26) also holds for $G(\tilde{Z}_1, \tilde{Z}_2, z)$ with $q = n^{-c_a+\varepsilon}$, and hence concludes the proof by (C.28). \square

Using Corollary C.8, we can complete the proof of Lemma A.2 and Lemma B.2.

Proof of Lemma A.2. In the setting of Lemma A.2, we can write

$$\mathcal{R} := (w^2 X_1^\top X_1 + X_2^\top X_2)^{-1} = n^{-1} \left(\tilde{\Sigma}_1^{1/2} Z_1^\top Z_1 \tilde{\Sigma}_1^{1/2} + \Sigma_2^{1/2} Z_2^\top Z_2 \Sigma_2^{1/2} \right)^{-1},$$

where $\tilde{\Sigma}_1 := w^2 \Sigma_1$, Σ_2 , Z_1 and Z_2 satisfy Assumption C.1. Here the extra n^{-1} is due to the choice of the variances—in the setting of Lemma A.2 the variances of the entries of $Z_{1,2}$ are equal to 1, while in (C.1) they are taken to be n^{-1} . As in (C.6), we assume that $M := \tilde{\Sigma}_1^{1/2} \Sigma_2^{-1/2}$ has singular value decomposition

$$M = U \Lambda V^\top, \quad \Lambda = \text{diag}(\sigma, \dots, \sigma_p). \quad (\text{C.30})$$

Then as in (C.9), we can write

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

Now by Corollary C.8, we obtain that for any small constant $\varepsilon > 0$, with probability $1 - o(1)$,

$$\max_{1 \leq i \leq p} |(\Sigma_2 \mathcal{R} - \Sigma_2^{1/2} V \Pi(0) V^\top \Sigma_2^{-1/2})_{ii}| \leq n^\varepsilon q, \quad q = n^{2/a-1/2}, \quad (\text{C.31})$$

where by (C.18),

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

with (b_2, b_3) satisfying (C.19). Thus from (C.31) we get that

$$n^{-1} \text{Tr}(\Sigma_2 \mathcal{R}) = n^{-1} \text{Tr}(r_1 b_2 M^\top M + r_2 b_3)^{-1} + O(n^\varepsilon q)$$

with probability $1 - o(1)$. This concludes (??) if we rename $r_1 b_2 \rightarrow a_1$ and $r_2 b_3 \rightarrow a_2$. For (??), it is a well-known result for inverse Wishart matrices (add some references). In fact, if we set $n_1 = 0$ and $n_2 = n$, then it is easy to check that $a_1 = 0$ and $a_2 = (n_2 - p)/n_2$ is the solution to (4.2). This gives (??) by (??). \square

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

$$\begin{aligned} \Delta &:= n^2 \left\| \Sigma_2^{1/2} (\hat{w}^2 X_1^\top X_1 + X_2^\top X_2)^{-1} \Sigma_1 (\beta_s - w \beta_t) \right\|^2 \\ &= (\beta_s - w \beta_t) \Sigma_1^{1/2} M (M^\top Z_1^\top Z_1 M + Z_2^\top Z_2)^{-2} M^\top \Sigma_1^{1/2} (\beta_s - w \beta_t), \end{aligned}$$

where $\tilde{\Sigma}_1 := w^2 \Sigma_1$, Σ_2 , Z_1 and Z_2 satisfy Assumption C.1 and $M := \tilde{\Sigma}_1^{1/2} \Sigma_2^{-1/2}$. Here again the n^2 factor disappears due to the choice of scaling. Again we assume that M has the singular value decomposition (C.30). Then we can write

$$\Delta := \mathbf{v}^\top (\mathcal{G}^2)(0) \mathbf{v}, \quad \mathbf{v} := V^\top M^\top \Sigma_1^{1/2} (\beta_s - w \beta_t).$$

Note that $\mathcal{G}^2(0) = \partial_z \mathcal{G}|_{z=0}$. Now using Cauchy's integral formula and Corollary C.8, we get that with probability $1 - o(1)$,

$$\mathbf{v}^\top \mathcal{G}^2(0) \mathbf{v} = \frac{1}{2\pi i} \oint_{\mathcal{C}} \frac{\mathbf{v}^\top \mathcal{G}(z) \mathbf{v}}{z^2} dz = \frac{1}{2\pi i} \oint_{\mathcal{C}} \frac{\mathbf{v}^\top \Pi(z) \mathbf{v}}{z^2} dz + O_\prec(q) = \mathbf{v}^\top \Pi'(0) \mathbf{v} + O_\prec(q), \quad (\text{C.32})$$

where \mathcal{C} is the contour $\{z \in \mathbb{C} : |z| \leq (\log n)^{-1}\}$ and we used (C.26) in the second step. Hence it remains to study the derivatives

$$\mathbf{v}^\top \Pi'(0) \mathbf{v} = \mathbf{v} \frac{1 + r_1 m'_{2c}(0) \Lambda^2 + r_2 m'_{3c}(0)}{(r_1 m_{2c}(0) \Lambda^2 + r_2 m_{3c}(0))^2} \mathbf{v}. \quad (\text{C.33})$$

It remains to calculate the derivatives $m'_{2c}(0)$ and $m'_{3c}(0)$.

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

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

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 equation becomes

$$\begin{aligned} \left(\frac{r_2}{a_2^2} - \frac{1}{n} \sum_{i=1}^p \frac{1}{(\sigma_i^2 a_1 + a_2)^2} \right) a_3 - \left(\frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^2}{(\sigma_i^2 a_1 + a_2)^2} \right) a_4 &= \frac{1}{n} \sum_{i=1}^p \frac{1}{(\sigma_i^2 a_1 + a_2)^2}, \\ \left(\frac{r_1}{a_1^2} - \frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^4}{(\sigma_i^2 a_1 + a_2)^2} \right) a_4 - \left(\frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^2}{(\sigma_i^2 a_1 + a_2)^2} \right) a_3 &= \frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^2}{(\sigma_i^2 a_1 + a_2)^2}. \end{aligned} \quad (\text{C.34})$$

Then by (C.32) and (C.33), we get

$$\begin{aligned} \Delta &= (\beta_s - w\beta_t)^\top \Sigma_1^{1/2} M V \frac{1 + a_3 + a_4 \Lambda^2}{(a_1 \Lambda^2 + a_2)} V^\top M^\top \Sigma_1^{1/2} (\beta_s - w\beta_t) \\ &= (\beta_s - w\beta_t)^\top \Sigma_1^{1/2} M \frac{1 + a_3 + a_4 M^\top M}{(a_1 M^\top M + a_2)} M^\top \Sigma_1^{1/2} (\beta_s - w\beta_t) \end{aligned}$$

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

D Proof of Theorem C.7

The main difficulty for the proof of Theorem C.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. [4, 22, 25]. To go from the Gaussian case to the general X case, we will adopt a continuous self-consistent comparison argument developed in [16].

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

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

and prove the following result.

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

In Section D.1, we will collect some a priori estimates and resolvent identities that will be used in the proof of Theorem C.7 and Proposition D.1. Then in Section D.2 we give the proof of Proposition D.1, which, as discussed above, concludes Theorem C.7 for i.i.d. Gaussian Z_1 and Z_2 . Finally, in Section D.3, we will describe how to extend the result in Theorem C.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 $z = E + i\eta$.

D.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 $1 - o(1)$ [3]. In our proof, we shall need a slightly stronger probability bound, which is given by the following lemma. Denote the nonzero eigenvalues of $Z_1^\top Z_1$ and $Z_2^\top Z_2$ by $\lambda_1(Z_1^\top Z_1) \geq \dots \geq \lambda_{p \wedge n_1}(Z_1^\top Z_1)$ and $\lambda_1(Z_2^\top Z_2) \geq \dots \geq \lambda_p(Z_2^\top Z_2)$.

Lemma D.2. *Suppose Assumption C.1 holds, and Z_1, Z_2 satisfy the bounded support condition (C.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,*

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

and

$$(1 - \sqrt{d_2})^2 - \varepsilon \leq \lambda_p(Z_2^\top Z_2) \leq \lambda_1(Z_2^\top Z_2) \leq (1 + \sqrt{d_2})^2 + \varepsilon. \quad (\text{D.3})$$

Proof. This lemma essentially follows from [4, Theorem 2.10], although the authors considered the case with $q \prec n^{-1/2}$ only. The results for larger q follows from [7, Lemma 3.12], but only the bounds for the largest eigenvalues are given there in order to avoid the issue with the smallest eigenvalue when d_2 is close to 1. However, under the assumption (C.3), the lower bound for the smallest eigenvalue follows from the exactly the same arguments as in [7]. Hence we omit the details. \square

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

Lemma D.3. *Suppose the assumptions of Lemma D.2 holds. Then there exists a constant $C > 0$ such that the following estimates hold uniformly in $z, z' \in \mathbf{D}$ with high probability:*

$$\|G(z)\| \leq C, \quad (\text{D.4})$$

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

$$|\mathbf{u}^\top [G(z) - G(z')] \mathbf{v}| \leq C|z - z'|. \quad (\text{D.5})$$

Proof. As in (C.15), we let $\{\lambda_k\}_{1 \leq k \leq p}$ be the eigenvalues of WW^\top . By Lemma D.2 and the assumption (C.3), we obtain that

$$\lambda_p \geq \lambda_p(Z_2^\top Z_2) \geq \varepsilon > 0 \quad (\text{D.6})$$

for some constant $\varepsilon > 0$. In particular, it implies that

$$\inf_{z \in \mathbf{D}} \min_{1 \leq k \leq p} |\lambda_k - z| \gtrsim 1.$$

Together with (C.15), it implies the estimates (D.4) and (D.5). \square

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

Definition D.4 (Minors). *For any $(p+n) \times (p+n)$ matrix \mathcal{A} and $\mathbb{T} \subseteq \mathcal{I}$, we define the minor $\mathcal{A}^{(\mathbb{T})} := (\mathcal{A}_{ab} : a, b \in \mathcal{I} \setminus \mathbb{T})$ as the $(p+n-|\mathbb{T}|) \times (p+n-|\mathbb{T}|)$ matrix obtained by removing all rows and columns indexed by \mathbb{T} . Note that we keep the names of indices when defining $\mathcal{A}^{(\mathbb{T})}$, i.e. $(\mathcal{A}^{(\mathbb{T})})_{ab} = \mathcal{A}_{ab}$ for $a, b \notin \mathbb{T}$. Correspondingly, we define the resolvent minor as (recall (C.13))*

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

and the partial traces (recall (C.14))

$$\begin{aligned} m^{(\mathbb{T})} &:= \frac{1}{p} \sum_{i \in \mathcal{I}_1} G_{ii}^{(\mathbb{T})}(z), & m_1^{(\mathbb{T})} &:= \frac{1}{p} \sum_{i \in \mathcal{I}_1}^{(\mathbb{T})} \sigma_i^2 G_{ii}^{(\mathbb{T})}(z), \\ m_2^{(\mathbb{T})}(z) &:= \frac{1}{n_1} \sum_{\mu \in \mathcal{I}_2} G_{\mu\mu}^{(\mathbb{T})}(z), & m_3^{(\mathbb{T})}(z) &:= \frac{1}{n_2} \sum_{\mu \in \mathcal{I}_3} G_{\mu\mu}^{(\mathbb{T})}(z), \end{aligned} \quad (\text{D.7})$$

where we abbreviated that $\sum_a^{(\mathbb{T})} := \sum_{a \notin \mathbb{T}}$. For convenience, 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 $b \in \mathbb{T}$. Moreover, we will abbreviate $(\{a\}) \equiv (a)$ and $(\{a, b\}) \equiv (ab)$.

Then we record the following resolvent identities.

Lemma D.5. (*Resolvent identities*).

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

$$\frac{1}{G_{ii}} = -z - \left(W G^{(i)} W^\top \right)_{ii}, \quad \frac{1}{G_{\mu\mu}} = -1 - \left(W^\top G^{(\mu)} W \right)_{\mu\mu}. \quad (\text{D.8})$$

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

$$G_{ij} = -G_{ii} \left(W G^{(i)} \right)_{ij}, \quad G_{\mu\nu} = -G_{\mu\mu} \left(W^\top G^{(\mu)} \right)_{\mu\nu}. \quad (\text{D.9})$$

For $i \in \mathcal{I}_1$ and $\mu \in \mathcal{I}_2$, we have

$$G_{i\mu} = -G_{ii} \left(W G^{(i)} \right)_{i\mu}, \quad G_{\mu i} = -G_{\mu\mu} \left(W^\top G^{(\mu)} \right)_{\mu i}. \quad (\text{D.10})$$

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

$$G_{bc}^{(a)} = G_{bc} - \frac{G_{ba} G_{ac}}{G_{aa}}, \quad \frac{1}{G_{bb}} = \frac{1}{G_{bb}^{(a)}} - \frac{G_{ba} G_{ab}}{G_{bb}^{(a)} G_{aa}^{(a)}}. \quad (\text{D.11})$$

Proof. All these identities can be proved directly using Schur's complement formula. The reader can also refer to, for example, [16, Lemma 4.4]. \square

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

Lemma D.6 (Lemma 3.8 of [11]). *Let $(x_i), (y_j)$ be independent families of centered and independent random variables, and $(A_i), (B_{ij})$ be families of deterministic complex numbers. Suppose the entries x_i, y_j have variance at most n^{-1} and satisfy the bounded support condition (C.8) with $q \leq n^{-\phi}$ for some constant $\phi > 0$. Then we have the following bound:*

$$\left| \sum_i A_i x_i \right| \prec q \max_i |A_i| + \frac{1}{\sqrt{n}} \left(\sum_i |A_i|^2 \right)^{1/2}, \quad \left| \sum_{i,j} x_i B_{ij} y_j \right| \prec q^2 B_d + q B_o + \frac{1}{n} \left(\sum_{i \neq j} |B_{ij}|^2 \right)^{1/2}, \quad (\text{D.12})$$

$$\left| \sum_i \bar{x}_i B_{ii} x_i - \sum_i (\mathbb{E} |x_i|^2) B_{ii} \right| \prec q B_d, \quad \left| \sum_{i \neq j} \bar{x}_i B_{ij} x_j \right| \prec q B_o + \frac{1}{n} \left(\sum_{i \neq j} |B_{ij}|^2 \right)^{1/2}, \quad (\text{D.13})$$

where $B_d := \max_i |B_{ii}|$ and $B_o := \max_{i \neq j} |B_{ij}|$.

D.2 Entrywise local law

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

$$W = (\Lambda Z_1^\top, Z_2^\top). \quad (\text{D.14})$$

Lemma D.7. *Suppose the assumptions in Proposition D.1 hold. Then the following estimate holds uniformly for $z \in \mathbf{D}$:*

$$\max_{\mathbf{a}, \mathbf{b}} |(G_0)_{\mathbf{a}\mathbf{b}}(z) - \Pi_{\mathbf{a}\mathbf{b}}(z)| \prec n^{-1/2}. \quad (\text{D.15})$$

Proof. The proof of Lemma D.7 is divided into three steps. For simplicity, we will still denote $G \equiv G_0$ in the following proof, while keeping in mind that W takes the form in (D.14).

Step 1: Large deviations estimates. In this step, we prove some (almost) optimal large deviations estimates on the off-diagonal entries of G , and on the following Z variables. In analogy to [11, Section 3] and [16, Section 5], we introduce the Z variables

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

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 column of H . By (D.8), we have

$$\begin{aligned} Z_i = (\mathbb{E}_i - 1) \left(W G^{(i)} W^\top \right)_{ii} &= \sigma_i^2 \sum_{\mu, \nu \in \mathcal{I}_2} G_{\mu\nu}^{(i)} \left(\frac{1}{n} \delta_{\mu\nu} - z_{\mu i} z_{\nu i} \right) \\ &+ \sum_{\mu, \nu \in \mathcal{I}_3} G_{\mu\nu}^{(i)} \left(\frac{1}{n} \delta_{\mu\nu} - z_{\mu i} z_{\nu i} \right), \quad i \in \mathcal{I}_1, \end{aligned} \quad (\text{D.16})$$

and

$$\begin{aligned} Z_\mu = (\mathbb{E}_\mu - 1) \left(W^\top G^{(\mu)} W \right)_{\mu\mu} &= \sum_{i, j \in \mathcal{I}_1} \sigma_i \sigma_j G_{ij}^{(\mu)} \left(\frac{1}{n} \delta_{ij} - z_{\mu i} z_{\mu j} \right), \quad \mu \in \mathcal{I}_2, \\ Z_\mu = (\mathbb{E}_\mu - 1) \left(W^\top G^{(\mu)} W \right)_{\mu\mu} &= \sum_{i, j \in \mathcal{I}_1} G_{ij}^{(\mu)} \left(\frac{1}{n} \delta_{ij} - z_{\mu i} z_{\mu j} \right), \quad \mu \in \mathcal{I}_3. \end{aligned} \quad (\text{D.17})$$

For simplicity, we introduce the following random error

$$\Lambda_o := \max_{\mathbf{a} \neq \mathbf{b}} |G_{\mathbf{a}\mathbf{a}}^{-1} G_{\mathbf{a}\mathbf{b}}|. \quad (\text{D.18})$$

The following lemma gives the desired large deviations estimates on the Λ_o and the Z variables.

Lemma D.8. *Suppose the assumptions in Proposition D.1 hold. Then the following estimates hold uniformly for all $z \in \mathbf{D}$:*

$$\Lambda_o + \max_{\mathbf{a} \in \mathcal{I}} |Z_{\mathbf{a}}| \prec n^{-1/2}. \quad (\text{D.19})$$

Proof. Note that for any $\mathbf{a} \in \mathcal{I}$, $H^{(\mathbf{a})}$ and $G^{(\mathbf{a})}$ also satisfies the assumptions for Lemma D.3. Hence (D.4) and (D.5) also hold for $G^{(\mathbf{a})}$. Now applying Lemma D.6 to (D.16) and (D.17), and using the a priori bound (D.4), we get that for any $i \in \mathcal{I}_1$,

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

where in the last step we used that for any μ ,

$$\sum_{\nu \in \mathcal{I}_2 \cup \mathcal{I}_3} |G_{\mu\nu}^{(i)}|^2 \leq \sum_{\mathbf{a} \in \mathcal{I}} |G_{\mu\mathbf{a}}^{(i)}|^2 = [G^{(i)} (G^{(i)})^*]_{\mu\mu} = O(1) \quad (\text{D.20})$$

by (D.4). Similarly, applying Lemma D.6 to Z_μ in (D.17) and using (D.4), we obtain the same bound.

Then we prove the off-diagonal estimates. For $i \neq j \in \mathcal{I}_1$ and $\mu \neq \nu \in \mathcal{I}_2 \cup \mathcal{I}_3$, using (D.9), Lemma D.6 and (D.4), we obtain that

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

and

$$|G_{\mu\mu}^{-1} G_{\mu\nu}| \prec n^{-1/2} + \frac{1}{\sqrt{n}} \left(\sum_{i \in \mathcal{I}_1} |G_{i\nu}^{(\mu)}|^2 \right)^{1/2} \prec n^{-1/2}.$$

For $i \in \mathcal{I}_1 \cup \mathcal{I}_2$ and $\mu \in \mathcal{I}_3$, using (D.10), Lemma D.6 and (D.4), we obtain that

$$|G_{ii}^{-1} G_{i\mu}| + |G_{\mu\mu}^{-1} G_{\mu i}| \prec n^{-1/2} + \frac{1}{\sqrt{n}} \left(\sum_{\nu \in \mathcal{I}_2 \cup \mathcal{I}_3} |G_{\nu\mu}^{(i)}|^2 \right)^{1/2} + \frac{1}{\sqrt{n}} \left(\sum_{j \in \mathcal{I}_1} |G_{ji}^{(\mu)}|^2 \right)^{1/2} \prec n^{-1/2}.$$

Thus we obtain that $\Lambda_o \prec n^{-1/2}$, which concludes (D.19). \square

Note that combining (D.4) and (D.19), we immediately conclude (D.15) for $\mathbf{a} \neq \mathbf{b}$.

Step 2: Self-consistent equations. This is the key step of the proof for Proposition D.7, which derives approximate self-consistent equations satisfied by $m_2(z)$ and $m_3(z)$. More precisely, we will show that $(m_2(z), m_3(z))$ satisfies (C.23) up to some small error $|\mathcal{E}_{2,3}| \prec n^{-1/2}$. Then applying Lemma C.6 shows that $(m_2(z), m_3(z))$ is close to $(m_{2c}(z), m_{3c}(z))$ —this will be discussed in Step 3.

We define the following z -dependent event

$$\Xi(z) := \left\{ |m_2(z) - m_{2c}(z)| + |m_3(z) - m_{3c}(z)| \leq (\log n)^{-1/2} \right\}. \quad (\text{D.21})$$

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

$$|m_{2c}| \sim 1, \quad |m_{3c}| \sim 1, \quad |z + \sigma_i^2 r_1 m_{2c} + r_2 m_{3c}| \sim 1, \quad |1 + \gamma_n m_c| \sim 1, \quad |1 + \gamma_n m_{1c}| \sim 1, \quad (\text{D.22})$$

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

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

Plugging (D.22) into (C.18), we get

$$|\Pi_{\mathbf{a}\mathbf{a}}(z)| \sim 1 \quad \text{uniformly in } z \in \mathbf{D}, \mathbf{a} \in \mathcal{I}. \quad (\text{D.23})$$

Then we claim the following result.

Lemma D.9. *Suppose the assumptions in Proposition D.1 hold. Then the following estimates hold uniformly in $z \in \mathbf{D}$:*

$$\begin{aligned} \mathbf{1}(\Xi) \left| \frac{1}{m_2} + 1 - \frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^2}{z + \sigma_i^2 r_1 m_2 + r_2 m_3} \right| &\prec n^{-1/2}, \\ \mathbf{1}(\Xi) \left| \frac{1}{m_3} + 1 - \frac{1}{n} \sum_{i=1}^p \frac{1}{z + \sigma_i^2 r_2 m_2 + r_2 m_3} \right| &\prec n^{-1/2}. \end{aligned} \quad (\text{D.24})$$

Proof. By (D.8), (D.16) and (D.17), we obtain that

$$\frac{1}{G_{ii}} = -z - \frac{\sigma_i^2}{n} \sum_{\mu \in \mathcal{I}_2} G_{\mu\mu}^{(i)} - \frac{1}{n} \sum_{\mu \in \mathcal{I}_3} G_{\mu\mu}^{(i)} + Z_i = -z - \sigma_i^2 r_1 m_2 - r_2 m_3 + \varepsilon_i, \quad i \in \mathcal{I}_1, \quad (\text{D.25})$$

$$\frac{1}{G_{\mu\mu}} = -1 - \frac{1}{n} \sum_{i \in \mathcal{I}_1} \sigma_i^2 G_{ii}^{(\mu)} + Z_\mu = -1 - \gamma_n m_1 + \varepsilon_\mu, \quad \mu \in \mathcal{I}_2, \quad (\text{D.26})$$

$$\frac{1}{G_{\nu\nu}} = -1 - \frac{1}{n} \sum_{i \in \mathcal{I}_1} G_{ii}^{(\nu)} + Z_\nu = -1 - \gamma_n m + \varepsilon_\nu, \quad \nu \in \mathcal{I}_3, \quad (\text{D.27})$$

where we recall (D.7), and

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

By (D.11) we can bound that

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

where we used (D.19) in the second step. Similarly, we can get that

$$|m - m^{(\mu)}| + |m_1 - m_1^{(\mu)}| + |m_2 - m_2^{(i)}| + |m_3 - m_3^{(i)}| \prec n^{-1} \quad (\text{D.28})$$

for any $i \in \mathcal{I}_1$ and $\mu \in \mathcal{I}_2 \cup \mathcal{I}_3$. Together with (D.19), we obtain that for all i and μ ,

$$|\varepsilon_i| + |\varepsilon_\mu| \prec n^{-1/2}. \quad (\text{D.29})$$

With (D.22) 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 (D.25), (D.29) and (D.19), we obtain that

$$\mathbf{1}(\Xi)G_{ii} = \mathbf{1}(\Xi) \left[-\frac{1}{z + \sigma_i^2 r_1 m_2 + r_2 m_3} + O_{\prec} \left(n^{-1/2} \right) \right]. \quad (\text{D.30})$$

Plugging it into the definitions of m and m_1 in (D.7), we get

$$\mathbf{1}(\Xi)m = \mathbf{1}(\Xi) \left[-\frac{1}{p} \sum_{i \in \mathcal{I}_1} \frac{1}{z + \sigma_i^2 r_1 m_2 + r_2 m_3} + O_{\prec} \left(n^{-1/2} \right) \right], \quad (\text{D.31})$$

$$\mathbf{1}(\Xi)m_1 = \mathbf{1}(\Xi) \left[-\frac{1}{p} \sum_{i \in \mathcal{I}_1} \frac{\sigma_i^2}{z + \sigma_i^2 r_1 m_2 + r_2 m_3} + O_{\prec} \left(n^{-1/2} \right) \right]. \quad (\text{D.32})$$

As a byproduct, we obtain from the two estimates that

$$\mathbf{1}(\Xi) (|m - m_c| + |m_1 - m_{1c}|) \lesssim (\log n)^{-1/2}, \quad \text{with high probability.} \quad (\text{D.33})$$

Together with (D.22), we get that

$$|1 + \gamma_n m_1| \sim 1, \quad |1 + \gamma_n m| \sim 1, \quad \text{with high probability on } \Xi. \quad (\text{D.34})$$

Now using (D.26), (D.27), (D.29), (D.19) and (D.34), we can obtain that with high probability,

$$\mathbf{1}(\Xi)G_{\mu\mu} = \mathbf{1}(\Xi) \left[-\frac{1}{1 + \gamma_n m_1} + O_{\prec} \left(n^{-1/2} \right) \right], \quad \mu \in \mathcal{I}_2, \quad (\text{D.35})$$

$$\mathbf{1}(\Xi)G_{\nu\nu} = \mathbf{1}(\Xi) \left[-\frac{1}{1 + \gamma_n m} + O_{\prec} \left(n^{-1/2} \right) \right], \quad \nu \in \mathcal{I}_3. \quad (\text{D.36})$$

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

$$\mathbf{1}(\Xi)m_2 = \mathbf{1}(\Xi) \left[-\frac{1}{1 + \gamma_n m_1} + O_{\prec} \left(n^{-1/2} \right) \right], \quad \mathbf{1}(\Xi)m_3 = \mathbf{1}(\Xi) \left[-\frac{1}{1 + \gamma_n m} + O_{\prec} \left(n^{-1/2} \right) \right], \quad (\text{D.37})$$

which further implies

$$\mathbf{1}(\Xi) \left(\frac{1}{m_2} + 1 + \gamma_n m_1 \right) \prec n^{-1/2}, \quad \mathbf{1}(\Xi) \left(\frac{1}{m_3} + 1 + \gamma_n m \right) \prec n^{-1/2}. \quad (\text{D.38})$$

Finally, plugging (D.31) and (D.32) into (D.38), we conclude (D.24). \square

Step 3: Ξ holds with high probability. In this step, we show that the event $\Xi(z)$ in fact holds with high probability for all $z \in \mathbf{D}$. Once we have proved this fact, then applying Lemma C.6 to (D.24) immediately shows that $(m_2(z), m_3(z))$ is equal to $(m_{2c}(z), m_{3c}(z))$ up to an error of order $n^{-1/2}$.

First we claim that it suffices to show that

$$|m_2(0) - m_{2c}(0)| + |m_3(0) - m_{3c}(0)| \prec n^{-1/2}. \quad (\text{D.39})$$

Once we know (D.39), then by (C.22) and (D.5), we know $\max_{\alpha=2}^3 |m_{\alpha c}(z) - m_{\alpha c}(0)| = O((\log n)^{-1})$ and $\max_{\alpha=2}^3 |m_{\alpha}(z) - m_{\alpha}(0)| = O((\log n)^{-1})$ with high probability for $z \in \mathbf{D}$. Together with (D.39), we obtain that

$$|m_2(z) - m_{2c}(z)| + |m_3(z) - m_{3c}(z)| \lesssim (\log n)^{-1} \quad \text{with high probability.} \quad (\text{D.40})$$

and

$$\sup_{z \in \mathbf{D}} (|m_2(z) - m_{2c}(0)| + |m_3(z) - m_{3c}(0)|) \lesssim (\log n)^{-1} \quad \text{with high probability.} \quad (\text{D.41})$$

The condition (D.40) shows that Ξ holds with high probability, and the condition (D.41) verifies the condition (C.21) of Lemma C.6. Hence applying Lemma C.6 to (D.24), we obtain that

$$|m_2(z) - m_{2c}(z)| + |m_3(z) - m_{3c}(z)| \prec n^{-1/2} \quad (\text{D.42})$$

for all $z \in \mathbf{D}$. Plugging (D.42) into (D.25)-(D.27), we get the diagonal estimate

$$\max_{a \in \mathcal{I}} |G_{aa}(z) - \Pi_{aa}(z)| \prec n^{-1/2}. \quad (\text{D.43})$$

Together with the off-diagonal estimate in (D.19), we conclude (D.15).

Lemma D.10. *Under the assumptions in Proposition D.1, the estimate (D.39) holds.*

Proof. By (C.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 (D.26), (D.27) and (D.29), we obtain that (D.37) and (D.38) hold at $z = 0$ even without the indicator function $\mathbf{1}(\Xi)$. Furthermore, it gives that

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

with high probability. Then using (D.25) and (D.29), we obtain that (D.31) and (D.32) hold at $z = 0$ even without the indicator function $\mathbf{1}(\Xi)$. Finally, plugging (D.31) and (D.32) into (D.38), we conclude (D.24) holds at $z = 0$, that is,

$$\begin{aligned} \left| \frac{1}{m_2(0)} + 1 - \frac{1}{n} \sum_{i=1}^p \frac{\sigma_i^2}{\sigma_i^2 r_1 m_2(0) + r_2 m_3(0)} \right| &\prec n^{-1/2}, \\ \left| \frac{1}{m_3(0)} + 1 - \frac{1}{n} \sum_{i=1}^p \frac{1}{\sigma_i^2 r_2 m_2(0) + r_2 m_3(0)} \right| &\prec n^{-1/2}. \end{aligned} \quad (\text{D.44})$$

Denoting $\omega_2 = -m_{2c}(0)$ and $\omega_3 = -m_{3c}(0)$. By (D.38), we have

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

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

$$c \leq \omega_2 \leq 1, \quad c \leq \omega_3 \leq 1, \quad \text{with high probability.} \quad (\text{D.45})$$

Moreover, one can verify from (D.44) that (ω_2, ω_3) satisfy approximately the same equations as in (C.19):

$$r_1 \omega_2 + r_2 \omega_3 = 1 - \gamma_n + O_{\prec}(n^{-1/2}), \quad f(\omega_2) = 1 + O_{\prec}(n^{-1/2}). \quad (\text{D.46})$$

The first equation and (D.45) together implies that $\omega_2 \in [0, r_1^{-1}(1 - \gamma_n)]$ with high probability. Since f is strictly increasing and has bounded derivatives on $[0, r_1^{-1}(1 - \gamma_n)]$, by basic calculus the second equation in (D.46) gives that $|\omega_2 - b_2| \prec n^{-1/2}$. Together with the first equation in (D.46), we get $|\omega_3 - b_3| \prec n^{-1/2}$. This concludes (D.39). \square

This lemma concludes (D.39), and as explained above, concludes the proof of Lemma D.7. \square

With Lemma D.7, we can conclude the proof of Proposition D.1.

Proof of Proposition D.1. With (D.15), one can repeat the polynomialization method in [4, Section 5] to get the anisotropic local law (C.26) for G_0 . The proof is exactly the same, except for some minor notation difference, so we omit the details. \square

D.3 Anisotropic local law

In this section, we finish the proof of Theorem C.7 for a general X satisfying the bounded support condition (C.8) with $q \leq n^{-\phi}$ for some constant $\phi > 0$. The proposition D.1 implies that (C.26) holds for Gaussian Z_1^{Gauss} and Z_2^{Gauss} . Thus the basic idea is to prove that for Z_1 and Z_2 satisfying the assumptions in Theorem C.7,

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

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 $Z^{Gauss} := \begin{pmatrix} Z_1^{Gauss} \\ Z_2^{Gauss} \end{pmatrix}$. We prove the above statement using a continuous comparison argument introduced in [16]. The proof is similar to the ones in Sections 7-8 of [16], so we only give an outline without writing down all the details.

Definition D.11 (Interpolating matrices). *We denote Introduce the notations $Z^0 := Z^{Gauss}$ and $Z^1 := Z$. Let $\rho_{\mu i}^0$ and $\rho_{\mu i}^1$ be the laws of $Z_{\mu i}^0$ and $Z_{\mu i}^1$, respectively. For $\theta \in [0, 1]$, we define the interpolated law*

$$\rho_{\mu i}^\theta := (1 - \theta)\rho_{\mu i}^0 + \theta\rho_{\mu i}^1.$$

We shall work on the probability space consisting of triples (Z^0, Z^θ, Z^1) of independent $n \times p$ random matrices, where the matrix $Z^\theta = (Z_{\mu i}^\theta)$ has law

$$\prod_{i \in \mathcal{I}_1} \prod_{\mu \in \mathcal{I}_2 \cup \mathcal{I}_3} \rho_{\mu i}^\theta(dZ_{\mu i}^\theta). \quad (\text{D.47})$$

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

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

We also introduce the matrices

$$G^\theta(z) := G(Z^\theta, z), \quad G_{(\mu i)}^{\theta, \lambda}(z) := G(Z_{(\mu i)}^{\theta, \lambda}, z).$$

We shall prove (C.26) through interpolation matrices Z^θ between Z^0 and Z^1 . We have see that (C.26) holds for Z^0 by Proposition D.1. Using (D.47) and fundamental calculus, we get the following basic interpolation formula: for $F : \mathbb{R}^{n \times p} \rightarrow \mathbb{C}$,

$$\frac{d}{d\theta} \mathbb{E} F(Z^\theta) = \sum_{i \in \mathcal{I}_1} \sum_{\mu \in \mathcal{I}_2 \cup \mathcal{I}_3} \left[\mathbb{E} F \left(Z_{(\mu i)}^{\theta, Z_{\mu i}^1} \right) - \mathbb{E} F \left(Z_{(\mu i)}^{\theta, Z_{\mu i}^0} \right) \right] \quad (\text{D.48})$$

provided all the expectations exist.

We shall apply (D.48) to $F(Z) := F_{\mathbf{u}\mathbf{v}}^p(Z, z)$ for (large) $p \in 2\mathbb{N}$ and $F_{\mathbf{v}}(Z, z)$ defined as

$$F_{\mathbf{u}\mathbf{v}}(Z, z) := |G_{\mathbf{u}\mathbf{v}}(Z, z) - \Pi_{\mathbf{u}\mathbf{v}}(z)|. \quad (\text{D.49})$$

Here for simplicity of notations, we introduce the following notation of generalized entries: for $\mathbf{u}, \mathbf{v} \in \mathbb{R}^{\mathcal{I}}$, we shall denote $G_{\mathbf{u}\mathbf{v}} := \mathbf{u}^\top G \mathbf{v}$. Moreover, we shall abbreviate $G_{\mathbf{u}\mathbf{a}} := G_{\mathbf{u}\mathbf{e}_a}$ for $\mathbf{a} \in \mathcal{I}$, where \mathbf{e}_a is the standard unit vector along \mathbf{a} -th axis. Given any vector $\mathbf{u} \in \mathbb{R}^{\mathcal{I}_{1,2,3}}$, we always identify it with its natural embedding in $\mathbb{R}^{\mathcal{I}}$. The exact meanings will be clear from the context. The main work is to show the following self-consistent estimate for the right-hand side of (D.48) for any fixed $p \in 2\mathbb{N}$ and constant $\varepsilon > 0$:

$$\sum_{i \in \mathcal{I}_1} \sum_{\mu \in \mathcal{I}_2 \cup \mathcal{I}_3} \left[\mathbb{E} F_{\mathbf{u}\mathbf{v}}^p \left(Z_{(\mu i)}^{\theta, Z_{\mu i}^1}, z \right) - \mathbb{E} F_{\mathbf{u}\mathbf{v}}^p \left(Z_{(\mu i)}^{\theta, Z_{\mu i}^0}, z \right) \right] = O((n^\varepsilon q)^p + \mathbb{E} F_{\mathbf{u}\mathbf{v}}^p(Z^\theta, z)) \quad (\text{D.50})$$

for all $\theta \in [0, 1]$. If (D.50) holds, then combining (D.48) with a Grönwall's argument we obtain that for any fixed $p \in 2\mathbb{N}$ and constant $\varepsilon > 0$:

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

Together with Markov's inequality, we conclude (C.26). In order to prove (D.50), we compare $Z_{(\mu i)}^{\theta, Z_{\mu i}^0}$ and $Z_{(\mu i)}^{\theta, Z_{\mu i}^1}$ via a common $Z_{(\mu i)}^{\theta, 0}$, i.e. we will prove that

$$\sum_{i \in \mathcal{I}_1} \sum_{\mu \in \mathcal{I}_2 \cup \mathcal{I}_3} \left[\mathbb{E} F_{\mathbf{u}\mathbf{v}}^p \left(Z_{(\mu i)}^{\theta, Z_{\mu i}^a}, z \right) - \mathbb{E} F_{\mathbf{v}}^p \left(Z_{(\mu i)}^{\theta, 0}, z \right) \right] = O \left((n^\varepsilon q)^p + \mathbb{E} F_{\mathbf{u}\mathbf{v}}^p (Z^\theta, z) \right) \quad (\text{D.51})$$

for all $a \in \{0, 1\}$ and $\theta \in [0, 1]$. Underlying the proof of (D.51) is an expansion approach. We define the $\mathcal{I} \times \mathcal{I}$ matrix $\Delta_{(\mu i)}^\lambda$ as

$$\Delta_{(\mu i)}^\lambda := \lambda \begin{pmatrix} 0 & \mathbf{u}_i^{(\mu)} \mathbf{e}_\mu^\top \\ \mathbf{e}_\mu (\mathbf{u}_i^{(\mu)})^\top & 0 \end{pmatrix}, \quad (\text{D.52})$$

where we denote $\mathbf{u}_i^{(\mu)} := \Lambda U \mathbf{e}_i$ if $\mu \in \mathcal{I}_2$ and $\mathbf{u}_i^{(\mu)} := V \mathbf{e}_i$ if $\mu \in \mathcal{I}_3$. Then by the definition of H in (C.11), we have for any $\lambda, \lambda' \in \mathbb{R}$ and $K \in \mathbb{N}$,

$$G_{(i\mu)}^{\theta, \lambda'} = G_{(\mu i)}^{\theta, \lambda} + \sum_{k=1}^K G_{(\mu i)}^{\theta, \lambda} \left(\Delta_{(\mu i)}^{\lambda - \lambda'} G_{(\mu i)}^{\theta, \lambda} \right)^k + G_{(\mu i)}^{\theta, \lambda'} \left(\Delta_{(\mu i)}^{\lambda - \lambda'} G_{(\mu i)}^{\theta, \lambda} \right)^{K+1}. \quad (\text{D.53})$$

Using this expansion and the a priori bound (D.4), it is easy to prove the following estimate: if y is a random variable satisfying $|y| \prec q$, then

$$G_{(\mu i)}^{\theta, y} = O(1), \quad i \in \mathcal{I}_1, \mu \in \mathcal{I}_2 \cup \mathcal{I}_3, \quad (\text{D.54})$$

with high probability.

In the following proof, for simplicity of notations, we introduce $f_{(\mu i)}(y) := F_{\mathbf{v}}^p(Z_{(\mu i)}^{\theta, y})$. We use $f_{(\mu i)}^{(r)}$ to denote the r -th derivative of $f_{(\mu i)}$. By (D.54), it is easy to see that for any fixed $r \in \mathbb{N}$, $f_{(\mu i)}^{(r)}(y) = O(1)$ with high probability for any random variable y satisfying $|y| \prec q$. Then the Taylor expansion of $f_{(\mu i)}$ gives

$$f_{(\mu i)}(y) = \sum_{r=0}^{p+4} \frac{y^r}{r!} f_{(\mu i)}^{(r)}(0) + O_{\prec}(q^{p+4}), \quad (\text{D.55})$$

Therefore we have for $a \in \{0, 1\}$,

$$\begin{aligned} \mathbb{E} F_{\mathbf{v}}^p \left(Z_{(\mu i)}^{\theta, Z_{\mu i}^a} \right) - \mathbb{E} F_{\mathbf{v}}^p \left(Z_{(\mu i)}^{\theta, 0} \right) &= \mathbb{E} [f_{(\mu i)}(Z_{i\mu}^a) - f_{(\mu i)}(0)] \\ &= \mathbb{E} f_{(\mu i)}(0) + \frac{1}{2n} \mathbb{E} f_{(\mu i)}^{(2)}(0) + \sum_{r=4}^{p+4} \frac{1}{r!} \mathbb{E} f_{(\mu i)}^{(r)}(0) \mathbb{E} (Z_{i\mu}^a)^r + O_{\prec}(q^{p+4}). \end{aligned} \quad (\text{D.56})$$

Here to illustrate the idea in a more concise way, we assume the extra condition

$$\mathbb{E}(Z_{\mu i}^1)^3 = 0, \quad 1 \leq \mu \leq n, \quad 1 \leq i \leq p. \quad (\text{D.57})$$

Hence the $r = 3$ term in the Taylor expansion vanishes. However, this is not necessary as we will explain at the end of the proof.

By (C.2) and the bounded support condition, we have

$$|\mathbb{E} (Z_{i\mu}^a)^r| \prec n^{-2} q^{r-4}, \quad r \geq 4. \quad (\text{D.58})$$

Thus to show (D.51), we only need to prove for $r = 4, 5, \dots, p+4$,

$$n^{-2} q^{r-4} \sum_{i \in \mathcal{I}_1} \sum_{\mu \in \mathcal{I}_2 \cup \mathcal{I}_3} \left| \mathbb{E} f_{(\mu i)}^{(r)}(0) \right| = O \left((n^\varepsilon q)^p + \mathbb{E} F_{\mathbf{u}\mathbf{v}}^p (Z^\theta, z) \right). \quad (\text{D.59})$$

In order to get a self-consistent estimate in terms of the matrix Z^θ on the right-hand side of (D.59), we want to replace $Z_{(\mu i)}^{\theta, 0}$ in $f_{(\mu i)}(0) = F_{\mathbf{u}\mathbf{v}}^p(Z_{(\mu i)}^{\theta, 0})$ with $Z^\theta = Z_{(\mu i)}^{\theta, Z_{\mu i}^0}$.

Lemma D.12. *Suppose that*

$$n^{-2}q^{r-4} \sum_{i \in \mathcal{I}_1} \sum_{\mu \in \mathcal{I}_2 \cup \mathcal{I}_3} \left| \mathbb{E} f_{(\mu)}^{(r)}(Z_{i\mu}^\theta) \right| = O((n^\varepsilon q)^p + \mathbb{E} F_{\mathbf{v}}^p(X^\theta, z)) \quad (\text{D.60})$$

holds for $r = 4, \dots, 4p + 4$. Then (D.59) holds for $r = 4, \dots, 4p + 4$.

Proof. The proof is the same as the one for [16, Lemma 7.16]. \square

What remains now is to prove (D.60). For simplicity of notations, we shall abbreviate $Z^\theta \equiv Z$. For any $k \in \mathbb{N}$, we denote

$$A_{\mu i}(k) := \left(\frac{\partial}{\partial Z_{\mu i}} \right)^k (G_{\mathbf{u}\mathbf{v}} - \Pi_{\mathbf{u}\mathbf{v}}).$$

The derivative on the right-hand side can be calculated using the expansion (D.53). In particular, it is easy to verify that it satisfies the following bound

$$|A_{\mu i}(k)| \prec \begin{cases} (\mathcal{R}_i^{(\mu)})^2 + \mathcal{R}_\mu^2, & \text{if } k \geq 2 \\ \mathcal{R}_i^{(\mu)} \mathcal{R}_\mu, & \text{if } k = 1 \end{cases}, \quad (\text{D.61})$$

where for $i \in \mathcal{I}_1$ and $\mu \in \mathcal{I}_2 \cup \mathcal{I}_3$, we denote

$$\mathcal{R}_i^{(\mu)} := |G_{\mathbf{u}\mathbf{u}_i^{(\mu)}}| + |G_{\mathbf{v}\mathbf{u}_i^{(\mu)}}|, \quad \mathcal{R}_\mu := |G_{\mathbf{u}\mu}| + |G_{\mathbf{v}\mu}|. \quad (\text{D.62})$$

Then we can calculate the derivative

$$\left(\frac{\partial}{\partial Z_{\mu i}} \right)^r F_{\mathbf{u}\mathbf{v}}^p(Z) = \sum_{k_1 + \dots + k_p = r} \prod_{t=1}^{p/2} \left(A_{\mu i}(k_t) \overline{A_{\mu i}(k_{t+p/2})} \right).$$

Then to prove (D.60), it suffices to show that

$$n^{-2}q^{r-4} \sum_{i \in \mathcal{I}_1} \sum_{\mu \in \mathcal{I}_2 \cup \mathcal{I}_3} \left| \mathbb{E} \prod_{t=1}^{p/2} A_{\mu i}(k_t) \overline{A_{\mu i}(k_{t+p/2})} \right| = O((n^\varepsilon q)^p + \mathbb{E} F_{\mathbf{u}\mathbf{v}}^p(Z, z)) \quad (\text{D.63})$$

for $4 \leq r \leq p + 4$ and $(k_1, \dots, k_p) \in \mathbb{N}^p$ satisfying $k_1 + \dots + k_p = r$. Treating zero k 's separately (note $A_{\mu i}(0) = (G_{\mathbf{u}\mathbf{v}} - \Pi_{\mathbf{u}\mathbf{v}})$ by definition), we find that it suffices to prove

$$n^{-2}q^{r-4} \sum_{i \in \mathcal{I}_1} \sum_{\mu \in \mathcal{I}_2 \cup \mathcal{I}_3} \mathbb{E} |A_{\mu i}(0)|^{p-l} \prod_{t=1}^l |A_{\mu i}(k_t)| = O((n^\varepsilon q)^p + \mathbb{E} F_{\mathbf{u}\mathbf{v}}^p(Z, z)) \quad (\text{D.64})$$

for $4 \leq r \leq p + 4$ and $1 \leq l \leq p$. Here without loss of generality, we assume that $k_t = 0$ for $l + 1 \leq t \leq p$, and $\sum_{t=1}^l k_t = r$ with $k_t \geq 1$ for $t \leq l$.

Now we first consider the case $r \leq 2l - 2$. Then by pigeonhole principle, there exist at least two k_t 's with $k_t = 1$. Therefore by (D.61) we have

$$\prod_{t=1}^l |A_{\mu i}(k_t)| \prec \mathbf{1}(r \geq 2l - 1) \left[(\mathcal{R}_i^{(\mu)})^2 + \mathcal{R}_\mu^2 \right] + \mathbf{1}(r \leq 2l - 2) (\mathcal{R}_i^{(\mu)})^2 \mathcal{R}_\mu^2. \quad (\text{D.65})$$

Using (D.4) and a similar argument as in (D.20), we get that

$$\sum_{i \in \mathcal{I}_1} (\mathcal{R}_i^{(\mu)})^2 = O(1), \quad \sum_{\mu \in \mathcal{I}_2 \cup \mathcal{I}_3} \mathcal{R}_\mu^2 = O(1), \quad \text{with high probability.} \quad (\text{D.66})$$

Using (D.66) and $n^{-1/2} \leq q$, we get that

$$\begin{aligned} n^{-2}q^{r-4} \sum_{i \in \mathcal{I}_1} \sum_{\mu \in \mathcal{I}_2 \cup \mathcal{I}_3} |A_{\mu i}(0)|^{p-l} \prod_{t=1}^l |A_{\mu i}(k_t)| &\prec q^{r-4} F_{\mathbf{u}\mathbf{v}}^{p-l}(Z) [\mathbf{1}(r \geq 2l - 1) n^{-1} + \mathbf{1}(r \leq 2l - 2) n^{-2}] \\ &\leq F_{\mathbf{u}\mathbf{v}}^{p-l}(Z) [\mathbf{1}(r \geq 2l - 1) q^{r-2} + \mathbf{1}(r \leq 2l - 2) q^r]. \end{aligned}$$

If $r \leq 2l - 2$, then we get $q^r \leq q^l$ using the trivial inequality $r \geq l$. On the other hand, if $r \geq 4$ and $r \geq 2l - 1$, then $r \geq l + 2$ and we get $q^r \leq q^{l+2}$. Therefore we conclude that

$$n^{-2} q^{r-4} \sum_{i \in \mathcal{I}_1} \sum_{\mu \in \mathcal{I}_2 \cup \mathcal{I}_3} |A_{\mu i}(0)|^{p-l} \prod_{t=1}^l |A_{\mu i}(k_t)| \prec F_{\mathbf{u}\mathbf{v}}^{p-l}(Z) q^l.$$

Now (D.64) follows from Hölder's inequality. This concludes the proof of (D.60), and hence of (D.51), and hence of (C.26).

Finally, if the condition (D.57) does not hold, then there is also an $r = 3$ term in the Taylor expansion (D.56):

$$\frac{1}{6} \mathbb{E} f_{(\mu i)}^{(3)}(0) \mathbb{E} (Z_{i\mu}^a)^3.$$

Note that $\mathbb{E} (Z_{i\mu}^a)^3$ is of order $n^{-3/2}$, while the sum over i and μ in (D.51) provides a factor n^2 . In fact, $\mathbb{E} f_{(\mu i)}^{(3)}(0)$ will provide an extra $n^{-1/2}$ to compensate the remaining $n^{1/2}$ factor. This follows from an improved self-consistent comparison argument for sample covariance matrices in [16, Section 8]. The argument for our case is almost the same except for some notational differences, so we omit the details.

D.4 Proof of Lemma C.5 and Lemma C.6

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

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

$$r_1 m_{2c} = -(1 - \gamma_n) - r_2 m_{3c} - z \left(\frac{1}{m_{3c}} + 1 \right), \quad g_z(m_{3c}(z)) = 1, \quad (\text{D.67})$$

where

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

We first show that there exists a unique solution $m_{3c}(z)$ to the equation $g_z(m_{3c}(z)) = r_2$ under the conditions in (C.21), and the solution satisfies (C.22). Now we abbreviate $\varepsilon(z) := m_{3c}(z) - m_{3c}(0)$, and from (D.67) we can obtain that

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

which implies

$$g'_z(m_{3c}(0))\varepsilon(z) = -[g_z(m_{3c}(0)) - g_0(m_{3c}(0))] - [g_z(m_{3c}(0) + \varepsilon(z)) - g_z(m_{3c}(0)) - g'_z(m_{3c}(0))\varepsilon(z)]. \quad (\text{D.68})$$

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

$$\varepsilon^{(k+1)} = -\frac{g_z(m_{3c}(0)) - g_0(m_{3c}(0))}{g'_z(m_{3c}(0))} - \frac{g_z(m_{3c}(0) + \varepsilon^{(k)}) - g_z(m_{3c}(0)) - g'_z(m_{3c}(0))\varepsilon^{(k)}}{g'_z(m_{3c}(0))}. \quad (\text{D.69})$$

In other words, the above equation 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

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

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

$$|[g'_z(m_{3c}(0))]^{-1}| \leq \tilde{C}, \quad \left| \frac{g_z(m_{3c}(0) + \varepsilon_1) - g_z(m_{3c}(0) + \varepsilon_2) - g'_z(m_{3c}(0))(\varepsilon_1 - \varepsilon_2)}{g'_z(m_{3c}(0))} \right| \leq \tilde{C} |\varepsilon_1 - \varepsilon_2|^2, \quad (\text{D.70})$$

and

$$\left| \frac{g_z(m_{3c}(0)) - g_0(m_{3c}(0))}{g'_z(m_{3c}(0))} \right| \leq \tilde{C}|z|, \quad (\text{D.71})$$

for all $|z| \leq \tilde{c}$ and $|\varepsilon_1| \leq \tilde{c}$, $|\varepsilon_2| \leq \tilde{c}$. Then with (D.70) and (D.71), 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

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

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 to prove that h restricted to B_r is a contraction, which then implies that $\varepsilon := \lim_{k \rightarrow \infty} \varepsilon^{(k)}$ exists and $m_{3c}(0) + \varepsilon$ is a unique solution to the second equation of (D.67) subject to the condition $\|\varepsilon\|_\infty \leq r$.

From the iteration relation (D.69), using (D.70) one can readily check that

$$\varepsilon^{(k+1)} - \varepsilon^{(k)} = h(\varepsilon^{(k)}) - h(\varepsilon^{(k-1)}) \leq \tilde{C}|\varepsilon^{(k)} - \varepsilon^{(k-1)}|^2. \quad (\text{D.72})$$

Hence as long as r is chosen to be sufficiently small such that $2r\tilde{C} \leq 1/2$, then h is indeed a contraction mapping on B_r , which proves both the existence and uniqueness of the solution $m_{3c}(z) = m_{3c}(0) + \varepsilon$, if we choose c_0 in (C.21) as $c_0 = \min\{c_\delta, r\}$. After obtaining $m_{3c}(z)$, we can then find $m_{2c}(z)$ using the first equation in (D.67).

Note that with (D.71) and $\varepsilon^{(0)} = 0$, we get from (D.69) that

$$|\varepsilon^{(1)}| \leq \tilde{C}|z|.$$

With the contraction mapping, we have the bound

$$|\varepsilon| \leq \sum_{k=0}^{\infty} \|\varepsilon^{(k+1)} - \varepsilon^{(k)}\|_\infty \leq 2\tilde{C}|z|. \quad (\text{D.73})$$

This gives the bound (C.22) for $m_{3c}(z)$. Using the first equation in (D.67), we immediately obtain the bound

$$r_1|m_{2c}(z) - m_{2c}(0)| \leq C|z|.$$

This gives (C.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 (C.16) and treat $m_{2c}(z)$ as the solution to the following equation:

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

We can calculate that

$$g'_z(m_{2c}(0)) = -1 + \frac{\gamma_n}{p} \sum_{i=1}^p \frac{\sigma_i^2(z + r_2 m_{3c}(z))}{(z + \sigma_i^2 r_1 m_{2c}(0) + r_2 m_{3c}(z))^2}.$$

At $z = 0$, we have

$$|g'_0(m_{2c}(0))| = \left| 1 + \frac{\gamma_n}{p} \sum_{i=1}^p \frac{\sigma_i^2 r_2 b_3}{(\sigma_i^2 r_1 b_2 + r_2 b_3)^2} \right| \geq 1,$$

where b_2 and b_3 satisfy (C.20). Thus under (C.21) we have $|g'_z(m_{2c}(0))| \sim 1$ as long as c_0 is taken sufficiently small. Then with the above arguments for $m_{3c}(z)$ between (D.67) and (D.73), we can conclude (C.22) for $m_{2c}(z)$. This concludes the proof of Lemma C.5. \square

Proof of Lemma C.6. Under (C.21), we can obtain equation (D.67) approximately up to some small error

$$r_1 m_{2c} = -(1 - \gamma_n) - r_2 m_{3c} - z \left(\frac{1}{m_{3c}} + 1 \right) + \mathcal{E}'_2(z), \quad g_z(m_{3c}(z)) = 1 + \mathcal{E}'_3(z), \quad (\text{D.74})$$

with $|\mathcal{E}'_2(z)| + |\mathcal{E}'_3(z)| = O(\delta(z))$. Then we subtract the equations (D.67) from (D.74), and consider the contraction principle for the functions $\varepsilon(z) := m_3(z) - m_{3c}(z)$. The rest of the proof is exactly the one for Lemma C.5, so we omit the details. \square