

000 BENIGN OVERFITTING IN OUT-OF-DISTRIBUTION 001 GENERALIZATION OF LINEAR MODELS 002

003 **Anonymous authors**
004

005 Paper under double-blind review
006

007 ABSTRACT 008

009 Benign overfitting refers to the phenomenon where an over-parameterized model
010 fits the training data perfectly, including noise in the data, but still generalizes well
011 to the unseen test data. While prior work provides some theoretical understanding
012 of this phenomenon under the in-distribution setup, modern machine learning often
013 operates in a more challenging Out-of-Distribution (OOD) regime, where the
014 target (test) distribution can be rather different from the source (training) distribution.
015 In this work, we take an initial step towards understanding benign overfitting
016 in the OOD regime by focusing on the basic setup of over-parameterized linear
017 models under covariate shift. We provide non-asymptotic guarantees proving that
018 benign overfitting occurs in standard ridge regression, even under the OOD regime
019 when the target covariance satisfies certain structural conditions. We identify several
020 vital quantities relating to source and target covariance, which govern the
021 performance of OOD generalization. Our result is sharp, which provably recovers
022 prior in-distribution benign overfitting guarantee (Tsigler & Bartlett, 2023), as
023 well as under-parameterized OOD guarantee (Ge et al., 2024) when specializing
024 to each setup. Moreover, we also present theoretical results for a more general
025 family of target covariance matrix, where standard ridge regression only achieves
026 a slow statistical rate of $\mathcal{O}(1/\sqrt{n})$ for the excess risk, while Principal Component
027 Regression (PCR) is guaranteed to achieve the fast rate $\mathcal{O}(1/n)$, where n is the
028 number of samples.

030 1 INTRODUCTION 031

032 In modern machine learning, distribution shift has become a ubiquitous challenge where models
033 trained on a source data distribution are tested on a different target distribution (Zou et al., 2018;
034 Hendrycks & Dietterich, 2019; Guan & Liu, 2021; Koh et al., 2021). Generalization under distribu-
035 tion shift, known as Out-of-Distribution (OOD) generalization, remains a fundamental issue in the
036 practical application of machine learning (Recht et al., 2019; Hendrycks et al., 2021; Miller et al.,
037 2021; Wenzel et al., 2022). While there has been extensive work on the theoretical understand-
038 ing of OOD generalization, most of it has focused on under-parameterized models (Shimodaira, 2000;
039 Lei et al., 2021; Ge et al., 2024; Zhang et al., 2022). However, over-parameterized models, such as
040 deep neural networks and large language models (LLMs) in the fine-tuning stage, which have more
041 parameters than training samples, are widely used in modern machine learning. Surprisingly, despite
042 the classic bias-variance tradeoff for under-parameterized models, over-parameterized models tend
043 to overfit the data while still achieving strong in-distribution generalization, a phenomenon known
044 as benign overfitting (Hastie et al., 2022; Shamir, 2023) or harmless interpolation (Muthukumar
045 et al., 2020). Therefore, it is crucial to theoretically understand how benign overfitting shapes OOD
046 generalization in over-parameterized models.

047 It is established in over-parameterized models that “benign overfitting” occurs when the data essen-
048 tially resides on a low-dimensional manifold. The manifold assumption (Belkin & Niyogi, 2003) is
049 widely applicable across image, speech, and language data, where although features are embedded
050 in a high-dimensional ambient space, their generation is governed by a few degrees of freedom im-
051 posed by physical constraints (Niyogi, 2013). Specifically, the covariance matrix of the data should
052 be characterized by several major directions corresponding to large eigenvalues, while the remaining
053 directions are high-dimensional but associated with small eigenvalues. In this setting, even though
the estimator may overfit the noise in the low-variance directions, it can still capture the signal along

the major directions while the noise is damped in the minor directions. Recent non-asymptotic analyses have provided upper bounds on the excess risk for the minimum-norm interpolant and over-parameterized ridge estimator under this framework (Bartlett et al., 2020; Hastie et al., 2022; Tsigler & Bartlett, 2023).

However, a theoretical understanding of OOD generalization in over-parameterized models remains elusive. In this paper, we take an initial step towards characterizing OOD generalization in over-parameterized models under *general* covariate shift, a standard assumption in the OOD regime (Ben-David et al., 2006), where the conditional distribution of the outcome given the covariates remains invariant. We derive the first vanishing, non-asymptotic excess risk bound for ridge regression and minimum-norm interpolation, assuming that the source covariance is dominated by a few major eigenvalues, which satisfies the in-distribution benign overfitting condition. Notably, we allow the target covariance to be arbitrary. This result stands in contrast to recent work that either addresses only a restrictive form of covariate shift (Hao et al., 2024; Mallinar et al., 2024) or provides excess risk bounds that asymptotically remain above a constant (Tripuraneni et al., 2021b; Hao et al., 2024).

In summary, our excess risk bound identifies several key quantities that relate to the source and target covariance, suggesting that “benign overfitting” occurs when these quantities are well controlled. In such cases, the target distribution data lies on the low-dimensional manifold of the source distribution. Otherwise, ridge regression may incur excess risk, lower bounded by the slow statistical rate of $\mathcal{O}(1/\sqrt{n})$. In contrast, we show that Principal Component Regression (PCR) achieves the fast rate of $\mathcal{O}(1/n)$ in such scenarios. **The non-asymptotic rates of both ridge regression and PCR are validated through simulation experiments on multivariate Gaussian data in Appendix A.**

Our contributions.

1. We provide a sharp, instance-dependent excess risk bound for ridge regression (Theorem 2). Our result applies to any target distribution, requiring only that the source covariance be dominated by a few major eigenvectors and that the minor components are high-dimensional. We show that ridge regression exhibits “benign overfitting,” achieving excess risk comparable to the in-distribution case, provided that certain key quantities relating to the source and target distributions are bounded. Importantly, this condition requires that the *overall magnitude* of the target covariance along the minor directions scales similarly to, or smaller than, that of the source, but it does not depend on the spectral structure of the target covariance. Our results recover the in-distribution bound from Tsigler & Bartlett (2023) when the source and target match and also recover the sharp bound from Ge et al. (2024) for under-parameterized linear regression under covariate shift when the minor components vanish.
2. We extend our analysis to examine the scenario where the target distribution exhibits significant variance in the minor directions. In this setting, ridge regression incurs a higher error rate compared to the in-distribution case, specifically the slow statistical rate of $\mathcal{O}(1/\sqrt{n})$ in certain instances (Theorem 4). However, we demonstrate that Principal Component Regression ensures a fast rate of $\mathcal{O}(1/n)$ in these cases, provided that the true signal primarily lies in the major directions of the source (Theorem 5). Additionally, PCR does not rely on the minor directions of the source distribution being high-dimensional, highlighting its advantage over ridge regression in such settings.

1.1 RELATED WORK

Over-parameterization. The success of over-parameterized models in machine learning has sparked significant research on their theoretical foundations. Harmless interpolation (Muthukumar et al., 2020) or benign overfitting (Shamir, 2023) describes cases where linear models interpolate noise yet still generalize well. Double descent in prediction error is also observed as the ambient dimension surpasses the number of training samples (Nakkiran, 2019; Xu & Hsu, 2019).

Research in this field can be divided into two categories based on assumptions about the spectral structure of the sample covariance. The first category assumes an almost isotropic sample covariance matrix with a bounded condition number or an isotropic prior distribution of parameters (Belkin et al., 2020). In this case, a limiting covariance spectral structure may emerge when $n \asymp d$ and both tend to infinity, allowing for asymptotic risk bounds (Dobriban & Wager, 2018; Richards et al., 2021). However, ridgeless regression is sub-optimal in this setting unless the signal-to-noise ratio

108 is infinite (Wu & Xu, 2020), and non-asymptotic error bounds are lacking. Our work falls into
 109 the second category, focusing on the covariance model where a small number of eigenvalues dom-
 110 inate the sample covariance, and the signal is concentrated in the subspace spanned by the leading
 111 eigenvectors (Bibas et al., 2019; Chinot & Lerasle, 2022; Hastie et al., 2022). Linear regression can
 112 be optimal without regularization under this covariance structure (Kobak et al., 2020), which is of
 113 practical interest because ridgeless regression is equivalent as gradient descent from zero initializa-
 114 tion (Zhou et al., 2020). Sharp non-asymptotic bounds for variance and bias in ridge and ridgeless
 115 regression have been derived (Bartlett et al., 2020; Tsigler & Bartlett, 2023).

116 Extending the analysis of ridgeless estimators (i.e., minimum norm interpolants), uniform conver-
 117 gence bounds for generalization error have been studied for all interpolants with arbitrary norms.
 118 However, uniformly bounding the difference between population and empirical errors generally
 119 fails to ensure a consistent predictor (Zhou et al., 2020), necessitating strong assumptions on dis-
 120 tributions (Koehler et al., 2021) or hypothesis classes (Negrea et al., 2020). Over-parameterization
 121 theory for linear models has also been applied to two-layer neural networks approximated via kernel
 122 ridge regression (Liang et al., 2020; Ghorbani et al., 2020; 2021; Bartlett et al., 2021; Mei & Monta-
 123 nari, 2022; Mei et al., 2022; Montanari & Zhong, 2022; Simon et al., 2023), though this lies beyond
 124 the scope of the present work.

125 **Out-of-Distribution generalization.** Out-of-distribution generalization is well studied for under-
 126 parameterized models, particularly in transfer learning between two distributions, where labeled
 127 source data is combined with unlabeled target data to train models. For covariate shift, importance
 128 weighting (Cortes et al., 2010; Agapiou et al., 2017) is asymptotically optimal when using den-
 129 sity ratio as weights (Shimodaira, 2000). More generally, the theoretical limits of transfer learning
 130 are explored through minimax lower bounds for distribution shifts bounded by divergence met-
 131 rics (Mousavi Kalan et al., 2020; Zhang et al., 2022). A number of algorithms are proposed to
 132 achieve matching upper bounds (Lei et al., 2021). However, Ge et al. (2024) shows that even without
 133 target data, vanilla MLE (Empirical Risk Minimization, ERM) is minimax optimal for well-specified
 134 models under covariate shift, with a sharp $1/n$ excess risk bound based on Fisher information.

135 Research on over-parameterized models under distribution shift has largely focused on covariate
 136 shifts in linear regression. Importance weighting for over-parameterized models (Chen et al., 2024)
 137 and general sample reweighting offer no advantage over ERM since both converge to the same esti-
 138 mator via gradient descent (Zhai et al., 2022). Consequently, much literature focuses on minimum-
 139 norm interpolation as the natural ERM solution. For isotropic signals, Tripuraneni et al. (2021a)
 140 prove that over-parameterization improves robustness to covariate shift, deriving an asymptotic gen-
 141 eralization bound decreasing with d/n . Under the covariance model dominated by several major
 142 eigenvectors, Hao et al. (2024) derive a non-asymptotic bound for a specific instance of covariate
 143 shift where features are translated by a constant, but the covariance matrix is preserved. However, a
 144 constant excess risk remains in their bound due to estimation variance. Kausik et al. (2024) study a
 145 linear model with additive noise on covariates when data strictly lies in a low-dimensional subspace,
 146 also showing a non-vanishing bound. Mallinar et al. (2024) investigate minimum-norm interpola-
 147 tion with independent covariates and simultaneously diagonal source and target covariance matrices,
 148 allowing them to directly extend in-distribution bounds of Bartlett et al. (2020); Tsigler & Bartlett
 149 (2023). Still, their estimation bias bound is looser than ours, as it exhibits a gap compared to Tsigler
 150 & Bartlett (2023)'s sharp bound even when the source and target distributions are aligned. In con-
 151 trast, our work achieves the first vanishing non-asymptotic error bound for general covariate shift,
 152 assuming only finite second moments for the target covariance matrix.

153 Another line of research considers non-parametric models under covariate shift (Kpotufe & Mar-
 154 tinet, 2018; Hanneke & Kpotufe, 2019; Pathak et al., 2022; Ma et al., 2023), presenting minimax
 155 results governed by a transfer-exponent that measures the similarity between source and target dis-
 156 tributions. However, this falls outside the scope of our work.

157 **Principal Component Regression.** Principal Component Regression (PCR) has been designed to
 158 treat multicollinearity in high-dimensional linear regression, where the covariates possess a latent,
 159 low-dimensional representation (Massy, 1965; Jeffers, 1967; Jolliffe, 1982; Jeffers, 1981). PCR has
 160 been widely used in statistics (Liu et al., 2003; Fan et al., 2021; Fan & Gu, 2023), econometrics
 161 (Stock & Watson, 2002; Bai & Ng, 2002; Fan et al., 2020), chemometrics (Næs & Martens, 1988;
 Sun, 1995; Vigneau et al., 1997; Depczynski et al., 2000; Keithley et al., 2009), construction man-

gement (Chan & Park, 2005), environmental science (Kumar & Goyal, 2011; Hidalgo et al., 2000), signal processing (Huang & Yang, 2012) and etc.

Regarding the theory of PCR, Hadi & Ling (1998) identify conditions under which PCR will fail. Bair et al. (2006) suggest selecting principal components based on their association with the outcome and provide corresponding asymptotic consistency results. Xu & Hsu (2019) establish asymptotic risk bounds for PCR with varying numbers of selected components k . They show that the “double descent” behavior also happens in PCR as k/d increases, where d represents the data dimension. Most relevant to our work, Agarwal et al. (2019) derive non-asymptotic error bounds for PCR, and show that the error decays at a rate of $\mathcal{O}(1/\sqrt{n})$ (n is the sample size), assuming all singular values of the data matrix are of similar magnitude. Agarwal et al. (2020) further improves this rate to $\mathcal{O}(1/n)$. However, both results consider a fixed design with strict low-rank assumptions, making them inapplicable to our setting of OOD generalization.

2 COVARIATE SHIFT SETUP UNDER OVER-PARAMETERIZATION

2.1 DATA WITH COVARIATE SHIFT

We address the out-of-distribution (OOD) generalization of over-parameterized models under covariate shift, where the covariates, denoted by a random vector $x \in \mathbb{R}^d$, follow different distributions during training and evaluation. Specifically, we assume that the training data is sampled from a source distribution \mathcal{P}_S , and the learned model is subsequently applied to data from an unknown target distribution \mathcal{P}_T . Let the covariates be zero-mean on the source distribution, and define the covariance matrix as $\Sigma_S := \mathbb{E}_{x \sim \mathcal{P}_S} [xx^T]$. Since we can always choose an orthonormal basis such that Σ_S becomes diagonal, we express $\Sigma_S = \text{diag}(\lambda_1, \dots, \lambda_d)$ without loss of generality, where the eigenvalues are arranged in non-increasing order: $\lambda_1 \geq \dots \geq \lambda_d \geq 0$. Moreover, we assume sub-gaussianity of the source covariates, i.e., $\Sigma_S^{-1/2}x$ is σ -sub-gaussian where the precise definition of the sub-Gaussian norm is given in section B. We consider a general covariate distribution for the target, assuming only that it has a finite second moment, denoted by $\Sigma_T := \mathbb{E}_{x \sim \mathcal{P}_T} [xx^T]$, which is not necessarily diagonal.

We consider a linear response model that remains consistent across the source and target distributions. The outcome follows $y = x^T \beta^* + \epsilon$, where $\beta^* \in \mathbb{R}^d$ represents the true parameter, and ϵ is an independent noise with zero-mean and variance v^2 .

2.2 LEARNING PROCEDURE AND EVALUATION

The learning procedure involves training a linear model with n i.i.d. samples $\{(x_i, y_i)\}_{i=1}^n$ drawn from the source distribution. Define $X := (x_1, \dots, x_n)^T \in \mathbb{R}^{n \times d}$, $Y := (y_1, \dots, y_n)^T$ and $\epsilon := (\epsilon_1, \dots, \epsilon_n)^T$. We focus on models $\hat{\beta}(Y)$ that are linear in Y , allowing us to write $\hat{\beta}(Y) = \hat{\beta}(X\beta^*) + \hat{\beta}(\epsilon)$. We consider ridge regression and Principal Component Regression to be two examples of such algorithms. With a regularization coefficient $\lambda \geq 0$, the ridge estimator is defined as

$$\hat{\beta}(Y) = X^T(XX^T + \lambda I_n)^{-1}Y.$$

The estimator is assessed on the target distribution by its excess risk relative to the true model, expressed as the following equation:

$$\mathcal{R}(\hat{\beta}(Y)) := \mathbb{E}_{(x,y) \sim \mathcal{P}_T} [(y - x^T \hat{\beta}(Y))^2 - (y - x^T \beta^*)^2] = \|\hat{\beta}(Y) - \beta^*\|_{\Sigma_T}^2,$$

where we define $\|x\|_A := \sqrt{x^T A x}$ for any positive semi-definite matrix A . The metric of interest is the expected excess risk with respect to the noise, given by $\mathbb{E}_\epsilon [\mathcal{R}(\hat{\beta}(Y))]$. Following from the linearity of the model, the expected excess risk can be decomposed into bias and variance components:

$$\mathbb{E}_\epsilon [\mathcal{R}(\hat{\beta}(Y))] = \mathbb{E}_\epsilon \|\hat{\beta}(\epsilon)\|_{\Sigma_T}^2 + \|\hat{\beta}(X\beta^*) - \beta^*\|_{\Sigma_T}^2,$$

where we define the variance as $V := \mathbb{E}_\epsilon \|\hat{\beta}(\epsilon)\|_{\Sigma_T}^2$ and the bias as $B := \|\hat{\beta}(X\beta^*) - \beta^*\|_{\Sigma_T}^2$.

216 2.3 THE STRUCTURE OF COVARIANCE IN BENIGN OVERFITTING
 217

218 Throughout this paper, we follow the convention of [Tsigler & Bartlett \(2023\)](#) and consider the
 219 *source* covariance matrix Σ_S as characterized by a few number of high-variance directions and a
 220 large number of low-variance directions of similar magnitude. We refer to the high-variance direc-
 221 tions as “major directions” and the low-variance directions as “minor directions”. We denote the
 222 number of major directions as k . For the remaining $d - k$ minor directions, we use the following
 223 notions of effective rank to approximate the number of directions with a similar scale. For the ridge
 224 regularization coefficient $\lambda \geq 0$, we define:

$$225 \quad r_k := \frac{\lambda + \sum_{j>k} \lambda_j}{\lambda_{k+1}}, \quad R_k := \frac{(\lambda + \sum_{j>k} \lambda_j)^2}{\sum_{j>k} \lambda_j^2}.$$

226 We have $1 \leq r_k \leq R_k$. When $\lambda = 0$, it further holds that $R_k \leq d - k$. We denote the first k
 227 columns of X as X_k and the remaining $d - k$ columns as X_{-k} . Correspondingly, we partition β^* into
 228 β_k^* and β_{-k}^* . The covariance matrix blocks along the diagonals are denoted by $\Sigma_{S,k}$, $\Sigma_{S,-k}$, $\Sigma_{T,k}$
 229 and $\Sigma_{T,-k}$. We define the following quantities to facilitate our presentation, which are crucial in our
 230 analysis.

$$231 \quad \mathcal{T} = \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}}, \quad \mathcal{U} = \Sigma_{S,-k} \Sigma_{T,-k}, \quad \mathcal{V} = \Sigma_{S,-k}^2. \quad (1)$$

232 3 OVER-PARAMETERIZED RIDGE REGRESSION
 233

234 In the context of in-distribution generalization, where $\Sigma_S = \Sigma_T$, for over-parameterized linear
 235 models, [Bartlett et al. \(2020\)](#) and [Tsigler & Bartlett \(2023\)](#) demonstrate that the ridge estimator
 236 (with the minimum-norm interpolant as a special case) can effectively learn the signal from the
 237 subspace of data spanned by the major eigenvectors, while benignly overfitting the noise in the
 238 minor directions under certain scenarios. They argue that benign overfitting occurs when the true
 239 signal predominantly lies in the major directions, and the minor directions have a small scale but
 240 highly effective rank. This section explores whether this mechanism still holds under covariate
 241 shift. We derive upper bounds (Theorem 2) for the excess risk of the over-parameterized ridge
 242 estimator in the context of OOD generalization, demonstrating that “benign overfitting” also happens
 243 under covariate shift, given that the target distribution’s covariance remains dominated by the first
 244 k dimensions. Specifically, we show that \mathcal{T} characterizes the shift in the major directions, while
 245 the *overall magnitude* of $\Sigma_{T,-k}$, which captures the shift in the minor directions, is crucial for
 246 benign overfitting. When the overall magnitude of $\Sigma_{T,-k}$ scales similarly to or smaller than those
 247 of the source, ridge regression achieves the same non-asymptotic error rate under covariate shift as
 248 in the in-distribution setting. Surprisingly, although a high effective rank in the minor directions
 249 of the source is essential for benign overfitting, only the overall magnitude matters for the target
 250 distribution.

251 3.1 WARM-UP: IN-DISTRIBUTION BENIGN OVERFITTING
 252

253 As a warm-up, we introduce [Tsigler & Bartlett \(2023\)](#)’s in-distribution result on benign overfitting
 254 in ridge regression. When the data dimension d exceeds the sample size n , the ridge estimator
 255 interpolates the training data, fitting the noise. In this case, the estimator $\hat{\beta}$ lies within the subspace
 256 spanned by the n samples. If d is much larger than n , a new test point is highly likely to be orthogonal
 257 to this subspace, preventing noise from affecting the prediction. Therefore, the minor components
 258 of the covariance matrix actually provide implicit regularization. [Tsigler & Bartlett \(2023\)](#) assume
 259 that the data lies in a space with k major directions and $d - k$ weak but essentially high-dimensional
 260 minor directions, allowing for benign overfitting. This intuition is formalized through an assumption
 261 that controls the condition number of the Gram matrix for the remaining $d - k$ dimensions.

262 **Assumption 1** (CondNum(k, δ, L) ([Tsigler & Bartlett, 2023](#))). Define a matrix $A_k = \lambda I_n +$
 263 $X_{-k} X_{-k}^T$. With probability at least $1 - \delta$, A_k is positive definite and has a condition number no
 264 greater than L , i.e.,

$$265 \quad \frac{\mu_1(A_k)}{\mu_n(A_k)} \leq L,$$

266 where the i -th largest eigenvalue of a matrix is denoted by $\mu_i(\cdot)$.

270 **Remark 1.** This assumption essentially posits that the minor directions of the source covariance
 271 have an effective rank significantly greater than n . As evidence, [Tsigler & Bartlett \(2023\)](#) prove that
 272 if CondNum holds, the effective rank r_k is lower bounded by n/L , up to a constant. Conversely, a
 273 lower bound on the effective rank r_k also implies an upper bound of the condition number of A_k .
 274 For more details, refer to [Tsigler & Bartlett \(2023\)](#), Lemma 3).

275 Assuming CondNum, [Tsigler & Bartlett \(2023\)](#) obtain sharp upper bounds for both the variance and
 276 bias of the ridge estimator, with matching lower bounds (see their Theorem 2). To facilitate the
 277 presentation, we define $\tilde{\lambda} := \lambda + \sum_{j>k} \lambda_j$ to represent the combined regularization term from both
 278 ridge and implicit regularization.

279 **Theorem 1** ([Tsigler & Bartlett \(2023\)](#)). There exists a constant c that only depends on σ, L , such
 280 that for any $n > ck$, if the assumption condNum(k, δ, L) (Assumption 1) is satisfied, then it holds
 281 that $n < cr_k$, and with probability at least $1 - \delta - ce^{-n/c}$,

$$\frac{V}{cv^2} \leq \frac{k}{n} + \frac{n}{R_k}, \quad \frac{B}{c} \leq B_{\text{ID}} := \|\beta_k^*\|_{\Sigma_{S,k}^{-1}}^2 \left(\frac{\tilde{\lambda}}{n} \right)^2 + \|\beta_{-k}^*\|_{\Sigma_{S,-k}}^2,$$

286 where v denotes the standard deviation of the noise ϵ .

287 The first variance term arises from estimating the k major signal dimensions, corresponding to the
 288 classic variance in k -dimensional ordinary least squares. The second variance term, n/R_k , vanishes
 289 when the minor directions are sufficiently high-dimensional, i.e., when $R_k \gg n$. However, the
 290 signal in the minor directions, $\|\beta_{-k}^*\|_{\Sigma_{S,-k}}^2$, is nearly lost when projected from the high-dimensional
 291 ambient space onto the low-dimensional sample space, contributing to the second bias term. Finally,
 292 the first bias term relates to the signal estimation in the first k dimensions and is introduced by the
 293 overall regularization from both ridge and implicit regularization imposed by the minor components.

296 3.2 OUT-OF-DISTRIBUTION BENIGN OVERFITTING

297 We now investigate the out-of-distribution performance of the ridge estimator. Intuitively, when the
 298 minor components vanish for both the source and target distributions, over-parameterized ridge
 299 regression essentially reduces to under-parameterized ridge regression in the major directions, achiev-
 300 ing a rate of $\tilde{\mathcal{O}}(\text{tr}[\mathcal{T}]/n)$, as demonstrated by [Ge et al. \(2024\)](#). When the minor components do
 301 not vanish, a high effective rank of the minor components in the source distribution is essential for
 302 “benign overfitting”, as shown by [Tsigler & Bartlett \(2023\)](#). However, we argue that for the target
 303 distribution, only the *overall magnitude* of the minor components is critical for benign overfitting.
 304 This is because when the source’s minor directions have an effective rank much larger than n , the
 305 n -dimensional subspace spanned by the training samples is already almost orthogonal to any test
 306 point with high probability. As a result, the spectral structure of the target becomes irrelevant—only
 307 a small overall magnitude of the target’s minor components is required.

308 We formalize those intuitive claims by deriving upper bounds for both the variance and bias of ridge
 309 regression under covariate shift, assuming a source distribution similar to the in-distribution case.
 310 Our upper bound is sharp and can be applied to any target distributions, reducing to [Tsigler & Bartlett](#)
 311 ([2023](#))’s bound (Theorem 1) when the target and source distributions are aligned. Additionally, we
 312 recover [Ge et al. \(2024\)](#)’s sharp bound for under-parameterized linear regression under a covariate
 313 shift when the high-dimensional minor components vanish.

314 **Theorem 2.** There exists a constant $c > 2$ depending only on σ, L , such that for any $cN < n < r_k$,
 315 if the assumption condNum(k, δ, L) (Assumption 1) is satisfied, then with probability at least $1 - 3\delta$,

$$\begin{aligned} \frac{V}{cv^2} &\leq \frac{k}{n} \cdot \frac{\text{tr}[\mathcal{T}]}{k} + \frac{n}{R_k} \cdot \frac{\text{tr}[\mathcal{U}]}{\text{tr}[\mathcal{V}]}. \\ \frac{B}{c} &\leq B_{\text{ID}} \cdot \left(\|\mathcal{T}\| + \frac{n}{r_k} \frac{\|\Sigma_{T,-k}\|}{\|\Sigma_{S,-k}\|} \right). \end{aligned}$$

322 where $\mathcal{T}, \mathcal{U}, \mathcal{V}$ are defined in Equation (1), $N = \text{Poly}(k + \ln(1/\delta), \lambda_1 \lambda_k^{-1}, 1 + \tilde{\lambda} \lambda_k^{-1})$, and $\text{Poly}(\cdot)$
 323 denotes a polynomial function.

Recall that B_{ID} is the bias upper bound from Theorem 1. Theorem 2 establishes an upper bound for the excess risk of ridge regression under general covariate shift, expressed in a multiplicative form based on Theorem 1. This formulation enables a straightforward comparison of the impact of covariate shifts on the bias and variance of ridge estimators relative to the in-distribution case. The first conclusion is that Theorem 2 well reduces to the corresponding result in Theorem 1 when no distribution shift occurs—i.e., $\Sigma_S = \Sigma_T$. This connection follows directly from the condition $n < r_k$.

The second conclusion is that covariate shift in the first k dimensions and last $d - k$ dimensions introduce multiplicative factors of $\frac{\text{tr}[\mathcal{T}]}{k}$, $\|\mathcal{T}\|$ and $\frac{\text{tr}[\mathcal{U}]}{\text{tr}[\mathcal{V}]}$, $nr_k^{-1} \frac{\|\Sigma_{T,-k}\|}{\|\Sigma_{S,-k}\|}$, respectively, on the excess risk. Therefore, as long as these factors are bounded by constants, over-parameterized ridge regression achieves the same non-asymptotic rate of excess risk under covariate shift as the in-distribution setting. This scenario, well addressed by ridge regression, occurs when the target distribution’s covariance structure remains dominated by the first k dimensions. In the following, we analyze the impact of the factors introduced by covariate shifts on both the major and minor directions.

1. \mathcal{T} characterizes the shift in the major directions. Under covariate shift within the first k dimensions, we obtain the same non-asymptotic error rate as in Theorem 1, only if $\|\mathcal{T}\|$ is bounded by a constant, as $\text{tr}[\mathcal{T}]/k \leq \|\mathcal{T}\|$. The matrix \mathcal{T} plays a central role in Theorem 2 to quantify covariate shift within the first k dimensions, matching our intuition. This echoes with Ge et al. (2024)’s finding that $\text{tr}[\mathcal{T}]$ captures the difficulty of covariate shift for under-parameterized ridgeless regression (MLE). They establish a sharp upper bound on excess risk using Fisher information (see their Theorem 3.1), which simplifies to a rate of $\tilde{\mathcal{O}}(\text{tr}[\mathcal{T}]/n)$ for linear models. Theorem 2 recovers this result when applied to a k -dimensional under-parameterized setting where all high-dimensional minor components vanish, specifically when $\Sigma_{S,-k} = \Sigma_{T,-k} = \mathbf{0}$. Under the same condition as Theorem 2, for a constant c depending only on σ, L , with high probability, the variance and bias terms are bounded by:

$$\frac{V}{cv^2} \leq \frac{\text{tr}[\mathcal{T}]}{n}, \quad \frac{B}{c} \leq \|\beta_k^*\|_{\Sigma_{S,k}^{-1}}^2 \left(\frac{\lambda}{n} \right)^2 \|\mathcal{T}\|.$$

The variance bound aligns with Ge et al. (2024)’s result while the bias vanishes as $\lambda \rightarrow 0$.

2. The overall magnitude of $\Sigma_{T,-k}$ is crucial for benign overfitting. Under covariate shift within the last $d - k$ dimensions, when both $\frac{\text{tr}[\mathcal{U}]}{\text{tr}[\mathcal{V}]}$ and $nr_k^{-1} \frac{\|\Sigma_{T,-k}\|}{\|\Sigma_{S,-k}\|}$ are bounded by constants, we achieve the same non-asymptotic error rate as in Theorem 1. Note that $\frac{\text{tr}[\mathcal{U}]}{\text{tr}[\mathcal{V}]} \leq \frac{\|\Sigma_{T,-k}\|_{\text{F}}}{\|\Sigma_{S,-k}\|_{\text{F}}}$. In other words, matching our intuition, if the *overall magnitude* of the minor components of target covariance scales similarly to or smaller than those of the source, in terms of the covariance norms, “benign overfitting” also happens under covariate shift. Importantly, this condition does not impose constraints on the internal spectral structure of the minor components of the target covariance. For example, we do not force each eigenvalue of $\Sigma_{T,-k}$ to scale with its corresponding eigenvalue of $\Sigma_{S,-k}$ in decreasing order, as assumed in prior work (Mallinar et al., 2024). Surprisingly, for benign overfitting to happen, the source distribution must have a high effective rank in the minor directions. However, for the target distribution, only the overall magnitude of the minor components is relevant.

Another observation is that the bias scales with $nr_k^{-1} \frac{\|\Sigma_{T,-k}\|}{\|\Sigma_{S,-k}\|}$, meaning that we only require $\frac{\|\Sigma_{T,-k}\|}{\|\Sigma_{S,-k}\|} = \mathcal{O}(r_k/n)$, which is a less restrictive condition for larger r_k . Thus, over-parameterization improves the robustness of the estimation bias against covariate shift in the minor direction.

Remark 2 (Sample complexity). We have assumed $n > cN$ in Theorem 2. The explicit formula for N is deferred to Theorem 25 and Remark 8. Here we summarize the sample complexity required for the bound to hold. The dependence on k varies between $\Omega(k)$ and $\Omega(k^3)$, depending on the degree of covariate shift. The optimal case, aligning with the sample complexity of classic linear regression, occurs when $\Sigma_{S,k} \approx \Sigma_{T,k}$. The worst case arises when there is a significant covariate shift in the first k dimensions, such as when the test data lies predominantly in the subspace of the first dimension. This variation in sample complexity under covariate shift parallels the analysis of Ge et al. (2024) (see their Theorem 4.2) for the under-parameterized setting. Additionally, we require $n \gg \lambda + \sum_{j>k} \lambda_j$, ensuring that the regularization is not too strong to introduce a bias

378 exceeding a constant (as reflected in the first term of B_{ID}). On the other hand, we assume $n < r_k$ in
 379 the theorem, consistent with the over-parameterized regime and Assumption 1, where the last $d - k$
 380 components are considered to be essentially high-dimensional.

381 **Remark 3** (Dependence on L). Theorem 2 does not explicitly show how the excess risk depends
 382 on the condition number L of A_k . However, we demonstrate in Theorem 25 that the upper bounds
 383 scale at most as L^2 . Notably, we maintain the same order of dependence on L in each term of the
 384 upper bounds as in the analysis by Tsigler & Bartlett (2023) (see their Theorem 5).

385 Finally, Theorem 2 suggests an $\mathcal{O}(1/n)$ vanishing error under several conditions that naturally fol-
 386 low from the previous discussions, which we now state rigorously. First, the covariate space decom-
 387 poses into subspaces spanned by low-dimensional major directions and high-dimensional minor
 388 directions, with $k = \mathcal{O}(1)$ and $R_k = \Omega(n^2)$. Second, the low-rank covariance structure is preserved
 389 after covariate shift, such that $\|\mathcal{T}\|, \frac{\text{tr}[\mathcal{U}]}{\text{tr}[\mathcal{V}]}, nr_k^{-1} \frac{\|\Sigma_{T,-k}\|}{\|\Sigma_{S,-k}\|} = \mathcal{O}(1)$. Third, the signal lies predom-
 390 inantly in the major directions, with $\|\beta_k^*\|_{\Sigma_{S,k}^{-1}} = \mathcal{O}(1)$ and $\|\beta_{-k}^*\|_{\Sigma_{S,-k}} = \mathcal{O}(1/\sqrt{n})$. Lastly, the
 391 regularization is not excessively strong to introduce a significant bias, with $\tilde{\lambda} = \lambda + \sum_{j>k} \lambda_j =$
 392 $\mathcal{O}(\sqrt{n})$.

395 4 LARGE SHIFT IN MINOR DIRECTIONS

396 In the previous section, we established an upper bound for over-parameterized ridge regression under
 397 covariate shift. We showed that when the shift in the minor directions is controlled—specifically,
 398 when the overall magnitude of $\Sigma_{T,-k}$ is small—“benign overfitting” also occurs under covariate
 399 shift. However, when the shift in minor directions is significant, meaning the target covariance ma-
 400 trix has many large eigenvalues with corresponding eigenvectors outside the major directions, the
 401 excess risk for ridge regression deteriorates. In this section, we further illustrate the limitations of
 402 ridge regression in such cases by providing a lower bound for its performance for large distribu-
 403 tion shifts in the minor directions. We show that, in certain instances, ridge regression can only
 404 achieve the slow statistical rate of $\mathcal{O}(1/\sqrt{n})$ for the excess risk. On the other hand, it is natural
 405 to consider alternative algorithms to ridge regression. We demonstrate that even with a large shift
 406 in the minor directions, Principal Component Regression (PCR) is guaranteed to achieve the fast
 407 statistical rate $\mathcal{O}(1/n)$ in the same instances, provided that the signal β^* lies primarily within the
 408 subspace spanned by the major directions. Moreover, PCR does not require the minor directions to
 409 have a high effective rank in the source distribution, highlighting its advantage over ridge regres-
 410 sion in such cases. **Throughout this section, we maintain the setup and source covariance structure**
 411 **described in Section 2. However, Assumption 1 is no longer required.**

413 4.1 SLOW RATE FOR RIDGE REGRESSION

414 In this subsection, we demonstrate the limitations of ridge regression when the overall magnitude of
 415 $\Sigma_{T,-k}$ is large. Consider an instance where Σ_S has its first k components as $\Theta(1)$, while the minor
 416 directions have eigenvalues of $o(1)$. If we set $\Sigma_T = I_d$, in contrast to the “benign overfitting” regime
 417 described in Theorem 2, ridge regression will have a large excess risk for this instance. Although
 418 the signal from the major directions is effectively captured, the signal in the minor directions is
 419 nearly lost. Unlike the case in Section 3, here, the estimation error in the minor directions is crucial
 420 because the target distribution has significant components in these directions. We formalize this
 421 intuitive example through the following theorems:

422 **Corollary 3.** For some absolute constants C_1, C_2 , consider the following instance of Σ_S :

$$423 \quad \lambda_1 = \dots = \lambda_k = 1, \quad \lambda_{k+1} = \dots = \lambda_{k+\lfloor \frac{\sqrt{n}}{C_2} \rfloor} = \frac{C_1}{\sqrt{n}}, \quad \lambda_{k+\lfloor \frac{\sqrt{n}}{C_2} \rfloor+1} = \dots = \lambda_d = 0.$$

424 Assume $\Sigma_{T,-k} = \mathbf{0}$, $\Sigma_{T,k} = I_k$, and $\beta_{-k}^* = 0$. By choosing $\lambda = \sqrt{n}$, under the same conditions of
 425 Theorem 2, we can bound the excess risk of the ridge estimator with probability at least $1 - 3\delta$:

$$426 \quad \mathbb{E}_\epsilon [\mathcal{R}(\hat{\beta}(Y))] \leq \mathcal{O}\left(\frac{v^2 k + \|\beta^*\|^2}{n}\right).$$

427 **Remark 4.** Corollary 3 is a direct application of Theorem 2.

Theorem 4. Consider the same instance of Σ_S as in Corollary 3. Assume $\Sigma_T = I_d$ and $\lambda = \sqrt{n}$. There exists an absolute constant $C > 0$, such that for some $0 < \delta < 1$, $N_2 > 0$ and for any $n > N_2$, with probability at least $1 - \delta$, we have $V \geq Cv^2$.

Furthermore, for any $\lambda > 0$, we can bound the excess risk of the ridge estimator with probability at least $1 - \delta$:

$$\mathbb{E}_\epsilon [\mathcal{R}(\hat{\beta}(Y))] \geq C \frac{\|\beta^*\|^2 \wedge v^2}{\sqrt{n}}.$$

From Theorem 4, we observe that when $\Sigma_T = I_d$, the performance of ridge regression deteriorates compared to the case where $\Sigma_{T,-k} = 0$. If we set $\lambda = \sqrt{n}$ as in Corollary 3, ridge regression incurs a constant excess risk under covariate shift while achieving an in-distribution error rate of $\mathcal{O}(1/n)$. Furthermore, Theorem 4 shows no matter how we choose the regularization parameter λ , the excess risk is always lower bounded by the slow statistical rate $\mathcal{O}(1/\sqrt{n})$, which is worse than the fast rate of $\mathcal{O}(1/n)$. However, as we will prove in the next subsection, Principal Component Regression (PCR) can achieve an excess risk of $\mathcal{O}(1/n)$ under this instance, even with $\Sigma_T = I_d$.

4.2 FAST RATE FOR PRINCIPAL COMPONENT REGRESSION

As discussed earlier, ridge regression faces significant limitations when there is a large shift in the minor directions. In Section 3.1, it was shown that the signal in the minor directions, β_{-k}^* , is nearly lost when projected from the high-dimensional ambient space onto the low-dimensional sample space. In other words, learning the true signal from the minor directions is essentially impossible. Therefore, in this subsection, we continue to focus on the scenario where the true signal β^* primarily resides in the major directions. In this case, Principal Component Regression (PCR) emerges as a natural algorithm that estimates the space spanned by the major directions and performs regression on that subspace.

Principal Component Regression (PCR).

- **Step 1: Obtain an estimator \hat{U} of the top- k subspace of Σ_S .** For simplicity, we assume a sample size of $2n$ and use the first half of the data to compute \hat{U} by principal component analysis (PCA) on the sample covariance matrix $\hat{\Sigma}_S := \frac{1}{n} X^T X$. Specifically, $\hat{U} = (\hat{u}_1, \dots, \hat{u}_k)$ where \hat{u}_i is the i -th eigenvector of $\hat{\Sigma}_S$.

- **Step 2: Use the data projected on \hat{U} to conduct linear regression.** With a little abuse of notation, we use $X \in \mathbb{R}^{n \times d}$ to denote the data matrix $(x_{n+1}, \dots, x_{2n})^T$, and $Y \in \mathbb{R}^n$ to denote $(y_{n+1}, \dots, y_{2n})^T$. If we let $Z := X\hat{U} \in \mathbb{R}^{n \times k}$ be the projected data matrix, the estimator $\hat{\beta}$ we obtained is given by

$$\hat{\beta} = \hat{U}(Z^T Z)^{-1} Z^T Y = \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T X^T Y.$$

Consider the scenario where the last $d - k$ components of the true signal β^* is exactly zero, i.e., $\beta_{-k}^* = 0$. In this case, if the subspace represented by \hat{U} perfectly matches the subspace represented by $U = \begin{pmatrix} I_k \\ 0 \end{pmatrix} \in \mathbb{R}^{d \times k}$, corresponding to the first k components, then PCR performs linear regression using only the first k components of the covariates. As a result, the excess risk would just be the usual variance of linear regression in the major directions. In this scenario, regardless of the norm $\|\Sigma_{T,-k}\|$, the PCR estimator assigns coefficients of zero to the last $d - k$ components, thus avoiding any large excess risk. Furthermore, if the distance between \hat{U} and U is nonzero, an additional term in the excess risk will arise due to the estimation error of \hat{U} . We formalize this intuition with the following upper bound for the excess risk of PCR. To facilitate the presentation, we introduce a measure of the estimation accuracy of \hat{U} . We define $\Delta = \text{dist}(\hat{U}, U) := \|UU^T - \hat{U}\hat{U}^T\|$, which represents the distance between the subspaces spanned by the columns of \hat{U} and U . Then we present the following theorem.

Theorem 5. Assume $\beta_{-k}^* = 0$. If $\Delta \leq \Theta$, for any $0 < \delta < 1$ and any $n \geq N_1$, we can bound the excess risk of PCR estimator $\widehat{\beta}$ with probability $1 - \delta$:

$$\mathbb{E}_\epsilon [\mathcal{R}(\widehat{\beta}(Y))] \leq \mathcal{O}\left(v^2 \frac{\text{tr}(\mathcal{T})}{n} + \|\beta^*\|^2 \left(\frac{\lambda_1}{\lambda_k}\right)^2 \|\Sigma_T\| \Delta^2\right),$$

where $\Theta^{-1} = \text{Poly}(\lambda_1 \lambda_k^{-1}, \|\Sigma_T\| \lambda_k^{-1}, k \text{tr}(\mathcal{T})^{-1})$ and $N_1 = \text{Poly}(\sigma, \lambda_1 \lambda_k^{-1}, \|\Sigma_T\| \lambda_k^{-1}, k \ln(1/\delta), k \text{tr}(\mathcal{T})^{-1})$.

Remark 5. Theorem 5 is a special case of Lemma 31. For explicit formulas of Θ and N_1 , as well as an upper bound for cases where $\beta_{-k}^* \neq 0$, refer to Lemma 31 for details.

The excess risk upper bound provided by Theorem 5 consists of two terms. The variance term $\text{tr}(\mathcal{T})/n$ is incurred by the nature of linear regression on the major directions and remains unavoidable even when the subspace estimation is exact (i.e., $\Delta = 0$). This term also appears as the first variance term in Theorem 2, and exactly matches the sharp rate $\text{tr}[\Sigma_S^{-1} \Sigma_T]/n$ for under-parameterized linear regression under covariate shift (Ge et al., 2024). The second term $\|\beta^*\|^2 (\frac{\lambda_1}{\lambda_k})^2 \|\Sigma_T\| \Delta^2$ represents the bias induced by the subspace estimation error in Step 1, which exhibits a quadratic dependence on Δ . By combining Theorem 5 with a bound on Δ , we can derive an end-to-end excess risk upper bound of PCR. We present the following lemma to control Δ .

Lemma 6. With probability at least $1 - \delta$, if $n \geq r + \ln(1/\delta)$, we have

$$\Delta \leq \mathcal{O}\left(\sigma^4 \frac{\lambda_1}{\lambda_k - \lambda_{k+1}} \sqrt{\frac{r + \ln \frac{1}{\delta}}{n}}\right),$$

where $r = \lambda_1^{-1} \sum_{i=1}^d \lambda_i$ is the effective rank of the entire covariance matrix Σ_S .

Remark 6. Lemma 6 shows that Δ depends on several quantities: the eigenvalue gap between the major and minor directions, i.e., $\lambda_k - \lambda_{k+1}$, and the effective rank r . We observe that Δ will be small if the major and minor directions are well separated, meaning $\lambda_k - \lambda_{k+1}$ is large, and the minor directions are relatively small compared to λ_1 .

Combining Theorem 5 with Lemma 6, an end-to-end error bound for PCR directly follows (for a detailed statement, refer to Theorem 29), suggesting that PCR will achieve a small excess risk as long as the major and minor directions are well separated, and the effective rank of the entire source covariance matrix is small. In contrast to ridge regression, PCR does not rely on the minor components having a high effective rank. This highlights the superiority of PCR over ridge regression in certain scenarios.

As an example, consider the instance in Theorem 4, where $k, \|\Sigma_T\|, \lambda_1, \lambda_k$ are all $\Theta(1)$. In this case, the variance term scales as $1/n$, and the bias term scales as $\mathcal{O}(\Delta^2)$. Since $r = \Theta(1)$ in this instance, we have $\Delta \leq \mathcal{O}(1/\sqrt{n})$. Consequently, we conclude that PCR achieves a $\mathcal{O}(1/n)$ rate in this instance, even when $\Sigma_T = I_d$. Compared with the excess risk for ridge regression, which is at least $1/\sqrt{n}$, PCR shows its superiority against ridge regression when there is a large shift in the minor directions.

5 CONCLUSION AND DISCUSSION

In conclusion, we provide an instance-dependent upper bound on the excess risk for ridge regression under general covariate shift. Our findings demonstrate that “benign overfitting” also occurs in OOD generalization when the shift in the minor directions is well controlled. We also investigate the regime with a large shift in the minor directions, where ridge regression may incur a large excess risk, whereas Principal Component Regression (PCR) exhibits superior performance.

Our work opens several directions for future research. First, while we have established a lower bound for ridge regression in certain instances, a key challenge remains in deriving a general lower bound that matches our upper bounds, offering a more precise characterization of the excess risk under covariate shift. Second, our analysis has focused on linear models as an initial step in understanding over-parameterized OOD problems. Extending this investigation to more complex, non-linear models would be a intriguing direction for future exploration.

540 REFERENCES
541

- 542 Sergios Agapiou, Omiros Papaspiliopoulos, Daniel Sanz-Alonso, and Andrew M Stuart. Importance
543 sampling: Intrinsic dimension and computational cost. *Statistical Science*, pp. 405–431, 2017.
- 544 Anish Agarwal, Devavrat Shah, Dennis Shen, and Dogyoon Song. On robustness of principal com-
545 ponent regression. *Advances in Neural Information Processing Systems*, 32, 2019.
- 546 Anish Agarwal, Devavrat Shah, and Dennis Shen. On model identification and out-of-sample pre-
547 diction of principal component regression: Applications to synthetic controls. *arXiv preprint arXiv:2010.14449*, 2020.
- 548 Jushan Bai and Serena Ng. Determining the number of factors in approximate factor models. *Econo-*
549 *metrica*, 70(1):191–221, 2002.
- 550 Eric Bair, Trevor Hastie, Debmashis Paul, and Robert Tibshirani. Prediction by supervised principal
551 components. *Journal of the American Statistical Association*, 101(473):119–137, 2006.
- 552 Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear
553 regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- 554 Peter L Bartlett, Andrea Montanari, and Alexander Rakhlin. Deep learning: a statistical viewpoint.
555 *Acta numerica*, 30:87–201, 2021.
- 556 Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps for dimensionality reduction and data
557 representation. *Neural computation*, 15(6):1373–1396, 2003.
- 558 Mikhail Belkin, Daniel Hsu, and Ji Xu. Two models of double descent for weak features. *SIAM
559 Journal on Mathematics of Data Science*, 2(4):1167–1180, 2020.
- 560 Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of representations
561 for domain adaptation. *Advances in neural information processing systems*, 19, 2006.
- 562 Koby Bibas, Yaniv Fogel, and Meir Feder. A new look at an old problem: A universal learning
563 approach to linear regression. In *2019 IEEE International Symposium on Information Theory
564 (ISIT)*, pp. 2304–2308. IEEE, 2019.
- 565 Swee Lean Chan and Moonseo Park. Project cost estimation using principal component regression.
566 *Construction Management and Economics*, 23(3):295–304, 2005.
- 567 Yihang Chen, Fanghui Liu, Taiji Suzuki, and Volkan Cevher. High-dimensional kernel methods
568 under covariate shift: Data-dependent implicit regularization. In *Forty-first International Con-
569 ference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net,
570 2024. URL <https://openreview.net/forum?id=bBzlapzeR1>.
- 571 Yuxin Chen, Yuejie Chi, Jianqing Fan, Cong Ma, et al. Spectral methods for data science: A
572 statistical perspective. *Foundations and Trends® in Machine Learning*, 14(5):566–806, 2021.
- 573 Geoffrey Chinot and Matthieu Lerasle. On the robustness of the minimim l2 interpolator. *Bernoulli*,
574 2022.
- 575 Corinna Cortes, Yishay Mansour, and Mehryar Mohri. Learning bounds for importance weighting.
576 *Advances in neural information processing systems*, 23, 2010.
- 577 Uwe Depczynski, VJ Frost, and K Molt. Genetic algorithms applied to the selection of factors in
578 principal component regression. *Analytica Chimica Acta*, 420(2):217–227, 2000.
- 579 Edgar Dobriban and Stefan Wager. High-dimensional asymptotics of prediction: Ridge regression
580 and classification. *The Annals of Statistics*, 46(1):247–279, 2018.
- 581 Jianqing Fan and Yihong Gu. Factor augmented sparse throughput deep relu neural networks for
582 high dimensional regression. *Journal of the American Statistical Association*, pp. 1–15, 2023.
- 583 Jianqing Fan, Yuan Ke, and Kaizheng Wang. Factor-adjusted regularized model selection. *Journal
584 of Econometrics*, 216(1):71–85, 2020.

- 594 Jianqing Fan, Kaizheng Wang, Yiqiao Zhong, and Ziwei Zhu. Robust high dimensional factor
 595 models with applications to statistical machine learning. *Statistical science: a review journal of*
 596 *the Institute of Mathematical Statistics*, 36(2):303, 2021.
- 597
- 598 Jiawei Ge, Shange Tang, Jianqing Fan, and Chi Jin. On the provable advantage of unsupervised
 599 pretraining. *arXiv preprint arXiv:2303.01566*, 2023.
- 600 Jiawei Ge, Shange Tang, Jianqing Fan, Cong Ma, and Chi Jin. Maximum likelihood estimation is all
 601 you need for well-specified covariate shift. In *The Twelfth International Conference on Learning*
 602 *Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL
 603 <https://openreview.net/forum?id=eoTCKK0gIs>.
- 604
- 605 Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. When do neural
 606 networks outperform kernel methods? *Advances in Neural Information Processing Systems*, 33:
 607 14820–14830, 2020.
- 608 Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Linearized two-layers
 609 neural networks in high dimension. *The Annals of Statistics*, 49(2), 2021.
- 610
- 611 Hao Guan and Mingxia Liu. Domain adaptation for medical image analysis: a survey. *IEEE Trans-*
 612 *actions on Biomedical Engineering*, 69(3):1173–1185, 2021.
- 613
- 614 Ali S Hadi and Robert F Ling. Some cautionary notes on the use of principal components regression.
 615 *The American Statistician*, 52(1):15–19, 1998.
- 616 Steve Hanneke and Samory Kpotufe. On the value of target data in transfer learning. *Advances in*
 617 *Neural Information Processing Systems*, 32, 2019.
- 618
- 619 Yifan Hao, Yong Lin, Difan Zou, and Tong Zhang. On the benefits of over-parameterization for out-
 620 of-distribution generalization. *CoRR*, abs/2403.17592, 2024. doi: 10.48550/ARXIV.2403.17592.
 621 URL <https://doi.org/10.48550/arXiv.2403.17592>.
- 622 Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-
 623 dimensional ridgeless least squares interpolation. *Annals of statistics*, 50(2):949, 2022.
- 624
- 625 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common cor-
 626 ruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019.
- 627
- 628 Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul
 629 Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A criti-
 630 cal analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF international*
 631 *conference on computer vision*, pp. 8340–8349, 2021.
- 632
- Hugo G Hidalgo, Thomas C Piechota, and John A Dracup. Alternative principal components re-
 633 gression procedures for dendrohydrologic reconstructions. *Water Resources Research*, 36(11):
 634 3241–3249, 2000.
- 635
- Shih-Ming Huang and Jar-Ferr Yang. Improved principal component regression for face recognition
 636 under illumination variations. *IEEE signal processing letters*, 19(4):179–182, 2012.
- 637
- JNR Jeffers. Investigation of alternative regressions: Some practical examples. *Journal of the Royal*
 639 *Statistical Society. Series D (The Statistician)*, 30(2):79–88, 1981.
- 640
- John NR Jeffers. Two case studies in the application of principal component analysis. *Journal of*
 641 *the Royal Statistical Society: Series C (Applied Statistics)*, 16(3):225–236, 1967.
- 643
- Ian T Jolliffe. A note on the use of principal components in regression. *Journal of the Royal*
 644 *Statistical Society Series C: Applied Statistics*, 31(3):300–303, 1982.
- 645
- Chinmaya Kausik, Kashvi Srivastava, and Rishi Sonthalia. Double descent and overfitting under
 646 noisy inputs and distribution shift for linear denoisers. *Trans. Mach. Learn. Res.*, 2024, 2024.
 647 URL <https://openreview.net/forum?id=HxfqTdLIRF>.

- 648 Richard B Keithley, R Mark Wightman, and Michael L Heien. Multivariate concentration determi-
 649 nation using principal component regression with residual analysis. *TrAC Trends in Analytical
 650 Chemistry*, 28(9):1127–1136, 2009.
- 651
 652 Dmitry Kobak, Jonathan Lomond, and Benoit Sanchez. The optimal ridge penalty for real-world
 653 high-dimensional data can be zero or negative due to the implicit ridge regularization. *Journal of
 654 Machine Learning Research*, 21(169):1–16, 2020.
- 655 Frederic Koehler, Lijia Zhou, Danica J Sutherland, and Nathan Srebro. Uniform convergence of
 656 interpolators: Gaussian width, norm bounds and benign overfitting. *Advances in Neural Informa-
 657 tion Processing Systems*, 34:20657–20668, 2021.
- 658
 659 Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Bal-
 660 subramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Irena Gao, et al. Wilds: A
 661 benchmark of in-the-wild distribution shifts. In *International conference on machine learning*, pp.
 662 5637–5664. PMLR, 2021.
- 663 Samory Kpotufe and Guillaume Martinet. Marginal singularity, and the benefits of labels in
 664 covariate-shift. In *Conference On Learning Theory*, pp. 1882–1886. PMLR, 2018.
- 665 Anikender Kumar and Pramila Goyal. Forecasting of air quality in delhi using principal component
 666 regression technique. *Atmospheric Pollution Research*, 2(4):436–444, 2011.
- 667
 668 Qi Lei, Wei Hu, and Jason Lee. Near-optimal linear regression under distribution shift. In *Interna-
 669 tional Conference on Machine Learning*, pp. 6164–6174. PMLR, 2021.
- 670 Tengyuan Liang, Alexander Rakhlin, and Xiyu Zhai. On the multiple descent of minimum-norm
 671 interpolants and restricted lower isometry of kernels. In *Conference on Learning Theory*, pp.
 672 2683–2711. PMLR, 2020.
- 673
 674 RX Liu, J Kuang, Qiong Gong, and XL Hou. Principal component regression analysis with spss.
 675 *Computer methods and programs in biomedicine*, 71(2):141–147, 2003.
- 676 Cong Ma, Reese Pathak, and Martin J Wainwright. Optimally tackling covariate shift in rkhs-based
 677 nonparametric regression. *The Annals of Statistics*, 51(2):738–761, 2023.
- 678
 679 Neil Mallinar, Austin Zane, Spencer Frei, and Bin Yu. Minimum-norm interpolation under covariate
 680 shift. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria,
 681 July 21-27, 2024*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=Zw7TcnTmHj>.
- 682
 683 William F Massy. Principal components regression in exploratory statistical research. *Journal of
 684 the American Statistical Association*, 60(309):234–256, 1965.
- 685
 686 Song Mei and Andrea Montanari. The generalization error of random features regression: Precise
 687 asymptotics and the double descent curve. *Communications on Pure and Applied Mathematics*,
 688 75(4):667–766, 2022.
- 689
 690 Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Generalization error of random feature
 691 and kernel methods: hypercontractivity and kernel matrix concentration. *Applied and Compu-
 692 tational Harmonic Analysis*, 59:3–84, 2022.
- 693 John P Miller, Rohan Taori, Aditi Raghunathan, Shiori Sagawa, Pang Wei Koh, Vaishaal Shankar,
 694 Percy Liang, Yair Carmon, and Ludwig Schmidt. Accuracy on the line: on the strong correlation
 695 between out-of-distribution and in-distribution generalization. In *International conference on
 696 machine learning*, pp. 7721–7735. PMLR, 2021.
- 697
 698 Andrea Montanari and Yiqiao Zhong. The interpolation phase transition in neural networks: Memo-
 699 rization and generalization under lazy training. *The Annals of Statistics*, 50(5):2816–2847, 2022.
- 700 Mohammadreza Mousavi Kalan, Zalan Fabian, Salman Avestimehr, and Mahdi Soltanolkotabi. Min-
 701 imax lower bounds for transfer learning with linear and one-hidden layer neural networks. *Ad-
 702 vances in Neural Information Processing Systems*, 33:1959–1969, 2020.

- 702 Vidya Muthukumar, Kailas Vodrahalli, Vignesh Subramanian, and Anant Sahai. Harmless interpolation
 703 of noisy data in regression. *IEEE Journal on Selected Areas in Information Theory*, 1(1):
 704 67–83, 2020.
- 705 Tormod Næs and Harald Martens. Principal component regression in nir analysis: viewpoints,
 706 background details and selection of components. *Journal of chemometrics*, 2(2):155–167, 1988.
- 708 Preetum Nakkiran. More data can hurt for linear regression: Sample-wise double descent. *arXiv
 709 preprint arXiv:1912.07242*, 2019.
- 710 Jeffrey Negrea, Gintare Karolina Dziugaite, and Daniel Roy. In defense of uniform convergence:
 711 Generalization via derandomization with an application to interpolating predictors. In *International
 712 Conference on Machine Learning*, pp. 7263–7272. PMLR, 2020.
- 714 Partha Niyogi. Manifold regularization and semi-supervised learning: Some theoretical analyses.
 715 *Journal of Machine Learning Research*, 14(5), 2013.
- 717 Reese Pathak, Cong Ma, and Martin Wainwright. A new similarity measure for covariate shift with
 718 applications to nonparametric regression. In *International Conference on Machine Learning*, pp.
 719 17517–17530. PMLR, 2022.
- 720 Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers
 721 generalize to imagenet? In *International conference on machine learning*, pp. 5389–5400. PMLR,
 722 2019.
- 723 Dominic Richards, Jaouad Mourtada, and Lorenzo Rosasco. Asymptotics of ridge (less) regres-
 724 sion under general source condition. In *International Conference on Artificial Intelligence and
 725 Statistics*, pp. 3889–3897. PMLR, 2021.
- 727 Ohad Shamir. The implicit bias of benign overfitting. *Journal of Machine Learning Research*, 24
 728 (113):1–40, 2023.
- 729 Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-
 730 likelihood function. *Journal of statistical planning and inference*, 90(2):227–244, 2000.
- 732 James B Simon, Dhruva Karkada, Nikhil Ghosh, and Mikhail Belkin. More is better in modern
 733 machine learning: when infinite overparameterization is optimal and overfitting is obligatory.
 734 *arXiv preprint arXiv:2311.14646*, 2023.
- 735 James H Stock and Mark W Watson. Forecasting using principal components from a large number
 736 of predictors. *Journal of the American statistical association*, 97(460):1167–1179, 2002.
- 738 Jianguo Sun. A correlation principal component regression analysis of nir data. *Journal of Chemo-
 739 metrics*, 9(1):21–29, 1995.
- 740 Nilesh Tripuraneni, Ben Adlam, and Jeffrey Pennington. Overparameterization improves ro-
 741 bustness to covariate shift in high dimensions. In Marc’Aurelio Ranzato, Alina Beygelz-
 742 imer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances
 743 in Neural Information Processing Systems 34: Annual Conference on Neural Infor-
 744 mation Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 13883–
 745 13897, 2021a. URL <https://proceedings.neurips.cc/paper/2021/hash/73fed7fd472e502d8908794430511f4d-Abstract.html>.
- 747 Nilesh Tripuraneni, Chi Jin, and Michael Jordan. Provable meta-learning of linear representations.
 748 In *International Conference on Machine Learning*, pp. 10434–10443. PMLR, 2021b.
- 750 Alexander Tsigler and Peter L. Bartlett. Benign overfitting in ridge regression. *J. Mach. Learn. Res.*,
 751 24:123:1–123:76, 2023. URL <http://jmlr.org/papers/v24/22-1398.html>.
- 752 J Leo van Hemmen and Tsuneya Ando. An inequality for trace ideals. *Communications in Mathe-
 753 matical Physics*, 76:143–148, 1980.
- 755 Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *CoRR*,
 abs/1011.3027, 2010. URL <http://arxiv.org/abs/1011.3027>.

- 756 Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*,
 757 volume 47. Cambridge university press, 2018.
 758
- 759 Evelyne Vigneau, MF Devaux, EM Qannari, and P Robert. Principal component regression, ridge re-
 760 gression and ridge principal component regression in spectroscopy calibration. *Journal of Chemo-*
 761 *metrics: A Journal of the Chemometrics Society*, 11(3):239–249, 1997.
- 762 Per-Ake Wedin. Perturbation theory for pseudo-inverses. *BIT Numerical Mathematics*, 13:217–232,
 763 1973.
- 764 Florian Wenzel, Andrea Dittadi, Peter Gehler, Carl-Johann Simon-Gabriel, Max Horn, Dominik
 765 Zietlow, David Kernert, Chris Russell, Thomas Brox, Bernt Schiele, et al. Assaying out-of-
 766 distribution generalization in transfer learning. *Advances in Neural Information Processing Sys-*
 767 *tems*, 35:7181–7198, 2022.
- 768
- 769 Denny Wu and Ji Xu. On the optimal weighted l2 regularization in overparameterized linear regres-
 770 sion. *Advances in Neural Information Processing Systems*, 33:10112–10123, 2020.
- 771
- 772 Ji Xu and Daniel J Hsu. On the number of variables to use in principal component regression.
 773 *Advances in neural information processing systems*, 32, 2019.
- 774 Runtian Zhai, Chen Dan, Zico Kolter, and Pradeep Ravikumar. Understanding why generalized
 775 reweighting does not improve over erm. *arXiv preprint arXiv:2201.12293*, 2022.
- 776
- 777 Xuhui Zhang, Jose Blanchet, Soumyadip Ghosh, and Mark S Squillante. A class of geometric
 778 structures in transfer learning: Minimax bounds and optimality. In *International Conference on*
 779 *Artificial Intelligence and Statistics*, pp. 3794–3820. PMLR, 2022.
- 780 Lijia Zhou, Danica J Sutherland, and Nati Srebro. On uniform convergence and low-norm interpo-
 781 lation learning. *Advances in Neural Information Processing Systems*, 33:6867–6877, 2020.
- 782
- 783 Yang Zou, Zhiding Yu, BVK Kumar, and Jinsong Wang. Unsupervised domain adaptation for se-
 784 mantic segmentation via class-balanced self-training. In *Proceedings of the European conference*
 785 *on computer vision (ECCV)*, pp. 289–305, 2018.
- 786
- 787
- 788
- 789
- 790
- 791
- 792
- 793
- 794
- 795
- 796
- 797
- 798
- 799
- 800
- 801
- 802
- 803
- 804
- 805
- 806
- 807
- 808
- 809

A SIMULATION

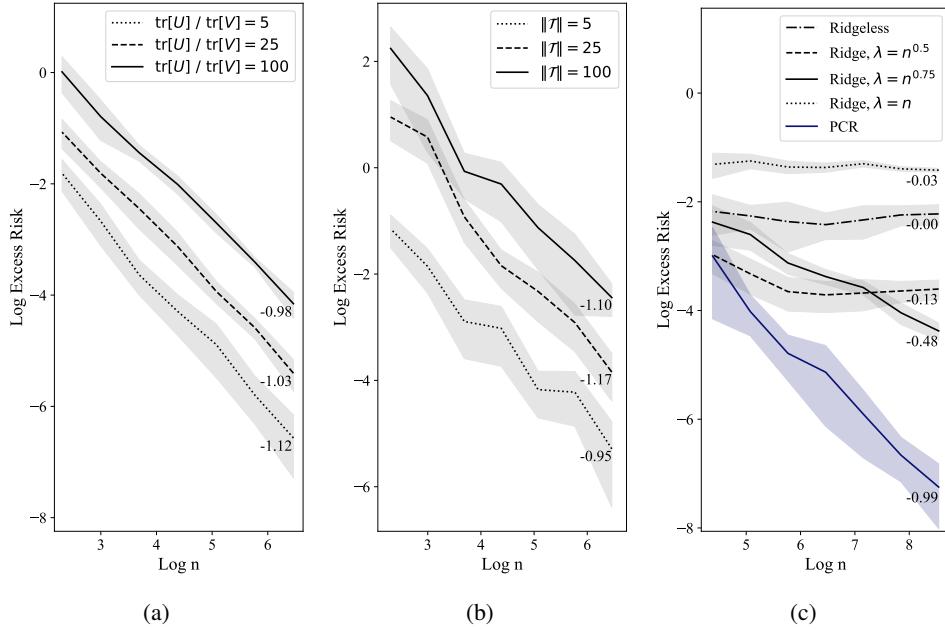


Figure 1: Simulation results for excess risks across varying training sample sizes. The shaded regions represent standard errors of 10 runs, using different samples of training and test sets. The slope of the fitted OLS model is marked along each curve. (a)(b) Minimum norm interpolation under distinct target covariance matrices with small shifts in minor directions. The source covariance matrix remains constant. (a) Various magnitudes of shifts in minor directions, with $\|\mathcal{T}\| = 1$. (b) Various magnitudes of shifts in major directions, with $\text{tr}[\mathcal{U}] / \text{tr}[\mathcal{V}] = 1$. (c) Ridge and PCR under large shifts in minor directions, following the setting of Theorem 4.

A.1 BENIGN-OVERFITTING: SMALL SHIFT IN MINOR DIRECTIONS

We simulate the covariate shift discussed in section 3, where the overall magnitude of the target covariance matrix's minor directions is comparable to that of the source. Theorem 2 establishes an $\mathcal{O}(1/n)$ excess risk rate for ridge regression under certain benign-overfitting conditions. Specifically, for the training data, we assume $k = \mathcal{O}(1)$, $R_k = \Omega(n^2)$, $\|\beta_k^*\|_{\Sigma_{S,k}^{-1}} = \mathcal{O}(1)$, $\|\beta_{-k}^*\|_{\Sigma_{S,-k}} = \mathcal{O}(1/\sqrt{n})$, and $\tilde{\lambda} = \mathcal{O}(\sqrt{n})$. For the test data, we assume $\|\mathcal{T}\|, \frac{\text{tr}[\mathcal{U}]}{\text{tr}[\mathcal{V}]}, nr_k^{-1} \frac{\|\Sigma_{T,-k}\|}{\|\Sigma_{S,-k}\|} = \mathcal{O}(1)$. In the experiment, data is generated according to these conditions.

$$y = x^T \beta^* + \epsilon,$$

where $\beta^* \in \mathbb{R}^{k+n^2}$, with $\beta_k^* = (\frac{1}{\sqrt{k}}, \dots, \frac{1}{\sqrt{k}})^T$, $\beta_{-k}^* = \mathbf{0}$ and $k = 10$. The noise ϵ follows a centered gaussian distribution with variance 0.1, and x is drawn from a multivariate normal distribution with zero mean and a source covariance matrix $\Sigma_S = \text{diag}(I_k, n^{-1.5} I_{n^2})$. The target covariance matrix is $\Sigma_T = \text{diag}(\Sigma_{T,k}, \Sigma_{T,-k})$ where $\Sigma_{T,k}$ and $\Sigma_{T,-k}$ are randomly generated and scaled to achieve specific values for $\|\mathcal{T}\|$ and $\frac{\text{tr}[\mathcal{U}]}{\text{tr}[\mathcal{V}]}$, respectively. At the same time, we do not explicitly control $nr_k^{-1} \frac{\|\Sigma_{T,-k}\|}{\|\Sigma_{S,-k}\|}$, because it equals $\frac{1}{n} \frac{\|\Sigma_{T,-k}\|}{\|\Sigma_{S,-k}\|}$ in this setting and is typically bounded for a randomly generated $\Sigma_{T,-k}$. We run minimum norm interpolation (ridgeless regression) with $\lambda = 0$.

The source covariance matrix Σ_S is fixed for all experiments while we vary the target covariance matrices. To study covariate shifts in major directions, we vary $\|\mathcal{T}\|$ among 5, 25, 100 and keep $\frac{\text{tr}[\mathcal{U}]}{\text{tr}[\mathcal{V}]} = 1$. To study covariate shifts in minor directions, we vary $\frac{\text{tr}[\mathcal{U}]}{\text{tr}[\mathcal{V}]} = 1$ among 5, 25, 100 and keep $\|\mathcal{T}\| = 1$. For each pair of $\|\mathcal{T}\|$ and $\frac{\text{tr}[\mathcal{U}]}{\text{tr}[\mathcal{V}]}$, we generate training samples of various sizes n and

864 1000 test samples. For each n , a target covariance matrix Σ_T is randomly generated to satisfy the
 865 specified $\|\mathcal{T}\|$ and $\frac{\text{tr}[\mathcal{U}]}{\text{tr}[\mathcal{V}]}$. We run 10 experiments for each set of (Σ_S, Σ_T, n) , using independently
 866 sampled training sets and test sets, and the mean and standard error of the excess risks are reported.
 867

868 The results are shown in Figure 1a, 1b. The fast rate $\mathcal{O}(1)$ of minimum norm interpolation is
 869 confirmed, as the log-log plot of excess risk versus n has a slope near -1 across all combinations
 870 of $\|\mathcal{T}\|$ and $\frac{\text{tr}[\mathcal{U}]}{\text{tr}[\mathcal{V}]}$. The excess risk increases with larger $\frac{\text{tr}[\mathcal{U}]}{\text{tr}[\mathcal{V}]}$, indicating a greater shift in minor
 871 directions. Similarly, the excess risk increases with larger $\|\mathcal{T}\|$, indicating a greater shift in major
 872 directions.

873

874 A.2 RIDGE v.s. PCR: LARGE SHIFT IN MINOR DIRECTIONS

875

876 Theorem 4 identifies a setting where large covariate shifts occur in minor directions of the covariance
 877 matrix, leading to a lower bound of $\mathcal{O}(1/\sqrt{n})$ on the excess risk for ridge regression, while Principal
 878 Component Regression (PCR) achieves the fast rate of $\mathcal{O}(1)$. We design a simulation experiment
 879 under the same instance of the source and target covariance matrices. Specifically, data is generated
 880 as:

881

$$y = x^T \beta^* + \epsilon,$$

882

883 where $\beta^* \in \mathbb{R}^{k+\lfloor \sqrt{n} \rfloor}$, with $\beta_k^* = (\frac{1}{\sqrt{k}}, \dots, \frac{1}{\sqrt{k}})^T$, $\beta_{-k}^* = \mathbf{0}$, and $k = 10$. The noise ϵ is drawn
 884 from a centered gaussian distribution with variance 0.1, and x follows a multivariate normal distri-
 885 bution with zero mean. The source covariance matrix is $\Sigma_S = \text{diag}(I_k, n^{-0.5} I_{\lfloor \sqrt{n} \rfloor})$, and the target
 886 covariance matrix is $\Sigma_T = I_{k+\lfloor \sqrt{n} \rfloor}$.

887

888 We evaluate ridge regression with various regularization strengths: $\lambda = 10^{-8}, n^{0.5}, n^{0.75}, n$. Here,
 889 we use $\lambda = 10^{-8}$ to approximate ridgeless regression, which has a singular solution under this setup.
 890 We compare PCR to ridge regression for different training sample sizes n . The test set contains 1000
 891 samples. For each n , 10 experiments are conducted with independently sampled training and test
 892 sets. We report the mean and standard error of the excess risks.

893

894 Figure 1c present the results. As expected, PCR nearly achieves the fast rate of $\mathcal{O}(1/n)$, with the
 895 log-log slope of excess risk versus n being -0.99. In contrast, the optimal rate of ridge regression
 896 is $\mathcal{O}(n^{-0.48})$, achieved with $\lambda = n^{0.75}$. This aligns with the lower bound of $\mathcal{O}(1/\sqrt{n})$ from The-
 897 orems 4, and its proof also suggests $\lambda = n^{0.75}$ as the optimal regularization strength. Additionally,
 898 ridge regression exhibits excess risks above a constant for certain choices of λ .

899

900 B RIDGE REGRESSION

901

902 Let $X = (x_1, \dots, x_n)^T \in \mathbb{R}^{n \times d}$, $Y = (y_1, \dots, y_n)^T \in \mathbb{R}^n$ and $\epsilon = (\epsilon_1, \dots, \epsilon_n)^T \in \mathbb{R}^n$. We denote
 903 the first k columns of X as X_k and the remaining $d - k$ columns as X_{-k} . Similarly, β_k^* and β_{-k}^*
 904 represent the corresponding components of β^* . $\Sigma_{S,k}$, $\Sigma_{S,-k}$ are the corresponding blocks on the
 905 diagonal of Σ_S . The i -th eigenvalue of a matrix is denoted by $\mu_i(\cdot)$. Define $Z = X \Sigma_S^{-1/2}$, where
 906 the rows of Z are i.i.d. centered isotropic random vectors. Additionally, we assume the rows of Z
 907 are σ -sub-gaussian, where the sub-gaussian norm is defined as follows.

908

909 For a random variable s , the sub-gaussian norm $\|s\|_{\psi_2}$ is given by:
 910

$$\|s\|_{\psi_2} = \inf \left\{ t > 0 : \mathbb{E} \left[\exp \frac{s^2}{t^2} \right] \leq 2 \right\}.$$

911

912 For a random vector S , the sub-gaussian norm $\|S\|_{\psi_2}$ is given by:
 913

$$\|S\|_{\psi_2} = \sup_{v \neq 0} \frac{\|\langle S, v \rangle\|_{\psi_2}}{\|v\|}.$$

914

915 For $\lambda \geq 0$, consider the ridge estimator:
 916

$$\hat{\beta}(Y) = X^T (X X^T + \lambda I_n)^{-1} Y$$

$$\begin{aligned}
&= X^T(XX^T + \lambda I_n)^{-1}X\beta^* + X^T(XX^T + \lambda I_n)^{-1}\epsilon \\
&= \widehat{\beta}(X\beta^*) + \widehat{\beta}(\epsilon),
\end{aligned}$$

where we define $\widehat{\beta}(X\beta^*) = X^T(XX^T + \lambda I_n)^{-1}X\beta^*$ and $\widehat{\beta}(\epsilon) = X^T(XX^T + \lambda I_n)^{-1}\epsilon$. Additionally, we define $\widehat{\Sigma}_S = \Sigma_S + \frac{\lambda}{n}I_d$. The effective rank of $\widehat{\Sigma}_{S,k}$ is defined as $r_k = \lambda_{k+1}^{-1}(\lambda + \sum_{j>k} \lambda_j)$.

Assumption 2 (CondNum(k, δ, L)). Define a matrix $A_k = \lambda I_n + X_{-k}X_{-k}^T$. With probability at least $1 - \delta$, A_k is positive definite and has a condition number no greater than L , i.e.,

$$\frac{\mu_1(A_k)}{\mu_n(A_k)} \leq L.$$

B.1 CONCENTRATION INEQUALITIES

Denote the element of a matrix X in the i -th row and the j -th column as $X[i, j]$, and the i -th row of the matrix X as $X[i, *]$.

Lemma 7 (Lemma 20 of Tsigler & Bartlett (2023)). Let z be a sub-gaussian vector in \mathbb{R}^p with $\|z\|_{\psi_2} \leq \sigma$, and consider $\Sigma = \text{diag}(\lambda_1, \dots, \lambda_p)$ where the sequence $\{\lambda_j\}_{j=1}^p$ is positive and non-increasing. Then there exists some absolute constant c , for any $t > 0$, with probability at least $1 - 2e^{-t/c}$:

$$\|\Sigma^{1/2}z\|^2 \leq c\sigma^2 \left(t\lambda_1 + \sum_{j=1}^p \lambda_j \right).$$

Lemma 8 (Lemma 23 of Tsigler & Bartlett (2023)). Let \mathring{A}_k represent the matrix $X_{-k}X_{-k}^T$ with its diagonal elements set to zero:

$$\mathring{A}_k[i, j] = (1 - \delta_{i,j})(X_{-k}X_{-k}^T)[i, j].$$

Then there exists some absolute constant c , for any $t > 0$, with probability at least $1 - 4e^{-t/c}$:

$$\|\mathring{A}_k\| \leq c\sigma^2 \sqrt{(t+n) \left(\lambda_{k+1}^2(t+n) + \sum_{j>k} \lambda_j^2 \right)}.$$

Lemma 9 (Lemma 21 of Tsigler & Bartlett (2023)). Suppose $\{z_i\}_{i=1}^n$ is a sequence of independent isotropic sub-gaussian random vectors, where $\|z_i\|_{\psi_2} \leq \sigma$. Let $\Sigma = \text{diag}(\lambda_1, \dots, \lambda_p)$ represent a diagonal matrix with a positive, non-increasing sequence $\{\lambda_i\}_{i=1}^p$. Then there exists some absolute constant c , for any $t \in (0, n)$, with probability at least $1 - 2e^{-ct}$:

$$(n - \sqrt{nt}\sigma^2) \sum_{j=1}^p \lambda_j \leq \sum_{i=1}^n \|\Sigma^{1/2}z_i\|^2 \leq (n + \sqrt{nt}\sigma^2) \sum_{j=1}^p \lambda_j.$$

Lemma 10. There exists a constant c_x , depending only on σ , such that for any n satisfying $n\lambda_{k+1} \leq (\lambda + \sum_{j>k} \lambda_j)$, under the assumption CondNum(k, δ, L) (Assumption 2), with probability at least $1 - \delta - c_x e^{-n/c_x}$:

$$\begin{aligned}
\frac{1}{c_x L} \left(\lambda + \sum_{j>k} \lambda_j \right) &\leq \mu_n(A_k) \leq \mu_1(A_k) \leq c_x \left(\lambda + \sum_{j>k} \lambda_j \right). \\
\mu_1(X_{-k}X_{-k}^T) &\leq c_x \left(n\lambda_{k+1} + \sum_{j>k} \lambda_j \right).
\end{aligned}$$

Proof. This result follows from the proof of Lemma 3 in Tsigler & Bartlett (2023), which establishes both upper and lower bounds of $\mu_1(A_k)$. By combining the lower bound with the assumption CondNum, we derive a lower bound of $\mu_n(A_k)$. For completeness, we restate the entire proof here.

According to lemma 7 and lemma 8, there exists an absolute constant c , such that for any $t > 0$:

972 1. for all $1 \leq i \leq n$, with probability at least $1 - 2e^{-t/c}$:

$$973 \\ 974 \|X_{-k}[i, *]\|^2 \leq c\sigma^2 \left(t\lambda_{k+1} + \sum_{j>k} \lambda_j \right). \\ 975 \\ 976$$

977 2. with probability at least $1 - 4e^{-t/c}$:

$$978 \\ 979 \|A_k\| \leq c\sigma^2 \sqrt{(t+n) \left(\lambda_{k+1}^2(t+n) + \sum_{j>k} \lambda_j^2 \right)}. \\ 980 \\ 981 \\ 982$$

983 Since $\mu_1(A_k) \leq \lambda + \|A_k\| + \max_i \|X_{-k}[i, *]\|^2$, by setting $t = n$, we have with probability at least
984 $1 - (2n+4)e^{-n/c}$:

$$985 \\ 986 \mu_1(A_k) \leq \lambda + c\sigma^2 \left(n\lambda_{k+1} + \sum_{j>k} \lambda_j + \sqrt{(2n\lambda_{k+1})^2 + 2n \sum_{j>k} \lambda_j^2} \right) \\ 987 \\ 988 \leq \lambda + c\sigma^2 \left(n\lambda_{k+1} + \sum_{j>k} \lambda_j + 2n\lambda_{k+1} + \sqrt{2n \sum_{j>k} \lambda_j^2} \right) \\ 989 \\ 990 \leq \lambda + c\sigma^2 \left(n\lambda_{k+1} + \sum_{j>k} \lambda_j + 2n\lambda_{k+1} + \sqrt{2n\lambda_{k+1} \sum_{j>k} \lambda_j} \right) \\ 991 \\ 992 \leq \lambda + c\sigma^2 \left(n\lambda_{k+1} + \sum_{j>k} \lambda_j + 2n\lambda_{k+1} + n\lambda_{k+1} + \frac{1}{2} \sum_{j>k} \lambda_j \right) \\ 993 \\ 994 \leq \lambda + 4c\sigma^2 \left(n\lambda_{k+1} + \sum_{j>k} \lambda_j \right) \\ 995 \\ 996 \leq \max \{1, 4c\sigma^2\} \left(\lambda + \sum_{j>k} \lambda_j + n\lambda_{k+1} \right) \\ 997 \\ 998 \leq 2 \max \{1, 4c\sigma^2\} \left(\lambda + \sum_{j>k} \lambda_j \right). \quad (2) \\ 999 \\ 1000 \\ 1001 \\ 1002 \\ 1003 \\ 1004 \\ 1005 \\ 1006 \\ 1007 \\ 1008 \\ 1009$$

1010 The last inequality follows from $n\lambda_{k+1} \leq (\lambda + \sum_{j>k} \lambda_j)$. Similarly,

$$1011 \\ 1012 \mu_1(X_{-k} X_{-k}^T) \leq 4c\sigma^2 \left(n\lambda_{k+1} + \sum_{j>k} \lambda_j \right). \quad (3) \\ 1013 \\ 1014$$

1015 On the other hand, by applying Lemma 9 with $t = \frac{n}{4\sigma^4}$, there exists an absolute constant c' , such
1016 that with probability at least $1 - 2 \exp \left\{ -\frac{c'}{4\sigma^4} n \right\}$:

$$1017 \\ 1018 \sum_{i=1}^n \|X_{-k}[i, *]\|^2 \geq \frac{1}{2} n \sum_{j>k} \lambda_j. \\ 1019 \\ 1020$$

1021 On this event,

$$1022 \\ 1023 \mu_1(A_k) \geq \lambda + \frac{1}{n} \text{tr}(X_{-k} X_{-k}^T) \\ 1024 \\ 1025 = \lambda + \frac{1}{n} \sum_{i=1}^n \|X_{-k}[i, *]\|^2$$

$$\begin{aligned}
&\geq \lambda + \frac{1}{2} \sum_{j>k} \lambda_j \\
&\geq \frac{1}{2} \left(\lambda + \sum_{j>k} \lambda_j \right).
\end{aligned}$$

By the assumption CondNum(k, δ, L), with probability at least $1 - \delta - 2 \exp\left\{-\frac{c'}{4\sigma^4}n\right\}$:

$$\mu_n(A_k) \geq \frac{1}{L} \mu_1(A_k) \geq \frac{1}{2L} \left(\lambda + \sum_{j>k} \lambda_j \right). \quad (4)$$

Combining Equation 2, 3 and 4, there exists a constant c_x depending only on σ , such that with probability at least $1 - \delta - c_x e^{-n/c_x}$:

$$\begin{aligned}
\frac{1}{c_x L} \left(\lambda + \sum_{j>k} \lambda_j \right) &\leq \mu_n(A_k) \leq \mu_1(A_k) \leq c_x \left(\lambda + \sum_{j>k} \lambda_j \right). \\
\mu_1(X_{-k} X_{-k}^T) &\leq c_x \left(n \lambda_{k+1} + \sum_{j>k} \lambda_j \right).
\end{aligned}$$

□

Lemma 11. There exists a constant c_x depending only on σ , such that with probability at least $1 - \delta$, if $n > k + \ln(1/\delta)$,

$$\left\| \frac{1}{n} X_k^T X_k - \Sigma_{S,k} \right\| \leq c_x \lambda_1 \sqrt{\frac{k + \ln \frac{1}{\delta}}{n}}.$$

Proof. This follows directly from Theorem 5.39 and Remark 5.40 of Vershynin (2010), which shows there exists a constant c'_x depending only on σ , such that for any $t \geq 0$, with probability at least $1 - 2 \exp\{-t^2/c'_x\}$:

$$\left\| \frac{1}{n} X_k^T X_k - \Sigma_{S,k} \right\| \leq \lambda_1 \max \left\{ c'_x \sqrt{\frac{k}{n}} + \frac{t}{\sqrt{n}}, \left(c'_x \sqrt{\frac{k}{n}} + \frac{t}{\sqrt{n}} \right)^2 \right\}.$$

Taking $t = \sqrt{c'_x \ln(2/\delta)}$ completes the proof. □

Corollary 12. Under the same conditions as in Lemma 11, and on the same event, the following holds:

$$\left\| (X_k^T X_k)^{\frac{1}{2}} - \sqrt{n} \Sigma_{S,k}^{\frac{1}{2}} \right\| \leq c_x \sqrt{k + \ln \frac{1}{\delta}} \lambda_1 \lambda_k^{-\frac{1}{2}}.$$

Proof. According to Proposition 3.2 of van Hemmen & Ando (1980), for any positive semi-definite matrix $A, B \in \mathbb{R}^k$, we have

$$\|A - B\| \geq \left(\mu_k \left(A^{\frac{1}{2}} \right) + \mu_k \left(B^{\frac{1}{2}} \right) \right) \|A^{\frac{1}{2}} - B^{\frac{1}{2}}\|.$$

Therefore,

$$\begin{aligned}
\left\| (X_k^T X_k)^{\frac{1}{2}} - \sqrt{n} \Sigma_{S,k}^{\frac{1}{2}} \right\| &\leq \frac{1}{\mu_k \left(\sqrt{n} \Sigma_{S,k}^{\frac{1}{2}} \right)} \|X_k^T X_k - n \Sigma_{S,k}\| \\
&= \sqrt{n} \lambda_k^{-\frac{1}{2}} \left\| \frac{1}{n} X_k^T X_k - \Sigma_{S,k} \right\|.
\end{aligned}$$

By applying Lemma 11, the proof is complete. □

1080
1081 **Lemma 13.** There exists a constant c_x depending only on σ , such that for any $n > c_x k$, with
1082 probability at least $1 - 2e^{-n/c_x}$:

$$1083 \frac{1}{c_x} n \leq \mu_k \left(\Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right) \leq \mu_1 \left(\Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right) \leq c_x n.$$

1085
1086 *Proof.* According to Theorem 5.39 of Vershynin (2010), there exists a constant c'_x depending only
1087 on σ , such that for any $t \geq 0$, with probability at least $1 - 2 \exp\{-t^2/c'_x\}$:

$$1088 \mu_k \left(\Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right) \geq \left(\sqrt{n} - c'_x \sqrt{k} - t \right)^2.$$

$$1089 \mu_1 \left(\Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right) \leq \left(\sqrt{n} + c'_x \sqrt{k} + t \right)^2.$$

1092 Let $t = \frac{1}{2}\sqrt{n}$. For $n > 16(c'_x)^2 k$, with probability at least $1 - 2 \exp\{-n/(4c'_x)\}$:

$$1094 \mu_k \left(\Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right) \geq \left(\sqrt{n} - \frac{1}{4}\sqrt{n} - \frac{1}{2}\sqrt{n} \right)^2 = \frac{1}{16} n.$$

$$1096 \mu_1 \left(\Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right) \leq \left(\sqrt{n} + \frac{1}{4}\sqrt{n} + \frac{1}{2}\sqrt{n} \right)^2 = \frac{49}{16} n.$$

1099 By taking $c_x = \max\{16(c'_x)^2, 4c'_x, 16\}$, the proof is complete. \square

1100 **Remark 7.** On the same event, the following inequalities also hold:

$$1102 \mu_1(X_k^T X_k) \leq \|\Sigma_{S,k}\| \left\| \Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right\| \leq c_x \lambda_1 n.$$

$$1104 \mu_k(X_k^T X_k) \geq \mu_k(\Sigma_{S,k}) \mu_k \left(\Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right) \geq \frac{1}{c_x} \lambda_k n.$$

1106 **Lemma 14.** There exists a constant c_x depending only on σ , with probability at least $1 - 2e^{-n/c_x}$:

$$1108 \text{tr}(X_{-k} \Sigma_{T,-k} X_{-k}^T) \leq c_x n \text{tr} \left(\Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}} \right).$$

1110 *Proof.* According to Hanson-Wright Inequality (Vershynin, 2018), there exists an absolute constant
1111 c , such that for any $1 \leq i \leq n$,

$$1113 \left\| Z_{-k}[i, *] \Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}} Z_{-k}[i, *]^T \right\|_{\psi_1} \leq c \sigma^2 \left\| \Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}} \right\|_{\text{F}}$$

$$1115 \leq c \sigma^2 \text{tr} \left(\Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}} \right).$$

1116 By Bernstein Inequality (Proposition 5.16 of Vershynin (2010)), there exists an absolute constant c' ,
1117 for any $t \geq 0$,

$$1119 \mathbb{P} \left\{ \frac{1}{n} \left| \sum_{i=1}^n \left[Z_{-k}[i, *] \Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}} Z_{-k}[i, *]^T - \text{tr} \left(\Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}} \right) \right] \right| \geq t \right\}$$

$$1123 \leq 2 \exp \left\{ -c'n \min \left\{ \frac{t^2}{K^2}, \frac{t}{K} \right\} \right\},$$

1124 where $K = \max_i \left\| Z_{-k}[i, *] \Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}} Z_{-k}[i, *]^T \right\|_{\psi_1}$.

1127 Let $t = c \sigma^2 \text{tr} \left(\Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}} \right)$. Then, with probability at least $1 - 2e^{-c'n}$:

$$1129 \text{tr}(X_{-k} \Sigma_{T,-k} X_{-k}^T) = \sum_{i=1}^n Z_{-k}[i, *] \Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}} Z_{-k}[i, *]^T$$

$$1131 \leq (1 + c \sigma^2) n \text{tr} \left(\Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}} \right).$$

1133 By taking $c_x = \max\{1 + c \sigma^2, \frac{1}{c'}\}$, the proof is complete. \square

1134 **Lemma 15.** There exists a constant c_x depending only on σ , with probability at least $1 - 2e^{-n/c_x}$:

$$1135 \quad (\beta_{-k}^*)^T X_{-k}^T X_{-k} \beta_{-k}^* \leq c_x n (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^*.$$

1136 *Proof.* The result follows from the proof of Lemma 3 in [Tsigler & Bartlett \(2023\)](#), which we restate
1137 here for completeness. Consider the isotropic vector $[(\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^*]^{-1/2} X_{-k} \beta_{-k}^*$. For the
1140 i -th component,

$$\begin{aligned} 1141 \quad \left\| \left[(\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^* \right]^{-\frac{1}{2}} X_{-k}[i,*] \beta_{-k}^* \right\|_{\psi_2} &= \left[(\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^* \right]^{-\frac{1}{2}} \left\| Z_{-k}[i,*] \Sigma_{S,-k}^{\frac{1}{2}} \beta_{-k}^* \right\|_{\psi_2} \\ 1142 \quad &\leq \left[(\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^* \right]^{-\frac{1}{2}} \sigma \left\| \Sigma_{S,-k}^{\frac{1}{2}} \beta_{-k}^* \right\| \\ 1143 \quad &= \sigma. \end{aligned}$$

1144 By applying Lemma 9 for the sequence $\left\{ \left[(\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^* \right]^{-1/2} X_{-k}[i,*] \beta_{-k}^* \right\}_{i=1}^n$, there exists
1145 an absolute constant c , for any $t \in (0, n)$, with probability at least $1 - 2e^{-ct}$:

$$1146 \quad \frac{(\beta_{-k}^*)^T X_{-k}^T X_{-k} \beta_{-k}^*}{(\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^*} \leq n + \sqrt{nt}\sigma^2.$$

1147 Let $t = n/4$, with probability at least $1 - 2e^{-cn/4}$:

$$1148 \quad (\beta_{-k}^*)^T X_{-k}^T X_{-k} \beta_{-k}^* \leq (1 + \frac{1}{2}\sigma^2)n \cdot (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^*.$$

1149 By taking $c_x = \max\{1 + \frac{1}{2}\sigma^2, \frac{4}{c}\}$, the proof is complete. \square

1150 B.2 BLOCK DECOMPOSITION OF $X_{-k} X_{-k}^T$

1151 Let $X_k = U \widetilde{M}^{\frac{1}{2}} V$, where $U \in \mathbb{R}^{n \times n}$ and $V \in \mathbb{R}^{d \times d}$ are orthogonal matrices representing the left
1152 and right singular vectors, respectively. The matrix $\widetilde{M}^{\frac{1}{2}}$ is defined as:

$$1153 \quad \widetilde{M}^{\frac{1}{2}} = \begin{pmatrix} m_1^{\frac{1}{2}} & & \\ & \ddots & \\ & & m_k^{\frac{1}{2}} \\ & \mathbf{0}^{(n-k) \times k} & \end{pmatrix} \in \mathbb{R}^{n \times k}.$$

1154 Therefore, we have $X_k X_k^T = U M U^T$, where $M = \text{diag}(m_1, \dots, m_k, 0, \dots, 0) \in \mathbb{R}^{n \times n}$. Similarly,
1155 $X_k^T X_k = V^T M_k V$, where $M_k = \text{diag}(m_1, \dots, m_k) \in \mathbb{R}^{k \times k}$.

1156 Let $\Delta = U^T X_{-k} X_{-k}^T U$, and write Δ in block matrix form as:

$$1157 \quad \Delta = \begin{pmatrix} \Delta_{11} & \Delta_{12} \\ \Delta_{12}^T & \Delta_{22} \end{pmatrix},$$

1158 where $\Delta_{11} \in \mathbb{R}^{k \times k}$, $\Delta_{12} \in \mathbb{R}^{k \times (n-k)}$, and $\Delta_{22} \in \mathbb{R}^{(n-k) \times (n-k)}$.

1159 We will repeatedly use the first k rows of $(M + \lambda I_n + \Delta)^{-1}$, which we compute here. Because
1160 $M + \lambda I_n + \Delta$ and $\lambda I_{n-k} + \Delta_{22}$ are invertible when A_k is positive definite, by block matrix inverse,

$$\begin{aligned} 1161 \quad &(M + \lambda I_n + \Delta)^{-1}[k,*] \\ 1162 \quad &= (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} (I_k, -\Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1}). \end{aligned} \tag{5}$$

1163 **Corollary 16** (Corollary of Lemma 10). There exists a constant depending only on σ , such that
1164 for any $n < \lambda_{k+1}^{-1}(\lambda + \sum_{j>k} \lambda_j)$, if the assumption $\text{condNum}(k, \delta, L)$ is satisfied, the following
1165 inequalities hold with probability at least $1 - \delta - c_x e^{-n/c_x}$, on the same event as in Lemma 10.

$$1166 \quad \|\Delta_{11}\|, \|\Delta_{12}\| \leq \|\Delta\| \leq c_x \left(\lambda + \sum_{j>k} \lambda_j \right).$$

$$\begin{aligned}
& \|(\lambda I_{n-k} + \Delta_{22})^{-1}\| \leq \|\Delta^{-1}\| \leq c_x L \left(\lambda + \sum_{j>k} \lambda_j \right)^{-1}. \\
& \|\Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-2}\Delta_{12}^T\| \leq c_x^4 L^2. \\
& \|\Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1}\Delta_{12}^T\| \leq c_x^3 L \left(\lambda + \sum_{j>k} \lambda_j \right). \\
& \|\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1}\Delta_{12}^T\| \leq c_x \left(\lambda + \sum_{j>k} \lambda_j \right).
\end{aligned}$$

Proof. 1. The first inequality.

$$\|\Delta_{11}\|, \|\Delta_{12}\| \leq \|\Delta\| = \|X_{-k} X_{-k}^T\| \leq \|A_k\| \leq c_x \left(\lambda + \sum_{j>k} \lambda_j \right).$$

2. The second inequality.

$$\|(\lambda I_{n-k} + \Delta_{22})^{-1}\| \leq \|(\lambda I_n + \Delta)^{-1}\| = \|A_k^{-1}\| \leq c_x L \left(\lambda + \sum_{j>k} \lambda_j \right)^{-1},$$

where the first inequality holds because $\lambda I_n + \Delta$ is positive definite.

3. The third inequality.

$$\|\Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-2}\Delta_{12}^T\| \leq \|\Delta_{12}\|^2 \|(\lambda I_{n-k} + \Delta_{22})^{-1}\|^2 \leq c_x^4 L^2.$$

4. The fourth inequality.

$$\|\Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1}\Delta_{12}^T\| \leq \|\Delta_{12}\|^2 \|(\lambda I_{n-k} + \Delta_{22})^{-1}\| \leq c_x^3 L \left(\lambda + \sum_{j>k} \lambda_j \right).$$

5. The last inequality.

$$\begin{aligned}
& \|\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1}\Delta_{12}^T\| \\
&= \|\Delta_{11} + \lambda I_k - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1}\Delta_{12}^T\| - \lambda \\
&\leq \|\Delta_{11} + \lambda I_k\| - \lambda \\
&= \|\Delta_{11}\| \\
&\leq c_x \left(\lambda + \sum_{j>k} \lambda_j \right).
\end{aligned}$$

The first inequality holds because $\Delta_{11} + \lambda I_k - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1}\Delta_{12}^T$ is the Schur complement of the block $\Delta_{11} + \lambda I_k$ of the matrix $\Delta + \lambda I_n$, which is positive definite. Therefore, we have

$$\Delta_{11} + \lambda I_k \succ \Delta_{11} + \lambda I_k - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1}\Delta_{12}^T.$$

□

Lemma 17. There exists a constant $c_x > 2$ depending only on σ , such that for any $N_1 < n < N_2$, if the assumption $\text{condNum}(k, \delta, L)$ is satisfied, the following holds with probability at least $1 - 2\delta - c_x e^{-n/c_x}$, on both events from Lemma 10 and Lemma 11,

$$\left\| [X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1}\Delta_{12}^T) V]^{-1} - (n \tilde{\Sigma}_{S,k})^{-1} \right\|$$

$$\leq \frac{c_x^2 \left(\sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 + c_x^2 L \left(\lambda + \sum_{j>k} \lambda_j \right) \right)}{(\lambda + n\lambda_k)^2}.$$

where

$$N_1 = \max \left\{ 4c_x^4 (k + \ln(1/\delta)) \frac{\lambda_1^2}{\lambda_k^2}, 2c_x^4 L \lambda_k^{-1} \left(\lambda + \sum_{j>k} \lambda_j \right) \right\}.$$

$$N_2 = \frac{1}{\lambda_{k+1}} \left(\lambda + \sum_{j>k} \lambda_j \right).$$

Proof.

$$\begin{aligned} & \left\| [X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V]^{-1} - (n \tilde{\Sigma}_{S,k})^{-1} \right\| \\ & \leq \left\| [X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V]^{-1} \right\| \\ & \quad \cdot \left\| [X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V] - (n \tilde{\Sigma}_{S,k}) \right\| \\ & \quad \cdot \left\| (n \tilde{\Sigma}_{S,k})^{-1} \right\| \\ & = \frac{1}{\lambda + n\lambda_k} \left\| [X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V]^{-1} \right\| \\ & \quad \cdot \|X_k^T X_k - n \Sigma_{S,k} + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V\|. \end{aligned}$$

According to Lemma 11, Corollary 16, there exists a constant $c_x > 2$ depending only on σ , such that for any $k + \ln(1/\delta) < N_1 < n < N_2 = \lambda_{k+1}^{-1} (\lambda + \sum_{j>k} \lambda_j)$, with probability at least $1 - 2\delta - c_x e^{-n/c_x}$, on both events in Lemma 10 and Lemma 11,

$$\left\| \frac{1}{n} X_k^T X_k - \Sigma_{S,k} \right\| \leq c_x \lambda_1 \sqrt{\frac{k + \ln \frac{1}{\delta}}{n}}.$$

$$\|\Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T\| \leq c_x^3 L \left(\lambda + \sum_{j>k} \lambda_j \right).$$

$$\begin{aligned} 1. & \|X_k^T X_k - n \Sigma_{S,k} + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V\| \\ & \quad \|X_k^T X_k - n \Sigma_{S,k} + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V\| \\ & \leq \|X_k^T X_k - n \Sigma_{S,k}\| + \|(\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)\| \\ & \leq c_x \sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 + c_x^3 L \left(\lambda + \sum_{j>k} \lambda_j \right). \end{aligned}$$

$$\begin{aligned} 2. & \left\| [X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V]^{-1} \right\| \\ & \quad \frac{1}{\lambda + n\lambda_k} \|X_k^T X_k - n \Sigma_{S,k} + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V\| \\ & \leq \frac{1}{\lambda + n\lambda_k} \left(c_x \sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 + c_x^3 L \left(\lambda + \sum_{j>k} \lambda_j \right) \right). \end{aligned}$$

Since $n > 4c_x^4 (k + \ln(1/\delta)) \frac{\lambda_1^2}{\lambda_k^2}$,

$$\frac{1}{\lambda + n\lambda_k} c_x \sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 \leq \frac{c_x \sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1}{n\lambda_k}$$

$$\begin{aligned}
&= \frac{c_x \sqrt{(k + \ln \frac{1}{\delta}) \lambda_1}}{\sqrt{n} \lambda_k} \\
&< \frac{1}{2c_x}.
\end{aligned}$$

Since $n > 2c_x^4 L \lambda_k^{-1} (\lambda + \sum_{j>k} \lambda_j)$,

$$\begin{aligned}
\frac{1}{\lambda + n\lambda_k} c_x^3 L \left(\lambda + \sum_{j>k} \lambda_j \right) &\leq \frac{c_x^3 L (\lambda + \sum_{j>k} \lambda_j)}{n\lambda_k} \\
&< \frac{1}{2c_x}.
\end{aligned}$$

Therefore, we have

$$\frac{1}{\lambda + n\lambda_k} \|X_k^T X_k - n\Sigma_{S,k} + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V\| < \frac{1}{c_x}.$$

Now we derive the upper bound for our target.

$$\begin{aligned}
&\left\| [X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V]^{-1} \right\| \\
&= \left\| [n\tilde{\Sigma}_{S,k} + X_k^T X_k - n\Sigma_{S,k} + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V]^{-1} \right\| \\
&\leq \left\| (n\tilde{\Sigma}_{S,k})^{-1} \right\| \left\| 1 - \left\| (n\tilde{\Sigma}_{S,k})^{-1} \right\| \cdot \left\| [X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V]^{-1} \right\| \right\|^{-1} \\
&\leq \frac{1}{\lambda + n\lambda_k} \left(1 - \frac{1}{c_x} \right)^{-1} \\
&\leq \frac{c_x}{\lambda + n\lambda_k}.
\end{aligned}$$

The first inequality follows from the result $\|(A + T)^{-1}\| \leq \|A^{-1}\| (1 - \|A^{-1}\| \|T\|)^{-1}$, provided that both A and $A + T$ are invertible and $\|A^{-1}\| \|T\| < 1$ (see Lemma 3.1 in [Wedin \(1973\)](#)).

Combining the above two inequalities,

$$\begin{aligned}
&\left\| [X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V]^{-1} - (n\tilde{\Sigma}_{S,k})^{-1} \right\| \\
&\leq \frac{1}{\lambda + n\lambda_k} \frac{c_x}{\lambda + n\lambda_k} \left(c_x \sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 + c_x^3 L \left(\lambda + \sum_{j>k} \lambda_j \right) \right) \\
&= \frac{c_x^2 \left(\sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 + c_x^2 L \left(\lambda + \sum_{j>k} \lambda_j \right) \right)}{(\lambda + n\lambda_k)^2}.
\end{aligned}$$

□

B.3 BIAS VARIANCE DECOMPOSITION

We consider the expectation of the excess risk $\mathcal{R}(\hat{\beta}(Y)) = \mathcal{R}(\hat{\beta}(X\beta^*) + \hat{\beta}(\epsilon))$ with respect to the distribution of the noise ϵ .

$$\mathbb{E}_\epsilon [\mathcal{R}(\hat{\beta}(Y))] = \mathbb{E}_\epsilon \left[(\hat{\beta}(Y) - \beta^*)^T \Sigma_T (\hat{\beta}(Y) - \beta^*) \right]$$

$$= \mathbb{E}_\epsilon \left[\widehat{\beta}(\epsilon)^T \Sigma_T \widehat{\beta}(\epsilon) \right] + \left(\widehat{\beta}(X\beta^*) - \beta^* \right)^T \Sigma_T \left(\widehat{\beta}(X\beta^*) - \beta^* \right).$$

We decompose the expected excess risk into variance and bias terms.

$$\begin{aligned} V &= \mathbb{E}_\epsilon \left[\widehat{\beta}(\epsilon)^T \Sigma_T \widehat{\beta}(\epsilon) \right] \\ &\leq 2\mathbb{E}_\epsilon \left[\widehat{\beta}(\epsilon)_k^T \Sigma_{T,k} \widehat{\beta}(\epsilon)_k \right] + 2\mathbb{E}_\epsilon \left[\widehat{\beta}(\epsilon)_{-k}^T \Sigma_{T,-k} \widehat{\beta}(\epsilon)_{-k} \right]. \\ B &= \left(\widehat{\beta}(X\beta^*) - \beta^* \right)^T \Sigma_T \left(\widehat{\beta}(X\beta^*) - \beta^* \right) \\ &\leq 2 \left(\widehat{\beta}(X\beta^*)_k - \beta_k^* \right)^T \Sigma_{T,k} \left(\widehat{\beta}(X\beta^*)_k - \beta_k^* \right) \\ &\quad + 2 \left(\widehat{\beta}(X\beta^*)_{-k} - \beta_{-k}^* \right)^T \Sigma_{T,-k} \left(\widehat{\beta}(X\beta^*)_{-k} - \beta_{-k}^* \right). \end{aligned}$$

The inequalities follow from the result for a positive definite block quadratic form:

$$(x_1^T, x_2^T) \begin{pmatrix} A & B \\ B^T & D \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = x_1^T Ax_1 + 2x_1^T Bx_2 + x_1^T Dx_1,$$

where the positive definiteness implies $x_1^T Ax_1 + x_1^T Dx_1 \geq 2x_1^T Bx_2$.

Lemma 18. There exists a constant $c_x > 2$ depending only on σ , such that for any $N_1 < n < N_2$, if the assumption condNum(k, δ, L) (Assumption 2) is satisfied, then with probability at least $1 - 2\delta - c_x e^{-n/c_x}$, the following inequalities hold simultaneously:

$$\begin{aligned} \mu_n(A_k) &\geq \frac{1}{c_x L} \left(\lambda + \sum_{j>k} \lambda_j \right). \\ \mu_1(A_k) &\leq c_x \left(\lambda + \sum_{j>k} \lambda_j \right). \\ \mu_1(X_{-k} X_{-k}^T) &\leq c_x \left(n\lambda_{k+1} + \sum_{j>k} \lambda_j \right). \\ \left\| \frac{1}{n} X_k^T X_k - \Sigma_{S,k} \right\| &\leq c_x \lambda_1 \sqrt{\frac{k + \ln \frac{1}{\delta}}{n}}. \\ \left\| (X_k^T X_k)^{\frac{1}{2}} - \sqrt{n} \Sigma_{S,k}^{\frac{1}{2}} \right\| &\leq c_x \sqrt{k + \ln \frac{1}{\delta}} \lambda_1 \lambda_k^{-\frac{1}{2}}. \\ \mu_k \left(\Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right) &\geq \frac{1}{c_x} n. \\ \mu_1 \left(\Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right) &\leq c_x n. \\ \mu_1(X_k^T X_k) &\leq c_x \lambda_1 n. \\ \mu_k(X_k^T X_k) &\geq \frac{1}{c_x} \lambda_k n. \\ \text{tr} \left(X_{-k} \Sigma_{T,-k} X_{-k}^T \right) &\leq c_x n \text{tr} \left(\Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}} \right). \\ (\beta_{-k}^*)^T X_{-k}^T X_{-k} \beta_{-k}^* &\leq c_x n (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^*. \\ \|\Delta_{11}\|, \|\Delta_{12}\|, \|\Delta\| &\leq c_x \left(\lambda + \sum_{j>k} \lambda_j \right). \\ \|(\lambda I_{n-k} + \Delta_{22})^{-1}\|, \|\Delta^{-1}\| &\leq c_x L \left(\lambda + \sum_{j>k} \lambda_j \right)^{-1}. \end{aligned}$$

$$\begin{aligned} & \|\Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-2} \Delta_{12}^T\| \leq c_x^4 L^2. \\ & \|\Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T\| \leq c_x^3 L \left(\lambda + \sum_{j>k} \lambda_j \right). \\ & \|\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T\| \leq c_x \left(\lambda + \sum_{j>k} \lambda_j \right). \end{aligned}$$

And,

$$\begin{aligned} & \left\| [X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V]^{-1} - (n \tilde{\Sigma}_{S,k})^{-1} \right\| \\ & \leq \frac{c_x^2 \left(\sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 + c_x^2 L \left(\lambda + \sum_{j>k} \lambda_j \right) \right)}{(\lambda + n \lambda_k)^2}. \end{aligned}$$

N_1 and N_2 are defined as follows:

$$\begin{aligned} N_1 &= \max \left\{ 4c_x^4 (k + \ln(1/\delta)) \frac{\lambda_1^2}{\lambda_k^2}, 2c_x^4 L \lambda_k^{-1} \left(\lambda + \sum_{j>k} \lambda_j \right) \right\}. \\ N_2 &= \frac{1}{\lambda_{k+1}} \left(\lambda + \sum_{j>k} \lambda_j \right). \end{aligned}$$

Proof. The lemma is a direct corollary from Lemma 10, Lemma 11, Corollary 12, Lemma 13, Lemma 14, Lemma 15, Corollary 16, Lemma 17. \square

B.3.1 VARIANCE IN THE FIRST k DIMENSIONS

Lemma 19. Under the same conditions as in Lemma 18, and on the same event, for any $N_1 < n < N_2$,

$$\mathbb{E}_\epsilon [\hat{\beta}(\epsilon)_k^T \Sigma_{T,k} \hat{\beta}(\epsilon)_k] \leq 16v^2 (1 + c_x^4 L^2) \frac{1}{n} \text{tr} [\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}}],$$

where

$$\begin{aligned} N_1 &= \max \left\{ 4c_x^4 \left(k + \ln \frac{1}{\delta} \right) \lambda_1^4 \lambda_k^{-4}, \right. \\ &\quad 2c_x^4 L \lambda_1 \lambda_k^{-2} \left(\lambda + \sum_{j>k} \lambda_j \right), \\ &\quad 4c_x^4 \left(k + \ln \frac{1}{\delta} \right) \lambda_1^6 \lambda_k^{-8} \|\Sigma_{T,k}\|^2 k^2 \left(\text{tr} [\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}}] \right)^{-2}, \\ &\quad \left. 2c_x^4 L \lambda_1^2 \lambda_k^{-4} \left(\lambda + \sum_{j>k} \lambda_j \right) \|\Sigma_{T,k}\| k \left(\text{tr} [\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}}] \right)^{-1} \right\}, \\ N_2 &= \frac{1}{\lambda_{k+1}} \left(\lambda + \sum_{j>k} \lambda_j \right). \end{aligned}$$

Proof.

$$\begin{aligned} & \mathbb{E}_\epsilon [\hat{\beta}(\epsilon)_k^T \Sigma_{T,k} \hat{\beta}(\epsilon)_k] \\ &= \mathbb{E}_\epsilon \text{tr} [\epsilon \epsilon^T (XX^T + \lambda I_n)^{-1} X_k \Sigma_{T,k} X_k^T (XX^T + \lambda I_n)^{-1}] \end{aligned}$$

$$\begin{aligned}
&= v^2 \operatorname{tr} [(XX^T + \lambda I_n)^{-1} X_k \Sigma_{T,k} X_k^T (XX^T + \lambda I_n)^{-1}] \\
&= v^2 \operatorname{tr} [(UMU^T + U\Delta U^T + \lambda I_n)^{-1} U \widetilde{M}^{\frac{1}{2}} V \Sigma_{T,k} \\
&\quad \cdot V^T (\widetilde{M}^{\frac{1}{2}})^T U^T (UMU^T + U\Delta U^T + \lambda I_n)^{-1}] \\
&= v^2 \operatorname{tr} [U(M + \Delta + \lambda I_n)^{-1} \widetilde{M}^{\frac{1}{2}} V \Sigma_{T,k} V^T (\widetilde{M}^{\frac{1}{2}})^T (M + \Delta + \lambda I_n)^{-1} U^T] \\
&= v^2 \operatorname{tr} [\left(\widetilde{M}^{\frac{1}{2}}\right)^T (M + \Delta + \lambda I_n)^{-1} (M + \Delta + \lambda I_n)^{-1} \widetilde{M}^{\frac{1}{2}} V \Sigma_{T,k} V^T] \\
&= v^2 \operatorname{tr} [M_k^{\frac{1}{2}} (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} (I_k, -\Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1}) \\
&\quad \cdot (I_k, -\Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1})^T (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} M_k^{\frac{1}{2}} \\
&\quad \cdot V \Sigma_{T,k} V^T] \\
&= v^2 \operatorname{tr} [M_k^{\frac{1}{2}} (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} \\
&\quad \cdot (I_k + \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-2} \Delta_{12}^T) (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} M_k^{\frac{1}{2}} \\
&\quad \cdot V \Sigma_{T,k} V^T] \\
&= v^2 \operatorname{tr} [(I_k + \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-2} \Delta_{12}^T) (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} \\
&\quad \cdot M_k^{\frac{1}{2}} V \Sigma_{T,k} V^T M_k^{\frac{1}{2}} (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1}] \\
&\leq v^2 \|I_k + \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-2} \Delta_{12}^T\| \operatorname{tr} [(M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} \\
&\quad \cdot M_k^{\frac{1}{2}} V \Sigma_{T,k} V^T M_k^{\frac{1}{2}} (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1}] \\
&\leq v^2 (1 + c_x^4 L^2) \operatorname{tr} [(M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} \\
&\quad \cdot M_k^{\frac{1}{2}} V \Sigma_{T,k} V^T M_k^{\frac{1}{2}} (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1}] \\
&= v^2 (1 + c_x^4 L^2) \operatorname{tr} [(V^T (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1} \\
&\quad \cdot V^T M_k^{\frac{1}{2}} V \cdot \Sigma_{T,k} \cdot V^T M_k^{\frac{1}{2}} V \cdot (V^T (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1}] \\
&= v^2 (1 + c_x^4 L^2) \operatorname{tr} [(X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1} \\
&\quad \cdot (X_k^T X_k)^{\frac{1}{2}} \Sigma_{T,k} (X_k^T X_k)^{\frac{1}{2}} \\
&\quad \cdot (X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1}].
\end{aligned}$$

The sixth equation follows from Equation 5. The first inequality follows from the result $\operatorname{tr}[AB] \leq \|A\| \operatorname{tr}[B]$ where the matrix B is positive semi-definite.

We define two quantities that represent concentration error terms:

$$\begin{aligned}
E_1 &= \left\| [X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V]^{-1} - (n \widetilde{\Sigma}_{S,k})^{-1} \right\|. \\
E_2 &= (X_k^T X_k)^{\frac{1}{2}} - (n \Sigma_{S,k})^{\frac{1}{2}}. \\
\text{Since } n &> 4c_x^4 (k + \ln \frac{1}{\delta}) \lambda_1^6 \lambda_k^{-8} \|\Sigma_{T,k}\|^2 k^2 \left(\operatorname{tr} [\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}}] \right)^{-2}, \\
\text{and } n &> 2c_x^4 L \left(\lambda + \sum_{j>k} \lambda_j \right) \lambda_1^2 \lambda_k^{-4} \|\Sigma_{T,k}\| k \left(\operatorname{tr} [\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}}] \right)^{-1}, \\
\|E_1\| \left\| n \widetilde{\Sigma}_{S,k} \right\| \left\| (n \widetilde{\Sigma}_{S,k})^{-1} \right\| \left\| (n \Sigma_{S,k})^{\frac{1}{2}} \right\| \left\| \Sigma_{T,k} \right\| \left\| (n \Sigma_{S,k})^{\frac{1}{2}} \right\| \left\| (n \widetilde{\Sigma}_{S,k})^{-1} \right\|
\end{aligned}$$

$$\begin{aligned}
&\leq \frac{c_x^2 \left(\sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 + c_x^2 L \left(\lambda + \sum_{j>k} \lambda_j \right) \right)}{(\lambda + n\lambda_k)^2} (\lambda + n\lambda_1) \frac{n\lambda_1}{(\lambda + n\lambda_k)^2} \|\Sigma_{T,k}\| \\
&\leq \frac{c_x^2 \left(\sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 + c_x^2 L \left(\lambda + \sum_{j>k} \lambda_j \right) \right)}{n^2} \frac{\lambda_1^2}{\lambda_k^4} \|\Sigma_{T,k}\| \\
&= \frac{c_x^2 \sqrt{(k + \ln \frac{1}{\delta})}}{n\sqrt{n}} \frac{\lambda_1^3}{\lambda_k^4} \|\Sigma_{T,k}\| + \frac{c_x^4 L \left(\lambda + \sum_{j>k} \lambda_j \right)}{n^2} \frac{\lambda_1^2}{\lambda_k^4} \|\Sigma_{T,k}\| \\
&< \frac{1}{2nk} \text{tr} \left[\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right] + \frac{1}{2nk} \text{tr} \left[\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right] \\
&= \frac{1}{nk} \text{tr} \left[\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right].
\end{aligned}$$

Since $n > 4c_x^4 (k + \ln \frac{1}{\delta}) \lambda_1^4 \lambda_k^{-4}$ and $n > 2c_x^4 L \left(\lambda + \sum_{j>k} \lambda_j \right) \lambda_1 \lambda_k^{-2}$,

$$\begin{aligned}
&\|E_1\| \left\| n \tilde{\Sigma}_{S,k} \right\| \\
&\leq \frac{c_x^2 \left(\sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 + c_x^2 L \left(\lambda + \sum_{j>k} \lambda_j \right) \right)}{(\lambda + n\lambda_k)^2} (\lambda + n\lambda_1) \\
&\leq \frac{c_x^2 \left(\sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 + c_x^2 L \left(\lambda + \sum_{j>k} \lambda_j \right) \right)}{n} \frac{\lambda_1}{\lambda_k^2} \\
&= \frac{c_x^2 \sqrt{(k + \ln \frac{1}{\delta})}}{\sqrt{n}} \frac{\lambda_1^2}{\lambda_k^2} + \frac{c_x^4 L \left(\lambda + \sum_{j>k} \lambda_j \right)}{n} \frac{\lambda_1}{\lambda_k^2} \\
&< \frac{1}{2} + \frac{1}{2} \\
&= 1.
\end{aligned} \tag{6}$$

Since $n > c_x^2 (k + \ln \frac{1}{\delta}) \lambda_1^4 \lambda_k^{-6} \|\Sigma_{T,k}\|^2 k^2 \left(\text{tr} \left[\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right] \right)^{-2}$,

$$\begin{aligned}
&\|E_2\| \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-\frac{1}{2}} \right\| \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right\| \left\| (n \Sigma_{S,k})^{\frac{1}{2}} \right\| \left\| \Sigma_{T,k} \right\| \left\| (n \Sigma_{S,k})^{\frac{1}{2}} \right\| \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right\| \\
&\leq c_x \sqrt{k + \ln \frac{1}{\delta}} \lambda_1 \lambda_k^{-\frac{1}{2}} (n \lambda_k)^{-\frac{1}{2}} \frac{n \lambda_1}{(\lambda + n \lambda_k)^2} \|\Sigma_{T,k}\| \\
&\leq \frac{c_x \sqrt{k + \ln \frac{1}{\delta}}}{n \sqrt{n}} \frac{\lambda_1^2}{\lambda_k^3} \|\Sigma_{T,k}\| \\
&\leq \frac{1}{nk} \text{tr} \left[\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right].
\end{aligned}$$

Since $n > c_x^2 (k + \ln \frac{1}{\delta}) \lambda_1^2 \lambda_k^{-2}$,

$$\begin{aligned}
&\|E_2\| \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-\frac{1}{2}} \right\| \\
&\leq c_x \sqrt{k + \ln \frac{1}{\delta}} \lambda_1 \lambda_k^{-\frac{1}{2}} (n \lambda_k)^{-\frac{1}{2}} \\
&= \frac{c_x \sqrt{k + \ln \frac{1}{\delta}}}{\sqrt{n}} \frac{\lambda_1}{\lambda_k} \\
&< 1.
\end{aligned} \tag{7}$$

Combing the above four inequalities, we have

$$\text{tr} \left[(X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12} (\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1} \right]$$

$$\begin{aligned}
& \cdot (X_k^T X_k)^{\frac{1}{2}} \Sigma_{T,k} (X_k^T X_k)^{\frac{1}{2}} \\
& \cdot (X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1} \\
= & \text{tr} \left[\left(n \tilde{\Sigma}_{S,k} \right)^{-1} (n \Sigma_{S,k})^{\frac{1}{2}} \Sigma_{T,k} (n \Sigma_{S,k})^{\frac{1}{2}} \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right] \\
& + 2 \text{tr} \left[E_1 (n \Sigma_{S,k})^{\frac{1}{2}} \Sigma_{T,k} (n \Sigma_{S,k})^{\frac{1}{2}} \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right] \\
& + 2 \text{tr} \left[\left(n \tilde{\Sigma}_{S,k} \right)^{-1} E_2 \Sigma_{T,k} (n \Sigma_{S,k})^{\frac{1}{2}} \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right] \\
& + \text{tr} \left[E_1 (n \Sigma_{S,k})^{\frac{1}{2}} \Sigma_{T,k} (n \Sigma_{S,k})^{\frac{1}{2}} E_1 \right] \\
& + \text{tr} \left[\left(n \tilde{\Sigma}_{S,k} \right)^{-1} E_2 \Sigma_{T,k} E_2 \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right] \\
& + 2 \text{tr} \left[E_1 (n \Sigma_{S,k})^{\frac{1}{2}} \Sigma_{T,k} E_2 \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right] \\
& + 2 \text{tr} \left[E_1 E_2 \Sigma_{T,k} (n \Sigma_{S,k})^{\frac{1}{2}} \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right] \\
& + 2 \text{tr} \left[E_1 E_2 \Sigma_{T,k} (n \Sigma_{S,k})^{\frac{1}{2}} E_1 \right] \\
& + 2 \text{tr} \left[E_1 E_2 \Sigma_{T,k} \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right] \\
& + \text{tr} [E_1 E_2 \Sigma_{T,k} E_2 E_1].
\end{aligned}$$

In particular,

$$\begin{aligned}
& \text{tr} \left[\left(n \tilde{\Sigma}_{S,k} \right)^{-1} (n \Sigma_{S,k})^{\frac{1}{2}} \Sigma_{T,k} (n \Sigma_{S,k})^{\frac{1}{2}} \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right] \\
& = \frac{1}{n} \text{tr} \left[\tilde{\Sigma}_{S,k}^{-1} \Sigma_{S,k}^{\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{\frac{1}{2}} \tilde{\Sigma}_{S,k}^{-1} \right] \\
& \leq \frac{1}{n} \text{tr} \left[\Sigma_{S,k}^{-1} \Sigma_{S,k}^{\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{\frac{1}{2}} \Sigma_{S,k}^{-1} \right] \\
& = \frac{1}{n} \text{tr} \left[\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right].
\end{aligned}$$

The inequality follows from the fact that $\text{tr}[BAB] = \text{tr}[A^{\frac{1}{2}}BA^{\frac{1}{2}}] \leq \text{tr}[A^{\frac{1}{2}}CA^{\frac{1}{2}}] = \text{tr}[CAC]$, where A, B, C are positive semi-definite matrices, and $C \succcurlyeq B$, which implies that $A^{\frac{1}{2}}CA^{\frac{1}{2}} \succcurlyeq A^{\frac{1}{2}}BA^{\frac{1}{2}}$.

$$\begin{aligned}
& \text{tr} [E_1 E_2 \Sigma_{T,k} E_2 E_1] \\
& = \text{tr} \left[E_1 n \tilde{\Sigma}_{S,k} \left(n \tilde{\Sigma}_{S,k} \right)^{-1} E_2 \left(n \tilde{\Sigma}_{S,k} \right)^{-\frac{1}{2}} \left(n \tilde{\Sigma}_{S,k} \right)^{\frac{1}{2}} \right. \\
& \quad \left. \cdot \Sigma_{T,k} E_2 \left(n \tilde{\Sigma}_{S,k} \right)^{-\frac{1}{2}} \left(n \tilde{\Sigma}_{S,k} \right)^{\frac{1}{2}} E_1 n \tilde{\Sigma}_{S,k} \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right] \\
& \leq k \left(\|E_1\| \left\| n \tilde{\Sigma}_{S,k} \right\| \right)^2 \left(\|E_2\| \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-\frac{1}{2}} \right\| \right) \\
& \quad \cdot \|E_2\| \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-\frac{1}{2}} \right\| \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right\| \left\| (n \Sigma_{S,k})^{\frac{1}{2}} \right\| \|\Sigma_{T,k}\| \left\| (n \Sigma_{S,k})^{\frac{1}{2}} \right\| \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right\| \\
& \leq \frac{1}{n} \text{tr} \left[\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right].
\end{aligned}$$

The other terms can be similarly bounded. Therefore,

$$\text{tr} \left[(X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1} \right]$$

$$\begin{aligned}
& \cdot (X_k^T X_k)^{\frac{1}{2}} \Sigma_{T,k} (X_k^T X_k)^{\frac{1}{2}} \\
& \cdot (X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1} \\
& \leq \frac{16}{n} \text{tr} [\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}}].
\end{aligned}$$

The proof is complete by combining all the inequalities above. \square

B.3.2 VARIANCE IN THE LAST $d - k$ DIMENSIONS

Lemma 20. Under the same conditions as in Lemma 18, and on the same event, for any $N_1 < n < N_2$,

$$\mathbb{E}_\epsilon [\widehat{\beta}(\epsilon)_{-k}^T \Sigma_{T,-k} \widehat{\beta}(\epsilon)_{-k}] \leq v^2 c_x^3 L^2 n \left(\lambda + \sum_{j>k} \lambda_j \right)^{-2} \text{tr} [\Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}}].$$

where N_1, N_2 are defined as in Lemma 18.

Proof.

$$\begin{aligned}
& \mathbb{E}_\epsilon [\widehat{\beta}(\epsilon)_{-k}^T \Sigma_{T,-k} \widehat{\beta}(\epsilon)_{-k}] \\
& = \mathbb{E}_\epsilon \text{tr} [\epsilon \epsilon^T (XX^T + \lambda I_n)^{-1} X_{-k} \Sigma_{T,-k} X_{-k}^T (XX^T + \lambda I_n)^{-1}] \\
& = v^2 \text{tr} [(XX^T + \lambda I_n)^{-1} X_{-k} \Sigma_{T,-k} X_{-k}^T (XX^T + \lambda I_n)^{-1}] \\
& \leq v^2 \| (XX^T + \lambda I_n)^{-1} \| \text{tr} [X_{-k} \Sigma_{T,-k} X_{-k}^T] \\
& \leq v^2 \| (X_{-k} X_{-k}^T + \lambda I_n)^{-1} \| \text{tr} [X_{-k} \Sigma_{T,-k} X_{-k}^T] \\
& \leq v^2 \left(\frac{1}{c_x L} \left(\lambda + \sum_{j>k} \lambda_j \right) \right)^{-2} c_x n \text{tr} [\Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}}] \\
& = v^2 c_x^3 L^2 n \left(\lambda + \sum_{j>k} \lambda_j \right)^{-2} \text{tr} [\Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}}].
\end{aligned}$$

The first inequality follows from the result $\text{tr}[ABA] = \text{tr}[A^2 B] \leq \|A^2\| \text{tr}[B]$ where the matrix B is positive semi-definite. The second inequality follows from $XX^T + \lambda I_n \succcurlyeq X_{-k} X_{-k}^T + \lambda I_n$. \square

B.3.3 BIAS IN THE FIRST k DIMENSIONS

The bias in the first k dimensions can be decomposed into two terms.

$$\begin{aligned}
& (\widehat{\beta}(X\beta^\star)_k - \beta_k^\star)^T \Sigma_{T,k} (\widehat{\beta}(X\beta^\star)_k - \beta_k^\star) \\
& = (X_k^T (XX^T + \lambda I_n)^{-1} X \beta^\star - \beta_k^\star)^T \Sigma_{T,k} (X_k^T (XX^T + \lambda I_n)^{-1} X \beta^\star - \beta_k^\star) \\
& \leq 2 (X_k^T (XX^T + \lambda I_n)^{-1} X_k \beta_k^\star - \beta_k^\star)^T \Sigma_{T,k} (X_k^T (XX^T + \lambda I_n)^{-1} X_k \beta_k^\star - \beta_k^\star) \\
& \quad + 2 (X_k^T (XX^T + \lambda I_n)^{-1} X_{-k} \beta_{-k}^\star)^T \Sigma_{T,k} (X_k^T (XX^T + \lambda I_n)^{-1} X_{-k} \beta_{-k}^\star).
\end{aligned}$$

The inequality follows from the result $x_1^T A x_1 + x_2^T A x_2 \geq 2x_1^T A x_2$ where A is positive semi-definite.

Lemma 21. Under the same conditions as in Lemma 18, and on the same event, for any $N_1 < n < N_2$,

$$\begin{aligned}
& (X_k^T (XX^T + \lambda I_n)^{-1} X_k \beta_k^\star - \beta_k^\star)^T \Sigma_{T,k} (X_k^T (XX^T + \lambda I_n)^{-1} X_k \beta_k^\star - \beta_k^\star) \\
& \leq \frac{16c_x^4}{n^2} \left(\lambda + \sum_{j>k} \lambda_j \right)^2 (\beta_k^\star)^T \Sigma_{S,k}^{-1} \beta_k^\star \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|.
\end{aligned}$$

1674 where

$$\begin{aligned}
 N_1 &= \max \left\{ 2c_x^3(\lambda + \sum_{j>k} \lambda_j) \lambda_1 \lambda_k^{-2}, \right. \\
 &\quad 4c_x^4(k + \ln(1/\delta)) \lambda_1^2 \lambda_k^{-2}, \\
 &\quad \left. 2c_x^4 L \lambda_k^{-1} \left(\lambda + \sum_{j>k} \lambda_j \right) \right\}, \\
 N_2 &= \frac{1}{\lambda_{k+1}} \left(\lambda + \sum_{j>k} \lambda_j \right).
 \end{aligned}$$

1686 *Proof.*

$$\begin{aligned}
 & (X_k^T (XX^T + \lambda I_n)^{-1} X_k \beta_k^* - \beta_k^*)^T \Sigma_{T,k} (X_k^T (XX^T + \lambda I_n)^{-1} X_k \beta_k^* - \beta_k^*) \\
 &= (\beta_k^*)^T (X_k^T (XX^T + \lambda I_n)^{-1} X_k - I_k)^T \Sigma_{T,k} (X_k^T (XX^T + \lambda I_n)^{-1} X_k - I_k) \beta_k^* \\
 &= (\beta_k^*)^T \Sigma_{S,k}^{-\frac{1}{2}} \left(\Sigma_{S,k}^{\frac{1}{2}} X_k^T (XX^T + \lambda I_n)^{-1} X_k \Sigma_{S,k}^{\frac{1}{2}} - \Sigma_{S,k} \right)^T \Sigma_{S,k}^{-\frac{1}{2}} \\
 &\quad \cdot \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \left(\Sigma_{S,k}^{\frac{1}{2}} X_k^T (XX^T + \lambda I_n)^{-1} X_k \Sigma_{S,k}^{\frac{1}{2}} - \Sigma_{S,k} \right)^T \Sigma_{S,k}^{-\frac{1}{2}} \beta_k^* \\
 &\leq (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^* \cdot \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| \\
 &\quad \cdot \left\| \Sigma_{S,k}^{\frac{1}{2}} X_k^T (XX^T + \lambda I_n)^{-1} X_k \Sigma_{S,k}^{\frac{1}{2}} - \Sigma_{S,k} \right\|^2.
 \end{aligned}$$

1699 Subsequently,

$$\begin{aligned}
 & \left\| \Sigma_{S,k}^{\frac{1}{2}} X_k^T (XX^T + \lambda I_n)^{-1} X_k \Sigma_{S,k}^{\frac{1}{2}} - \Sigma_{S,k} \right\| \\
 &= \left\| \Sigma_{S,k}^{\frac{1}{2}} V^T \left(\widetilde{M}^{\frac{1}{2}} \right)^T U^T U (M + \lambda I_n + \Delta)^{-1} U^T U \widetilde{M}^{\frac{1}{2}} V \Sigma_{S,k}^{\frac{1}{2}} - \Sigma_{S,k} \right\| \\
 &= \left\| \Sigma_{S,k}^{\frac{1}{2}} V^T M_k^{\frac{1}{2}} (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} M_k^{\frac{1}{2}} V \Sigma_{S,k}^{\frac{1}{2}} - \Sigma_{S,k} \right\| \\
 &= \left\| \Sigma_{S,k}^{\frac{1}{2}} \left(V^T M_k^{-\frac{1}{2}} (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} V \right)^{-1} \Sigma_{S,k}^{\frac{1}{2}} \right. \\
 &\quad \left. - \Sigma_{S,k} \right\| \\
 &= \left\| \Sigma_{S,k}^{\frac{1}{2}} \left(I_k + V^T M_k^{-\frac{1}{2}} (\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} V \right)^{-1} \Sigma_{S,k}^{\frac{1}{2}} \right. \\
 &\quad \left. - \Sigma_{S,k} \right\| \\
 &= \left\| \Sigma_{S,k}^{\frac{1}{2}} \left(\left(I_k + V^T M_k^{-\frac{1}{2}} (\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} V \right)^{-1} - I_k \right) \right. \\
 &\quad \left. \cdot \Sigma_{S,k}^{\frac{1}{2}} \right\| \\
 &= \left\| \Sigma_{S,k}^{\frac{1}{2}} V^T M_k^{-\frac{1}{2}} (\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} V \right. \\
 &\quad \left. \cdot \left(I_k + V^T M_k^{-\frac{1}{2}} (\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} V \right)^{-1} \Sigma_{S,k}^{\frac{1}{2}} \right\| \\
 &\leq \left\| \Sigma_{S,k}^{\frac{1}{2}} V^T M_k^{-\frac{1}{2}} (\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} V \Sigma_{S,k}^{\frac{1}{2}} \right\| \\
 &\quad + \left\| \Sigma_{S,k}^{\frac{1}{2}} V^T M_k^{-\frac{1}{2}} (\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} V \right. \\
 &\quad \left. \cdot \left[\left(I_k + V^T M_k^{-\frac{1}{2}} (\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} V \right)^{-1} - I_k \right] \Sigma_{S,k}^{\frac{1}{2}} \right\|.
 \end{aligned}$$

1728 The second equation follows from Equation 5.
 1729

1730 We will derive upper bounds for both terms in the last equation above.
 1731

1732 1. The first term.

$$\begin{aligned}
 & \left\| \Sigma_{S,k}^{\frac{1}{2}} V^T M_k^{-\frac{1}{2}} (\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} V \Sigma_{S,k}^{\frac{1}{2}} \right\| \\
 & \leq \left\| \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T \right\| \left\| \Sigma_{S,k}^{\frac{1}{2}} V^T M_k^{-1} V \Sigma_{S,k}^{\frac{1}{2}} \right\| \\
 & = \left\| \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T \right\| \left\| \left(\Sigma_{S,k}^{-\frac{1}{2}} (X_k^T X_k)^{-1} \Sigma_{S,k}^{-\frac{1}{2}} \right)^{-1} \right\| \\
 & \leq \left(\lambda + c_x \left(\lambda + \sum_{j>k} \lambda_j \right) \right) \frac{c_x}{n} \\
 & \leq \frac{2c_x^2}{n} \left(\lambda + \sum_{j>k} \lambda_j \right).
 \end{aligned}$$

1746 The inequality follows from $c_x > 2$.
 1747

1748 2. The second term.

1749 Since $n > 2c_x^3(\lambda + \sum_{j>k} \lambda_j)\lambda_k^{-1}$,
 1750

$$\begin{aligned}
 & \left\| M_k^{-\frac{1}{2}} (\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} \right\| \\
 & \leq \|M_k^{-1}\| \left\| \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T \right\| \\
 & \leq \frac{c_x}{n\lambda_k} \cdot 2c_x \left(\lambda + \sum_{j>k} \lambda_j \right) \\
 & < \frac{1}{c_x}.
 \end{aligned}$$

1760 Therefore,

$$\begin{aligned}
 & \left\| \left(I_k + V^T M_k^{-\frac{1}{2}} (\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} V \right)^{-1} - I_k \right\| \\
 & \leq \left\| \left(I_k + V^T M_k^{-\frac{1}{2}} (\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} V \right)^{-1} \right\| \\
 & \quad \cdot \left\| V^T M_k^{-\frac{1}{2}} (\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} V \right\| \\
 & \leq \left(1 - \frac{1}{c_x} \right)^{-1} \left\| V^T M_k^{-\frac{1}{2}} (\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} V \right\| \\
 & \leq c_x \cdot \frac{c_x}{n\lambda_k} \cdot 2c_x \left(\lambda + \sum_{j>k} \lambda_j \right) \\
 & = \frac{2c_x^3}{n} \frac{\lambda + \sum_{j>k} \lambda_j}{\lambda_k}.
 \end{aligned}$$

1777 The second inequality follows from $\|(A + T)^{-1}\| \leq \|A^{-1}\| (1 - \|A^{-1}\| \|T\|)^{-1}$, where
 1778 both A and $A + T$ are invertible and $\|A^{-1}\| \|T\| < 1$. Note that $c_x > 2$.
 1779

1780 Since $n > 2c_x^3(\lambda + \sum_{j>k} \lambda_j)\lambda_1\lambda_k^{-2}$,

$$\left\| \Sigma_{S,k}^{\frac{1}{2}} V^T M_k^{-\frac{1}{2}} (\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} V \right\|$$

$$\begin{aligned}
& \cdot \left[\left(I_k + V^T M_k^{-\frac{1}{2}} (\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} V \right)^{-1} - I_k \right] \Sigma_{S,k}^{\frac{1}{2}} \Big\| \\
& \leq \|\Sigma_{S,k}\| \|M_k^{-1}\| \|\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T\| \\
& \quad \cdot \left\| \left(I_k + V^T M_k^{-\frac{1}{2}} (\lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) M_k^{-\frac{1}{2}} V \right)^{-1} - I_k \right\| \\
& \leq \lambda_1 \cdot \frac{c_x}{n \lambda_k} \cdot 2c_x \left(\lambda + \sum_{j>k} \lambda_j \right) \cdot \frac{2c_x^3}{n} \frac{\lambda + \sum_{j>k} \lambda_j}{\lambda_k} \\
& = \frac{1}{n} \cdot \frac{4c_x^5}{n} \lambda_1 \lambda_k^{-2} \left(\lambda + \sum_{j>k} \lambda_j \right)^2 \\
& < \frac{2c_x^2}{n} \left(\lambda + \sum_{j>k} \lambda_j \right).
\end{aligned}$$

Combining both terms above, we have

$$\left\| \Sigma_{S,k}^{\frac{1}{2}} X_k^T (XX^T + \lambda I_n)^{-1} X_k \Sigma_{S,k}^{\frac{1}{2}} - \Sigma_{S,k} \right\| \leq \frac{4c_x^2}{n} \left(\lambda + \sum_{j>k} \lambda_j \right).$$

Therefore,

$$\begin{aligned}
& (X_k^T (XX^T + \lambda I_n)^{-1} X_k \beta_k^* - \beta_k^*)^T \Sigma_{T,k} (X_k^T (XX^T + \lambda I_n)^{-1} X_k \beta_k^* - \beta_k^*) \\
& \leq (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^* \cdot \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| \\
& \quad \cdot \left\| \Sigma_{S,k}^{\frac{1}{2}} X_k^T (XX^T + \lambda I_n)^{-1} X_k \Sigma_{S,k}^{\frac{1}{2}} - \Sigma_{S,k} \right\|^2 \\
& \leq \frac{16c_x^4}{n^2} \left(\lambda + \sum_{j>k} \lambda_j \right)^2 (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^* \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|.
\end{aligned}$$

□

Lemma 22. Under the same conditions as in Lemma 18, and on the same event, for any $N_1 < n < N_2$,

$$\begin{aligned}
& (X_k^T (XX^T + \lambda I_n)^{-1} X_{-k} \beta_{-k}^*)^T \Sigma_{T,k} (X_k^T (XX^T + \lambda I_n)^{-1} X_{-k} \beta_{-k}^*) \\
& \leq 16c_x (1 + c_x^4 L^2) \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^*.
\end{aligned}$$

where

$$\begin{aligned}
N_1 &= \max \left\{ 4c_x^4 \left(k + \ln \frac{1}{\delta} \right) \lambda_1^4 \lambda_k^{-4}, \right. \\
&\quad 2c_x^4 L \lambda_1 \lambda_k^{-2} \left(\lambda + \sum_{j>k} \lambda_j \right), \\
&\quad 4c_x^4 \left(k + \ln \frac{1}{\delta} \right) \lambda_1^6 \lambda_k^{-8} \|\Sigma_{T,k}\|^2 \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|^{-2}, \\
&\quad \left. 2c_x^4 L \left(\lambda + \sum_{j>k} \lambda_j \right) \lambda_1^2 \lambda_k^{-4} \|\Sigma_{T,k}\| \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|^{-1} \right\}, \\
N_2 &= \frac{1}{\lambda_{k+1}} \left(\lambda + \sum_{j>k} \lambda_j \right).
\end{aligned}$$

1836 *Proof.*

$$1838 \quad (X_k^T (XX^T + \lambda I_n)^{-1} X_{-k} \beta_{-k}^*)^T \Sigma_{T,k} (X_k^T (XX^T + \lambda I_n)^{-1} X_{-k} \beta_{-k}^*) \\ 1839 \leq \| (XX^T + \lambda I_n)^{-1} X_k \Sigma_{T,k} X_k^T (XX^T + \lambda I_n)^{-1} \| \cdot (\beta_{-k}^*)^T X_{-k}^T X_{-k} \beta_{-k}^*.$$

1840 From Lemma 18,

$$1842 \quad (\beta_{-k}^*)^T X_{-k}^T X_{-k} \beta_{-k}^* \leq c_x n (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^*.$$

1843 In the following, we derive an upper bound for the other term.

$$1844 \quad \| (XX^T + \lambda I_n)^{-1} X_k \Sigma_{T,k} X_k^T (XX^T + \lambda I_n)^{-1} \| \\ 1845 = \left\| (M + \lambda I_n + \Delta)^{-1} \tilde{M}^{\frac{1}{2}} V \Sigma_{T,k} V^T \left(\tilde{M}^{\frac{1}{2}} \right)^T (M + \lambda I_n + \Delta)^{-1} \right\| \\ 1846 = \left\| \Sigma_{T,k}^{\frac{1}{2}} V^T \left(\tilde{M}^{\frac{1}{2}} \right)^T (M + \lambda I_n + \Delta)^{-2} \tilde{M}^{\frac{1}{2}} V \Sigma_{T,k}^{\frac{1}{2}} \right\| \\ 1847 = \left\| \Sigma_{T,k}^{\frac{1}{2}} V^T M_k^{\frac{1}{2}} (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} \right. \\ 1848 \quad \cdot (I_k + \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-2} \Delta_{12}^T) (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} \\ 1849 \quad \cdot M_k^{\frac{1}{2}} V \Sigma_{T,k}^{\frac{1}{2}} \left. \right\| \\ 1850 \leq \| I_k + \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-2} \Delta_{12}^T \| \\ 1851 \quad \cdot \| (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} M_k^{\frac{1}{2}} V \Sigma_{T,k} \\ 1852 \quad \cdot V^T M_k^{\frac{1}{2}} (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} \| \\ 1853 \leq (1 + c_x^4 L^2) \| (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} M_k^{\frac{1}{2}} V \Sigma_{T,k} \\ 1854 \quad \cdot V^T M_k^{\frac{1}{2}} (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} \| \\ 1855 = (1 + c_x^4 L^2) \| (V^T (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1} \\ 1856 \quad \cdot V^T M_k^{\frac{1}{2}} V \Sigma_{T,k} V^T M_k^{\frac{1}{2}} V \\ 1857 \quad \cdot (V^T (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1} \| \\ 1858 \leq (1 + c_x^4 L^2) \| (X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1} \\ 1859 \quad \cdot (X_k^T X_k)^{\frac{1}{2}} \Sigma_{T,k} (X_k^T X_k)^{\frac{1}{2}} \\ 1860 \quad \cdot (X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1} \|.$$

1861 The third equation follows from Equation 5.

1862 We define two quantities that represent concentration error terms:

$$1863 \quad E_1 = \left\| [X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V]^{-1} - (n \tilde{\Sigma}_{S,k})^{-1} \right\|. \\ 1864 \quad E_2 = (X_k^T X_k)^{\frac{1}{2}} - (n \Sigma_{S,k})^{\frac{1}{2}}.$$

1865 Since $n > 4c_x^4 (k + \ln \frac{1}{\delta}) \lambda_1^6 \lambda_k^{-8} \|\Sigma_{T,k}\|^2 \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|^{-2}$,

1866 and $n > 2c_x^4 L \left(\lambda + \sum_{j>k} \lambda_j \right) \lambda_1^2 \lambda_k^{-4} \|\Sigma_{T,k}\| \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|^{-1}$,

$$1867 \quad \|E_1\| \left\| n \tilde{\Sigma}_{S,k} \right\| \left\| (n \tilde{\Sigma}_{S,k})^{-1} \right\| \left\| (n \Sigma_{S,k})^{\frac{1}{2}} \right\| \left\| \Sigma_{T,k} \right\| \left\| (n \Sigma_{S,k})^{\frac{1}{2}} \right\| \left\| (n \tilde{\Sigma}_{S,k})^{-1} \right\| \\ 1868 \leq \frac{c_x^2 \left(\sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 + c_x^2 L \left(\lambda + \sum_{j>k} \lambda_j \right) \right)}{(\lambda + n \lambda_k)^2} (\lambda + n \lambda_1) \frac{n \lambda_1}{(\lambda + n \lambda_k)^2} \|\Sigma_{T,k}\|$$

$$\begin{aligned}
&\leq \frac{c_x^2 \left(\sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 + c_x^2 L \left(\lambda + \sum_{j>k} \lambda_j \right) \right)}{n^2} \frac{\lambda_1^2}{\lambda_k^4} \|\Sigma_{T,k}\| \\
&= \frac{c_x^2 \sqrt{(k + \ln \frac{1}{\delta})}}{n\sqrt{n}} \frac{\lambda_1^3}{\lambda_k^4} \|\Sigma_{T,k}\| + \frac{c_x^4 L \left(\lambda + \sum_{j>k} \lambda_j \right)}{n^2} \frac{\lambda_1^2}{\lambda_k^4} \|\Sigma_{T,k}\| \\
&< \frac{1}{2n} \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| + \frac{1}{2n} \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| \\
&= \frac{1}{n} \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|.
\end{aligned}$$

Similar to Equation 6, since $n > 4c_x^4 (k + \ln \frac{1}{\delta}) \lambda_1^4 \lambda_k^{-4}$ and $n > 2c_x^4 L \left(\lambda + \sum_{j>k} \lambda_j \right) \lambda_1 \lambda_k^{-2}$,

$$\|E_1\| \left\| n \tilde{\Sigma}_{S,k} \right\| < 1.$$

Since $n > c_x^2 (k + \ln \frac{1}{\delta}) \lambda_1^4 \lambda_k^{-6} \|\Sigma_{T,k}\|^2 \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|^{-2}$,

$$\begin{aligned}
&\|E_2\| \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-\frac{1}{2}} \right\| \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right\| \left\| (n \Sigma_{S,k})^{\frac{1}{2}} \right\| \left\| \Sigma_{T,k} \right\| \left\| (n \Sigma_{S,k})^{\frac{1}{2}} \right\| \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right\| \\
&\leq c_x \sqrt{k + \ln \frac{1}{\delta}} \lambda_1 \lambda_k^{-\frac{1}{2}} (n \lambda_k)^{-\frac{1}{2}} \frac{n \lambda_1}{(\lambda + n \lambda_k)^2} \|\Sigma_{T,k}\| \\
&\leq \frac{c_x \sqrt{k + \ln \frac{1}{\delta}}}{n \sqrt{n}} \frac{\lambda_1^2}{\lambda_k^3} \|\Sigma_{T,k}\| \\
&\leq \frac{1}{n} \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|.
\end{aligned}$$

Similar to Equation 7, since $n > c_x^2 (k + \ln \frac{1}{\delta}) \lambda_1^2 \lambda_k^{-2}$,

$$\|E_2\| \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-\frac{1}{2}} \right\| < 1.$$

Combining the four inequalities above,

$$\begin{aligned}
&\left\| \left(X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12} (\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V \right)^{-1} \right. \\
&\quad \cdot \left. \left(X_k^T X_k \right)^{\frac{1}{2}} \Sigma_{T,k} \left(X_k^T X_k \right)^{\frac{1}{2}} \right. \\
&\quad \cdot \left. \left(X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12} (\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V \right)^{-1} \right\| \\
&\leq \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-1} (n \Sigma_{S,k})^{\frac{1}{2}} \Sigma_{T,k} (n \Sigma_{S,k})^{\frac{1}{2}} \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right\| \\
&\quad + 2 \left\| E_1 (n \Sigma_{S,k})^{\frac{1}{2}} \Sigma_{T,k} (n \Sigma_{S,k})^{\frac{1}{2}} \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right\| \\
&\quad + 2 \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-1} E_2 \Sigma_{T,k} (n \Sigma_{S,k})^{\frac{1}{2}} \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right\| \\
&\quad + \left\| E_1 (n \Sigma_{S,k})^{\frac{1}{2}} \Sigma_{T,k} (n \Sigma_{S,k})^{\frac{1}{2}} E_1 \right\| \\
&\quad + \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-1} E_2 \Sigma_{T,k} E_2 \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right\| \\
&\quad + 2 \left\| E_1 (n \Sigma_{S,k})^{\frac{1}{2}} \Sigma_{T,k} E_2 \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right\| \\
&\quad + 2 \left\| E_1 E_2 \Sigma_{T,k} (n \Sigma_{S,k})^{\frac{1}{2}} \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right\| \\
&\quad + 2 \left\| E_1 E_2 \Sigma_{T,k} (n \Sigma_{S,k})^{\frac{1}{2}} E_1 \right\|
\end{aligned}$$

$$\begin{aligned} & + 2 \left\| E_1 E_2 \Sigma_{T,k} E_2 \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right\| \\ & + \| E_1 E_2 \Sigma_{T,k} E_2 E_1 \| . \end{aligned}$$

In particular,

$$\begin{aligned} & \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-1} (n \Sigma_{S,k})^{\frac{1}{2}} \Sigma_{T,k} (n \Sigma_{S,k})^{\frac{1}{2}} \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right\| \\ & = \frac{1}{n} \left\| \tilde{\Sigma}_{S,k}^{-1} \Sigma_{S,k}^{\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{\frac{1}{2}} \tilde{\Sigma}_{S,k}^{-1} \right\| \\ & \leq \frac{1}{n} \left\| \Sigma_{S,k}^{-1} \Sigma_{S,k}^{\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{\frac{1}{2}} \Sigma_{S,k}^{-1} \right\| \\ & = \frac{1}{n} \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|. \end{aligned}$$

The inequality follows from the fact that $\|B A B\| = \|A^{\frac{1}{2}} B A^{\frac{1}{2}}\| \leq \|A^{\frac{1}{2}} C A^{\frac{1}{2}}\| = \|C A C\|$, where A, B, C are positive semi-definite matrices, and $C \succcurlyeq B$, which implies that $A^{\frac{1}{2}} C A^{\frac{1}{2}} \succcurlyeq A^{\frac{1}{2}} B A^{\frac{1}{2}}$.

$$\begin{aligned} & \|E_1 E_2 \Sigma_{T,k} E_2 E_1\| \\ & = \left\| E_1 n \tilde{\Sigma}_{S,k} \left(n \tilde{\Sigma}_{S,k} \right)^{-1} E_2 \left(n \tilde{\Sigma}_{S,k} \right)^{-\frac{1}{2}} \left(n \tilde{\Sigma}_{S,k} \right)^{\frac{1}{2}} \right. \\ & \quad \cdot \Sigma_{T,k} E_2 \left(n \tilde{\Sigma}_{S,k} \right)^{-\frac{1}{2}} \left(n \tilde{\Sigma}_{S,k} \right)^{\frac{1}{2}} E_1 n \tilde{\Sigma}_{S,k} \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \Big\| \\ & \leq \left(\|E_1\| \left\| n \tilde{\Sigma}_{S,k} \right\| \right)^2 \left(\|E_2\| \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-\frac{1}{2}} \right\| \right) \\ & \quad \cdot \|E_2\| \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-\frac{1}{2}} \right\| \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right\| \left\| (n \Sigma_{S,k})^{\frac{1}{2}} \right\| \left\| \Sigma_{T,k} \right\| \left\| (n \Sigma_{S,k})^{\frac{1}{2}} \right\| \left\| \left(n \tilde{\Sigma}_{S,k} \right)^{-1} \right\| \\ & \leq \frac{1}{n} \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|. \end{aligned}$$

The other terms can be similarly bounded. Therefore,

$$\begin{aligned} & \left\| (X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12} (\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1} \right. \\ & \quad \cdot (X_k^T X_k)^{\frac{1}{2}} \Sigma_{T,k} (X_k^T X_k)^{\frac{1}{2}} \\ & \quad \cdot (X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12} (\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1} \Big\| \\ & \leq \frac{16}{n} \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|. \end{aligned}$$

□

B.3.4 BIAS IN THE LAST $d - k$ DIMENSIONS

The upper bound for the bias in the last $d - k$ dimensions is extended from [Tsigler & Bartlett \(2023\)](#)'s Lemma 28. The bias can be decomposed into three terms.

$$\begin{aligned} & \left(\widehat{\beta}(X\beta^*)_{-k} - \beta_{-k}^* \right)^T \Sigma_{T,-k} \left(\widehat{\beta}(X\beta^*)_{-k} - \beta_{-k}^* \right) \\ & \leq 3(\beta_{-k}^*)^T \Sigma_{T,-k} \beta_{-k}^* \\ & \quad + 3(\beta_{-k}^*)^T X_{-k}^T (XX^T + \lambda I_n)^{-1} X_{-k} \Sigma_{T,-k} X_{-k}^T (XX^T + \lambda I_n)^{-1} X_{-k} \beta_{-k}^* \\ & \quad + 3(\beta_k^*)^T X_k^T (XX^T + \lambda I_n)^{-1} X_{-k} \Sigma_{T,-k} X_{-k}^T (XX^T + \lambda I_n)^{-1} X_k \beta_k^*. \end{aligned}$$

Lemma 23. Under the same conditions as in Lemma 18, and on the same event, for any $N_1 < n < N_2$,

$$(\beta_{-k}^*)^T X_{-k}^T (XX^T + \lambda I_n)^{-1} X_{-k} \Sigma_{T,-k} X_{-k}^T (XX^T + \lambda I_n)^{-1} X_{-k} \beta_{-k}^*$$

$$\begin{aligned} & \leq c_x^2 L \left(\lambda + \sum_j \lambda_j \right)^{-1} n \|\Sigma_{T,-k}\| (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}. \end{aligned}$$

where N_1, N_2 are defined as in Lemma 18.

Proof.

$$\begin{aligned} & (\beta_{-k}^*)^T X_{-k}^T (XX^T + \lambda I_n)^{-1} X_{-k} \Sigma_{T,-k} X_{-k}^T (XX^T + \lambda I_n)^{-1} X_{-k} \beta_{-k}^* \\ & \leq \|\Sigma_{T,-k}\| (\beta_{-k}^*)^T X_{-k}^T (XX^T + \lambda I_n)^{-1} X_{-k} X_{-k}^T (XX^T + \lambda I_n)^{-1} X_{-k} \beta_{-k}^* \\ & \leq \|\Sigma_{T,-k}\| (\beta_{-k}^*)^T X_{-k}^T (XX^T + \lambda I_n)^{-1} (XX^T + \lambda I_n) (XX^T + \lambda I_n)^{-1} X_{-k} \beta_{-k}^* \\ & \leq \|\Sigma_{T,-k}\| \|(XX^T + \lambda I_n)^{-1}\| (\beta_{-k}^*)^T X_{-k}^T X_{-k} \beta_{-k}^* \\ & \leq \|\Sigma_{T,-k}\| \|(X_{-k} X_{-k}^T + \lambda I_n)^{-1}\| (\beta_{-k}^*)^T X_{-k}^T X_{-k} \beta_{-k}^* \\ & \leq \|\Sigma_{T,-k}\| \left(\frac{1}{c_x L} \left(\lambda + \sum_j \lambda_j \right) \right)^{-1} c_x n (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^* \\ & = c_x^2 L \left(\lambda + \sum_j \lambda_j \right)^{-1} n \|\Sigma_{T,-k}\| (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^*. \end{aligned}$$

The fourth inequality follows from $XX^T + \lambda I_n \succcurlyeq X_{-k} X_{-k}^T + \lambda I_n$. \square

Lemma 24. Under the same conditions as in Lemma 18, and on the same event, for any $N_1 < n < N_2$,

$$\begin{aligned} & (\beta_k^*)^T X_k^T (XX^T + \lambda I_n)^{-1} X_{-k} \Sigma_{T,-k} X_{-k}^T (XX^T + \lambda I_n)^{-1} X_k \beta_k^* \\ & \leq \frac{c_x^6}{n} L \left(\lambda + \sum_{j>k} \lambda_j \right) \|\Sigma_{T,-k}\| (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^*. \end{aligned}$$

where N_1, N_2 are defined as in Lemma 18.

Proof. It can be verified by Woodbury matrix identity that:

$$(XX^T + \lambda I_n)^{-1} X_k = (X_{-k} X_{-k}^T + \lambda I_n)^{-1} X_k (I_k + X_k^T (X_{-k} X_{-k}^T + \lambda I_n)^{-1} X_k)^{-1}.$$

Therefore,

$$\begin{aligned} & (\beta_k^*)^T X_k^T (XX^T + \lambda I_n)^{-1} X_{-k} \Sigma_{T,-k} X_{-k}^T (XX^T + \lambda I_n)^{-1} X_k \beta_k^* \\ & = \left\| \Sigma_{T,-k}^{\frac{1}{2}} X_{-k}^T (X_{-k} X_{-k}^T + \lambda I_n)^{-1} X_k (I_k + X_k^T (X_{-k} X_{-k}^T + \lambda I_n)^{-1} X_k)^{-1} \beta_k^* \right\|^2 \\ & \leq \|\Sigma_{T,-k}\| \|(X_{-k} X_{-k}^T + \lambda I_n)^{-1} X_{-k} X_{-k}^T (X_{-k} X_{-k}^T + \lambda I_n)^{-1}\| \\ & \quad \cdot \left\| X_k (I_k + X_k^T (X_{-k} X_{-k}^T + \lambda I_n)^{-1} X_k)^{-1} \beta_k^* \right\|^2 \\ & = \|\Sigma_{T,-k}\| \|(X_{-k} X_{-k}^T + \lambda I_n)^{-1} X_{-k} X_{-k}^T (X_{-k} X_{-k}^T + \lambda I_n)^{-1}\| \\ & \quad \cdot \left\| X_k \Sigma_{S,k}^{-\frac{1}{2}} \left(\Sigma_{S,k}^{-1} + \Sigma_{S,k}^{-\frac{1}{2}} X_k^T (X_{-k} X_{-k}^T + \lambda I_n)^{-1} X_k \Sigma_{S,k}^{-\frac{1}{2}} \right)^{-1} \Sigma_{S,k}^{-\frac{1}{2}} \beta_k^* \right\|^2 \\ & \leq \|\Sigma_{T,-k}\| \|(X_{-k} X_{-k}^T + \lambda I_n)^{-1}\| \left\| \Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right\| \\ & \quad \cdot \left\| \left(\Sigma_{S,k}^{-1} + \Sigma_{S,k}^{-\frac{1}{2}} X_k^T (X_{-k} X_{-k}^T + \lambda I_n)^{-1} X_k \Sigma_{S,k}^{-\frac{1}{2}} \right)^{-2} \right\| (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^*. \end{aligned}$$

In particular,

$$\left\| \left(\Sigma_{S,k}^{-1} + \Sigma_{S,k}^{-\frac{1}{2}} X_k^T (X_{-k} X_{-k}^T + \lambda I_n)^{-1} X_k \Sigma_{S,k}^{-\frac{1}{2}} \right)^{-2} \right\|$$

$$\begin{aligned}
&\leq \left\| \left(\Sigma_{S,k}^{-\frac{1}{2}} X_k^T (X_{-k} X_{-k}^T + \lambda I_n)^{-1} X_k \Sigma_{S,k}^{-\frac{1}{2}} \right)^{-1} \right\| \\
&\leq \|X_{-k} X_{-k}^T + \lambda I_n\| \left\| \left(\Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right)^{-1} \right\| \\
&\leq c_x \left(\lambda + \sum_{j>k} \lambda_j \right) \frac{c_x}{n} \\
&= \frac{c_x^2}{n} \left(\lambda + \sum_{j>k} \lambda_j \right).
\end{aligned}$$

The second inequality follows from $\mu_{\min}(ABA^T) \geq \mu_{\min}(B)\mu_{\min}(AA^T)$ where the matrix B is positive definite.

Therefore,

$$\begin{aligned}
&(\beta_k^*)^T X_k^T (XX^T + \lambda I_n)^{-1} X_{-k} \Sigma_{T,-k} X_{-k}^T (XX^T + \lambda I_n)^{-1} X_k \beta_k^* \\
&\leq \|\Sigma_{T,-k}\| \| (X_{-k} X_{-k}^T + \lambda I_n)^{-1} \| \left\| \Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right\| \\
&\quad \cdot \left\| \left(\Sigma_{S,k}^{-1} + \Sigma_{S,k}^{-\frac{1}{2}} X_k^T (X_{-k} X_{-k}^T + \lambda I_n)^{-1} X_k \Sigma_{S,k}^{-\frac{1}{2}} \right)^{-2} \right\| (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^* \\
&\leq \|\Sigma_{T,-k}\| \cdot \left(\frac{1}{c_x L} \left(\lambda + \sum_{j>k} \lambda_j \right) \right)^{-1} \cdot c_x n \cdot \frac{c_x^4}{n^2} \left(\lambda + \sum_{j>k} \lambda_j \right)^2 \cdot (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^* \\
&= \frac{c_x^6}{n} L \left(\lambda + \sum_{j>k} \lambda_j \right) \|\Sigma_{T,-k}\| (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^*.
\end{aligned}$$

□

B.4 MAIN RESULTS

Theorem 25. Let $\mathcal{T} = \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}}$ and $\mathcal{U} = \Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}}$. There exists a constant $c > 2$ depending only on σ , such that for any $cN < n < r_k$, if the assumption $\text{condNum}(k, \delta, L)$ (Assumption 2) is satisfied, then with probability at least $1 - 2\delta - ce^{-n/c}$,

$$\begin{aligned}
\frac{V}{cv^2} &\leq L^2 \frac{\text{tr}[\mathcal{T}]}{n} + L^2 \frac{n \text{tr}[\mathcal{U}]}{(\lambda + \sum_{j>k} \lambda_j)^2}. \\
\frac{B}{c} &\leq \|\beta_k^*\|_{\Sigma_{S,k}^{-1}}^2 \left(\frac{\lambda + \sum_{j>k} \lambda_j}{n} \right)^2 \left[\|\mathcal{T}\| + L \frac{n \|\Sigma_{T,-k}\|}{\lambda + \sum_{j>k} \lambda_j} \right] \\
&\quad + \|\beta_{-k}^*\|_{\Sigma_{S,-k}}^2 \left[L^2 \|\mathcal{T}\| + L \frac{n \|\Sigma_{T,-k}\|}{\lambda + \sum_{j>k} \lambda_j} \right].
\end{aligned}$$

N is defined as follows:

$$\begin{aligned}
N = \max \left\{ &\left(k + \ln \frac{1}{\delta} \right) \lambda_1^6 \lambda_k^{-8} \|\Sigma_{T,k}\|^2 k^2 (\text{tr}[\mathcal{T}])^{-2}, \right. \\
&\left. L \lambda_1^2 \lambda_k^{-4} (\lambda + \sum_{j>k} \lambda_j) \|\Sigma_{T,k}\| k (\text{tr}[\mathcal{T}])^{-1} \right\}.
\end{aligned}$$

Remark 8 (Sample complexity). We have assumed $n > cN$ in the theorem. The first condition on N indicates $n \gg k$. From the inequality $\lambda_k^2 \leq \|\Sigma_{T,k}\|^2 k^2 (\text{tr}[\mathcal{T}])^{-2} \leq k^2 \lambda_1^2$, it follows that $n = \Omega(k)$ in the best case, consistent with the sample complexity of classic linear regression.

This optimal case occurs when $\Sigma_{S,k} \approx \Sigma_{T,k}$. In the worst case, $n = \Omega(k^3)$ where covariate shift is significant in the first k dimensions—e.g., when the test data lies predominantly in the subspace of the first dimension. This shift in sample complexity under varying degrees of covariate shift parallels the analysis of Ge et al. (2024) (see their Theorem 4.2) for the under-parameterized setting. The second condition implies $n \gg \lambda + \sum_{j>k} \lambda_j$, such that the regularization is not too strong to introduce a bias greater than a constant (as shown in the first bias term). On the other hand, we assume $n < r_k$ in the theorem, which is consistent with the over-parameterized regime and Assumption 1, where the last $d - k$ components are considered to be essentially high-dimensional.

Proof. The theorem follows from Lemma 18, Lemma 19, Lemma 20, Lemma 21, Lemma 22, Lemma 23 and Lemma 24. For a constant $c'_x > 2$ depending only on σ , these lemmas hold for values of n that satisfy the following inequalities:

$$\begin{aligned} n &> 4c'^4_x(k + \ln(1/\delta))\lambda_1^2\lambda_k^{-2}, \\ n &> 2c'^4_x L \lambda_k^{-1} \left(\lambda + \sum_{j>k} \lambda_j \right), \\ n &> 4c'^4_x \left(k + \ln \frac{1}{\delta} \right) \lambda_1^4 \lambda_k^{-4}, \\ n &> 2c'^4_x L \lambda_1 \lambda_k^{-2} \left(\lambda + \sum_{j>k} \lambda_j \right), \\ n &> 4c'^4_x \left(k + \ln \frac{1}{\delta} \right) \lambda_1^6 \lambda_k^{-8} \|\Sigma_{T,k}\|^2 k^2 \left(\text{tr} \left[\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right] \right)^{-2}, \\ n &> 2c'^4_x L \lambda_1^2 \lambda_k^{-4} \left(\lambda + \sum_{j>k} \lambda_j \right) \|\Sigma_{T,k}\| k \left(\text{tr} \left[\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right] \right)^{-1}, \\ n &> 2c'^3_x (\lambda + \sum_{j>k} \lambda_j) \lambda_1 \lambda_k^{-2}, \\ n &> 4c'^4_x \left(k + \ln \frac{1}{\delta} \right) \lambda_1^6 \lambda_k^{-8} \|\Sigma_{T,k}\|^2 \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|^{-2}, \\ n &> 2c'^4_x L \left(\lambda + \sum_{j>k} \lambda_j \right) \lambda_1^2 \lambda_k^{-4} \|\Sigma_{T,k}\| \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|^{-1}, \\ n &< \lambda_{k+1}^{-1} \left(\lambda + \sum_{j>k} \lambda_j \right). \end{aligned}$$

A sufficient condition for all the inequalities above is given by $4c'^4_x N_1 < n < r_k$. This follows from the following facts:

$$\begin{aligned} \lambda_1 \lambda_k^{-1} &\geq 1, \\ c'_x &> 2, \\ L &\geq 1, \\ k \left(\text{tr} \left[\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right] \right)^{-1} &\geq \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|^{-1}, \\ k \|\Sigma_{T,k}\| \left(\text{tr} \left[\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right] \right)^{-1} &\geq \lambda_k. \end{aligned}$$

Then, with probability at least $1 - 2\delta - c'_x e^{-n/c'_x}$:

$$V/2 \leq 16v^2(1 + c'^4_x L^2) \frac{1}{n} \text{tr} \left[\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right]$$

$$\begin{aligned}
& + v^2 c_x'^3 L^2 n \left(\lambda + \sum_{j>k} \lambda_j \right)^{-2} \text{tr} \left[\Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}} \right] \\
& \leq 32v^2 c_x'^4 L^2 \frac{1}{n} \text{tr} \left[\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right] \\
& \quad + v^2 c_x'^3 L^2 n \left(\lambda + \sum_{j>k} \lambda_j \right)^{-2} \text{tr} \left[\Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}} \right], \\
B/2 & \leq \frac{16c_x'^4}{n^2} \left(\lambda + \sum_{j>k} \lambda_j \right)^2 (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^* \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| \\
& \quad + 32c_x'(1 + c_x'^4 L^2) \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^* \\
& \quad + 3c_x'^2 L \left(\lambda + \sum_j \lambda_j \right)^{-1} n \|\Sigma_{T,-k}\| (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^* \\
& \quad + 3 \frac{c_x'^6}{n} L \left(\lambda + \sum_{j>k} \lambda_j \right) \|\Sigma_{T,-k}\| (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^* \\
& \quad + 3(\beta_{-k}^*)^T \Sigma_{T,-k} \beta_{-k}^* \\
& \leq 16c_x'^4 \frac{1}{n^2} \left(\lambda + \sum_{j>k} \lambda_j \right)^2 \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^* \\
& \quad + 64c_x'^5 L^2 \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^* \\
& \quad + 3c_x'^2 L n \left(\lambda + \sum_j \lambda_j \right)^{-1} \|\Sigma_{T,-k}\| (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^* \\
& \quad + 3c_x'^6 L \frac{1}{n} \left(\lambda + \sum_{j>k} \lambda_j \right) \|\Sigma_{T,-k}\| (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^* \\
& \quad + 3(\beta_{-k}^*)^T \Sigma_{T,-k} \beta_{-k}^* \\
& \leq 16c_x'^4 \frac{1}{n^2} \left(\lambda + \sum_{j>k} \lambda_j \right)^2 \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^* \\
& \quad + 64c_x'^5 L^2 \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^* \\
& \quad + 3c_x'^2 L n \left(\lambda + \sum_{j>k} \lambda_j \right)^{-1} \|\Sigma_{T,-k}\| (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^* \\
& \quad + 3c_x'^6 L \frac{1}{n} \left(\lambda + \sum_{j>k} \lambda_j \right) \|\Sigma_{T,-k}\| (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^* \\
& \quad + 3c_x'^5 L^2 \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^*.
\end{aligned}$$

The last inequality follows from:

$$\begin{aligned}
(\beta_{-k}^*)^T \Sigma_{T,-k} \beta_{-k}^* & = (\beta_{-k}^*)^T \Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{S,-k}^{-\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{-\frac{1}{2}} \Sigma_{S,-k}^{\frac{1}{2}} \beta_{-k}^* \\
& \leq \left\| \Sigma_{S,-k}^{-\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{-\frac{1}{2}} \right\| (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^*.
\end{aligned}$$

2214 By taking $c = 134c_x'^6$, the proof is complete. \square
 2215

2216 **Corollary 26** (Restatement of Theorem 2). Let $\mathcal{T} = \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}}$, $\mathcal{U} = \Sigma_{S,-k} \Sigma_{T,-k}$ and $\mathcal{V} =$
 2217 $\Sigma_{S,-k}^2$. There exists a constant $c > 2$ depending only on σ, L , such that for any $cN < n < r_k$, if the
 2218 assumption $\text{condNum}(k, \delta, L)$ (Assumption 2) is satisfied, then with probability at least $1 - 3\delta$,
 2219

$$\begin{aligned} 2220 \frac{V}{cv^2} &\leq \frac{k}{n} \frac{\text{tr}[\mathcal{T}]}{k} + \frac{n}{R_k} \frac{\text{tr}[\mathcal{U}]}{\text{tr}[\mathcal{V}]} \\ 2221 \\ 2222 \frac{B}{c} &\leq \left(\|\beta_k^*\|_{\Sigma_{S,k}^{-1}}^2 \left(\frac{\lambda + \sum_{j>k} \lambda_j}{n} \right)^2 + \|\beta_{-k}^*\|_{\Sigma_{S,-k}}^2 \right) \left[\|\mathcal{T}\| + \frac{n}{r_k} \frac{\|\Sigma_{T,-k}\|}{\|\Sigma_{S,-k}\|} \right]. \end{aligned}$$

2225 N is a polynomial function of $k + \ln(1/\delta), \lambda_1 \lambda_k^{-1}, 1 + (\lambda + \sum_{j>k} \lambda_j) \lambda_k^{-1}$.
 2226

2227 *Proof.* The first variance term follows directly from Theorem 25.

2228 For the second variance term, by plugging in the definition of R_k ,

$$\begin{aligned} 2230 L^2 \frac{n \text{tr}[\mathcal{U}]}{\left(\lambda + \sum_{j>k} \lambda_j \right)^2} &= L^2 \frac{n}{R_k} \frac{\text{tr}[\Sigma_{S,-k} \Sigma_{T,-k}]}{\sum_{j>k} \lambda_j^2} \\ 2231 \\ 2233 &= L^2 \frac{n}{R_k} \frac{\text{tr}[\mathcal{U}]}{\text{tr}[\mathcal{V}]} . \end{aligned}$$

2235 For the first bias term, by plugging in the definition of r_k ,

$$\begin{aligned} 2237 \|\beta_k^*\|_{\Sigma_{S,k}^{-1}}^2 \left(\frac{\lambda + \sum_{j>k} \lambda_j}{n} \right)^2 \left[\|\mathcal{T}\| + L \frac{n \|\Sigma_{T,-k}\|}{\lambda + \sum_{j>k} \lambda_j} \right] \\ 2238 \\ 2240 &= \|\beta_k^*\|_{\Sigma_{S,k}^{-1}}^2 \left(\frac{\lambda + \sum_{j>k} \lambda_j}{n} \right)^2 \left[\|\mathcal{T}\| + L \frac{n}{r_k} \frac{\|\Sigma_{T,-k}\|}{\lambda_{k+1}} \right]. \end{aligned}$$

2242 Similarly, the second bias term can be transformed into:

$$2244 \|\beta_{-k}^*\|_{\Sigma_{S,-k}}^2 \left[L^2 \|\mathcal{T}\| + L \frac{n \|\Sigma_{T,-k}\|}{\lambda + \sum_{j>k} \lambda_j} \right] = \|\beta_{-k}^*\|_{\Sigma_{S,-k}}^2 \left[L^2 \|\mathcal{T}\| + L \frac{n}{r_k} \frac{\|\Sigma_{T,-k}\|}{\lambda_{k+1}} \right].$$

2246 Since the statement of Theorem 25 holds with probability at least $1 - 2\delta - ce^{-n/c}$, we only require
 2247 $ce^{-n/c} < \delta$, which is equivalent as $n > c \ln c + c \ln(1/\delta)$. Combining the lower bounds of n in
 2248 Theorem 25, we should have:
 2249

$$\begin{aligned} 2250 n &> \max \left\{ c \ln c + c \ln \frac{1}{\delta}, \right. \\ 2251 \\ 2253 &\quad c \left(k + \ln \frac{1}{\delta} \right) \lambda_1^6 \lambda_k^{-8} \|\Sigma_{T,k}\|^2 k^2 (\text{tr}[\mathcal{T}])^{-2}, \\ 2255 \\ 2257 &\quad \left. c L \lambda_1^2 \lambda_k^{-4} \left(\lambda + \sum_{j>k} \lambda_j \right) \|\Sigma_{T,k}\| k (\text{tr}[\mathcal{T}])^{-1} \right\}. \end{aligned}$$

2258 For the first term in the maximum argument,

$$\begin{aligned} 2259 c \ln c + c \ln \frac{1}{\delta} &\leq c^2 + c \ln \frac{1}{\delta} \\ 2260 \\ 2263 &\leq c^2 \left(k + \ln \frac{1}{\delta} \right). \end{aligned}$$

2264 The second term:

$$\begin{aligned} 2265 c \left(k + \ln \frac{1}{\delta} \right) \lambda_1^6 \lambda_k^{-8} \|\Sigma_{T,k}\|^2 k^2 (\text{tr}[\mathcal{T}])^{-2} \\ 2266 \\ 2267 &\leq c \left(k + \ln \frac{1}{\delta} \right) \lambda_1^6 \lambda_k^{-8} \|\Sigma_{T,k}\|^2 k^2 \left(\mu_k(\Sigma_{S,k}^{-1}) \text{tr}[\Sigma_{T,k}] \right)^{-2} \end{aligned}$$

$$\begin{aligned}
&\leq c \left(k + \ln \frac{1}{\delta} \right) \lambda_1^8 \lambda_k^{-8} \|\Sigma_{T,k}\|^2 k^2 \|\Sigma_{T,k}\|^{-2} \\
&= c \left(k + \ln \frac{1}{\delta} \right)^3 \lambda_1^8 \lambda_k^{-8}.
\end{aligned}$$

The first inequality follows from $\text{tr}[MN] \geq \mu_{\min}(M) \text{tr}[N]$ for positive semi-definite matrices M, N .

Similar, for the third term:

$$\begin{aligned}
&cL\lambda_1^2\lambda_k^{-4}(\lambda + \sum_{j>k}\lambda_j)\|\Sigma_{T,k}\|k(\text{tr}[\mathcal{T}])^{-1} \\
&\leq cL\lambda_1^2\lambda_k^{-4}(\lambda + \sum_{j>k}\lambda_j)\|\Sigma_{T,k}\|k\lambda_1\|\Sigma_{T,k}\|^{-1} \\
&\leq cL(k + \ln \frac{1}{\delta})\lambda_1^3\lambda_k^{-4}(\lambda + \sum_{j>k}\lambda_j).
\end{aligned}$$

The proof is complete by taking c as c^2L^2 and $N = (k + \ln \frac{1}{\delta})^3 (\lambda_1 \lambda_k^{-1})^8 [1 + (\lambda + \sum_{j>k} \lambda_j) \lambda_k^{-1}]$. \square

C LARGE SHIFT IN MINOR DIRECTIONS

In this section, we consider the scenario where the signal β^* mainly concentrate on the first k components (here we choose the basis to be the eigenvectors of Σ_S), but the target covariance Σ_T may not be small on the last $d - k$ components.

C.1 LOWER BOUND FOR RIDGE REGRESSION

In this subsection, we will show that the original ridge regression algorithm will not work under this scenario.

Recall our model:

$$y = \beta^{*T} x + \epsilon, \quad (8)$$

We can write our data as

$$Y = X\beta^* + \epsilon, \quad (9)$$

where $Y = (y_1, \dots, y_n)^T \in \mathbb{R}^{n \times 1}$, $X = (x_1, \dots, x_n)^T \in \mathbb{R}^{n \times d}$, $\epsilon = (\epsilon_1, \dots, \epsilon_n)^T \in \mathbb{R}^{n \times 1}$. We denote by $\widehat{\Sigma}_S := \frac{1}{n} X^T X$ the sample covariance matrix.

Assume the same assumptions as in our previous section still holds. We let $\Sigma_S = \mathbb{E}[x_i x_i^T]$ be the following: its eigenvalues $\lambda_1, \dots, \lambda_d$ satisfies $\lambda_1 = \dots = \lambda_k = 1$, $\lambda_{k+1} = \dots = \lambda_{k+\lfloor \sqrt{n}/C_2 \rfloor} = C_1/\sqrt{n}$ for sufficiently large constants C_1, C_2 , and the remaining eigenvalues are all set to zero. We let $\Sigma_T = I_d$. Then the excess risk is $\mathbb{E}_\epsilon[(\widehat{\beta} - \beta^*)^T \Sigma_T (\widehat{\beta} - \beta^*)] = \mathbb{E}_\epsilon \|\widehat{\beta} - \beta^*\|^2$. We will show that under this scenario, ridge regression can not obtain an error rate of $\mathcal{O}(\frac{1}{n})$. To see this, we explicitly write out the ridge solution:

$$\begin{aligned}
\widehat{\beta} &= (X^T X + \lambda I_d)^{-1} X^T Y \\
&= (\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-1} (\frac{1}{n} X^T Y) \\
&= (\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-1} (\frac{1}{n} X^T (X\beta^* + \epsilon)) \\
&= (\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-1} (\frac{1}{n} X^T X\beta^* + \frac{1}{n} X^T \epsilon) \\
&= (\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-1} (\widehat{\Sigma}_S \beta^* + \frac{1}{n} X^T \epsilon)
\end{aligned}$$

$$= (\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-1} \widehat{\Sigma}_S \beta^* + (\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-1} \frac{1}{n} X^T \epsilon. \quad (10)$$

Therefore

$$\begin{aligned} \widehat{\beta} - \beta^* &= (\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-1} \widehat{\Sigma}_S \beta^* - \beta^* + (\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-1} \frac{1}{n} X^T \epsilon \\ &= (\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-1} \widehat{\Sigma}_S \beta^* - (\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-1} (\widehat{\Sigma}_S + \frac{\lambda}{n} I_d) \beta^* + (\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-1} \frac{1}{n} X^T \epsilon \\ &= -\frac{\lambda}{n} (\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-1} \beta^* + (\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-1} \frac{1}{n} X^T \epsilon \end{aligned}$$

Taking expectation with respect to ϵ ,

$$\begin{aligned} \mathbb{E}_\epsilon \|\widehat{\beta} - \beta^*\|^2 &= \frac{\lambda^2}{n^2} \|(\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-1} \beta^*\|^2 + \frac{1}{n^2} \text{tr}(\epsilon^T X (\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-2} X^T \epsilon) \\ &= \frac{\lambda^2}{n^2} \|(\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-1} \beta^*\|^2 + v^2 \frac{1}{n} \text{tr}((\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-2} \widehat{\Sigma}_S) \\ &:= B + V \end{aligned} \quad (11)$$

where $B = \frac{\lambda^2}{n^2} \|(\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-1} \beta^*\|^2$ is the bias, $V = \frac{v^2}{n} \text{tr}((\widehat{\Sigma}_S + \frac{\lambda}{n} I_d)^{-2} \widehat{\Sigma}_S)$ is the variance. We state the formal version of Theorem 4 in the following:

Theorem 27. Under the instance we consider, namely $\lambda_1, \dots, \lambda_d$ satisfies $\lambda_1 = \dots = \lambda_k = 1$, $\lambda_{k+1} = \dots = \lambda_{k+\lfloor \sqrt{n}/C_2 \rfloor} = C_1/\sqrt{n}$, $\lambda_{k+\lfloor \sqrt{n}/C_2 \rfloor + 1} = \dots = \lambda_d = 0$. WLOG assume $\sigma = 1$, $C_2 \geq C_1((\frac{C_1}{4C})^2 - k - \log \frac{1}{\delta})^{-1}$ for some absolute constant C , and $n \geq (\frac{3C_1}{2})^4$. With probability $1 - \delta$, when $\lambda = c\sqrt{n}$, we have $\frac{V}{v^2} \geq C'$, where $C' > 0$ is some absolute constant. When $\lambda \leq n^{3/4}$, we have $\frac{V}{v^2} \geq C' \frac{1}{\sqrt{n}}$. When $\lambda \geq n^{3/4}$, $B \geq \frac{\|\beta^*\|^2}{9\sqrt{n}}$.

Proof. We will use the following concentration lemma modified from (Vershynin, 2018, Exercise 9.2.5):

Lemma 28. Let $\{x_i\}_{i=1}^n$ be i.i.d. d -dimensional random vectors, satisfying: x_i is mean zero, $\mathbb{E}[xx^T] = \Sigma$ and is $\sigma^2 \Sigma$ -sub-gaussian, in the sense that

$$\mathbb{E}[\exp(v^T x_i)] \leq \exp\left(\frac{\|\sigma \Sigma^{1/2} v\|^2}{2}\right).$$

$X = (x_1, \dots, x_n)^T \in \mathbb{R}^{n \times d}$. Then with probability $1 - \delta$,

$$\|\widehat{\Sigma} - \Sigma\| \leq C \sigma^4 \left(\sqrt{\frac{r + \log \frac{1}{\delta}}{n}} + \frac{r + \log \frac{1}{\delta}}{n} \right) \|\Sigma\|$$

where $r := \text{tr}(\Sigma)/\|\Sigma\|$ is the stable rank of Σ , C is an absolute constant.

Applying Lemma 28, we have

$$\|\widehat{\Sigma}_S - \Sigma_S\| \leq C \left(\sqrt{\frac{r + \log \frac{1}{\delta}}{n}} + \frac{r + \log \frac{1}{\delta}}{n} \right)$$

where $r = \sum_{i=1}^d \lambda_i = k + \lfloor \sqrt{n}/C_2 \rfloor \frac{C_1}{\sqrt{n}} \leq k + C_1/C_2$. When $n \geq C_1/C_2 + k + \log \frac{1}{\delta}$, we have

$$\|\widehat{\Sigma}_S - \Sigma_S\| \leq 2C \sqrt{\frac{C_1/C_2 + k + \log \frac{1}{\delta}}{n}}.$$

We denote by $\widehat{\lambda}_1 \geq \dots \geq \widehat{\lambda}_d$ the eigenvalues of $\widehat{\Sigma}_S$. Then by Weyl's inequality (Chen et al., 2021, Lemma 2.2), $\|\widehat{\lambda}_i - \lambda_i\| \leq \|\widehat{\Sigma}_S - \Sigma_S\|$. Combining with previous inequalities, we have $1 -$

2376 $2C\sqrt{\frac{C_1/C_2+k+\log \frac{1}{\delta}}{n}} \leq \hat{\lambda}_i \leq 1 + 2C\sqrt{\frac{C_1/C_2+k+\log \frac{1}{\delta}}{n}}$ for $1 \leq i \leq k$, $\frac{C_1}{\sqrt{n}} - 2C\sqrt{\frac{C_1/C_2+k+\log \frac{1}{\delta}}{n}} \leq$
 2377
 2378 $\hat{\lambda}_i \leq \frac{C_1}{\sqrt{n}} + 2C\sqrt{\frac{C_1/C_2+k+\log \frac{1}{\delta}}{n}}$ for $k+1 \leq i \leq k + \lfloor \sqrt{n}/C_2 \rfloor$. If we take $C_2 \geq C_1((\frac{C_1}{4C})^2 - k -$
 2379
 2380 $\log \frac{1}{\delta})^{-1}$ then $2C\sqrt{\frac{C_1/C_2+k+\log \frac{1}{\delta}}{n}} \leq \frac{C_1}{2\sqrt{n}}$. Therefore we have $\frac{C_1}{2\sqrt{n}} \leq \hat{\lambda}_i \leq \frac{3C_1}{2\sqrt{n}}$ for $k+1 \leq i \leq$
 2381 $k + \lfloor \sqrt{n}/C_2 \rfloor$. When $\lambda = c\sqrt{n}$, we have

$$\begin{aligned}
 \frac{V}{v^2} &= \frac{1}{n} \text{tr}((\hat{\Sigma}_S + \frac{\lambda}{n}I_d)^{-2}\hat{\Sigma}_S) \\
 &= \frac{1}{n} \sum_{i=1}^d (\hat{\lambda}_i + \frac{\lambda}{n})^{-2}\hat{\lambda}_i \\
 &\geq \frac{1}{n} \sum_{i=k+1}^{k+\lfloor \sqrt{n}/C_2 \rfloor} (\hat{\lambda}_i + \frac{\lambda}{n})^{-2}\hat{\lambda}_i \\
 &= \frac{1}{n} \sum_{i=k+1}^{k+\lfloor \sqrt{n}/C_2 \rfloor} (\hat{\lambda}_i + \frac{c}{\sqrt{n}})^{-2}\hat{\lambda}_i \\
 &\geq \frac{1}{n} \sum_{i=k+1}^{k+\lfloor \sqrt{n}/C_2 \rfloor} (\frac{3C_1}{2\sqrt{n}} + \frac{c}{\sqrt{n}})^{-2}\frac{C_1}{2\sqrt{n}} \\
 &= \frac{1}{n} \lfloor \sqrt{n}/C_2 \rfloor \frac{C_1}{2} (\frac{3C_1}{2} + c)^{-2}\sqrt{n} \\
 &\geq \frac{C_1}{4C_2} (\frac{3C_1}{2} + c)^{-2}.
 \end{aligned} \tag{12}$$

Similarly, if $\lambda \leq n^{3/4}$,

$$\begin{aligned}
 \frac{V}{v^2} &\geq \frac{1}{n} \sum_{i=k+1}^{k+\lfloor \sqrt{n}/C_2 \rfloor} (\hat{\lambda}_i + \frac{\lambda}{n})^{-2}\hat{\lambda}_i \\
 &\geq \frac{1}{n} \sum_{i=k+1}^{k+\lfloor \sqrt{n}/C_2 \rfloor} (\hat{\lambda}_i + n^{-1/4})^{-2}\hat{\lambda}_i \\
 &\geq \frac{1}{n} \sum_{i=k+1}^{k+\lfloor \sqrt{n}/C_2 \rfloor} (\frac{3C_1}{2\sqrt{n}} + n^{-1/4})^{-2}\frac{C_1}{2\sqrt{n}} \\
 &= \frac{1}{n} \lfloor \sqrt{n}/C_2 \rfloor \frac{C_1}{2} (\frac{3C_1}{2} + n^{1/4})^{-2}\sqrt{n} \\
 &\geq \frac{C_1}{16C_2} n^{-1/2},
 \end{aligned} \tag{13}$$

when $n \geq (\frac{3C_1}{2})^4$.

As for the bias term, assume $\lambda \geq n^{3/4}$. Using the same concentration argument, we have $2 > \hat{\lambda}_i > 1/2$, for $1 \leq i \leq k$. When $\lambda \leq n$, $\lambda_{\max}(\hat{\Sigma}_S + \frac{\lambda}{n}I_d) \leq 2 + \lambda/n \leq 3$, therefore $\lambda_{\min}((\hat{\Sigma}_S + \frac{\lambda}{n}I_d)^{-1}) \geq \frac{1}{3}$. This implies

$$\begin{aligned}
 B &= \frac{\lambda^2}{n^2} \|(\hat{\Sigma}_S + \frac{\lambda}{n}I_d)^{-1}\beta^*\|^2 \\
 &\geq \frac{n^{3/2}}{n^2} \|(\hat{\Sigma}_S + \frac{\lambda}{n}I_d)^{-1}\beta^*\|^2 \\
 &\geq \frac{1}{\sqrt{n}} \lambda_{\min}^2((\hat{\Sigma}_S + \frac{\lambda}{n}I_d)^{-1}) \|\beta^*\|^2 \\
 &\geq \frac{\|\beta^*\|^2}{9\sqrt{n}}.
 \end{aligned}$$

When $\lambda > n$, $\lambda_{\max}(\widehat{\Sigma}_S + \frac{\lambda}{n}I_d) \leq 2 + \lambda/n \leq \frac{3\lambda}{n}$, which means $\lambda_{\min}((\widehat{\Sigma}_S + \frac{\lambda}{n}I_d)^{-1}) \geq \frac{n}{3\lambda}$. This implies

$$\begin{aligned} B &= \frac{\lambda^2}{n^2} \|(\widehat{\Sigma}_S + \frac{\lambda}{n}I_d)^{-1}\beta^*\|^2 \\ &\geq \frac{\lambda^2}{n^2} \lambda_{\min}^2((\widehat{\Sigma}_S + \frac{\lambda}{n}I_d)^{-1}) \|\beta^*\|^2 \\ &\geq \frac{\lambda^2}{n^2} \frac{n^2}{9\lambda^2} \|\beta^*\|^2 \\ &\geq \frac{\|\beta^*\|^2}{9}. \end{aligned}$$

□

C.2 UPPER BOUND FOR PCR

In this subsection, we will give the following upper bound for Principal Component Regression.

Theorem 29. When $n \gtrsim \sigma^8(r + \log \frac{1}{\delta})(\frac{\lambda_1}{\lambda_k - \lambda_{k+1}})^2 \frac{\lambda_1^2 k^2 \|\Sigma_T\|^2}{\lambda_k^4 \text{tr}((\Sigma_{S,k})^{-1} \Sigma_{T,k})^2}$,

$$\begin{aligned} \mathbb{E}_{\epsilon} \|\widehat{\beta} - \beta^*\|_{\Sigma_T}^2 &\leq \mathcal{O}(\sigma^8(\frac{\lambda_1}{\lambda_k - \lambda_{k+1}})^2 (\frac{\lambda_1}{\lambda_k})^2 \|\Sigma_T\| (\frac{r + \log \frac{1}{\delta}}{n}) \|\beta_k^*\|^2 + \frac{1}{n} v^2 \text{tr}((\Sigma_{S,k})^{-1} \Sigma_{T,k}) \\ &\quad + \frac{\|\Sigma_{T,k}\| \|\beta_{-k}^*\|^2 \|\Sigma_{S,-k}\|}{\lambda_k} + \beta_{-k}^{*T} \Sigma_{T,-k} \beta_{-k}^*) \end{aligned}$$

where $r = \frac{\sum_{i=1}^d \lambda_i}{\lambda_1}$.

Proof. For simplicity, we assume we have a sample size of $2n$, and in the first step we obtain an estimator $\widehat{U} \in \mathbb{R}^{d \times k}$ of the top-k subspace $U = \begin{pmatrix} I_k \\ 0 \end{pmatrix} \in \mathbb{R}^{d \times k}$, by using principal component analysis on the sample covariance matrix $\widehat{\Sigma}_S := \frac{1}{n} X^T X = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$, namely $\widehat{U} = (\widehat{u}_1, \dots, \widehat{u}_k)$ where \widehat{u}_i is the i -th eigenvector of $\widehat{\Sigma}_S$. We denote the distance between the estimated subspace and the original one by $\Delta := \text{dist}(U, \widehat{U}) = \|UU^T - \widehat{U}\widehat{U}^T\|$. For controlling Δ , we have the following lemma (Lemma 6):

Lemma 30. With probability at least $1 - \delta$,

$$\Delta \leq C\sigma^4 \left(\sqrt{\frac{r + \log \frac{1}{\delta}}{n}} + \frac{r + \log \frac{1}{\delta}}{n} \right) \frac{\lambda_1}{\lambda_k - \lambda_{k+1}}$$

where $r = \frac{\sum_{i=1}^n \lambda_i}{\lambda_1}$.

In the second step, we do linear regression on the projected (second half) data. With a little abuse of notation, we still use $X \in \mathbb{R}^{n \times d}$ to denote the data matrix indexed from $n + 1$ to $2n$. The data here is independent from the data in step 1, and therefore independent of Δ . If we let $Z := X\widehat{U} \in \mathbb{R}^{n \times k}$ be the projected data matrix, the estimator $\widehat{\beta}$ we obtained is given by

$$\begin{aligned} \widehat{\beta} &= \widehat{U}(Z^T Z)^{-1} Z^T Y \\ &= \widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T Y. \end{aligned} \tag{14}$$

We aim to bound the excess risk on target, which is given by $\|\widehat{\beta} - \beta^*\|_{\Sigma_T}^2 := \|\Sigma_T^{\frac{1}{2}}(\widehat{\beta} - \beta^*)\|^2$. We introduce the following notations: suppose $\beta^* = (\beta_1^*, \dots, \beta_d^*)^T$. We let $\beta_U^* := (\beta_1^*, \dots, \beta_k^*, 0, \dots, 0)^T$, $\beta_{\perp}^* := (0, \dots, 0, \beta_{k+1}^*, \dots, \beta_d^*)^T = \beta^* - \beta_U^*$. Here we present an intermediate result for bounding the excess risk:

2484 Assume $\Delta \leq \frac{\lambda_k^2 \text{tr}((\Sigma_{S,k})^{-1}\Sigma_{T,k})}{4\lambda_1 k \|\Sigma_T\|}$. When $n \gtrsim \frac{\sigma^4 \lambda_1^2 \|\Sigma_T\|^2 k^3 \log(1/\delta)}{\lambda_k^4 \text{tr}((\Sigma_{S,k})^{-1}\Sigma_{T,k})^2}$, then with probability 1 - δ ,

$$\begin{aligned} \mathbb{E}_\epsilon \|\hat{\beta} - \beta^*\|_{\Sigma_T}^2 &\leq \mathcal{O}(\|\beta_U^*\|^2 \Delta^2 (\frac{\lambda_1}{\lambda_k})^2 \|\Sigma_T\| + \frac{1}{n} v^2 \text{tr}((\Sigma_{S,k})^{-1}\Sigma_{T,k})) \\ &\quad + \frac{\|\Sigma_{T,k}\| \|\beta_{-k}^*\|^2 \|\Sigma_{S,-k}\|}{\lambda_k} + \beta_{-k}^{*T} \Sigma_{T,-k} \beta_{-k}^* \end{aligned}$$

If further $n \gtrsim \sigma^4 \Delta^{-2} k \log(1/\delta)$,

$$\begin{aligned} \mathbb{E}_\epsilon \|\hat{\beta} - \beta^*\|_{\Sigma_T}^2 &\leq \mathcal{O}(\|\beta_U^*\|^2 (\Delta^4 (\frac{\lambda_1}{\lambda_k})^2 \|\Sigma_T\| + \Delta^2 \|\Sigma_{T,-k}\| + \Delta^3 \|\Sigma_T\|) \\ &\quad + \frac{1}{n} v^2 \text{tr}((\Sigma_{S,k})^{-1}\Sigma_{T,k}) + \frac{\|\Sigma_{T,k}\| \|\beta_{-k}^*\|^2 \|\Sigma_{S,-k}\|}{\lambda_k} + \beta_{-k}^{*T} \Sigma_{T,-k} \beta_{-k}^*) \end{aligned}$$

From Lemma 30, when $n \geq r + \log \frac{1}{\delta} = \frac{\sum_{i=1}^n \lambda_i}{\lambda_1} + \log \frac{1}{\delta}$, we have

$$\Delta \leq 2C \frac{\lambda_1}{\lambda_k - \lambda_{k+1}} \sigma^4 \sqrt{\frac{r + \log \frac{1}{\delta}}{n}}$$

Therefore when $n \gtrsim (r + \log \frac{1}{\delta}) \sigma^8 (\frac{\lambda_1}{\lambda_k - \lambda_{k+1}})^2 \frac{\lambda_1^2 k^2 \|\Sigma_T\|^2}{\lambda_k^4 \text{tr}((\Sigma_{S,k})^{-1}\Sigma_{T,k})^2}$, the assumption for Δ and n in Lemma 31 will be both satisfied. We can thus apply Lemma 31 to get

$$\begin{aligned} \mathbb{E}_\epsilon \|\hat{\beta} - \beta^*\|_{\Sigma_T}^2 &\leq \mathcal{O}(\sigma^8 (\frac{\lambda_1}{\lambda_k - \lambda_{k+1}})^2 (\frac{\lambda_1}{\lambda_k})^2 \|\Sigma_T\| \frac{r + \log \frac{1}{\delta}}{n} \|\beta_U^*\|^2 + \frac{1}{n} v^2 \text{tr}((\Sigma_{S,k})^{-1}\Sigma_{T,k}) \\ &\quad + \frac{\|\Sigma_{T,k}\| \|\beta_{-k}^*\|^2 \|\Sigma_{S,-k}\|}{\lambda_k} + \beta_{-k}^{*T} \Sigma_{T,-k} \beta_{-k}^*) \end{aligned}$$

where $r = \frac{\sum_{i=1}^d \lambda_i}{\lambda_1}$.

□

C.3 PROOFS FOR LEMMA 31

In the following we will prove Lemma 31.

Proof for Lemma 31. The proof idea is similar to (Ge et al., 2023, Theorem 4.4) and (Tripuraneni et al., 2021b, Theorem 4).

We can decompose $\hat{\beta} - \beta^*$ as

$$\begin{aligned} \hat{\beta} - \beta^* &= \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T X^T Y - \beta^* \\ &= \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T X^T (X \beta^* + \epsilon) - \beta^* \\ &= \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T X^T (X \beta_U^* + X \beta_\perp^* + \epsilon) - (\beta_U^* + \beta_\perp^*) \\ &= A_1 + A_2 + A_3 - \beta_\perp^*, \end{aligned}$$

where $A_1 := \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T X^T X \beta_U^* - \beta_U^*$, $A_2 := \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T X^T X \beta_\perp^*$, $A_3 := \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T X^T \epsilon$. Therefore

$$\|\hat{\beta} - \beta^*\|_{\Sigma_T}^2 \leq \|A_1\|_{\Sigma_T}^2 + \|A_2\|_{\Sigma_T}^2 + \|A_3\|_{\Sigma_T}^2 + \|\beta_\perp^*\|_{\Sigma_T}^2 \quad (15)$$

We give three lemmas for bounding the related terms. The first lemma considers the bias term A_1 :

Lemma 32. If $\Delta \leq \frac{\lambda_k}{4\lambda_1}$ and $n \gtrsim \max\{\sigma^4 (\frac{\lambda_1}{\lambda_k})^2 k \log(1/\delta), \sigma^4 k \log(1/\delta)\}$, then with probability at least $1 - \delta$,

2538

2539

$$2540 \quad \|A_1\|_{\Sigma_T}^2 \leq \mathcal{O}(\|\beta_U^*\|^2 \Delta^2 (\frac{\lambda_1}{\lambda_k})^2 \|\Sigma_T\|)$$

$$2541$$

2542 If we further have $n \gtrsim \sigma^4 \Delta^{-2} k \log(1/\delta)$, then with probability at least $1 - \delta$,

2543

2544

$$2545 \quad \|A_1\|_{\Sigma_T}^2 \leq \mathcal{O}(\|\beta_U^*\|^2 (\Delta^4 (\frac{\lambda_1}{\lambda_k})^2 \|\Sigma_T\| + \Delta^2 \|\Sigma_{T,-k}\| + \Delta^3 \|\Sigma_T\|)) \leq \mathcal{O}(\|\beta_U^*\|^2 \Delta^2 \|\Sigma_T\|)$$

$$2546$$

$$2547$$

2548 The second lemma considers the variance term A_3 :

2549

2550 **Lemma 33.** If $\Delta \leq \frac{\lambda_k^2 \text{tr}((\Sigma_{S,k})^{-1} \Sigma_{T,k})}{4\lambda_1 k \|\Sigma_T\|}$ and $n \gtrsim \frac{\sigma^4 \|\Sigma_S\|^2 \|\Sigma_T\|^2 k^3 \log(1/\delta)}{\lambda_k^4 \text{tr}((\Sigma_{S,k})^{-1} \Sigma_{T,k})^2}$, then with probability at

$$2551 \quad \text{least } 1 - \delta,$$

2552

2553

2554

$$\mathbb{E}_\epsilon[\|A_3\|_{\Sigma_T}^2] \leq \mathcal{O}(\frac{1}{n} v^2 \text{tr}((\Sigma_{S,k})^{-1} \Sigma_{T,k})).$$

2555

2556

For bounding A_2 , we actually have a similar result to bounding A_3 :

2557

2558

Lemma 34. If $n \gtrsim \sigma^4 (\frac{\lambda_1}{\lambda_k})^2 k \log(1/\delta)$ and $\Delta \leq \min\{\frac{\|\Sigma_{T,k}\|}{2\|\Sigma_T\|}, \frac{\lambda_k}{4\lambda_1}\}$, then with probability at least $1 - \delta$

2559

2560

2561

$$\|A_2\|_{\Sigma_T}^2 \leq \mathcal{O}(\frac{\|\Sigma_{T,k}\| \|\beta_{-k}^*\|^2 \|\Sigma_{S,-k}\|}{\lambda_k}) \quad (16)$$

2562

2563

2564

By Lemma 32, 33, 34, together with the decomposition (15), we have with probability $1 - \delta$, when $n \gtrsim N_1$,

2565

2566

$$\mathbb{E}_\epsilon \|\hat{\beta} - \beta^*\|_{\Sigma_T}^2 \leq \mathcal{O}(\|\beta_U^*\|^2 \Delta^2 (\frac{\lambda_1}{\lambda_k})^2 \|\Sigma_T\| + \frac{1}{n} v^2 \text{tr}((\Sigma_{S,k})^{-1} \Sigma_{T,k})) \quad (17)$$

2567

2568

2569

$$+ \frac{\|\Sigma_{T,k}\| \|\beta_{-k}^*\|^2 \|\Sigma_{S,-k}\|}{\lambda_k} + \beta_{-k}^{*T} \Sigma_{T,-k} \beta_{-k}^*) \quad (18)$$

2570

2571

If further $n \gtrsim \sigma^4 \Delta^{-2} k \log(1/\delta)$,

2572

2573

$$\mathbb{E}_\epsilon \|\hat{\beta} - \beta^*\|_{\Sigma_T}^2 \leq \mathcal{O}(\|\beta_U^*\|^2 (\Delta^4 (\frac{\lambda_1}{\lambda_k})^2 \|\Sigma_T\| + \Delta^2 \|\Sigma_{T,-k}\| + \Delta^3 \|\Sigma_T\|)) \quad (19)$$

2574

2575

2576

$$+ \frac{1}{n} v^2 \text{tr}((\Sigma_{S,k})^{-1} \Sigma_{T,k}) + \frac{\|\Sigma_{T,k}\| \|\beta_{-k}^*\|^2 \|\Sigma_{S,-k}\|}{\lambda_k} + \beta_{-k}^{*T} \Sigma_{T,-k} \beta_{-k}^*) \quad (20)$$

2577

2578

C.4 TECHNICAL PROOFS

2579

2580

2581

2582

2583

2584

2585

In the sequel, we give the proofs of Lemma 32, 33, 34 and 30. We first prove some additional technical lemmas. The following lemma, which is a simple corollary of (Tripuraneni et al., 2021b, Lemma 20), shows the concentration property of empirical covariance matrix.

2586

2587

2588

Lemma 35. Let $\{x_i\}_{i=1}^n$ be i.i.d. d -dimensional random vectors, satisfying: x_i is mean zero, $\mathbb{E}[xx^T] = \Sigma$ such that $\sigma_{\max}(\Sigma) \leq C_{\max}$ and is $\sigma^2 \Sigma$ -sub-gaussian, in the sense that

2589

2590

2591

$$\mathbb{E}[\exp(v^T x_i)] \leq \exp\left(\frac{\|\sigma \Sigma^{1/2} v\|^2}{2}\right).$$

2592

2593

$X = (x_1, \dots, x_n)^T \in \mathbb{R}^{n \times d}$. Then for any $A, B \in \mathbb{R}^{d \times k}$, we have with probability at least $1 - \delta$

2594

2595

2596

2597

2598

2599

2600

$$\|A^T (\frac{X^T X}{n}) B - A^T \Sigma B\|_2 \leq \mathcal{O}(\sigma^2 \|A\| \|B\| \|\Sigma\| (\sqrt{\frac{k}{n}} + \frac{k}{n} + \sqrt{\frac{\log(1/\delta)}{n}} + \frac{\log(1/\delta)}{n})). \quad (21)$$

Proof. We write the SVD of A and B : $A = U_1 \Lambda_1 V_1^T$, $B = U_2 \Lambda_2 V_2^T$, where $U_1, U_2 \in \mathbb{R}^{d \times k}$, $\Lambda_1, \Lambda_2, V_1, V_2 \in \mathbb{R}^{k \times k}$. Then

$$\begin{aligned} \|A^T(\frac{X^T X}{n})B - A^T \Sigma B\|_2 &= \|V_1 \Lambda_1 U_1^T(\frac{X^T X}{n})U_2 \Lambda_2 V_2^T - V_1 \Lambda_1 U_1^T \Sigma U_2 \Lambda_2 V_2^T\|_2 \\ &\leq \|V_1 \Lambda_1\| \|U_1^T(\frac{X^T X}{n})U_2 - U_1^T \Sigma U_2\| \|\Lambda_2 V_2^T\| \\ &\leq \|A\| \|B\| \|U_1^T(\frac{X^T X}{n})U_2 - U_1^T \Sigma U_2\|. \end{aligned} \quad (22)$$

Now since $U_1, U_2 \in \mathbb{R}^{d \times k}$ are projection matrices, we can apply Tripuraneni et al. (2021b) Lemma 20, therefore

$$\|U_1^T(\frac{X^T X}{n})U_2 - U_1^T \Sigma U_2\| \leq \mathcal{O}(\sigma^2 \|\Sigma\| (\sqrt{\frac{k}{n}} + \frac{k}{n} + \sqrt{\frac{\log(1/\delta)}{n}} + \frac{\log(1/\delta)}{n})) \quad (23)$$

which gives what we want. \square

The following lemma is a basic matrix perturbation result (see Tripuraneni et al. (2021b) Lemma 25).

Lemma 36. Let A be a positive definite matrix and E another matrix which satisfies $\|EA^{-1}\| \leq \frac{1}{4}$, then $F := (A + E)^{-1} - A^{-1}$ satisfies $\|F\| \leq \frac{4}{3} \|A^{-1}\| \|EA^{-1}\|$.

With these two technical lemmas, we are able to prove Lemma 32, 33.

Proof of Lemma 32. Notice that by the definition of U and β_U^* , we have $UU^T \beta_U^* = \beta_U^*$. We denote $\alpha^* := U^T \beta_U^*$, then we also have $\beta_U^* = U \alpha^*$. Therefore

$$\begin{aligned} A_1 &= \widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \beta_U^* - \beta_U^* \\ &= \widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X U \alpha^* - U \alpha^* \\ &= (\widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X U - U) \alpha^* \end{aligned}$$

We consider $\widehat{U} \in \mathbb{R}^{d \times k}$ and $\widehat{U}_\perp^T \in \mathbb{R}^{d \times (d-k)}$ be orthonormal projection matrices spanning orthogonal subspaces which are rank k and rank $d - k$ respectively, so that $\text{range}(\widehat{U}) \oplus \text{range}(\widehat{U}_\perp) = \mathbb{R}^d$. Then $\Delta = \text{dist}(\widehat{U}, U^*) = \|\widehat{U}_\perp^T U^*\|_2$. Notice that $I_d = \widehat{U} \widehat{U}^T + \widehat{U}_\perp \widehat{U}_\perp^T$, we have

$$\begin{aligned} &\widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X U^* - U^* \\ &= \widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X (\widehat{U} \widehat{U}^T + \widehat{U}_\perp \widehat{U}_\perp^T) U^* - U^* \\ &= \widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \widehat{U} \widehat{U}^T U^* + \widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \widehat{U}_\perp \widehat{U}_\perp^T U^* - U^* \\ &= \widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \widehat{U}_\perp \widehat{U}_\perp^T U^* + \widehat{U} \widehat{U}^T U^* - U^* \\ &= \widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \widehat{U}_\perp \widehat{U}_\perp^T U^* - \widehat{U}_\perp \widehat{U}_\perp^T U^* \end{aligned} \quad (24)$$

Thus

$$\begin{aligned} \|A_1\|_{\Sigma_T}^2 &= A_1^T \Sigma_T A_1 \\ &= \alpha^{*T} (\widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X U - U)^T \Sigma_T (\widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X U - U) \alpha^* \\ &= \alpha^{*T} (\widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \widehat{U}_\perp \widehat{U}_\perp^T U^* - \widehat{U}_\perp \widehat{U}_\perp^T U^*)^T \Sigma_T \\ &\quad (\widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \widehat{U}_\perp \widehat{U}_\perp^T U^* - \widehat{U}_\perp \widehat{U}_\perp^T U^*) \alpha^* \\ &\leq \|\alpha^*\|^2 \|\widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \widehat{U}_\perp \widehat{U}_\perp^T U^* - \widehat{U}_\perp \widehat{U}_\perp^T U^*\|_{\Sigma_T}^2 \\ &\leq \|\alpha^*\|^2 (\|\widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \widehat{U}_\perp \widehat{U}_\perp^T U^*\|_{\Sigma_T}^2 + \|\widehat{U}_\perp \widehat{U}_\perp^T U^*\|_{\Sigma_T}^2). \end{aligned} \quad (25)$$

Here we use the notation $\|M\|_{\Sigma_T} := \sqrt{\|M^T \Sigma_T M\|}$ for matrix M .

2646 For the second term,

$$2647 \quad \|\widehat{U}_\perp \widehat{U}_\perp^T U^*\|_{\Sigma_T}^2 \leq \|\widehat{U}_\perp^T \Sigma_T \widehat{U}_\perp\| \|\widehat{U}_\perp^T U^*\|^2 \leq \Delta^2 \|\widehat{U}_\perp^T \Sigma_T \widehat{U}_\perp\|. \quad (26)$$

2649 For the first term,

$$\begin{aligned} 2651 \quad & \|\widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \widehat{U}_\perp \widehat{U}_\perp^T U^*\|_{\Sigma_T}^2 \\ 2652 \quad & = \|\widehat{U}(\widehat{U}^T \frac{X^T X}{n} \widehat{U})^{-1} \widehat{U}^T \frac{X^T X}{n} \widehat{U}_\perp \widehat{U}_\perp^T U^*\|_{\Sigma_T}^2 \\ 2653 \quad & = \|\widehat{U}((\widehat{U}^T \Sigma_S \widehat{U})^{-1} + F)(\widehat{U}^T \Sigma_S \widehat{U}_\perp \widehat{U}_\perp^T U^* + E_1)\|_{\Sigma_T}^2 \\ 2654 \quad & = \|(\widehat{U}^T \Sigma_S \widehat{U}_\perp \widehat{U}_\perp^T U^* + E_1)^T ((\widehat{U}^T \Sigma_S \widehat{U})^{-1} + F)^T \widehat{U}^T \Sigma_T \widehat{U}((\widehat{U}^T \Sigma_S \widehat{U})^{-1} + F)(\widehat{U}^T \Sigma_S \widehat{U}_\perp \widehat{U}_\perp^T U^* + E_1)\| \\ 2655 \quad & \leq \|\widehat{U}^T \Sigma_S \widehat{U}_\perp \widehat{U}_\perp^T U^* + E_1\|^2 \|(\widehat{U}^T \Sigma_S \widehat{U})^{-1} + F\|^2 \|\widehat{U}^T \Sigma_T \widehat{U}\| \\ 2656 \quad & \leq (\|\widehat{U}^T \Sigma_S \widehat{U}_\perp \widehat{U}_\perp^T U^*\| + \|E_1\|)^2 (\|(\widehat{U}^T \Sigma_S \widehat{U})^{-1}\| + \|F\|)^2 \|\widehat{U}^T \Sigma_T \widehat{U}\| \end{aligned} \quad (27)$$

2660 where $E_1 = \widehat{U}^T \frac{X^T X}{n} \widehat{U}_\perp \widehat{U}_\perp^T U^* - \widehat{U}^T \Sigma_S \widehat{U}_\perp \widehat{U}_\perp^T U^*$, $F = (\widehat{U}^T \frac{X^T X}{n} \widehat{U})^{-1} - (\widehat{U}^T \Sigma_S \widehat{U})^{-1}$. We
2661 aim to show that $\|E_1\| \leq \|\widehat{U}^T \Sigma_S \widehat{U}_\perp \widehat{U}_\perp^T U^*\|$ and $\|F\| \leq \|(\widehat{U}^T \Sigma_S \widehat{U})^{-1}\| = C_{\min}^{-1}$ for suffi-
2662 ciently large n , therefore the term in (27) can be bounded well. First we need a careful analysis
2663 of $\|\widehat{U}^T \Sigma_S \widehat{U}_\perp \widehat{U}_\perp^T U^*\|$. It is obvious that

$$2664 \quad \|\widehat{U}^T \Sigma_S \widehat{U}_\perp \widehat{U}_\perp^T U^*\| \leq \|\widehat{U}^T \Sigma_S \widehat{U}_\perp\| \|\widehat{U}_\perp^T U^*\| \leq \Delta \|\widehat{U}^T \Sigma_S \widehat{U}_\perp\|. \quad (28)$$

2665 As for $\|\widehat{U}^T \Sigma_S \widehat{U}_\perp\|$, notice that if without the "hat", we have $U^T \Sigma_S U_\perp = 0$ by the definition of U
2666 and Σ_S is diagonal. By definition of distance between two subspaces, there exist $R \in \mathcal{O}^{k \times k}$ and
2667 $Q \in \mathcal{O}^{(d-k) \times (d-k)}$, such that $\|\widehat{U}R - U\| = \Delta = \|\widehat{U}_\perp Q - U_\perp\|$. Then we have

$$\begin{aligned} 2668 \quad \|\widehat{U}^T \Sigma_S \widehat{U}_\perp\| &= \|R^T \widehat{U}^T \Sigma_S \widehat{U}_\perp Q\| \\ 2669 \quad &= \|U^T \Sigma_S U_\perp + R^T \widehat{U}^T \Sigma_S \widehat{U}_\perp Q - U^T \Sigma_S U_\perp\| \\ 2670 \quad &= \|R^T \widehat{U}^T \Sigma_S \widehat{U}_\perp Q - U^T \Sigma_S U_\perp\| \\ 2671 \quad &= \|R^T \widehat{U}^T \Sigma_S \widehat{U}_\perp Q - U^T \Sigma_S \widehat{U}_\perp Q + U^T \Sigma_S \widehat{U}_\perp Q - U^T \Sigma_S U_\perp\| \\ 2672 \quad &\leq \|R^T \widehat{U}^T \Sigma_S \widehat{U}_\perp Q - U^T \Sigma_S \widehat{U}_\perp Q\| + \|U^T \Sigma_S \widehat{U}_\perp Q - U^T \Sigma_S U_\perp\| \\ 2673 \quad &\leq \|R^T \widehat{U}^T - U^T\| \|\Sigma_S \widehat{U}_\perp Q\| + \|U^T \Sigma_S\| \|\widehat{U}_\perp Q - U_\perp\| \\ 2674 \quad &\leq 2\Delta \|\Sigma_S\|. \end{aligned} \quad (29)$$

2675 Combine (28) and (29), we have

$$2676 \quad \|\widehat{U}^T \Sigma_S \widehat{U}_\perp \widehat{U}_\perp^T U^*\| \leq \mathcal{O}(\Delta^2 \|\Sigma_S\|) \quad (30)$$

2677 In order to bound $\|F\|$, let $E = \widehat{U}^T \frac{X^T X}{n} \widehat{U} - \widehat{U}^T \Sigma_S \widehat{U}$, then by Lemma 35, with probability at least
2678 $1 - \delta$,

$$2679 \quad \|E\| \leq \mathcal{O}(\sigma^2 \|\Sigma_S\| (\sqrt{\frac{k}{n}} + \frac{k}{n} + \sqrt{\frac{\log(1/\delta)}{n}} + \frac{\log(1/\delta)}{n})). \quad (31)$$

2680 Therefore,

$$\begin{aligned} 2681 \quad \|E(\widehat{U}^T \Sigma_S \widehat{U})^{-1}\| &\leq \|E\| \|(\widehat{U}^T \Sigma_S \widehat{U})^{-1}\| \\ 2682 \quad &\leq \|E\| C_{\min}^{-1} \\ 2683 \quad &\leq \mathcal{O}(\sigma^2 C_{\min}^{-1} \|\Sigma_S\| (\sqrt{\frac{k}{n}} + \frac{k}{n} + \sqrt{\frac{\log(1/\delta)}{n}} + \frac{\log(1/\delta)}{n})), \end{aligned} \quad (32)$$

2684 where $C_{\min} := \lambda_{\min}(\widehat{U}^T \Sigma_S \widehat{U})$. Notice that $n \gtrsim \sigma^4 C_{\min}^{-2} \|\Sigma_S\|^2 k \log(1/\delta)$ implies $\sqrt{\frac{k}{n}} + \frac{k}{n} +$
2685 $\sqrt{\frac{\log(1/\delta)}{n}} + \frac{\log(1/\delta)}{n} \lesssim \sigma^{-2} C_{\min} \|\Sigma_S\|^{-1}$. Thus, we show that when n is large enough, we have
2686 $\|E(\widehat{U}^T \Sigma_S \widehat{U})^{-1}\| \leq \frac{1}{4}$. Therefore we can apply Lemma 36, which gives

$$2687 \quad \|F\| \leq \frac{4}{3} \|E(\widehat{U}^T \Sigma_S \widehat{U})^{-1}\| \|(\widehat{U}^T \Sigma_S \widehat{U})^{-1}\|$$

$$\begin{aligned}
&\leq \frac{4}{3} \times \frac{1}{4} \|(\widehat{U}^T \Sigma_S \widehat{U})^{-1}\| \\
&\leq \frac{1}{3} C_{\min}^{-1}.
\end{aligned} \tag{33}$$

As for $\|E_1\|$, directly applying Lemma 35, when $n \gtrsim \sigma^4 \Delta^{-2} k \log(1/\delta)$ we get

$$\begin{aligned}
\|E_1\| &\leq \mathcal{O}(\sigma^2 \|\Sigma_S\| \|\widehat{U}_\perp \widehat{U}_\perp^T U^*\| (\sqrt{\frac{k}{n}} + \frac{k}{n} + \sqrt{\frac{\log(1/\delta)}{n}} + \frac{\log(1/\delta)}{n})) \\
&\leq \mathcal{O}(\sigma^2 \|\Sigma_S\| \Delta (\sqrt{\frac{k}{n}} + \frac{k}{n} + \sqrt{\frac{\log(1/\delta)}{n}} + \frac{\log(1/\delta)}{n}))
\end{aligned} \tag{34}$$

when $n \gtrsim \sigma^4 k \log(1/\delta)$ we have

$$\|E_1\| \leq \mathcal{O}(\Delta \|\Sigma_S\|) \tag{35}$$

, if further we have $n \gtrsim \sigma^4 \Delta^{-2} k \log(1/\delta)$, then

$$\|E_1\| \leq \mathcal{O}(\Delta^2 \|\Sigma_S\|). \tag{36}$$

Combining (27), (30), (33) and (36), we have

$$\begin{aligned}
&\|\widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \widehat{U}_\perp \widehat{U}_\perp^T U^*\|_{\Sigma_T}^2 \\
&\leq (\|\widehat{U}^T \Sigma_S \widehat{U}_\perp \widehat{U}_\perp^T U^*\| + \|E_1\|)^2 (\|(\widehat{U}^T \Sigma_S \widehat{U})^{-1}\| + \|F\|)^2 \|\widehat{U}^T \Sigma_T \widehat{U}\| \\
&\leq \mathcal{O}(\Delta^4 \|\Sigma_S\|^2 C_{\min}^{-2} \|\widehat{U}^T \Sigma_T \widehat{U}\|) \\
&\leq \mathcal{O}(\Delta^4 \|\Sigma_S\|^2 C_{\min}^{-2} \|\Sigma_T\|)
\end{aligned} \tag{37}$$

Combining (25), (26) and (37), we get

$$\begin{aligned}
\|A_1\|_{\Sigma_T}^2 &\leq \|\alpha^*\|^2 (\|\widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \widehat{U}_\perp \widehat{U}_\perp^T U^*\|_{\Sigma_T}^2 + \|\widehat{U}_\perp \widehat{U}_\perp^T U^*\|_{\Sigma_T}^2) \\
&\leq \mathcal{O}(\|\alpha^*\|^2 (\Delta^4 \|\Sigma_S\|^2 C_{\min}^{-2} \|\Sigma_T\| + \Delta^2 \|\widehat{U}_\perp^T \Sigma_T \widehat{U}_\perp\|))
\end{aligned} \tag{38}$$

with probability at least $1 - \delta$. Also, similar to (29), we have

$$\begin{aligned}
\|\widehat{U}_\perp^T \Sigma_T \widehat{U}_\perp\| &= \|Q^T \widehat{U}_\perp^T \Sigma_T \widehat{U}_\perp Q\| \\
&\leq \|U_\perp^T \Sigma_T U_\perp\| + \|Q^T \widehat{U}_\perp^T \Sigma_T \widehat{U}_\perp Q - U_\perp^T \Sigma_T U_\perp\| \\
&\leq \|U_\perp^T \Sigma_T U_\perp\| + 2\Delta \|\Sigma_T\|
\end{aligned} \tag{39}$$

Similarly, we can further know that C_{\min} is close to λ_k :

$$\begin{aligned}
C_{\min} &= \lambda_k(\widehat{U}^T \Sigma_S \widehat{U}) \\
&= \lambda_k(R^T \widehat{U}^T \Sigma_S \widehat{U} R) \\
&= \lambda_k(U^T \Sigma_S U + R^T \widehat{U}^T \Sigma_S \widehat{U} R - U^T \Sigma_S U) \\
&\geq \lambda_k(U^T \Sigma_S U) - \|R^T \widehat{U}^T \Sigma_S \widehat{U} R - U^T \Sigma_S U\| \\
&\geq \lambda_k(U^T \Sigma_S U) 2\Delta \|\Sigma_S\| \\
&\geq \lambda_k - 2\lambda_1 \Delta \\
&\geq \frac{1}{2} \lambda_k,
\end{aligned} \tag{40}$$

where the last inequality holds when $\Delta \leq \frac{\lambda_k}{4\lambda_1}$. Finally, combining (38), (39), (40), we have

$$\begin{aligned}
\|A_1\|_{\Sigma_T}^2 &\leq \mathcal{O}(\|\alpha^*\|^2 (\Delta^4 (\frac{\lambda_1}{\lambda_k})^2 \|\Sigma_T\| + \Delta^2 \|U_\perp^T \Sigma_T U_\perp\| + \Delta^3 \|\Sigma_T\|)) \\
&\leq \mathcal{O}(\|\beta_U^*\|^2 (\Delta^4 (\frac{\lambda_1}{\lambda_k})^2 \|\Sigma_T\| + \Delta^2 \|U_\perp^T \Sigma_T U_\perp\| + \Delta^3 \|\Sigma_T\|))
\end{aligned} \tag{41}$$

when $\Delta \leq \frac{\lambda_k}{4\lambda_1}$ and $n \gtrsim \max\{\sigma^4(\frac{\lambda_1}{\lambda_k})^2 k \log(1/\delta), \sigma^4 \Delta^{-2} k \log(1/\delta)\}$. If in the previous proofs we replace (36) by (35), we have

$$\|A_1\|_{\Sigma_T}^2 \leq \mathcal{O}(\|\beta_U^\star\|^2 (\Delta^2 (\frac{\lambda_1}{\lambda_k})^2 \|\Sigma_T\| + \Delta^2 \|U_\perp^T \Sigma_T U_\perp\| + \Delta^3 \|\Sigma_T\|)) \quad (42)$$

$$\leq \mathcal{O}(\|\beta_U^\star\|^2 \Delta^2 (\frac{\lambda_1}{\lambda_k})^2 \|\Sigma_T\|) \quad (43)$$

when $\Delta \leq \frac{\lambda_k}{4\lambda_1}$ and $n \gtrsim \max\{\sigma^4(\frac{\lambda_1}{\lambda_k})^2 k \log(1/\delta), \sigma^4 k \log(1/\delta)\}$. Notice that by definition of U , $U_\perp^T \Sigma_T U_\perp = \Sigma_{T,-k}$, therefore the result is exactly what we want. \square

Proof of Lemma 33. Recall $A_3 := \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T X^T \epsilon$. Therefore

$$\begin{aligned} \|A_3\|_{\Sigma_T}^2 &= \epsilon^T X \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T \Sigma_T \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T X^T \epsilon \\ &= \text{tr}(\epsilon^T X \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T \Sigma_T \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T X^T \epsilon) \\ &= \text{tr}(\epsilon \epsilon^T X \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T \Sigma_T \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T X^T) \end{aligned}$$

Taking expectation with respect to ϵ , using $\mathbb{E}[\epsilon \epsilon^T] = v^2 I_n$, we have

$$\begin{aligned} \mathbb{E}_\epsilon[\|A_3\|_{\Sigma_T}^2] &= \mathbb{E}[\text{tr}(\epsilon \epsilon^T X \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T \Sigma_T \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T X^T)] \\ &= v^2 \text{tr}(X \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T \Sigma_T \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T X^T) \\ &= v^2 \text{tr}((\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T \Sigma_T \hat{U}(\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T X^T X \hat{U}) \\ &= v^2 \text{tr}((\hat{U}^T X^T X \hat{U})^{-1} \hat{U}^T \Sigma_T \hat{U}) \\ &= \frac{1}{n} v^2 \text{tr}(((\hat{U}^T \Sigma_S \hat{U})^{-1} + F) \hat{U}^T \Sigma_T \hat{U}) \end{aligned} \quad (44)$$

Here we actually need a bound stronger than (33) for $\|F\|$: recall (32), we have with probability $1 - \delta$

$$\|E(\hat{U}^T \Sigma_S \hat{U})^{-1}\| \leq \mathcal{O}(\sigma^2 C_{\min}^{-1} \|\Sigma_S\| (\sqrt{\frac{k}{n}} + \frac{k}{n} + \sqrt{\frac{\log(1/\delta)}{n}} + \frac{\log(1/\delta)}{n})). \quad (45)$$

Applying Lemma 36, which gives

$$\begin{aligned} \|F\| &\leq \frac{4}{3} \|E(\hat{U}^T \Sigma_S \hat{U})^{-1}\| \|(\hat{U}^T \Sigma_S \hat{U})^{-1}\| \\ &\leq \mathcal{O}(\sigma^2 C_{\min}^{-2} \|\Sigma_S\| (\sqrt{\frac{k}{n}} + \frac{k}{n} + \sqrt{\frac{\log(1/\delta)}{n}} + \frac{\log(1/\delta)}{n})) \\ &\leq \mathcal{O}(\frac{1}{k \|\Sigma_T\|} \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U)) \end{aligned} \quad (46)$$

when $n \gtrsim \sigma^4 C_{\min}^{-4} \|\Sigma_S\|^2 \|\Sigma_T\|^2 \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U)^{-2} k^3 \log(1/\delta)$. Therefore we have

$$\begin{aligned} \mathbb{E}_\epsilon[\|A_3\|_{\Sigma_T}^2] &= \frac{1}{n} v^2 \text{tr}(((\hat{U}^T \Sigma_S \hat{U})^{-1} + F) \hat{U}^T \Sigma_T \hat{U}) \\ &= \frac{1}{n} v^2 (\text{tr}((\hat{U}^T \Sigma_S \hat{U})^{-1} \hat{U}^T \Sigma_T \hat{U}) + \text{tr}(F \hat{U}^T \Sigma_T \hat{U})) \\ &\leq \frac{1}{n} v^2 (\text{tr}((\hat{U}^T \Sigma_S \hat{U})^{-1} \hat{U}^T \Sigma_T \hat{U})) + \frac{1}{n} v^2 \|F\| \text{tr}(\hat{U}^T \Sigma_T \hat{U}) \\ &\leq \frac{1}{n} v^2 (\text{tr}((\hat{U}^T \Sigma_S \hat{U})^{-1} \hat{U}^T \Sigma_T \hat{U})) + \frac{1}{n} v^2 k \|F\| \|\Sigma_T\| \\ &\leq \frac{1}{n} v^2 (\text{tr}((\hat{U}^T \Sigma_S \hat{U})^{-1} \hat{U}^T \Sigma_T \hat{U})) + \frac{1}{n} v^2 \mathcal{O}(\text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U)) \end{aligned} \quad (47)$$

The remaining thing is to show that indeed $\text{tr}((\hat{U}^T \Sigma_S \hat{U})^{-1} \hat{U}^T \Sigma_T \hat{U})$ is close to $\text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U)$. In fact, $\text{tr}((\hat{U}^T \Sigma_S \hat{U})^{-1} \hat{U}^T \Sigma_T \hat{U}) = \text{tr}((R^T \hat{U}^T \Sigma_S \hat{U} R)^{-1} R^T \hat{U}^T \Sigma_T R \hat{U})$. Notice that

2808
2809
2810

$$\|R^T \widehat{U}^T \Sigma_T \widehat{U} R - U^T \Sigma_T U\| \leq 2\|\Delta\|\|\Sigma_T\|,$$

2811 we have

$$\begin{aligned} & \text{tr}((R^T \widehat{U}^T \Sigma_S \widehat{U} R)^{-1} R^T \widehat{U}^T \Sigma_T \widehat{U} R) \\ & \leq \text{tr}((R^T \widehat{U}^T \Sigma_S \widehat{U} R)^{-1} U^T \Sigma_T U) + \|R^T \widehat{U}^T \Sigma_T \widehat{U} R - U^T \Sigma_T U\| \text{tr}((\widehat{U}^T \Sigma_S \widehat{U})^{-1}) \\ & \leq \text{tr}((R^T \widehat{U}^T \Sigma_S \widehat{U} R)^{-1} U^T \Sigma_T U) + 2\|\Delta\|\|\Sigma_T\| \text{tr}((\widehat{U}^T \Sigma_S \widehat{U})^{-1}) \\ & \leq \text{tr}((R^T \widehat{U}^T \Sigma_S \widehat{U} R)^{-1} U^T \Sigma_T U) + 2\|\Delta\|\|\Sigma_T\| k C_{\min}^{-1} \\ & \leq \text{tr}((R^T \widehat{U}^T \Sigma_S \widehat{U} R)^{-1} U^T \Sigma_T U) + \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U) \end{aligned} \quad (48)$$

$$\begin{aligned} & \text{when } \Delta \leq \frac{\lambda_k \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U)}{4k\|\Sigma_T\|}. \text{ Also, we have} \\ & \| (R^T \widehat{U}^T \Sigma_S \widehat{U} R)^{-1} - (U^T \Sigma_S U)^{-1} \| \leq \| (R^T \widehat{U}^T \Sigma_S \widehat{U} R)^{-1} \| \| (U^T \Sigma_S U)^{-1} \| \| R^T \widehat{U}^T \Sigma_S \widehat{U} R - U^T \Sigma_S U \| \\ & \leq 4\lambda_k^{-2} \lambda_1 \Delta, \end{aligned}$$

2825 therefore

$$\begin{aligned} & \text{tr}((R^T \widehat{U}^T \Sigma_S \widehat{U} R)^{-1} U^T \Sigma_T U) \leq \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U) + \| (R^T \widehat{U}^T \Sigma_S \widehat{U} R)^{-1} - (U^T \Sigma_S U)^{-1} \| \text{tr}(U^T \Sigma_T U) \\ & \leq \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U) + 4\lambda_k^{-2} \lambda_1 \Delta \text{tr}(U^T \Sigma_T U) \\ & \leq 2 \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U), \end{aligned} \quad (50)$$

2830 if $\Delta \leq \frac{\lambda_k^2 \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U)}{4\lambda_1 \text{tr}(U^T \Sigma_T U)}$. Combining (47), (48) and (50) we have

$$\mathbb{E}_\epsilon[\|A_3\|_{\Sigma_T}^2] \leq \mathcal{O}\left(\frac{1}{n} v^2 \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U)\right),$$

2835 whenever $\Delta \leq \frac{\lambda_k^2 \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U)}{4\lambda_1 k\|\Sigma_T\|} \leq \min\left\{\frac{\lambda_k^2 \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U)}{4\lambda_1 \text{tr}(U^T \Sigma_T U)}, \frac{\lambda_k \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U)}{4k\|\Sigma_T\|}\right\}$
2836 and $n \gtrsim \sigma^4 C_{\min}^{-4} \|\Sigma_S\|^2 \|\Sigma_T\|^2 \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U)^{-2} k^3 \log(1/\delta)$, with probability at least
2837 $1 - \delta$. Notice that $U^T \Sigma_S U = \Sigma_{S,k}$ and $U^T \Sigma_T U = \Sigma_{T,k}$, therefore the result is exactly what we
2838 want. \square

2839
2840 *Proof of Lemma 34.* Recall $A_2 := \widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \beta_\perp^*$. Also we have

$$\begin{aligned} \|\widehat{U}^T \Sigma_T \widehat{U}\| &= \|R^T \widehat{U}^T \Sigma_T \widehat{U} R\| \\ &\leq \|U^T \Sigma_T U\| + \|R^T \widehat{U}^T \Sigma_T \widehat{U} R - U^T \Sigma_T U\| \\ &\leq \|U^T \Sigma_T U\| + 2\Delta\|\Sigma_T\| \end{aligned} \quad (51)$$

2846 Therefore

$$\begin{aligned} \|A_2\|_{\Sigma_T}^2 &= \|\beta_\perp^{*\top} X^T X \widehat{U} (\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T \Sigma_T \widehat{U} (\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \beta_\perp^*\| \\ &\leq \|X \widehat{U} (\widehat{U}^T X^T X \widehat{U})^{-1} (\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T\| \|\widehat{U}^T \Sigma_T \widehat{U}\| \|X \beta_\perp^*\|^2 \\ &\leq \|A\| (\|U^T \Sigma_T U\| + 2\Delta\|\Sigma_T\|) \|X \beta_\perp^*\|^2 \\ &\leq 2\|A\| \|U^T \Sigma_T U\| \|X \beta_\perp^*\|^2 \end{aligned} \quad (52)$$

2853 when $\Delta \leq \frac{\|U^T \Sigma_T U\|}{2\|\Sigma_T\|}$, where we let $A = \frac{1}{n} \frac{X \widehat{U}}{\sqrt{n}} (\widehat{U}^T \frac{X^T X}{n} \widehat{U})^{-2} \frac{\widehat{U}^T X^T}{\sqrt{n}}$. If we define $B = \frac{X \widehat{U}}{\sqrt{n}} \in$
2854 $\mathbb{R}^{n \times r}$, then $A = \frac{1}{n} B(B^T B)^{-2} B^T$. Let the SVD of B be $B = P M O^T$, where $P \in \mathbb{R}^{n \times k}$,
2855 $M, O \in \mathbb{R}^{k \times k}$, then

$$\begin{aligned} \|A\|_2 &= \frac{1}{n} \|B(B^T B)^{-2} B^T\|_2 \\ &= \frac{1}{n} \|P M O^T (O M^2 O^T)^{-2} O M P^T\|_2 \\ &= \frac{1}{n} \|P M^{-2} P^T\|_2 \end{aligned}$$

$$\begin{aligned}
&\leq \frac{1}{n} \|M^{-2}\|_2 \\
&= \frac{1}{n} \|(B^T B)^{-1}\|_2
\end{aligned} \tag{53}$$

Let $F = (\widehat{U}^T \frac{X^T X}{n} \widehat{U})^{-1} - (\widehat{U}^T \Sigma \widehat{U})^{-1}$. Recall (33), which states that with probability at least $1 - \delta$, we have $\|F\| \leq \frac{1}{3} C_{\min}^{-1} \leq \frac{2}{3} \lambda_k^{-1}$ when $n \gtrsim \sigma^4 C_{\min}^{-2} \|\Sigma_S\|^2 k \log(1/\delta)$ and $\Delta \leq \frac{\lambda_k}{4\lambda_1}$. Therefore

$$\begin{aligned}
\|A\| &\leq \frac{1}{n} \|(\widehat{U}^T \frac{X^T X}{n} \widehat{U})^{-1}\| \\
&= \|(\widehat{U}^T \Sigma_S \widehat{U})^{-1} + F\| \\
&\leq \frac{1}{n} \|(\widehat{U}^T \Sigma_S \widehat{U})^{-1}\| + \|F\| \\
&\leq \mathcal{O}\left(\frac{1}{n} \lambda_k^{-1}\right).
\end{aligned} \tag{54}$$

Thus $\|A\| \leq \mathcal{O}(\lambda_k^{-1})$. As for $\|X\beta_{\perp}^*\|^2$, notice that the first- k entries of β_{\perp}^* are zero, therefore $X\beta_{\perp}^* = X_{-k}\beta_{-k}^*$. by Lemma 35,

$$\|\beta_{-k}^{*T} \left(\frac{X_{-k}^T X_{-k}}{n}\right) \beta_{-k}^* - \beta_{-k}^{*T} \Sigma_{S,-k} \beta_{-k}^*\| \leq \mathcal{O}(\sigma^2 \|\beta_{-k}^*\|^2 \|\Sigma_{S,-k}\|) \left(\sqrt{\frac{1}{n}} + \frac{1}{n} + \sqrt{\frac{\log(1/\delta)}{n}} + \frac{\log(1/\delta)}{n}\right). \tag{55}$$

Therefore we have

$$\begin{aligned}
\|X\beta_{\perp}^*\|^2 &= n \beta_{-k}^{*T} \left(\frac{X_{-k}^T X_{-k}}{n}\right) \beta_{-k}^* \\
&\leq n (\beta_{-k}^{*T} \Sigma_{S,-k} \beta_{-k}^* + \|\beta_{-k}^{*T} \left(\frac{X_{-k}^T X_{-k}}{n}\right) \beta_{-k}^* - \beta_{-k}^{*T} \Sigma_{S,-k} \beta_{-k}^*\|) \\
&\leq \mathcal{O}(n \|\beta_{-k}^*\|^2 \|\Sigma_{S,-k}\|).
\end{aligned} \tag{56}$$

Combining (52)(54) and (56), we have

$$\|A_2\|_{\Sigma_T}^2 \leq \mathcal{O}\left(\frac{\|U^T \Sigma_T U\| \|\beta_{-k}^*\|^2 \|\Sigma_{S,-k}\|}{\lambda_k}\right) \tag{57}$$

when $n \gtrsim \sigma^4 C_{\min}^{-2} \|\Sigma_S\|^2 k \log(1/\delta)$ and $\Delta \leq \min\{\frac{\|U^T \Sigma_T U\|}{2\|\Sigma_T\|}, \frac{\lambda_k}{4\lambda_1}\}$. \square

Finally we prove Lemma 30 in the following.

Proof of Lemma 30. In the first step, we obtain $\widehat{U} \in \mathbb{R}^{d \times k}$ by selecting the top- k eigenvectors of the sample covariance matrix $\widehat{\Sigma}_S := \frac{1}{n} XX^T = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$ using PCA. Then by Davis-Kahan theorem (Chen et al., 2021, Corollary 2.8),

$$\Delta \leq \frac{2\|\widehat{\Sigma}_S - \Sigma_S\|}{\lambda_k - \lambda_{k+1}}. \tag{58}$$

Therefore it remains to bound $\|\widehat{\Sigma}_S - \Sigma_S\|$. Applying Lemma 28, we immediately have

$$\|\widehat{\Sigma}_S - \Sigma_S\| \leq C \sigma^4 \left(\sqrt{\frac{r + \log \frac{1}{\delta}}{n}} + \frac{r + \log \frac{1}{\delta}}{n} \right) \lambda_1$$

where $r = \frac{\sum_{i=1}^n \lambda_i}{\lambda_1}$. Together with (58), we have with probability at least $1 - \delta$,

$$\Delta \leq C \sigma^4 \left(\sqrt{\frac{r + \log \frac{1}{\delta}}{n}} + \frac{r + \log \frac{1}{\delta}}{n} \right) \frac{\lambda_1}{\lambda_k - \lambda_{k+1}}.$$

\square