
Benign Overfitting in Out-of-Distribution Generalization of Linear Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Benign overfitting refers to the phenomenon where a over-paramterized model fits
2 the training data perfectly, including noise in the data, but still generalizes well to
3 the unseen test data. While prior work provide a solid theoretical understanding
4 of this phenomenon under the in-distribution setup, modern machine learning of-
5 ten operates in a more challenging Out-of-Distribution (OOD) regime, where the
6 target (test) distribution can be rather different from the source (training) distribu-
7 tion. In this work, we take an initial step towards understanding benign overfitting
8 in the OOD regime by focusing on the basic setup of over-parameterized linear
9 models under covariate shift. We provide non-asymptotic guarantees proving that,
10 when the target covariance satisfies certain structural conditions, benign overfit-
11 ting occurs in standard ridge regression even under the OOD regime. We identify
12 a number of key quantities relating source and target covariance, which govern the
13 performance of OOD generalization. Our result is sharp, which provably recovers
14 prior in-distribution benign overfitting guarantee (Tsigler & Bartlett, 2023), as
15 well as under-parameterized OOD guarantee (Ge et al., 2024) when specializing
16 to each setup. Moreover, we also present theoretical results for a more general
17 family of target covariance matrix, where standard ridge regression only achieves
18 a slow statistical rate of $\mathcal{O}(1/\sqrt{n})$ for the excess risk, while Principal Component
19 Regression (PCR) is guaranteed to achieve the fast rate $\mathcal{O}(1/n)$, where n is the
20 number of samples.

21 1 Introduction

22 In modern machine learning, distribution shift has become a ubiquitous challenge where models
23 trained on a source data distribution are tested on a different target distribution (Zou et al., 2018;
24 Hendrycks & Dietterich, 2019; Guan & Liu, 2021; Koh et al., 2021). Generalization under distribu-
25 tion shift, known as Out-of-Distribution (OOD) generalization, remains a fundamental issue in the
26 practical application of machine learning (Recht et al., 2019; Hendrycks et al., 2021; Miller et al.,
27 2021; Wenzel et al., 2022). While there has been extensive work on the theoretical understanding
28 of OOD generalization, most of it has focused on under-parameterized models (Shimodaira, 2000;
29 Lei et al., 2021; Ge et al., 2024; Zhang et al., 2022). However, over-parameterized models, such
30 as deep neural networks and large language models (LLMs), which have more parameters than
31 training samples, are widely used in modern machine learning. Surprisingly, despite the classic
32 bias-variance tradeoff for under-parameterized models, over-parameterized models tend to overfit
33 the data while still achieving strong in-distribution generalization, a phenomenon known as benign
34 overfitting (Hastie et al., 2022; Shamir, 2023) or harmless interpolation (Muthukumar et al., 2020).
35 Therefore, it is crucial to theoretically understand how benign overfitting shapes OOD generalization
36 in over-parameterized models.

37 It is established in overparameterized models that “benign overfitting” occurs when the data essen-
 38 tially resides on a low-dimensional manifold. The manifold assumption (Belkin & Niyogi, 2003) is
 39 widely applicable across image, speech and language data, where although features are embedded
 40 in a high-dimensional ambient space, their generation is governed by a few degrees of freedom im-
 41 posed by physical constraints (Niyogi, 2013). Specifically, the covariance matrix of the data should
 42 be characterized by several major directions with large eigenvalues while the remaining directions
 43 are high-dimensional but have smaller scale. In this setting, even though the estimator may over-
 44 fit the noise, it can still capture the signal in the major directions while the noise is dampened in
 45 the minor directions. Recent non-asymptotic analyses have provided upper bounds on the excess
 46 risk for the minimum-norm interpolant and over-parameterized ridge estimator under this frame-
 47 work (Bartlett et al., 2020; Hastie et al., 2022; Tsigler & Bartlett, 2023).

48 However, theoretical characterization of OOD generalization in over-parameterized models remains
 49 elusive. In this paper, we take an initial step toward characterizing OOD generalization in over-
 50 parameterized models under *general* covariate shift, a standard assumption for OOD generaliza-
 51 tion (Ben-David et al., 2006), where the conditional distribution of the outcome given the covariates
 52 remains invariant. We derive the first vanishing, non-asymptotic excess risk bound for ridge regres-
 53 sion and minimum-norm interpolation, assuming that the source covariance is dominated by a few
 54 major eigenvalues, which satisfies the benign overfitting condition. But we allow the target covari-
 55 ance to be arbitrary. This result contrasts with recent work that either addresses only a restrictive
 56 form of covariate shift (Hao et al., 2024; Mallinar et al., 2024) or provides excess risk bounds that
 57 asymptotically remain above a constant (Tripuraneni et al., 2021b; Hao et al., 2024).

58 In summary, our excess risk bound identifies several key quantities that relate the source and target
 59 covariance, suggesting that “benign overfitting” occurs when these quantities are well controlled. In
 60 such cases, the target distribution data lies on the low-dimensional manifold of the source distribu-
 61 tion. Otherwise, ridge regression may incur excess risk, lower bounded by the slow statistical rate
 62 of $\mathcal{O}(1/\sqrt{n})$. In contrast, we show that principal component regression (PCR) achieves the fast rate
 63 of $\mathcal{O}(1/n)$ in such scenarios.

64 **Our contributions.**

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

84 **1.1 Related work**

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

90 Research in this field can be divided into two categories based on assumptions about the spectral
91 structure of the sample covariance. The first category assumes an almost isotropic sample covari-
92 ance matrix with a bounded condition number or an isotropic prior distribution of parameters (Belkin
93 et al., 2020). In this case, a limiting covariance spectral structure may emerge when $n \asymp d$ and both
94 tend to infinity, allowing for asymptotic risk bounds (Dobriban & Wager, 2018; Richards et al.,
95 2021). However, ridgeless regression is sub-optimal in this setting unless the signal-to-noise ratio
96 is infinite (Wu & Xu, 2020), and non-asymptotic error bounds are lacking. Our work falls into the
97 second category, focusing on covariance model where a small number of eigenvalues dominate the
98 sample covariance, and the signal is concentrated in the subspace spanned by the leading eigen-
99 vectors (Bibas et al., 2019; Chinot & Lerasle, 2022; Hastie et al., 2022). Linear regression can
100 be optimal without regularization under this covariance structure (Kobak et al., 2020), which is of
101 practical interest because ridgeless regression is equivalent as gradient descent from zero initializa-
102 tion (Zhou et al., 2020). Sharp non-asymptotic bounds for variance and bias in ridge regression have
103 been derived (Bartlett et al., 2020; Tsigler & Bartlett, 2023).

104 Extending the analysis of ridgeless estimators (i.e., minimum norm interpolants), uniform conver-
105 gence bounds for generalization error have been studied for all interpolants with arbitrary norms.
106 However, uniformly bounding the difference between population and empirical errors generally
107 fails to ensure a consistent predictor (Zhou et al., 2020), necessitating strong assumptions on dis-
108 tributions (Koehler et al., 2021) or hypothesis classes (Negrea et al., 2020). Over-parameterization
109 theory for linear models has also been applied to two-layer neural networks approximated via kernel
110 ridge regression (Liang et al., 2020; Ghorbani et al., 2020, 2021; Bartlett et al., 2021; Mei & Montanari,
111 2022; Mei et al., 2022; Montanari & Zhong, 2022; Simon et al., 2023), though this lies beyond
112 the scope of the present work.

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

123 Research on over-parameterized models under distribution shift has largely focused on covariate
124 shift in linear regression. Importance weighting for over-parameterized models (Chen et al., 2024)
125 and general sample reweighting offer no advantage over ERM since both converge to the same esti-
126 mator via gradient descent (Zhai et al., 2022). Consequently, much literature focuses on minimum-
127 norm interpolation as the natural ERM solution. For isotropic signals, Tripuraneni et al. (2021a)
128 prove that over-parameterization improves robustness to covariate shift, deriving an asymptotic gen-
129 eralization bound decreasing with d/n . Under the essentially low-rank covariance model, Hao et al.
130 (2024) derive a non-asymptotic bound for a specific covariate shift where features are translated by
131 a constant but the covariance matrix is preserved. However, a constant excess risk remains in their
132 bound due to estimation variance. Kausik et al. (2024) study a linear model with additive noise
133 on covariates when data strictly lies in a low-dimensional subspace, also showing a non-vanishing
134 bound. Mallinar et al. (2024) investigate minimum-norm interpolation with independent covariates
135 and simultaneously diagonalizable source and target covariance matrices, allowing them to directly
136 extend in-distribution bounds of Bartlett et al. (2020); Tsigler & Bartlett (2023). Still, their esti-
137 mation bias bound is looser than ours due to a gap compared to Tsigler & Bartlett (2023)’s sharp
138 bound even when the source matches the target. In contrast, our work achieves the first vanishing
139 non-asymptotic error bound for general covariate shift, assuming only finite second moments for the
140 target covariance matrix.

141 There also exist a line of work that considers non-parametric models under covariate shift (Kpotufe
142 & Martinet, 2018; Hanneke & Kpotufe, 2019; Pathak et al., 2022; Ma et al., 2023), presenting
143 minimax results controlling by a transfer-exponent that measures the similarity between source and
144 target, though this lies beyond the scope of our work.

145 **Principal component regression.** Principal component regression (PCR) has been designed as
 146 a method of treating multicollinearity problems in high-dimensional linear regression, where the
 147 covariates have a latent, low-dimensional representation (Massy, 1965; Jeffers, 1967; Jolliffe, 1982;
 148 Jeffers, 1981). PCR has been widely used in statistics (Liu et al., 2003), chemometrics (Næs &
 149 Martens, 1988; Sun, 1995; Vigneau et al., 1997; Depczynski et al., 2000; Keithley et al., 2009),
 150 construction management (Chan & Park, 2005), environmental science (Kumar & Goyal, 2011;
 151 Hidalgo et al., 2000), signal processing (Huang & Yang, 2012) and etc.

152 Regarding the theory for PCR, Hadi & Ling (1998) give conditions under which PCR will fail. Bair
 153 et al. (2006) suggest selecting principal components based on their association with the outcome,
 154 and provide corresponding asymptotic consistency results. Xu & Hsu (2019) give asymptotic risk
 155 bounds for PCR, under different number of selected components k . They show that the “double
 156 descent” behaviour also happens in PCR when k/d grows, where d is the data dimension. Most
 157 related to our work, Agarwal et al. (2019) provide non-asymptotic error bounds of PCR, and show
 158 that the error will decay as $\mathcal{O}(1/\sqrt{n})$ (n is the sample size) given that all the singular values of the
 159 data matrix are of the same order. Agarwal et al. (2020) further improves the rate to $\mathcal{O}(1/n)$. How-
 160 ever, the aforementioned two results both consider fixed design with strict low-rank assumptions,
 161 therefore not applicable to our setting of OOD-generalization.

162 2 Covariate shift setup under over-parameterization

163 2.1 Data with covariate shift

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

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

179 2.2 Learning procedure and evaluation

180 The learning procedure involves training a linear model with n i.i.d. samples $\{(x_i, y_i)\}_{i=1}^n$ drawn
 181 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 $\epsilon :=$
 182 $(\epsilon_1, \dots, \epsilon_n)^T$. We focus on models $\hat{\beta}(Y)$ that are linear in Y , allowing us to write $\hat{\beta}(Y) = \hat{\beta}(X\beta^*) +$
 183 $\hat{\beta}(\epsilon)$. We consider ridge regression and principal component regression as two instances of such
 184 algorithms. With a regularization coefficient $\lambda \geq 0$, the ridge estimator in the over-parameterized
 185 setting, where $n < d$, is defined as:

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

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

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

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

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

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

192 2.3 The structure of covariance in benign overfitting

193 Throughout this paper, we follow the convention of [Tsigler & Bartlett \(2023\)](#), consider the *source*
 194 covariance matrix Σ_S that has only a few number of high variance directions but a very large number
 195 of low variance directions with similar magnitude. We will also refer to those high variance direc-
 196 tions of the source as “major directions”, and those low variance directions as “minor directions”.
 197 We denote the number of major directions as k . For remaining $d - k$ minor directions, we use the
 198 following notions of effective ranks to approximately capture the number of directions that have a
 199 similar scale. Let the ridge regularization coefficient be $\lambda \geq 0$, we define:

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

200 We have $1 \leq r_k \leq R_k$. When $\lambda = 0$, we further have $R_k \leq d - k$. We denote the first k columns
 201 of X as X_k and the remaining $d - k$ columns as X_{-k} . Correspondingly, we partition β^* into β_k^*
 202 and β_{-k}^* . The covariance matrix blocks along the diagonals are denoted by $\Sigma_{S,k}$, $\Sigma_{S,-k}$, $\Sigma_{T,k}$ and
 203 $\Sigma_{T,-k}$. To facilitate our presentation, we define

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

204 These quantities turn out to be crucial in the analysis.

205 3 Over-parameterized ridge regression

206 In the context of in-distribution generalization for overparameterized linear models, [Bartlett et al.](#)
 207 [\(2020\)](#) and [Tsigler & Bartlett \(2023\)](#) demonstrate that the ridge estimator (minimum-norm interpo-
 208 late estimator as a special case) can effectively learn the signal from the subspace of data spanned
 209 by the major eigenvectors, while benignly overfitting noise from the minor directions under cer-
 210 tain scenarios. They argue that, when the true signal mainly lies in the major directions, and the
 211 minor directions have small scale but high effective rank, benign overfitting is possible. In this sec-
 212 tion, we explore whether this mechanism still holds under covariate shift. We derive upper bounds
 213 (Theorem 2) for the excess risk of the ridge estimator in the context of overparameterized OOD-
 214 generation, demonstrating that “benign overfitting” also happens under covariate shift, given that
 215 the target distribution’s covariance structure remains dominated by the first k dimensions. To be
 216 specific, we show that \mathcal{T} characterizes the shift in the major directions; the *overall magnitude* of
 217 $\Sigma_{T,-k}$, which characterizes the shift in the minor directions, is crucial for benign overfitting. When
 218 the *overall magnitude* of $\Sigma_{T,-k}$, scales similarly to or smaller than those of the source, ridge
 219 regression achieves the same non-asymptotic error rate under covariate shift as in the in-distribution
 220 setting. Surprisingly, although high effective rank in the minor directions of source is essential for
 221 benign overfitting, for target distribution only the overall magnitude matters.

222 3.1 Warm-up: in-distribution benign overfitting

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

232 **Assumption 1** (CondNum(k, δ, L), [\(Tsigler & Bartlett, 2023\)](#)). Define a matrix $A_k = \lambda I_n +$
 233 $X_{-k} X_{-k}^T$. With probability at least $1 - \delta$, A_k is positive definite and has a condition number no
 234 greater than L , i.e.,

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

235 **Remark 1.** This assumption is essentially assuming the minor directions have effective rank signif-
 236 icantly larger than n . As an evidence, Tsigler & Bartlett (2023) prove that if CondNum holds, then
 237 the effective rank r_k is lower bounded by n/L . On the other hand, a lower bound on the effective
 238 rank r_k can also imply an upper bound of the condition number of A_k . See Tsigler & Bartlett (2023,
 239 Lemma 3) for further detail.

240 Assuming CondNum, Tsigler & Bartlett (2023) obtain sharp upper bounds for both the variance and
 241 bias of the ridge estimator, with matching lower bounds (see their Theorem 2). To facilitate the
 242 presentation, we use $\tilde{\lambda} := \lambda + \sum_{j>k} \lambda_j$ to denote the overall regularization term.

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

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

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

253 3.2 Out-of-Distribution benign overfitting

254 We now investigate the out-of-distribution performance of ridge estimator. Intuitively, when all the
 255 minor components vanish (both on the source and the target), the over-parameterized ridge regres-
 256 sion is actually reduced to the usual ridge regression on the major directions, thus achieving a rate
 257 of $\tilde{O}(\text{tr}[\mathcal{T}]/n)$ as Ge et al. (2024) demonstrate. When the minor components do not vanish, high
 258 effective rank of minor components on the source is essential for “benign overfitting”, as Tsigler &
 259 Bartlett (2023) demonstrate. However, we argue that, regarding the target distribution, only the *over-*
 260 *all magnitude* of those minor components is crucial for benign overfitting. The reason is that, when
 261 the minor directions of source have effective rank much larger than n , the n -dimensional subspace
 262 spanned by training samples is already almost orthogonal to any test point, with a high probability.
 263 Therefore, no special spectral structure of the target is needed for benign overfitting. Only small
 264 overall magnitude of those minor components on target is required.

265 We formalize those intuitive claims, by deriving upper bounds for both the variance and bias of
 266 ridge regression under covariate shift, assuming a source distribution similar to the in-distribution
 267 case. Our upper bound is sharp, and can be applied to any target distributions, reducing to Tsigler
 268 & Bartlett (2023)’s bound (Theorem 1) when the target and source distributions are aligned. Addi-
 269 tionally, we recover Ge et al. (2024)’s sharp bound for under-parameterized linear regression under
 270 covariate shift when the high-dimensional minor components vanish.

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

$$\frac{V}{cv^2} \leq \frac{k}{n} \cdot \frac{\text{tr}[\mathcal{T}]}{k} + \frac{n}{R_k} \cdot \frac{\text{tr}[\mathcal{U}]}{\text{tr}[\mathcal{V}]},$$

$$\frac{B}{c} \leq B_{\text{ID}} \cdot \left(\|\mathcal{T}\| + \frac{n}{r_k} \frac{\|\Sigma_{T,-k}\|}{\|\Sigma_{S,-k}\|} \right).$$

273 where $\mathcal{T}, \mathcal{U}, \mathcal{V}$ are defined in Equation (1).

274 $N = \text{Poly}(k + \ln(1/\delta), \lambda_1 \lambda_k^{-1}, 1 + \tilde{\lambda} \lambda_k^{-1})$. $\text{Poly}(\cdot)$ denotes a polynomial function.

275 Recall B_{ID} is the upper bound for bias given by Theorem 1, we can see that Theorem 2 establishes
 276 an upper bound for the excess risk of ridge regression under general covariate shift, expressed as a
 277 multiplicative form of Theorem 1’s results. This formulation enables a direct analysis of the impact

278 of covariate shift on the bias and variance of ridge estimators, compared to the in-distribution case.
 279 The first conclusion is that Theorem 2 well reduces to the corresponding result in Theorem 1 when
 280 no distribution shift occurs—i.e., $\Sigma_S = \Sigma_T$. This connection follows directly from the condition
 281 $n < r_k$.

282 The second conclusion is that covariate shift in the first k dimensions and last $d - k$ dimensions
 283 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
 284 risk. Therefore, as long as these factors are bounded by constants, over-parameterized ridge regres-
 285 sion achieves the same non-asymptotic rate of excess risk under covariate shift as the in-distribution
 286 setting. This scenario, well addressed by ridge regression, occurs when the target distribution’s co-
 287 variance structure remains dominated by the first k dimensions. In the following, we analyze the
 288 impact of the factors introduced by covariate shift on both the major and minor directions.

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

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

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

301 **2. The overall magnitude of $\Sigma_{T,-k}$ is crucial for benign overfitting.** Under covariate shift within
 302 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
 303 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
 304 other words, matching our intuition, if the *overall magnitude* of the minor components of tar-
 305 get covariance scales similarly to or smaller than those of the source, in terms of the covariance
 306 norms, “benign overfitting” also happens under covariate shift. Importantly, this condition does
 307 not impose constraints on the internal spectral structure of the minor components of target co-
 308 variance. For example, we do not force each eigenvalue of $\Sigma_{T,-k}$ to scale with its corresponding
 309 eigenvalue of $\Sigma_{S,-k}$ in decreasing order, as assumed in prior work (Mallinar et al., 2024). Sur-
 310 prisingly, for benign overfitting to happen, it is essential for the source distribution to have high
 311 effective rank in the minor directions; however for target distribution, only the overall magnitude
 312 matters.

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

317 **Remark 2 (Sample complexity).** We have assumed $n \geq c_x N$ in Theorem 2. The explicit formula
 318 for N is deferred to Theorem 25 and Remark 8. Here we summarize the sample complexity required
 319 for the bound to hold. The dependence on k varies between $\Omega(k)$ and $\Omega(k^3)$, depending on the
 320 degree of covariate shift. The optimal case, aligning with the sample complexity of classic linear
 321 regression, occurs when $\Sigma_{S,k} \approx \Sigma_{T,k}$. The worst case arises when there is significant covariate shift
 322 in the first k dimensions, such as when the test data lies predominantly in the subspace of the first
 323 dimension. This variation in sample complexity under covariate shift parallels the analysis of Ge
 324 et al. (2024) (see their Theorem 4.2) for the under-parameterized setting. Additionally, we require
 325 $n \gg \lambda + \sum_{j>k} \lambda_j$, ensuring that the regularization is not too strong to introduce a bias exceeding a
 326 constant (as reflected in the first bias term). On the other hand, we assume $n < r_k$ in the theorem,
 327 consistent with the over-parameterized regime and Assumption 1, where the last $d - k$ components
 328 are considered to be essentially high-dimensional.

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

333 Finally, Theorem 2 suggests an $\mathcal{O}(1/n)$ vanishing error under several conditions that naturally fol-
 334 low from the previous discussions, which we now state rigorously. First, the covariate space decom-
 335 poses into subspaces spanned by low-dimensional major directions and high-dimensional minor
 336 directions, with $k = \mathcal{O}(1)$ and $R_k = \Omega(n^2)$. Second, the low-rank covariance structure is preserved
 337 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 predominantly
 338 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 regular-
 339 ization is not excessively strong to introduce a significant bias, with $\tilde{\lambda} = \lambda + \sum_{j>k} \lambda_j = \mathcal{O}(\sqrt{n})$.

340 4 Large shift in minor directions

341 In the previous section, we established an upper bound for overparameterized ridge regression under
 342 covariate shift. We showed that when the shift in the minor directions is controlled—specifically,
 343 when the overall magnitude of $\Sigma_{T,-k}$ is small—“benign overfitting” also occurs under covariate
 344 shift. However, when the shift in minor directions is significant, meaning the target covariance
 345 matrix has many large eigenvalues with corresponding eigenvectors outside the major directions, the
 346 excess risk for ridge regression deteriorates. In this section, we further illustrate the limitations of
 347 ridge regression in such cases by providing a lower bound for its performance for large distribution
 348 shift in the minor directions, showing that it can only achieve the slow rate of $\mathcal{O}(1/\sqrt{n})$ for the
 349 excess risk. On the other hand, it is natural to consider alternative algorithms to ridge regression
 350 in this scenario. We demonstrate that even with a large shift in the minor directions, principal
 351 component regression (PCR) is guaranteed to achieve the fast rate $\mathcal{O}(1/n)$, provided that the signal
 352 β^* lies primarily within the subspace spanned by the major directions. Moreover, PCR does not
 353 require the minor directions to have a high effective rank in the source distribution, highlighting its
 354 advantage over ridge regression in such cases.

355 4.1 Slow rate for ridge regression

356 In this subsection, we demonstrate the limitations of ridge regression when the overall magnitude of
 357 $\Sigma_{T,-k}$ is large. Consider an instance where Σ_S has its first k components as $\Theta(1)$, while the minor
 358 directions have eigenvalues of $o(1)$. If we set $\Sigma_T = I_d$, in contrast to the “benign overfitting” regime
 359 described in Theorem 2, ridge regression will have a large excess risk for this instance. Although the
 360 signal from the major directions is effectively captured, the signal in the minor directions is nearly
 361 lost. Unlike the case in Section 3, here the estimation error in the minor directions is crucial because
 362 the target distribution has significant components in these directions. We formalize this intuitive
 363 example through the following theorems:

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

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

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

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

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

368 **Theorem 4.** Consider the same instance of Σ_S as in Corollary 3. Assume $\Sigma_T = I_d$ and $\lambda = \sqrt{n}$.
 369 There exists an absolute constant $C > 0$, such that for some $0 < \delta < 1$, $N_2 > 0$ and for any
 370 $n > N_2$, with probability at least $1 - \delta$, we have $V \geq Cv^2$.

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

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

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

380 4.2 Fast rate for principal component regression

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

389 Principal Component Regression (PCR).

- 390 • **Step 1: Obtain an estimator \hat{U} of the top- k subspace of Σ_S .** For simplicity, we assume a sample
 391 size of $2n$ and use the first half of the data to compute \hat{U} by principal component analysis (PCA)
 392 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
 393 i -th eigenvector of $\hat{\Sigma}_S$.
- 394 • **Step 2: Use the data projected on \hat{U} to conduct linear regression.** With a little abuse of
 395 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
 396 $(y_{n+1}, \dots, y_{2n})^T$. If we let $Z := X\hat{U} \in \mathbb{R}^{n \times k}$ be the projected data matrix, the estimator $\hat{\beta}$ we
 397 obtained is given by

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

398 Consider the scenario where the last $d - k$ components of the true signal β^* is exactly zero, namely
 399 $\beta_{-k}^* = 0$. We can imagine that if the subspace represented by \hat{U} is exactly the same as the subspace
 400 represented by $U = \begin{pmatrix} I_k \\ 0 \end{pmatrix} \in \mathbb{R}^{d \times k}$, (i.e., the first k components), then PCR is actually doing linear
 401 regression using only the first k components of the samples, therefore will only have a excess risk
 402 induced by the usual variance of linear regression in the major directions. Under this scenario, no
 403 matter how large $\|\Sigma_{T,-k}\|$ is, the PCR estimator have zero estimates on the last $d - k$ components,
 404 therefore avoid inducing large excess risk. Further more, if the distance between \hat{U} and U is not
 405 zero, there should be another term in the excess risk induced by the estimation error of \hat{U} . We
 406 formalize this intuitive claim as the following upper bound for the excess risk of PCR. To facilitate
 407 the presentation, we introduce the following quantity for measuring the estimation accuracy of \hat{U} .
 408 We define $\Delta = \text{dist}(\hat{U}, U) := \|U U^T - \hat{U} \hat{U}^T\|$, the distance between the subspace spanned by the
 409 columns of \hat{U} and U . Then we have the following theorem:

410 **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
 411 excess risk of PCR estimator $\hat{\beta}$ with probability $1 - \delta$:

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

412 where Θ, N_1 is defined as follows:

$$413 \Theta^{-1} = \text{Poly}(\lambda_1 \lambda_k^{-1}, \|\Sigma_T\| \lambda_k^{-1}, k \text{tr}(\mathcal{T})^{-1}),$$

$$414 N_1 = \text{Poly}(\sigma, \lambda_1 \lambda_k^{-1}, \|\Sigma_T\| \lambda_k^{-1}, k \ln(1/\delta), k \text{tr}(\mathcal{T})^{-1}).$$

415 **Remark 5.** Theorem 5 is a special case of Lemma 31. For detailed characterization of Θ and N_1 ,
 416 as well as an upper bound for cases where $\beta_{-k}^* \neq 0$, one can refer to Lemma 31 for detail.

417 The excess risk upper bound given by Theorem 5 consists of two terms. The variance term $\frac{\text{tr}(\mathcal{T})}{n}$
 418 is incurred by the nature of linear regression on the major directions, and is unavoidable even if the
 419 subspace estimation is accurate (i.e., $\Delta = 0$). This term also appears in the first term of variance in
 420 Theorem 2, and exactly matches the sharp rate $\text{tr}[\Sigma_S^{-1}\Sigma_T]/n$ for under-parameterized linear regres-
 421 sion under covariate shift (Ge et al., 2024). The second term $\|\beta^*\|^2(\frac{\lambda_1}{\lambda_k})^2\|\Sigma_T\|\Delta^2$ is the bias term
 422 induced by the estimation error of the subspace in the first step. We can see that it has a quadratic
 423 dependence on Δ . If we combine Theorem 5 with a control of Δ , we can get the end-to-end excess
 424 risk upper bound of PCR. For controlling Δ , we have the following lemma:

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

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

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

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

431 Combining Theorem 5 and Lemma 6, an end-to-end error bound for PCR can be derived (for a
 432 detailed theorem, one can refer to Theorem 29), suggesting that PCR will achieve a small excess risk,
 433 as long as the major directions and the minor directions are well separated, and the effective rank of
 434 the entire source covariance matrix is small. Contrast to ridge regression, PCR does not require the
 435 minor components to have high-effective rank. This shows the superiority of PCR compared with
 436 ridge regression under certain scenarios.

437 As an example, consider the instance in Theorem 4, where $k, \|\Sigma_T\|, \lambda_1, \lambda_k$ are all $\Theta(1)$. In this
 438 case, the variance term will scale as $1/n$, and the bias term scales as $\mathcal{O}(\Delta^2)$. Notice that in this
 439 instance, $r = \Theta(1)$, therefore $\Delta \leq \mathcal{O}(1/\sqrt{n})$. We conclude that in this instance, PCR will achieve
 440 a $\mathcal{O}(1/n)$ rate even when $\Sigma_T = I_d$. Comparing with the excess risk for ridge regression, which is
 441 at least $1/\sqrt{n}$, PCR shows its superiority against ridge regression under the scenario where the shift
 442 in minor directions is large.

443 5 Conclusion and discussion

444 In conclusion, we provide an instance-dependent characterization of the excess risk for ridge regres-
 445 sion under general covariate shift. Our findings demonstrate that “benign overfitting” also happens
 446 in OOD generalization when the shift in the minor directions is well controlled. We also explore the
 447 “large shift in the minor directions” regime, under which ridge regression may incur a large excess
 448 risk, whereas principal component regression (PCR) exhibits superior performance.

449 Our work opens up several future research directions. First, while we have established a lower bound
 450 for ridge regression in certain instances, a key challenge remains in deriving a general lower bound
 451 that matches our upper bounds, offering a precise characterization of the excess risk under covari-
 452 ate shift. Second, our analysis has been focused on linear models as a first step in understanding
 453 overparameterized OOD problems. Extending this investigation to more complex, nonlinear models
 454 would be an interesting direction for future exploration.

455 References

- 456 Sergios Agapiou, Omiros Papaspiliopoulos, Daniel Sanz-Alonso, and Andrew M Stuart. Importance
 457 sampling: Intrinsic dimension and computational cost. *Statistical Science*, pp. 405–431, 2017.
- 458 Anish Agarwal, Devavrat Shah, Dennis Shen, and Dogyoon Song. On robustness of principal com-
 459 ponent regression. *Advances in Neural Information Processing Systems*, 32, 2019.
- 460 Anish Agarwal, Devavrat Shah, and Dennis Shen. On model identification and out-of-sample pre-
 461 diction of principal component regression: Applications to synthetic controls. *arXiv preprint*
 462 *arXiv:2010.14449*, 2020.

- 463 Eric Bair, Trevor Hastie, Debashis Paul, and Robert Tibshirani. Prediction by supervised principal
464 components. *Journal of the American Statistical Association*, 101(473):119–137, 2006.
- 465 Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear
466 regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- 467 Peter L Bartlett, Andrea Montanari, and Alexander Rakhlin. Deep learning: a statistical viewpoint.
468 *Acta numerica*, 30:87–201, 2021.
- 469 Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps for dimensionality reduction and data
470 representation. *Neural computation*, 15(6):1373–1396, 2003.
- 471 Mikhail Belkin, Daniel Hsu, and Ji Xu. Two models of double descent for weak features. *SIAM*
472 *Journal on Mathematics of Data Science*, 2(4):1167–1180, 2020.
- 473 Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of representations
474 for domain adaptation. *Advances in neural information processing systems*, 19, 2006.
- 475 Koby Bibas, Yaniv Fogel, and Meir Feder. A new look at an old problem: A universal learning
476 approach to linear regression. In *2019 IEEE International Symposium on Information Theory*
477 *(ISIT)*, pp. 2304–2308. IEEE, 2019.
- 478 Swee Lean Chan and Moonseo Park. Project cost estimation using principal component regression.
479 *Construction Management and Economics*, 23(3):295–304, 2005.
- 480 Yihang Chen, Fanghui Liu, Taiji Suzuki, and Volkan Cevher. High-dimensional kernel methods
481 under covariate shift: Data-dependent implicit regularization. In *Forty-first International Con-*
482 *ference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net,
483 2024. URL <https://openreview.net/forum?id=bBzlapzeR1>.
- 484 Yuxin Chen, Yuejie Chi, Jianqing Fan, Cong Ma, et al. Spectral methods for data science: A
485 statistical perspective. *Foundations and Trends® in Machine Learning*, 14(5):566–806, 2021.
- 486 Geoffrey Chinot and Matthieu Lerasle. On the robustness of the minimim l2 interpolator. *Bernoulli*,
487 2022.
- 488 Corinna Cortes, Yishay Mansour, and Mehryar Mohri. Learning bounds for importance weighting.
489 *Advances in neural information processing systems*, 23, 2010.
- 490 Uwe Depczynski, VJ Frost, and K Molt. Genetic algorithms applied to the selection of factors in
491 principal component regression. *Analytica Chimica Acta*, 420(2):217–227, 2000.
- 492 Edgar Dobriban and Stefan Wager. High-dimensional asymptotics of prediction: Ridge regression
493 and classification. *The Annals of Statistics*, 46(1):247–279, 2018.
- 494 Jiawei Ge, Shange Tang, Jianqing Fan, and Chi Jin. On the provable advantage of unsupervised
495 pretraining. *arXiv preprint arXiv:2303.01566*, 2023.
- 496 Jiawei Ge, Shange Tang, Jianqing Fan, Cong Ma, and Chi Jin. Maximum likelihood estimation is all
497 you need for well-specified covariate shift. In *The Twelfth International Conference on Learning*
498 *Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL
499 <https://openreview.net/forum?id=eOTCKK0gIs>.
- 500 Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. When do neural
501 networks outperform kernel methods? *Advances in Neural Information Processing Systems*, 33:
502 14820–14830, 2020.
- 503 Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Linearized two-layers
504 neural networks in high dimension. *The Annals of Statistics*, 49(2), 2021.
- 505 Hao Guan and Mingxia Liu. Domain adaptation for medical image analysis: a survey. *IEEE Trans-*
506 *actions on Biomedical Engineering*, 69(3):1173–1185, 2021.
- 507 Ali S Hadi and Robert F Ling. Some cautionary notes on the use of principal components regression.
508 *The American Statistician*, 52(1):15–19, 1998.

- 509 Steve Hanneke and Samory Kpotufe. On the value of target data in transfer learning. *Advances in*
510 *Neural Information Processing Systems*, 32, 2019.
- 511 Yifan Hao, Yong Lin, Difan Zou, and Tong Zhang. On the benefits of over-parameterization for out-
512 of-distribution generalization. *CoRR*, abs/2403.17592, 2024. doi: 10.48550/ARXIV.2403.17592.
513 URL <https://doi.org/10.48550/arXiv.2403.17592>.
- 514 Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-
515 dimensional ridgeless least squares interpolation. *Annals of statistics*, 50(2):949, 2022.
- 516 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common cor-
517 ruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019.
- 518 Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul
519 Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A criti-
520 cal analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF international*
521 *conference on computer vision*, pp. 8340–8349, 2021.
- 522 Hugo G Hidalgo, Thomas C Piechota, and John A Dracup. Alternative principal components re-
523 gression procedures for dendrohydrologic reconstructions. *Water Resources Research*, 36(11):
524 3241–3249, 2000.
- 525 Shih-Ming Huang and Jar-Ferr Yang. Improved principal component regression for face recognition
526 under illumination variations. *IEEE signal processing letters*, 19(4):179–182, 2012.
- 527 JNR Jeffers. Investigation of alternative regressions: Some practical examples. *Journal of the Royal*
528 *Statistical Society. Series D (The Statistician)*, 30(2):79–88, 1981.
- 529 John NR Jeffers. Two case studies in the application of principal component analysis. *Journal of*
530 *the Royal Statistical Society: Series C (Applied Statistics)*, 16(3):225–236, 1967.
- 531 Ian T Jolliffe. A note on the use of principal components in regression. *Journal of the Royal*
532 *Statistical Society Series C: Applied Statistics*, 31(3):300–303, 1982.
- 533 Chinmaya Kausik, Kashvi Srivastava, and Rishi Sonthalia. Double descent and overfitting under
534 noisy inputs and distribution shift for linear denoisers. *Trans. Mach. Learn. Res.*, 2024, 2024.
535 URL <https://openreview.net/forum?id=HxfqTdLIRF>.
- 536 Richard B Keithley, R Mark Wightman, and Michael L Heien. Multivariate concentration determi-
537 nation using principal component regression with residual analysis. *TrAC Trends in Analytical*
538 *Chemistry*, 28(9):1127–1136, 2009.
- 539 Dmitry Kobak, Jonathan Lomond, and Benoit Sanchez. The optimal ridge penalty for real-world
540 high-dimensional data can be zero or negative due to the implicit ridge regularization. *Journal of*
541 *Machine Learning Research*, 21(169):1–16, 2020.
- 542 Frederic Koehler, Lijia Zhou, Danica J Sutherland, and Nathan Srebro. Uniform convergence of
543 interpolators: Gaussian width, norm bounds and benign overfitting. *Advances in Neural Informa-*
544 *tion Processing Systems*, 34:20657–20668, 2021.
- 545 Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Bal-
546 subramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds: A
547 benchmark of in-the-wild distribution shifts. In *International conference on machine learning*,
548 pp. 5637–5664. PMLR, 2021.
- 549 Samory Kpotufe and Guillaume Martinet. Marginal singularity, and the benefits of labels in
550 covariate-shift. In *Conference On Learning Theory*, pp. 1882–1886. PMLR, 2018.
- 551 Anikender Kumar and Pramila Goyal. Forecasting of air quality in delhi using principal component
552 regression technique. *Atmospheric Pollution Research*, 2(4):436–444, 2011.
- 553 Qi Lei, Wei Hu, and Jason Lee. Near-optimal linear regression under distribution shift. In *Interna-*
554 *tional Conference on Machine Learning*, pp. 6164–6174. PMLR, 2021.

- 555 Tengyuan Liang, Alexander Rakhlin, and Xiyu Zhai. On the multiple descent of minimum-norm
556 interpolants and restricted lower isometry of kernels. In *Conference on Learning Theory*, pp.
557 2683–2711. PMLR, 2020.
- 558 RX Liu, J Kuang, Qiong Gong, and XL Hou. Principal component regression analysis with spss.
559 *Computer methods and programs in biomedicine*, 71(2):141–147, 2003.
- 560 Cong Ma, Reese Pathak, and Martin J Wainwright. Optimally tackling covariate shift in rkhs-based
561 nonparametric regression. *The Annals of Statistics*, 51(2):738–761, 2023.
- 562 Neil Mallinar, Austin Zane, Spencer Frei, and Bin Yu. Minimum-norm interpolation under co-
563 variate shift. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna,*
564 *Austria, July 21-27, 2024*. OpenReview.net, 2024. URL [https://openreview.net/forum?](https://openreview.net/forum?id=Zw7TcnTmHj)
565 [id=Zw7TcnTmHj](https://openreview.net/forum?id=Zw7TcnTmHj).
- 566 William F Massy. Principal components regression in exploratory statistical research. *Journal of*
567 *the American Statistical Association*, 60(309):234–256, 1965.
- 568 Song Mei and Andrea Montanari. The generalization error of random features regression: Precise
569 asymptotics and the double descent curve. *Communications on Pure and Applied Mathematics*,
570 75(4):667–766, 2022.
- 571 Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Generalization error of random feature
572 and kernel methods: hypercontractivity and kernel matrix concentration. *Applied and Computa-*
573 *tional Harmonic Analysis*, 59:3–84, 2022.
- 574 John P Miller, Rohan Taori, Aditi Raghunathan, Shiori Sagawa, Pang Wei Koh, Vaishaal Shankar,
575 Percy Liang, Yair Carmon, and Ludwig Schmidt. Accuracy on the line: on the strong correlation
576 between out-of-distribution and in-distribution generalization. In *International conference on*
577 *machine learning*, pp. 7721–7735. PMLR, 2021.
- 578 Andrea Montanari and Yiqiao Zhong. The interpolation phase transition in neural networks: Memo-
579 rization and generalization under lazy training. *The Annals of Statistics*, 50(5):2816–2847, 2022.
- 580 Mohammadreza Mousavi Kalan, Zalan Fabian, Salman Avestimehr, and Mahdi Soltanolkotabi. Min-
581 imax lower bounds for transfer learning with linear and one-hidden layer neural networks. *Ad-*
582 *vances in Neural Information Processing Systems*, 33:1959–1969, 2020.
- 583 Vidya Muthukumar, Kailas Vodrahalli, Vignesh Subramanian, and Anant Sahai. Harmless interpo-
584 lation of noisy data in regression. *IEEE Journal on Selected Areas in Information Theory*, 1(1):
585 67–83, 2020.
- 586 Tormod Næs and Harald Martens. Principal component regression in nir analysis: viewpoints,
587 background details and selection of components. *Journal of chemometrics*, 2(2):155–167, 1988.
- 588 Preetum Nakkiran. More data can hurt for linear regression: Sample-wise double descent. *arXiv*
589 *preprint arXiv:1912.07242*, 2019.
- 590 Jeffrey Negrea, Gintare Karolina Dziugaite, and Daniel Roy. In defense of uniform convergence:
591 Generalization via derandomization with an application to interpolating predictors. In *Internat-*
592 *ional Conference on Machine Learning*, pp. 7263–7272. PMLR, 2020.
- 593 Partha Niyogi. Manifold regularization and semi-supervised learning: Some theoretical analyses.
594 *Journal of Machine Learning Research*, 14(5), 2013.
- 595 Reese Pathak, Cong Ma, and Martin Wainwright. A new similarity measure for covariate shift with
596 applications to nonparametric regression. In *International Conference on Machine Learning*, pp.
597 17517–17530. PMLR, 2022.
- 598 Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers
599 generalize to imagenet? In *International conference on machine learning*, pp. 5389–5400. PMLR,
600 2019.

- 601 Dominic Richards, Jaouad Mourtada, and Lorenzo Rosasco. Asymptotics of ridge (less) regres-
602 sion under general source condition. In *International Conference on Artificial Intelligence and*
603 *Statistics*, pp. 3889–3897. PMLR, 2021.
- 604 Ohad Shamir. The implicit bias of benign overfitting. *Journal of Machine Learning Research*, 24
605 (113):1–40, 2023.
- 606 Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-
607 likelihood function. *Journal of statistical planning and inference*, 90(2):227–244, 2000.
- 608 James B Simon, Dhruva Karkada, Nikhil Ghosh, and Mikhail Belkin. More is better in modern
609 machine learning: when infinite overparameterization is optimal and overfitting is obligatory.
610 *arXiv preprint arXiv:2311.14646*, 2023.
- 611 Jianguo Sun. A correlation principal component regression analysis of nir data. *Journal of Chemo-*
612 *metrics*, 9(1):21–29, 1995.
- 613 Nilesh Tripuraneni, Ben Adlam, and Jeffrey Pennington. Overparameterization improves
614 robustness to covariate shift in high dimensions. In Marc’Aurelio Ranzato, Alina
615 Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.),
616 *Advances in Neural Information Processing Systems 34: Annual Conference on Neu-*
617 *ral Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*,
618 pp. 13883–13897, 2021a. URL [https://proceedings.neurips.cc/paper/2021/hash/
619 73fed7fd472e502d8908794430511f4d-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/73fed7fd472e502d8908794430511f4d-Abstract.html).
- 620 Nilesh Tripuraneni, Chi Jin, and Michael Jordan. Provable meta-learning of linear representations.
621 In *International Conference on Machine Learning*, pp. 10434–10443. PMLR, 2021b.
- 622 Alexander Tsigler and Peter L. Bartlett. Benign overfitting in ridge regression. *J. Mach. Learn. Res.*,
623 24:123:1–123:76, 2023. URL <http://jmlr.org/papers/v24/22-1398.html>.
- 624 J Leo van Hemmen and Tsuneya Ando. An inequality for trace ideals. *Communications in Mathe-*
625 *matical Physics*, 76:143–148, 1980.
- 626 Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *CoRR*,
627 abs/1011.3027, 2010. URL <http://arxiv.org/abs/1011.3027>.
- 628 Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*,
629 volume 47. Cambridge university press, 2018.
- 630 Evelyne Vigneau, MF Devaux, EM Qannari, and P Robert. Principal component regression, ridge re-
631 gression and ridge principal component regression in spectroscopy calibration. *Journal of Chemo-*
632 *metrics: A Journal of the Chemometrics Society*, 11(3):239–249, 1997.
- 633 Per-Ake Wedin. Perturbation theory for pseudo-inverses. *BIT Numerical Mathematics*, 13:217–232,
634 1973.
- 635 Florian Wenzel, Andrea Dittadi, Peter Gehler, Carl-Johann Simon-Gabriel, Max Horn, Dominik
636 Zietlow, David Kernert, Chris Russell, Thomas Brox, Bernt Schiele, et al. Assaying out-of-
637 distribution generalization in transfer learning. *Advances in Neural Information Processing Sys-*
638 *tems*, 35:7181–7198, 2022.
- 639 Denny Wu and Ji Xu. On the optimal weighted l2 regularization in overparameterized linear regres-
640 sion. *Advances in Neural Information Processing Systems*, 33:10112–10123, 2020.
- 641 Ji Xu and Daniel J Hsu. On the number of variables to use in principal component regression.
642 *Advances in neural information processing systems*, 32, 2019.
- 643 Runtian Zhai, Chen Dan, Zico Kolter, and Pradeep Ravikumar. Understanding why generalized
644 reweighting does not improve over erm. *arXiv preprint arXiv:2201.12293*, 2022.
- 645 Xuhui Zhang, Jose Blanchet, Soumyadip Ghosh, and Mark S Squillante. A class of geometric
646 structures in transfer learning: Minimax bounds and optimality. In *International Conference on*
647 *Artificial Intelligence and Statistics*, pp. 3794–3820. PMLR, 2022.

- 648 Lijia Zhou, Danica J Sutherland, and Nati Srebro. On uniform convergence and low-norm interpo-
649 lation learning. *Advances in Neural Information Processing Systems*, 33:6867–6877, 2020.
- 650 Yang Zou, Zhiding Yu, BVK Kumar, and Jinsong Wang. Unsupervised domain adaptation for se-
651 mantic segmentation via class-balanced self-training. In *Proceedings of the European conference*
652 *on computer vision (ECCV)*, pp. 289–305, 2018.

653 **A Ridge regression**

654 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
 655 the first k columns of X as X_k and the remaining $d - k$ columns as X_{-k} . Similarly, β_k^* and β_{-k}^*
 656 represent the corresponding components of β^* . $\Sigma_{S,k}$, $\Sigma_{S,-k}$ are the corresponding blocks on the
 657 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
 658 the rows of Z are i.i.d. centered isotropic random vectors. Additionally, we assume the rows of Z
 659 are σ -sub-gaussian, where the sub-gaussian norm is defined as follows.

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

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

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

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

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

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

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

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

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

667 **A.1 Concentration inequalities**

668 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
 669 the matrix X as $X[i, *]$.

670 **Lemma 7** (Lemma 20 of Tsigler & Bartlett (2023)). Let z be a sub-gaussian vector in \mathbb{R}^p with
 671 $\|z\|_{\psi_2} \leq \sigma$, and consider $\Sigma = \text{diag}(\lambda_1, \dots, \lambda_p)$ where the sequence $\{\lambda_j\}_{j=1}^p$ is positive and non-
 672 increasing. Then there exists some absolute constant c , for any $t > 0$, with probability at least
 673 $1 - 2e^{-t/c}$:

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

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

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

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

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

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

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

681 **Lemma 10.** There exists a constant c_x , depending only on σ , such that for any n satisfying $n\lambda_{k+1} \leq$
682 $(\lambda + \sum_{j>k} \lambda_j)$, under the assumption $\text{CondNum}(k, \delta, L)$ (Assumption 2), with probability at least
683 $1 - \delta - c_x e^{-n/c_x}$:

$$\frac{1}{c_x L} \left(\lambda + \sum_{j>k} \lambda_j \right) \leq \mu_n(A_k) \leq \mu_1(A_k) \leq c_x \left(\lambda + \sum_{j>k} \lambda_j \right).$$

$$\mu_1(X_{-k} X_{-k}^T) \leq c_x \left(n\lambda_{k+1} + \sum_{j>k} \lambda_j \right).$$

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

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

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

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

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

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

690 Since $\mu_1(A_k) \leq \lambda + \|\mathring{A}_k\| + \max_i \|X_{-k}[i, *]\|^2$, by setting $t = n$, we have with probability at least
691 $1 - (2n+4)e^{-n/c}$:

$$\begin{aligned} \mu_1(A_k) &\leq \lambda + c\sigma^2 \left(n\lambda_{k+1} + \sum_{j>k} \lambda_j + \sqrt{(2n\lambda_{k+1})^2 + 2n \sum_{j>k} \lambda_j^2} \right) \\ &\leq \lambda + c\sigma^2 \left(n\lambda_{k+1} + \sum_{j>k} \lambda_j + 2n\lambda_{k+1} + \sqrt{2n \sum_{j>k} \lambda_j^2} \right) \\ &\leq \lambda + c\sigma^2 \left(n\lambda_{k+1} + \sum_{j>k} \lambda_j + 2n\lambda_{k+1} + \sqrt{2n\lambda_{k+1} \sum_{j>k} \lambda_j} \right) \\ &\leq \lambda + c\sigma^2 \left(n\lambda_{k+1} + \sum_{j>k} \lambda_j + 2n\lambda_{k+1} + n\lambda_{k+1} + \frac{1}{2} \sum_{j>k} \lambda_j \right) \\ &\leq \lambda + 4c\sigma^2 \left(n\lambda_{k+1} + \sum_{j>k} \lambda_j \right) \end{aligned}$$

$$\begin{aligned}
&\leq \max\{1, 4c\sigma^2\} \left(\lambda + \sum_{j>k} \lambda_j + n\lambda_{k+1} \right) \\
&\leq 2 \max\{1, 4c\sigma^2\} \left(\lambda + \sum_{j>k} \lambda_j \right). \tag{2}
\end{aligned}$$

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

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

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

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

695 On this event,

$$\begin{aligned}
\mu_1(A_k) &\geq \lambda + \frac{1}{n} \text{tr}(X_{-k}X_{-k}^T) \\
&= \lambda + \frac{1}{n} \sum_{i=1}^n \|X_{-k}[i, *]\|^2 \\
&\geq \lambda + \frac{1}{2} \sum_{j>k} \lambda_j \\
&\geq \frac{1}{2} \left(\lambda + \sum_{j>k} \lambda_j \right).
\end{aligned}$$

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

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

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

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

699

□

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

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

702 *Proof.* This follows directly from Theorem 5.39 and Remark 5.40 of [Vershynin \(2010\)](#), which shows
 703 there exists a constant c'_x depending only on σ , such that for any $t \geq 0$, with probability at least
 704 $1 - 2 \exp\{-t^2/c'_x\}$:

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

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

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

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

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

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

710 Therefore,

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

711 By applying Lemma 11, the proof is complete. \square

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

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

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

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

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

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

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

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

$$\begin{aligned} \mu_1(X_k^T X_k) &\leq \|\Sigma_{S,k}\| \left\| \Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right\| \leq c_x \lambda_1 n. \\ \mu_k(X_k^T X_k) &\geq \mu_k(\Sigma_{S,k}) \mu_k \left(\Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right) \geq \frac{1}{c_x} \lambda_k n. \end{aligned}$$

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

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

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

$$\begin{aligned} \left\| Z_{-k}[i, *] \Sigma_{S, -k}^{\frac{1}{2}} \Sigma_{T, -k} \Sigma_{S, -k}^{\frac{1}{2}} Z_{-k}[i, *]^T \right\|_{\psi_1} &\leq c\sigma^2 \left\| \Sigma_{S, -k}^{\frac{1}{2}} \Sigma_{T, -k} \Sigma_{S, -k}^{\frac{1}{2}} \right\|_{\text{F}} \\ &\leq c\sigma^2 \text{tr} \left(\Sigma_{S, -k}^{\frac{1}{2}} \Sigma_{T, -k} \Sigma_{S, -k}^{\frac{1}{2}} \right). \end{aligned}$$

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

$$\begin{aligned} &\mathbb{P} \left\{ \frac{1}{n} \left| \sum_{i=1}^n \left[Z_{-k}[i, *] \Sigma_{S, -k}^{\frac{1}{2}} \Sigma_{T, -k} \Sigma_{S, -k}^{\frac{1}{2}} Z_{-k}[i, *]^T - \text{tr} \left(\Sigma_{S, -k}^{\frac{1}{2}} \Sigma_{T, -k} \Sigma_{S, -k}^{\frac{1}{2}} \right) \right] \right| \geq t \right\} \\ &\leq 2 \exp \left\{ -c'n \min \left\{ \frac{t^2}{K^2}, \frac{t}{K} \right\} \right\}, \end{aligned}$$

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

725 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 $1 - 2e^{-c'n}$:

$$\begin{aligned} \text{tr} \left(X_{-k} \Sigma_{T, -k} X_{-k}^T \right) &= \sum_{i=1}^n Z_{-k}[i, *] \Sigma_{S, -k}^{\frac{1}{2}} \Sigma_{T, -k} \Sigma_{S, -k}^{\frac{1}{2}} Z_{-k}[i, *]^T \\ &\leq (1 + c\sigma^2) n \text{tr} \left(\Sigma_{S, -k}^{\frac{1}{2}} \Sigma_{T, -k} \Sigma_{S, -k}^{\frac{1}{2}} \right). \end{aligned}$$

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

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

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

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

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

731 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
 732 an absolute constant c , for any $t \in (0, n)$, with probability at least $1 - 2e^{-ct}$:

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

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

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

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

735 **A.2 Block decomposition of $X_{-k}X_{-k}^T$**

736 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
737 and right singular vectors, respectively. The matrix $\widetilde{M}^{\frac{1}{2}}$ is defined as:

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

738 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,
739 $X_k^T X_k = V^T M_k V$, where $M_k = \text{diag}(m_1, \dots, m_k) \in \mathbb{R}^{k \times k}$.

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

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

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

742 We will repeatedly use the first k rows of $(M + \lambda I_n + \Delta)^{-1}$, which we compute here. Because
743 $M + \lambda I_n + \Delta$ and $\lambda I_{n-k} + \Delta_{22}$ are invertible when A_k is positive definite, by block matrix inverse,

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

744 **Corollary 16** (Corollary of Lemma 10). There exists a constant depending only on σ , such that
745 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
746 inequalities hold with probability at least $1 - \delta - c_x e^{-n/c_x}$, on the same event as in Lemma 10.

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

747 *Proof.* 1. The first inequality.

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

748 2. The second inequality.

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

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

750

3. The third inequality.

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

751

4. The fourth inequality.

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

752

5. The last inequality.

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

753

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

754

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

755

756

□

757

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,

758

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

760 where

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

Proof.

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

$$\cdot \left\| X_k^T X_k - n\Sigma_{S,k} + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V \right\|.$$

761 According to Lemma 11, Corollary 16, there exists a constant $c_x > 2$ depending only on σ , such
 762 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
 763 $1 - 2\delta - c_x e^{-n/c_x}$, on both events in Lemma 10 and Lemma 11,

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

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

764 1. $\left\| X_k^T X_k - n\Sigma_{S,k} + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V \right\|$

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

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

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

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

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

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

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

768 Therefore, we have

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

769 Now we derive the upper bound for our target.

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

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

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

773 Combining the above two inequalities,

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

774

□

775 A.3 Bias variance decomposition

776 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
777 distribution of the noise ϵ .

$$\begin{aligned}
\mathbb{E}_\epsilon \left[\mathcal{R}(\hat{\beta}(Y)) \right] &= \mathbb{E}_\epsilon \left[\left(\hat{\beta}(Y) - \beta^* \right)^T \Sigma_T \left(\hat{\beta}(Y) - \beta^* \right) \right] \\
&= \mathbb{E}_\epsilon \left[\hat{\beta}(\epsilon)^T \Sigma_T \hat{\beta}(\epsilon) \right] + \left(\hat{\beta}(X\beta^*) - \beta^* \right)^T \Sigma_T \left(\hat{\beta}(X\beta^*) - \beta^* \right).
\end{aligned}$$

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

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

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

$$\begin{pmatrix} x_1^T & x_2^T \end{pmatrix} \begin{pmatrix} A & B \\ B^T & D \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = x_1^T A x_1 + 2x_1^T B x_2 + x_2^T D x_2,$$

780 where the positive definiteness implies $x_1^T A x_1 + x_2^T D x_2 \geq 2x_1^T B x_2$.

781 **Lemma 18.** There exists a constant $c_x > 2$ depending only on σ , such that for any $N_1 < n <$
782 N_2 , if the assumption $\text{condNum}(k, \delta, L)$ (Assumption 2) is satisfied, then with probability at least
783 $1 - 2\delta - c_x e^{-n/c_x}$, the following inequalities hold simultaneously:

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

$$\begin{aligned}
\mu_1(A_k) &\leq c_x \left(\lambda + \sum_{j>k} \lambda_j \right). \\
\mu_1(X_{-k} X_{-k}^T) &\leq c_x \left(n\lambda_{k+1} + \sum_{j>k} \lambda_j \right). \\
\left\| \frac{1}{n} X_k^T X_k - \Sigma_{S,k} \right\| &\leq c_x \lambda_1 \sqrt{\frac{k + \ln \frac{1}{\delta}}{n}}. \\
\left\| (X_k^T X_k)^{\frac{1}{2}} - \sqrt{n} \Sigma_{S,k}^{\frac{1}{2}} \right\| &\leq c_x \sqrt{k + \ln \frac{1}{\delta}} \lambda_1 \lambda_k^{-\frac{1}{2}}. \\
\mu_k \left(\Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right) &\geq \frac{1}{c_x} n. \\
\mu_1 \left(\Sigma_{S,k}^{-\frac{1}{2}} X_k^T X_k \Sigma_{S,k}^{-\frac{1}{2}} \right) &\leq c_x n. \\
\mu_1(X_k^T X_k) &\leq c_x \lambda_1 n. \\
\mu_k(X_k^T X_k) &\geq \frac{1}{c_x} \lambda_k n. \\
\text{tr} \left(X_{-k} \Sigma_{T,-k} X_{-k}^T \right) &\leq c_x n \text{tr} \left(\Sigma_{S,-k}^{\frac{1}{2}} \Sigma_{T,-k} \Sigma_{S,-k}^{\frac{1}{2}} \right). \\
(\beta_{-k}^*)^T X_{-k}^T X_{-k} \beta_{-k}^* &\leq c_x n (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^*. \\
\|\Delta_{11}\|, \|\Delta_{12}\|, \|\Delta\| &\leq c_x \left(\lambda + \sum_{j>k} \lambda_j \right). \\
\|(\lambda I_{n-k} + \Delta_{22})^{-1}\|, \|\Delta^{-1}\| &\leq c_x L \left(\lambda + \sum_{j>k} \lambda_j \right)^{-1}. \\
\|\Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-2} \Delta_{12}^T\| &\leq c_x^4 L^2. \\
\|\Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T\| &\leq c_x^3 L \left(\lambda + \sum_{j>k} \lambda_j \right). \\
\|\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T\| &\leq c_x \left(\lambda + \sum_{j>k} \lambda_j \right).
\end{aligned}$$

784 And,

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

785 N_1 and N_2 are defined as follows:

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

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

788 **A.3.1 Variance in the first k dimensions**

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

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

791 where

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

Proof.

$$\begin{aligned} & \mathbb{E}_\epsilon \left[\widehat{\beta}(\epsilon)_k^T \Sigma_{T,k} \widehat{\beta}(\epsilon)_k \right] \\ &= \mathbb{E}_\epsilon \operatorname{tr} \left[\epsilon \epsilon^T (X X^T + \lambda I_n)^{-1} X_k \Sigma_{T,k} X_k^T (X X^T + \lambda I_n)^{-1} \right] \\ &= v^2 \operatorname{tr} \left[(X X^T + \lambda I_n)^{-1} X_k \Sigma_{T,k} X_k^T (X X^T + \lambda I_n)^{-1} \right] \\ &= v^2 \operatorname{tr} \left[(U M U^T + U \Delta U^T + \lambda I_n)^{-1} U \widetilde{M}^{\frac{1}{2}} V \Sigma_{T,k} \right. \\ & \quad \left. \cdot V^T \left(\widetilde{M}^{\frac{1}{2}} \right)^T U^T (U M U^T + U \Delta U^T + \lambda I_n)^{-1} \right] \\ &= v^2 \operatorname{tr} \left[U (M + \Delta + \lambda I_n)^{-1} \widetilde{M}^{\frac{1}{2}} V \Sigma_{T,k} V^T \left(\widetilde{M}^{\frac{1}{2}} \right)^T (M + \Delta + \lambda I_n)^{-1} U^T \right] \\ &= v^2 \operatorname{tr} \left[\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 \right] \\ &= v^2 \operatorname{tr} \left[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}) \right. \\ & \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 \left. \cdot V \Sigma_{T,k} V^T \right] \\ &= v^2 \operatorname{tr} \left[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. \\ & \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 \left. \cdot V \Sigma_{T,k} V^T \right] \\ &= v^2 \operatorname{tr} \left[(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} \right. \\ & \quad \left. \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} \right] \\ &\leq v^2 \left\| I_k + \Delta_{12} (\lambda I_{n-k} + \Delta_{22})^{-2} \Delta_{12}^T \right\| \operatorname{tr} \left[(M_k + \lambda I_k + \Delta_{11} - \Delta_{12} (\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} \right. \\ & \quad \left. \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} \right] \end{aligned}$$

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

792 The sixth equation follows from Equation 5. The first inequality follows from the result $\operatorname{tr}[AB] \leq$
793 $\|A\| \operatorname{tr}[B]$ where the matrix B is positive semi-definite.

794 We define two quantities that represent concentration error terms:

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

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

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

$$\begin{aligned}
&\|E_1\| \left\| n\tilde{\Sigma}_{S,k} \right\| \left\| (n\tilde{\Sigma}_{S,k})^{-1} \right\| \left\| (n\Sigma_{S,k})^{\frac{1}{2}} \right\| \|\Sigma_{T,k}\| \left\| (n\Sigma_{S,k})^{\frac{1}{2}} \right\| \left\| (n\tilde{\Sigma}_{S,k})^{-1} \right\| \\
&\leq \frac{c_x^2 \left(\sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 + c_x^2 L \left(\lambda + \sum_{j>k} \lambda_j \right) \right)}{(\lambda + n\lambda_k)^2} (\lambda + n\lambda_1) \frac{n\lambda_1}{(\lambda + n\lambda_k)^2} \|\Sigma_{T,k}\| \\
&\leq \frac{c_x^2 \left(\sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 + c_x^2 L \left(\lambda + \sum_{j>k} \lambda_j \right) \right)}{n^2} \frac{\lambda_1^2}{\lambda_k^4} \|\Sigma_{T,k}\| \\
&= \frac{c_x^2 \sqrt{(k + \ln \frac{1}{\delta})} \lambda_1^3}{n\sqrt{n}} \frac{\lambda_1^3}{\lambda_k^4} \|\Sigma_{T,k}\| + \frac{c_x^4 L \left(\lambda + \sum_{j>k} \lambda_j \right)}{n^2} \frac{\lambda_1^2}{\lambda_k^4} \|\Sigma_{T,k}\| \\
&< \frac{1}{2nk} \operatorname{tr} \left[\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right] + \frac{1}{2nk} \operatorname{tr} \left[\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right] \\
&= \frac{1}{nk} \operatorname{tr} \left[\Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right].
\end{aligned}$$

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

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

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

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

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

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

800 Combing the above four inequalities, we have

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

801 In particular,

$$\text{tr} \left[\left(n\tilde{\Sigma}_{S,k} \right)^{-1} \left(n\Sigma_{S,k} \right)^{\frac{1}{2}} \Sigma_{T,k} \left(n\Sigma_{S,k} \right)^{\frac{1}{2}} \left(n\tilde{\Sigma}_{S,k} \right)^{-1} \right]$$

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

802 The inequality follows from the fact that $\operatorname{tr}[BAB] = \operatorname{tr}[A^{\frac{1}{2}}BA^{\frac{1}{2}}] \leq \operatorname{tr}[A^{\frac{1}{2}}CA^{\frac{1}{2}}] = \operatorname{tr}[CAC]$,
803 where A, B, C are positive semi-definite matrices, and $C \succcurlyeq B$, which implies that $A^{\frac{1}{2}}CA^{\frac{1}{2}} \succcurlyeq$
804 $A^{\frac{1}{2}}BA^{\frac{1}{2}}$.

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

805 The other terms can be similarly bounded. Therefore,

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

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

807 A.3.2 Variance in the last $d - k$ dimensions

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

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

810 where N_1, N_2 are defined as in Lemma 18.

Proof.

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

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

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

813 A.3.3 Bias in the first k dimensions

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

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

815 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-
816 definite.

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

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

819 where

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

Proof.

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

820 Subsequently,

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

$$\begin{aligned}
&= \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}$$

821 The second equation follows from Equation 5.

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

823 1. The first term.

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

824 The inequality follows from $c_x > 2$.

825 2. The second term.

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

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

$$< \frac{1}{c_x}.$$

827

Therefore,

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

828

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

829

830

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

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

831 Combining both terms above, we have

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

832 Therefore,

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

$$\leq \frac{16c_x^4}{n^2} \left(\lambda + \sum_{j>k} \lambda_j \right)^2 (\beta_{-k}^*)^T \Sigma_{S,k}^{-1} \beta_{-k}^* \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|.$$

833

□

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

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

836 where

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

Proof.

$$\begin{aligned} & (X_k^T (XX^T + \lambda I_n)^{-1} X_{-k} \beta_{-k}^*)^T \Sigma_{T,k} (X_k^T (XX^T + \lambda I_n)^{-1} X_{-k} \beta_{-k}^*) \\ & \leq \left\| (XX^T + \lambda I_n)^{-1} X_k \Sigma_{T,k} X_k^T (XX^T + \lambda I_n)^{-1} \right\| \cdot (\beta_{-k}^*)^T X_{-k}^T X_{-k} \beta_{-k}^*. \end{aligned}$$

837 From Lemma 18,

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

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

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

$$\begin{aligned}
& \cdot V^T M_k^{\frac{1}{2}} (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T)^{-1} \Big\| \\
& = (1 + c_x^4 L^2) \Big\| (V^T (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1} \\
& \quad \cdot V^T M_k^{\frac{1}{2}} V \Sigma_{T,k} V^T M_k^{\frac{1}{2}} V \\
& \quad \cdot (V^T (M_k + \lambda I_k + \Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1} \Big\| \\
& = (1 + c_x^4 L^2) \Big\| (X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1} \\
& \quad \cdot (X_k^T X_k)^{\frac{1}{2}} \Sigma_{T,k} (X_k^T X_k)^{\frac{1}{2}} \\
& \quad \cdot (X_k^T X_k + \lambda I_k + V^T (\Delta_{11} - \Delta_{12}(\lambda I_{n-k} + \Delta_{22})^{-1} \Delta_{12}^T) V)^{-1} \Big\|.
\end{aligned}$$

839 The third equation follows from Equation 5.

840 We define two quantities that represent concentration error terms:

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

$$E_2 = (X_k^T X_k)^{\frac{1}{2}} - (n \Sigma_{S,k})^{\frac{1}{2}}.$$

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

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

$$\begin{aligned}
& \|E_1\| \left\| n \tilde{\Sigma}_{S,k} \right\| \left\| (n \tilde{\Sigma}_{S,k})^{-1} \right\| \left\| (n \Sigma_{S,k})^{\frac{1}{2}} \right\| \|\Sigma_{T,k}\| \left\| (n \Sigma_{S,k})^{\frac{1}{2}} \right\| \left\| (n \tilde{\Sigma}_{S,k})^{-1} \right\| \\
& \leq \frac{c_x^2 \left(\sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 + c_x^2 L \left(\lambda + \sum_{j>k} \lambda_j \right) \right)}{(\lambda + n \lambda_k)^2} (\lambda + n \lambda_1) \frac{n \lambda_1}{(\lambda + n \lambda_k)^2} \|\Sigma_{T,k}\| \\
& \leq \frac{c_x^2 \left(\sqrt{n(k + \ln \frac{1}{\delta})} \lambda_1 + c_x^2 L \left(\lambda + \sum_{j>k} \lambda_j \right) \right)}{n^2} \frac{\lambda_1^2}{\lambda_k^4} \|\Sigma_{T,k}\| \\
& = \frac{c_x^2 \sqrt{(k + \ln \frac{1}{\delta})} \lambda_1^3}{n \sqrt{n}} \frac{\lambda_1^3}{\lambda_k^4} \|\Sigma_{T,k}\| + \frac{c_x^4 L \left(\lambda + \sum_{j>k} \lambda_j \right)}{n^2} \frac{\lambda_1^2}{\lambda_k^4} \|\Sigma_{T,k}\| \\
& < \frac{1}{2n} \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| + \frac{1}{2n} \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| \\
& = \frac{1}{n} \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|.
\end{aligned}$$

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

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

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

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

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

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

846 Combining the four inequalities above,

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

847 In particular,

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

848 The inequality follows from the fact that $\|BAB\| = \|A^{\frac{1}{2}}BA^{\frac{1}{2}}\| \leq \|A^{\frac{1}{2}}CA^{\frac{1}{2}}\| = \|CAC\|$, where

849 A, B, C are positive semi-definite matrices, and $C \succcurlyeq B$, which implies that $A^{\frac{1}{2}}CA^{\frac{1}{2}} \succcurlyeq A^{\frac{1}{2}}BA^{\frac{1}{2}}$.

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

$$\leq \frac{1}{n} \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\|.$$

850 The other terms can be similarly bounded. Therefore,

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

851

□

852 A.3.4 Bias in the last $d - k$ dimensions

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

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

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

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

857 where N_1, N_2 are defined as in Lemma 18.

Proof.

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

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

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

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

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

861 where N_1, N_2 are defined as in Lemma 18.

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

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

863 Therefore,

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

864 In particular,

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

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

867 Therefore,

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

869 A.4 Main results

870 **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 $c > 2$
871 depending only on σ , such that for any $cN < n < r_k$, if the assumption $\text{condNum}(k, \delta, L)$ (As-
872 sumption 2) is satisfied, then with probability at least $1 - 2\delta - ce^{-n/c}$,

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

873 N is defined as follows:

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

874 **Remark 8** (Sample complexity). We have assumed $n \geq c_x N$ in the theorem. The first condition
875 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
876 $n = \Omega(k)$ in the best case, consistent with the sample complexity of classic linear regression.
877 This optimal case occurs when $\Sigma_{S,k} \approx \Sigma_{T,k}$. In the worst case, $n = \Omega(k^3)$ where covariate shift is
878 significant in the first k dimensions—e.g., when the test data lies predominantly in the subspace of the
879 first dimension. This shift in sample complexity under varying degrees of covariate shift parallels the
880 analysis of Ge et al. (2024) (see their Theorem 4.2) for the under-parameterized setting. The second
881 condition implies $n \gg \lambda + \sum_{j>k} \lambda_j$, such that the regularization is not too strong to introduce a
882 bias greater than a constant (as shown in the first bias term). On the other hand, we assume $n < r_k$
883 in the theorem, which is consistent with the over-parameterized regime and Assumption 1, where
884 the last $d - k$ components are considered to be essentially high-dimensional.

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

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

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

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

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

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

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

$$\begin{aligned}
& + 3c'_x{}^6 L \frac{1}{n} \left(\lambda + \sum_{j>k} \lambda_j \right) \|\Sigma_{T,-k}\| (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^* \\
& + 3(\beta_{-k}^*)^T \Sigma_{T,-k} \beta_{-k}^* \\
& \leq 16c'_x{}^4 \frac{1}{n^2} \left(\lambda + \sum_{j>k} \lambda_j \right)^2 \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^* \\
& + 64c'_x{}^5 L^2 \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^* \\
& + 3c'_x{}^2 L n \left(\lambda + \sum_{j>k} \lambda_j \right)^{-1} \|\Sigma_{T,-k}\| (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^* \\
& + 3c'_x{}^6 L \frac{1}{n} \left(\lambda + \sum_{j>k} \lambda_j \right) \|\Sigma_{T,-k}\| (\beta_k^*)^T \Sigma_{S,k}^{-1} \beta_k^* \\
& + 3c'_x{}^5 L^2 \left\| \Sigma_{S,k}^{-\frac{1}{2}} \Sigma_{T,k} \Sigma_{S,k}^{-\frac{1}{2}} \right\| (\beta_{-k}^*)^T \Sigma_{S,-k} \beta_{-k}^*.
\end{aligned}$$

891 The last inequality follows from:

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

892 By taking $c = 134c'_x{}^6$, the proof is complete. \square

893 **Corollary 26.** 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 $\mathcal{V} = \Sigma_{S,-k}^2$. There exists a
894 constant $c > 2$ depending only on σ, L , such that for any $cN < n < r_k$, if the assumption
895 $\text{condNum}(k, \delta, L)$ (Assumption 2) is satisfied, then with probability at least $1 - 3\delta$,

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

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

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

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

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

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

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

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

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

901 Since the statement of Theorem 25 holds with probability at least $1 - 2\delta - ce^{-n/c}$, we only require
 902 $ce^{-n/c} < \delta$, which is equivalent as $n > c \ln c + c \ln(1/\delta)$. Combining the lower bounds of n in
 903 Theorem 25, we should have:

$$n > \max \left\{ c \ln c + c \ln \frac{1}{\delta}, \right. \\
 c \left(k + \ln \frac{1}{\delta} \right) \lambda_1^6 \lambda_k^{-8} \|\Sigma_{T,k}\|^2 k^2 (\text{tr}[\mathcal{T}])^{-2}, \\
 \left. cL \lambda_1^2 \lambda_k^{-4} \left(\lambda + \sum_{j>k} \lambda_j \right) \|\Sigma_{T,k}\| k (\text{tr}[\mathcal{T}])^{-1} \right\}.$$

904 For the first term in the maximum argument,

$$c \ln c + c \ln \frac{1}{\delta} \leq c^2 + c \ln \frac{1}{\delta} \\
 \leq c^2 \left(k + \ln \frac{1}{\delta} \right).$$

905 The second term:

$$c \left(k + \ln \frac{1}{\delta} \right) \lambda_1^6 \lambda_k^{-8} \|\Sigma_{T,k}\|^2 k^2 (\text{tr}[\mathcal{T}])^{-2} \\
 \leq c \left(k + \ln \frac{1}{\delta} \right) \lambda_1^6 \lambda_k^{-8} \|\Sigma_{T,k}\|^2 k^2 \left(\mu_k(\Sigma_{S,k}^{-1}) \text{tr}[\Sigma_{T,k}] \right)^{-2} \\
 \leq c \left(k + \ln \frac{1}{\delta} \right) \lambda_1^8 \lambda_k^{-8} \|\Sigma_{T,k}\|^2 k^2 \|\Sigma_{T,k}\|^{-2} \\
 = c \left(k + \ln \frac{1}{\delta} \right)^3 \lambda_1^8 \lambda_k^{-8}.$$

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

908 Similar, for the third term:

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

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

910 □

911 B Large shift in minor directions

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

915 B.1 Lower bound for ridge regression

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

918 Recall our model:

$$y = \beta^{*T} x + \epsilon, \tag{8}$$

919 We can write our data as

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

920 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
921 denote by $\widehat{\Sigma}_S := \frac{1}{n}X^T X$ the sample covariance matrix.

922 Assume the same assumptions as in our previous section still holds. We let $\Sigma_S = \mathbb{E}[x_i x_i^T]$ be the
923 following: its eigenvalues $\lambda_1, \dots, \lambda_d$ satisfies $\lambda_1 = \dots = \lambda_k = 1$, $\lambda_{k+1} = \dots = \lambda_{k+\lfloor \sqrt{n}/C_2 \rfloor} =$
924 C_1/\sqrt{n} for sufficiently large constants C_1, C_2 , and the remaining eigenvalues are all set to zero. We
925 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
926 under this scenario, ridge regression can not obtain an error rate of $\mathcal{O}(\frac{1}{n})$. To see this, we explicitly
927 write out the ridge solution:

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

928 Therefore

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

929 Taking expectation with respect to ϵ ,

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

930 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
931 state the formal version of Theorem 4 in the following:

932 **Theorem 27.** Under the instance we consider, namely $\lambda_1, \dots, \lambda_d$ satisfies $\lambda_1 = \dots = \lambda_k = 1$,
933 $\lambda_{k+1} = \dots = \lambda_{k+\lfloor \sqrt{n}/C_2 \rfloor} = C_1/\sqrt{n}$, $\lambda_{k+\lfloor \sqrt{n}/C_2 \rfloor+1} = \dots = \lambda_d = 0$. WLOG assume $\sigma = 1$,
934 $C_2 \geq C_1((\frac{C_1}{4C})^2 - k - \log \frac{1}{\delta})^{-1}$ for some absolute constant C , and $n \geq (\frac{3C_1}{2})^4$. With probability
935 $1 - \delta$, when $\lambda = c\sqrt{n}$, we have $\frac{V}{v^2} \geq C'$, where $C' > 0$ is some absolute constant. When $\lambda \leq n^{3/4}$,
936 we have $\frac{V}{v^2} \geq C' \frac{1}{\sqrt{n}}$. When $\lambda \geq n^{3/4}$, $B \geq \frac{\|\beta^*\|^2}{9\sqrt{n}}$.

937 *Proof.* We will use the following concentration lemma modified from (Vershynin, 2018, Exercise
938 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,
 $\mathbb{E}[x x^T] = \Sigma$ and is $\sigma^2 \Sigma$ -sub-gaussian, in the sense that

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

939 $X = (x_1, \dots, x_n)^T \in \mathbb{R}^{n \times d}$. Then with probability $1 - \delta$,

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

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

941 Applying Lemma 28, we have

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

942 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

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

943 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.,
 944 2021, Lemma 2.2), $\|\widehat{\lambda}_i - \lambda_i\| \leq \|\widehat{\Sigma}_S - \Sigma_S\|$. Combining with previous inequalities, we have $1 -$
 945 $2C \sqrt{\frac{C_1/C_2 + k + \log \frac{1}{\delta}}{n}} \leq \widehat{\lambda}_i \leq 1 + 2C \sqrt{\frac{C_1/C_2 + k + \log \frac{1}{\delta}}{n}}$ for $1 \leq i \leq k$, $\frac{C_1}{\sqrt{n}} - 2C \sqrt{\frac{C_1/C_2 + k + \log \frac{1}{\delta}}{n}} \leq$
 946 $\widehat{\lambda}_i \leq \frac{C_1}{\sqrt{n}} + 2C \sqrt{\frac{C_1/C_2 + k + \log \frac{1}{\delta}}{n}}$ for $k + 1 \leq i \leq k + \lfloor \sqrt{n}/C_2 \rfloor$. If we take $C_2 \geq C_1((\frac{C_1}{4C})^2 - k -$
 947 $\log \frac{1}{\delta})^{-1}$ then $2C \sqrt{\frac{C_1/C_2 + k + \log \frac{1}{\delta}}{n}} \leq \frac{C_1}{2\sqrt{n}}$. Therefore we have $\frac{C_1}{2\sqrt{n}} \leq \widehat{\lambda}_i \leq \frac{3C_1}{2\sqrt{n}}$ for $k + 1 \leq i \leq$
 948 $k + \lfloor \sqrt{n}/C_2 \rfloor$. When $\lambda = c\sqrt{n}$, we have

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

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

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

$$\begin{aligned}
&= \frac{1}{n} \lfloor \sqrt{n}/C_2 \rfloor \frac{C_1}{2} \left(\frac{3C_1}{2} + n^{1/4} \right)^{-2} \sqrt{n} \\
&\geq \frac{C_1}{16C_2} n^{-1/2},
\end{aligned} \tag{13}$$

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

951 As for the bias term, assume $\lambda \geq n^{3/4}$. Using the same concentration argument, we have $2 > \widehat{\lambda}_i >$
952 $1/2$, for $1 \leq i \leq k$. When $\lambda \leq n$, $\lambda_{\max}(\widehat{\Sigma}_S + \frac{\lambda}{n}I_d) \leq 2 + \lambda/n \leq 3$, therefore $\lambda_{\min}((\widehat{\Sigma}_S +$
953 $\frac{\lambda}{n}I_d)^{-1}) \geq \frac{1}{3}$. This implies

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

954 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
955 implies

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

956

□

957 B.2 Upper bound for PCR

958 In this subsection, we will give the following upper bound for principal component regression.

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

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

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

961 *Proof.* For simplicity, we assume we have a sample size of $2n$, and in the first step we obtain an
962 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 anal-
963 ysis 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 $\widehat{U} = (\widehat{u}_1, \dots, \widehat{u}_k)$
964 where \widehat{u}_i is the i -th eigenvector of $\widehat{\Sigma}_S$. We denote the distance between the estimated subspace and
965 the original one by $\Delta := \text{dist}(U, \widehat{U}) = \|U U^T - \widehat{U} \widehat{U}^T\|$. For controlling Δ , we have the following
966 lemma (Lemma 6):

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

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

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

969 In the second step, we do linear regression on the projected (second half) data. With a little abuse of
 970 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
 971 is independent from the data in step 1, and therefore independent of Δ . If we let $Z := X\hat{U} \in \mathbb{R}^{n \times k}$
 972 be the projected data matrix, the estimator $\hat{\beta}$ we obtained is given by

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

973 We aim to bound the excess risk on target, which is given by $\|\hat{\beta} - \beta^*\|_{\Sigma_T}^2 := \|\Sigma_T^{\frac{1}{2}}(\hat{\beta} - \beta^*)\|^2$. 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:

977 **Lemma 31.** 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
 978 $1 - \delta$,

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

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

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

980 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

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

981 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
 982 Lemma 31 will be both satisfied. We can thus apply Lemma 31 to get

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

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

984 □

985 **B.3 Proofs for Lemma 31**

986 In the following we will prove Lemma 31.

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

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

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

990 where $A_1 := \widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \beta_U^* - \beta_U^*$, $A_2 := \widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \beta_{\perp}^*$, $A_3 :=$
991 $\widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T \epsilon$. Therefore

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

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

993 **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
994 least $1 - \delta$,

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

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

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

996 The second lemma considers the variance term A_3 :

997 **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
998 least $1 - \delta$,

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

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

1000 **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
1001 $1 - \delta$

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

1002 By Lemma 32, 33, 34, together with the decomposition (15), we have with probability $1 - \delta$, when
1003 $n \gtrsim N_1$,

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

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

1004 If further $n \gtrsim \sigma^4 \Delta^{-2} k \log(1/\delta)$,

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

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

1005 \square

1006 **B.4 Technical proofs**

1007 In the sequel, we give the proofs of Lemma 32, 33, 34 and 30. We first prove some additional
 1008 technical lemmas. The following lemma, which is a simple corollary of (Tripuraneni et al., 2021b,
 1009 Lemma 20), shows the concentration property of empirical covariance matrix.

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

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

1010 $X = (x_1, \dots, x_n)^T \in \mathbb{R}^{n \times d}$. Then for any $A, B \in \mathbb{R}^{d \times k}$, we have with probability at least $1 - \delta$

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

1011 *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 $U_1, U_2 \in \mathbb{R}^{d \times k}$,
 1012 $\Lambda_1, \Lambda_2, V_1, V_2 \in \mathbb{R}^{k \times k}$. Then

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

1013 Now since $U_1, U_2 \in \mathbb{R}^{d \times k}$ are projection matrices, we can apply Tripuraneni et al. (2021b) Lemma
 1014 20, therefore

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

1015 which gives what we want. \square

1016 The following lemma is a basic matrix perturbation result (see Tripuraneni et al. (2021b) Lemma
 1017 25).

1018 **Lemma 36.** Let A be a positive definite matrix and E another matrix which satisfies $\|EA^{-1}\| \leq \frac{1}{4}$,
 1019 then $F := (A + E)^{-1} - A^{-1}$ satisfies $\|F\| \leq \frac{4}{3} \|A^{-1}\| \|EA^{-1}\|$.

1020 With these two technical lemmas, we are able to prove Lemma 32, 33.

1021 *Proof of Lemma 32.* Notice that by the definition of U and β_U^* , we have $UU^T \beta_U^* = \beta_U^*$. We denote
 1022 $\alpha^* := U^T \beta_U^*$, then we also have $\beta_U^* = U \alpha^*$. Therefore

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

1023 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 orthogo-
 1024 nal subspaces which are rank k and rank $d - k$ respectively, so that $\text{range}(\widehat{U}) \oplus \text{range}(\widehat{U}_\perp) = \mathbb{R}^d$.
 1025 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

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

$$= \widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T X \widehat{U}_\perp \widehat{U}_\perp^T U^* - \widehat{U}_\perp \widehat{U}_\perp^T U^* \quad (24)$$

1026 Thus

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

1027 Here we use the notation $\|M\|_{\Sigma_T} := \sqrt{\|M^T \Sigma_T M\|}$ for matrix M .

1028 For the second term,

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

1029 For the first term,

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

1030 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
1031 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-
1032 ciently large n , therefore the term in (27) can be bounded well. First we need a careful analysis
1033 of $\|\widehat{U}^T \Sigma_S \widehat{U}_\perp \widehat{U}_\perp^T U^*\|$. It is obvious that

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

1034 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
1035 and Σ_S is diagonal. By definition of distance between two subspaces, there exist $R \in \mathcal{O}^{k \times k}$ and
1036 $Q \in \mathcal{O}^{(d-k) \times (d-k)}$, such that $\|\widehat{U} R - U\| = \Delta = \|\widehat{U}_\perp Q - U_\perp\|$. Then we have

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

1037 Combine (28) and (29), we have

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

1038 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
1039 $1 - \delta$,

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

1040 Therefore,

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

1041 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 $\sqrt{\frac{k}{n}} + \frac{k}{n} +$
1042 $\sqrt{\frac{\log(1/\delta)}{n}} + \frac{\log(1/\delta)}{n} \lesssim \sigma^{-2} C_{\min} \|\Sigma_S\|^{-1}$. Thus, we show that when n is large enough, we have
1043 $\|E(\widehat{U}^T \Sigma_S \widehat{U})^{-1}\| \leq \frac{1}{4}$. Therefore we can apply Lemma 36, which gives

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

1044 As for $\|E_1\|$, directly applying Lemma 35, when $n \gtrsim \sigma^4 \Delta^{-2} k \log(1/\delta)$ we get

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

1045 when $n \gtrsim \sigma^4 k \log(1/\delta)$ we have

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

1046 , if further we have $n \gtrsim \sigma^4 \Delta^{-2} k \log(1/\delta)$, then

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

1047 Combining (27), (30), (33) and (36), we have

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

1048 Combining (25), (26) and (37), we get

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

1049 with probability at least $1 - \delta$. Also, similar to (29), we have

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

1050 Similarly, we can further know that C_{\min} is close to λ_k :

$$\begin{aligned}
C_{\min} &= \lambda_k(\widehat{U}^T \Sigma_S \widehat{U}) \\
&= \lambda_k(R^T \widehat{U}^T \Sigma_S \widehat{U} R) \\
&= \lambda_k(U^T \Sigma_S U + R^T \widehat{U}^T \Sigma_S \widehat{U} R - U^T \Sigma_S U)
\end{aligned}$$

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

1051 where the last inequality holds when $\Delta \leq \frac{\lambda_k}{4\lambda_1}$. Finally, combining (38), (39), (40), we have

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

1052 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
1053 replace (36) by (35), we have

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

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

1054 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 ,
1055 $U_{\perp}^T \Sigma_T U_{\perp} = \Sigma_{T,-k}$, therefore the result is exactly what we want. \square

1056 *Proof of Lemma 33.* Recall $A_3 := \widehat{U}(\widehat{U}^T X^T X \widehat{U})^{-1} \widehat{U}^T X^T \epsilon$. Therefore

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

1057 Taking expectation with respect to ϵ , using $\mathbb{E}[\epsilon \epsilon^T] = v^2 I_n$, we have

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

1058 Here we actually need a bound stronger than (33) for $\|F\|$: recall (32), we have with probability
1059 $1 - \delta$

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

1060 Applying Lemma 36, which gives

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

1061 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

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

1062 The remaining thing is to show that indeed $\text{tr}((\widehat{U}^T \Sigma_S \widehat{U})^{-1} \widehat{U}^T \Sigma_T \widehat{U})$ is
1063 close to $\text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U)$. In fact, $\text{tr}((\widehat{U}^T \Sigma_S \widehat{U})^{-1} \widehat{U}^T \Sigma_T \widehat{U}) =$
1064 $\text{tr}((R^T \widehat{U}^T \Sigma_S \widehat{U} R)^{-1} R^T \widehat{U}^T \Sigma_T \widehat{U} R)$. Notice that

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

1065 we have

$$\begin{aligned}
&\text{tr}((R^T \widehat{U}^T \Sigma_S \widehat{U} R)^{-1} R^T \widehat{U}^T \Sigma_T \widehat{U} R) \quad (48) \\
&\leq \text{tr}((R^T \widehat{U}^T \Sigma_S \widehat{U} R)^{-1} U^T \Sigma_T U) + \|R^T \widehat{U}^T \Sigma_T \widehat{U} R - U^T \Sigma_T U\| \text{tr}((\widehat{U}^T \Sigma_S \widehat{U})^{-1}) \\
&\leq \text{tr}((R^T \widehat{U}^T \Sigma_S \widehat{U} R)^{-1} U^T \Sigma_T U) + 2\|\Delta\| \|\Sigma_T\| \text{tr}((\widehat{U}^T \Sigma_S \widehat{U})^{-1}) \\
&\leq \text{tr}((R^T \widehat{U}^T \Sigma_S \widehat{U} R)^{-1} U^T \Sigma_T U) + 2\|\Delta\| \|\Sigma_T\| k C_{\min}^{-1} \\
&\leq \text{tr}((R^T \widehat{U}^T \Sigma_S \widehat{U} R)^{-1} U^T \Sigma_T U) + \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U) \quad (49)
\end{aligned}$$

1066 when $\Delta \leq \frac{\lambda_k \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U)}{4k \|\Sigma_T\|}$. Also, we have

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

1067 therefore

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

1068 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

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

1069 whenever $\Delta \leq \frac{\lambda_k^2 \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U)}{4\lambda_1 k \|\Sigma_T\|} \leq \min\left\{\frac{\lambda_k^2 \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U)}{4\lambda_1 \text{tr}(U^T \Sigma_T U)}, \frac{\lambda_k \text{tr}((U^T \Sigma_S U)^{-1} U^T \Sigma_T U)}{4k \|\Sigma_T\|}\right\}$

1070 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
1071 $1 - \delta$. 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
1072 want. \square

1073 *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

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

1074 Therefore

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

1075 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 $B = \frac{X \widehat{U}}{\sqrt{n}} \in$
1076 $\mathbb{R}^{n \times r}$, then $A = \frac{1}{n} B (B^T B)^{-2} B^T$. Let the SVD of B be $B = P M O^T$, where $P \in \mathbb{R}^{n \times k}$,
1077 $M, O \in \mathbb{R}^{k \times k}$, then

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

1078 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 $1 - \delta$,
1079 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

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

1080 Thus $\|A\| \leq \mathcal{O}(\lambda_k^{-1})$. As for $\|X \beta_{\perp}^{\star}\|^2$, notice that the first- k entries of β_{\perp}^{\star} are zero, therefore
1081 $X \beta_{\perp}^{\star} = X_{-k} \beta_{-k}^{\star}$. by Lemma 35,

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

1082 Therefore we have

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

1083 Combining (52)(54) and (56), we have

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

1084 when $n \gtrsim \sigma^4 C_{\min}^{-2} \|\Sigma_S\|^2 k \log(1/\delta)$ and $\Delta \leq \min\left\{\frac{\|U^T \Sigma_T U\|}{2\|\Sigma_T\|}, \frac{\lambda_k}{4\lambda_1}\right\}$. \square

1085 Finally we prove Lemma 30 in the following.

1086 *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
 1087 the sample covariance matrix $\widehat{\Sigma}_S := \frac{1}{n} X X^T = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$ using PCA. Then by Davis-Kahan
 1088 theorem (Chen et al., 2021, Corollary 2.8),

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

1089 Therefore it remains to bound $\|\widehat{\Sigma}_S - \Sigma_S\|$. Applying Lemma 28, we immediately have

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

1090 where $r = \frac{\sum_{i=1}^n \lambda_i}{\lambda_1}$. Together with (58), we have with probability at least $1 - \delta$,

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

1091

□