

# Can Machines Learn Weak Signals?\*

Zhouyu Shen<sup>†</sup>

Dacheng Xiu<sup>‡</sup>

Booth School of Business

Booth School of Business

University of Chicago

University of Chicago and NBER

January 21, 2025

## Abstract

In high-dimensional regressions with low signal-to-noise ratios, we assess the predictive performance of several prevalent machine learning methods. Theoretical insights show Ridge regression’s superiority in exploiting weak signals, surpassing a zero benchmark. In contrast, Lasso fails to exceed this baseline, indicating its learning limitations. Simulations reveal that Random Forest generally outperforms Gradient Boosted Regression Trees when signals are weak. Moreover, Neural Networks with  $\ell_2$ -regularization excel in capturing nonlinear functions of weak signals. Our empirical analysis across six economic datasets suggests that the weakness of signals, not necessarily the absence of sparsity, may be Lasso’s major limitation in economic predictions.

**Keywords:** Weak Signals, Precise Error, Machine Learning, Bayes Risk

## 1 Introduction

In regression analysis, covariates with non-zero coefficients are identified as true signals, while those with zero coefficients are considered false signals. In a population model, this

---

\*We benefited tremendously from discussions with Gustavo Greire, Ulrich Muller, Mikkel Plagborg-Moller, Alberto Quaini, Pragma Sur, as well as seminar and conference participants at Aarhus University, Duke University, Tsinghua University, Bates White LLC, Econometrics Society North American Winter Meeting, ESIF Economics and AI+ML Meeting, Triangle Econometrics Conference, Applied Machine Learning, Economics, and Data Science Webinar, the Stevanovich Center Conference on Big Data and Machine Learning in Econometrics, Finance, and Statistics, the fifth International Workshop in Financial Econometrics, and Young Econometricians in Asia-Pacific Annual Meeting.

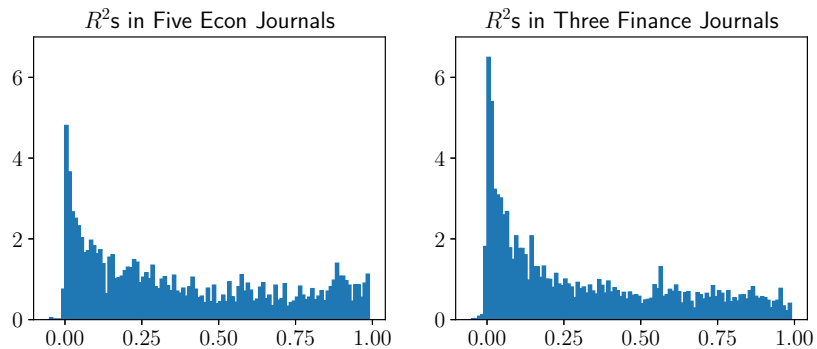
<sup>†</sup>Address: 5807 S Woodlawn Avenue, Chicago, IL 60637 USA. Email: [zshen10@chicagobooth.edu](mailto:zshen10@chicagobooth.edu).

<sup>‡</sup>Address: 5807 S Woodlawn Avenue, Chicago, IL 60637, USA. Email: [dacheng.xiu@chicagobooth.edu](mailto:dacheng.xiu@chicagobooth.edu).

distinction is clear-cut, resembling a “black and white” scenario. However, in finite samples, the presence of minuscule non-zero coefficients introduces a “gray” area, blurring the boundary between true and false signals.<sup>1</sup> This gray area represents weak signals—covariates that, individually, exert negligible influence on the outcome variable.

The investigation of weak signals holds tangible implications for economic and financial decision-making. Often, it is the collective impact of these weak signals that drives the outcomes in these fields. Supporting this, Figure 1 provides empirical evidence by presenting  $R^2$  values gathered from a compendium of Economics and Finance journal articles published in 2022. The 25% quantiles of these  $R^2$  values stand at 9.7% for economics and 5.8% for finance, suggesting that models in these disciplines frequently rely on covariates with modest explanatory power. Furthermore, Figure 1 is based exclusively on published studies, which are likely biased toward higher  $R^2$  values due to selection effects. This suggests that the presence of weak signals may be even more widespread than the data here indicates.

Figure 1: Histograms of  $R^2$ s in Selected Economics and Finance Journals



Note: The histograms depict  $R^2$ s manually collected from published papers in a selection of Economics and Finance journals in 2022. This collection comprises data from five Economics journals (left) and three Finance journals (right). Specifically, there are a total of 411 papers published in the American Economic Review, Econometrica, Journal of Political Economy, Quarterly Journal of Economics, and Review of Economic Studies, resulting in 8,129  $R^2$  observations. In addition, there are 380 papers from the Journal of Finance, Journal of Financial Economics, and Review of Financial Studies, contributing 12,198  $R^2$  observations.

The decision to incorporate weak signals into a regression model carries the risk of overfitting, potentially undermining predictive accuracy. Overfitting arises when the increased variance from estimating the coefficients of these weak signals outweighs the bias reduction.

<sup>1</sup> The comparison of regression coefficients’ magnitudes is meaningful only when predictor variables are normalized, an assumption implicitly adopted in the subsequent discussion.

tion gained from their inclusion. This trade-off becomes even more pronounced in high-dimensional settings, where the limited sample size relative to the number of covariates amplifies prediction errors.

Machine learning methods, renowned for their focus on variable selection and dimension reduction, have proven effective in mitigating overfitting and detecting true signals from false ones, particularly when the true signals are strong. These methods employ regularization techniques, such as penalizing the  $\ell_1$  or  $\ell_2$  norms of model parameters, to achieve this objective. However, a pivotal question emerges: Can machines learn weak signals, or in other words, can they surpass the naive zero predictor? The zero predictor, which disregards all covariates and always predicts zero, serves as a passive baseline in the context of weak signals. Surpassing this baseline indicates that a method has successfully extracted valuable signals. Conversely, failing to exceed it highlights a limitation in its learning capacity. In view of these considerations, we focus on evaluating the relative performance of regularized predictors, Ridge and Lasso, against the zero predictor in high-dimensional regressions, where the prediction target is driven by predictors that exhibit weak correlations with it.

In scenarios with sufficiently strong signals, both Lasso and Ridge are expected to outperform zero by effectively capturing and utilizing at least some of these signals. Hence, accurately defining the notion of “weak” signals is crucial at the outset of our investigation. This definition serves a dual purpose: it prevents the scenario from defaulting to trivial comparisons akin to strong signal cases and ensures practical relevance to finite sample scenarios in economics and finance. We characterize a weak signal scenario as one in which the zero predictor achieves the minimal Bayes prediction risk asymptotically. This setting turns out to encompass a wide class of data generating processes (DGPs), and it approximates a finite sample reality in which zero serves as a competitive benchmark.

In the defined weak signal scenarios, conventional error-bound analyses are insufficient in distinguishing the performance of different predictors. Considering the optimal Lasso and Ridge—each tuned to minimize prediction error—both can perform no worse than the zero predictor. Intuitively, as the tuning parameters for Ridge and Lasso approach infinity, Ridge converges to zero, while Lasso effectively becomes equivalent to zero. Consequently, all three predictors—optimal Lasso, optimal Ridge, and zero—can achieve the optimal Bayes prediction error, making them indistinguishable.

However, a more refined analysis reveals that, across a broad range of tuning parameter choices, Ridge outperforms zero. In contrast, the optimal Lasso effectively reduces to zero, meaning Lasso does not exceed zero’s performance regardless of its tuning parameter choice.

This conclusion is based on a precise error approach, which enables us to zoom in and explicitly characterize the *relative* prediction errors of Ridge and Lasso compared to the zero predictor. These results highlight an important distinction between shrinkage and selection methods. Shrinkage methods like Ridge are more effective in environments with more homogenous signal strength. On the other hand, selection methods like Lasso are preferable in scenarios where there is a clear distinction between true and false signals. In weak signal contexts where this distinction is blurred, the advantage of Lasso tends to wane.

Our study further demonstrates the effectiveness of the cross-validation (CV) algorithm in selecting the optimal tuning parameter for Ridge, even in weak signal contexts. This underscores CV’s robustness as a model-tuning tool in these settings. Moreover, the out-of-sample  $R^2$  from the optimal Ridge—a metric frequently used for assessing predictor performance on unseen data—proves a relevant indicator of the signal-to-noise ratio in the DGP, despite a notable gap between its asymptotic limit and the population  $R^2$ .

In the final aspect of our theoretical analysis, we expand our framework to include models featuring a mix of signal strengths. This section specifically addresses scenarios in which a benchmark model contains potentially strong signals. Our focus then shifts to evaluating the benefits of leveraging the predictive power of the remaining weak signals. To this end, we derive ordinary least squares (OLS) residuals, from which the impact of potentially strong covariates in the benchmark model has been removed. Consistent with our earlier findings, applying Ridge regression to these residuals, using the remaining covariates, enhances predictive performance compared to a baseline predictor that ignores these additional covariates.

Our simulation analysis corroborates our theoretical findings: Ridge surpasses zero, which in turn edges out Lasso, especially in DGPs characterized by low  $R^2$  values. When exploring more sophisticated machine learning techniques, we observe that Random Forest (RF), which produces dense models by including nearly all variables, outperforms the zero predictor. The latter, in turn, surpasses Gradient Boosted Regression Trees (GBRT), as GBRT, similar to Lasso, tends to generate sparse models. Furthermore, Neural Networks (NNs), when paired with the  $\ell_2$ -norm regularization, can deliver superior predictions. In contrast, applying an  $\ell_1$ -penalty in these networks fails to achieve comparable results.

Our empirical analysis spans six datasets across macroeconomics, microeconomics, and finance. Five align with [Giannone et al. \(2022\)](#), and one is sourced from [Gu et al. \(2020\)](#). Our finance examples delve into predicting market returns using financial and economic indicators, as well as firm-level return prediction based on their specific characteristics. In the macroeconomic context, we examine time-series predictions of industrial production

using macroeconomic indicators, and a cross-country GDP growth prediction, utilizing socioeconomic, institutional, and geographical factors. Our microeconomic studies focus on crime rate predictions and pro-plaintiff appellate decisions in takings law rulings.

The relevance of weak signals in datasets is contingent on the choice of benchmark models. For instance, when compared to a benchmark with an intercept alone, weak signals are revealed in four out of six datasets. Further benchmarking against covariates informed by economic theory reveals weak signals across all datasets, making them particularly well-suited for the application of our asymptotic theory. Drawing from their empirical analysis of these datasets, [Giannone et al. \(2022\)](#) argue that sparsity may be an illusion, as optimal predictive models often rely on a large number of covariates. Our collective theoretical and empirical evidence points to signal weakness as a key factor in the underperformance of Lasso. As our results suggest, even in cases where the majority of signals have zero coefficients in the true DGP, Ridge may still outperform Lasso if the true signals are weak. Their comparative performance thus does not necessarily offer insights into the sparsity of the DGP itself.

In light of these findings, we recommend a cautious approach to employing Lasso in economic and financial settings. Despite its popularity as a modern counterpart to OLS, Lasso’s effectiveness may be compromised in scenarios characterized by weak signals. Our study complements the findings of [Kolesár et al. \(2024\)](#), who highlight issues with sparsity-based estimators, such as their lack of invariance to reparametrization and sensitivity to normalizations that are otherwise innocuous to OLS.

Our paper is closely related to the literature on the theoretical performance of Ridge and Lasso, with two main threads being particularly relevant. The first focuses on error-bound analysis. For Ridge, [Hoerl and Kennard \(1970\)](#) show that the prediction error decreases at a rate of  $p/n$ , where  $p$  is the number of covariates and  $n$  the sample size, with its magnitude tied to the eigen-structure of the design matrix. For Lasso, the prediction error vanishes if  $s \log p/n \rightarrow 0$ , where  $s$  is the number of non-zero parameters (see, e.g., [Wainwright \(2019\)](#)). However, we consider an asymptotic setting where these error bounds fail to distinguish Ridge and Lasso from the zero predictor, as their leading-order prediction errors are identical. This motivates a more granular, higher-order analysis of prediction errors.

The second, more recent strand of research focuses on determining the precise probability limit of the prediction error for Ridge and Lasso.<sup>2</sup> [Bayati and Montanari \(2012\)](#)

---

<sup>2</sup>The precise error analysis has provided valuable insights into various machine learning methods. For example, [Liang and Sur \(2022\)](#) examine the properties of minimum  $\ell_1$ -norm interpolation and boosting in linear models. [Miolane and Montanari \(2021\)](#) explore cross-validation for Lasso, while [Hastie et al. \(2022\)](#) investigate minimum  $\ell_2$ -norm interpolation, shedding light on the double-descent phenomenon in neural

employ approximate message passing algorithms to link them with Lasso and derive its error limit. Alternatively, [Thrapoulidis et al. \(2015\)](#) use the Convex Gaussian Minimax Theory (CGMT) to simplify Lasso’s optimization problem, enabling precise error derivation. For Ridge, [Dicker \(2016\)](#) provides analogous insights into its prediction error. However, these precise error analyses often rely on stringent parametric assumptions, such as independently Gaussian-distributed design matrix elements. [Dobriban and Wager \(2018\)](#) extend [Dicker \(2016\)](#)’s work by accommodating dependent covariates and non-Gaussian predictors, leveraging universality results from random matrix theory.

This paper is organized as follows. Section 2 presents the main theoretical results regarding Ridge and Lasso. Section 3 conducts simulations to illuminate our theoretical predictions while also expanding the analysis to assess the performance of advanced machine learning methods under weak signals. Lastly, Section 4 provides empirical results. The appendix contains mathematical proofs of main results, while the online appendix provides additional theoretical results, technical lemmas, and their proofs.

**Notation:** For any  $x \in \mathbb{R}$ , we refer to  $\max(x, 0)$  as  $x_+$ . For any vector  $x$ ,  $\|x\|_0$ ,  $\|x\|_1$ ,  $\|x\|$  and  $\|x\|_\infty$  represent its  $\ell_0$ ,  $\ell_1$ ,  $\ell_2$  and  $\ell_\infty$  norms, respectively. For a real matrix  $A$ , we use  $\|A\|$  and  $\|A\|_F$  to denote its spectral norm (or  $\ell_2$  norm), and the Frobenius norm, that is,  $\sqrt{\lambda_{\max}(A^\top A)}$ , and  $\sqrt{\text{Tr}(A^\top A)}$ , respectively. In the case where  $A$  is a  $p \times p$  matrix,  $\lambda_i(A)$  denotes its  $i$ -th largest eigenvalue, for  $1 \leq i \leq p$ . We use the notation  $x_n \lesssim y_n$  when there exists a constant  $C$  such that  $x_n \leq C y_n$  holds for sufficiently large  $n$ . Similarly, we use  $x_n \lesssim_P y_n$  to denote  $x_n = O_P(y_n)$ . If  $x_n \lesssim y_n$  and  $y_n \lesssim x_n$ , we write  $x_n \asymp y_n$  for short. Similarly, we use  $x_n \asymp_P y_n$  if  $x_n \lesssim_P y_n$  and  $y_n \lesssim_P x_n$ .

## 2 Theoretical Results

### 2.1 Model Setup

We start with the following linear regression model:

$$y = X\beta_0 + \varepsilon, \tag{1}$$

where  $y \in \mathbb{R}^n$ ,  $X \in \mathbb{R}^{n \times p}$ ,  $\beta_0 \in \mathbb{R}^p$  and  $\varepsilon \in \mathbb{R}^n$ . Throughout our discussion,  $X$ ,  $\beta_0$ , and  $\varepsilon$  are treated as random variables, mutually independent of one another. Central to our analysis

---

networks. Regarding variable selection, [Su et al. \(2017\)](#) study the false discovery rate of the Lasso path, and [Wang et al. \(2020\)](#) compare the variable selection properties of bridge estimators.

is the calculation of the probability limit of the prediction errors, which necessitates assuming a (prior) probability distribution on the coefficients. This setting aligns with the literature on precise errors, which also connects our analysis of prediction error with Bayes risk.

Our objective is to investigate the predictive performance of machine learning techniques in the presence of weak signals.<sup>3</sup> To accomplish this, we focus on a high-dimensional regression setting characterized by an increasing number of predictors, that is,  $p \rightarrow \infty$ . In such a context, regularization techniques become not just relevant but often necessary. Moreover, our specific focus is on situations where the signals are weak, characterized by the condition:  $\|\beta_0\|^2 \asymp_P \tau \rightarrow 0$ . The choice to use  $\|\beta_0\|$  as the metric for characterizing weak signals is due to its close relationship with the widely-adopted  $R^2$  metric in regression analysis, which provides a familiar and intuitive understanding of signal strength. Our investigation then progresses to an asymptotic analysis under these two conditions. We will detail the specific requirements for  $p$ ,  $\tau$ , and sample size  $n$  after introducing the baseline predictor.

Now, we proceed to present the assumptions governing the DGP of  $X$ :

**Assumption 1.** *The covariates  $X \in \mathbb{R}^{n \times p}$  are generated as  $X = \Sigma_1^{1/2} Z \Sigma_2^{1/2}$  for an  $n \times p$  matrix  $Z$  with i.i.d. standard Gaussian entries, deterministic  $n \times n$  and  $p \times p$  positive definite matrices  $\Sigma_1$  and  $\Sigma_2$ . In addition, there exist positive constants  $c_1, C_1, c_2, C_2$  such that  $c_1 \leq \lambda_i(\Sigma_1) \leq C_1$ ,  $i = 1, 2, \dots, n$  and  $c_2 \leq \lambda_i(\Sigma_2) \leq C_2$ ,  $i = 1, 2, \dots, p$ .*

This assumption accommodates time series dependencies via  $\Sigma_1$  and cross-sectional correlations via  $\Sigma_2$ . The eigenvalue constraints serve two purposes: upper bounds limit excessive dependencies, while lower bounds prevent multicollinearity and ensure observations are not linearly dependent across time. Moreover, since we focus on prediction rather than variable selection, the dependence structure among  $X$  does not adversely impact Lasso’s predictive performance (see Section 7.4 of [Wainwright \(2019\)](#)), even though its variable selection properties are sensitive to strong dependence among covariates.

While the Gaussian assumption for  $X$  is integral to our use of Gordon’s inequality ([Gordon \(1988\)](#)) for Gaussian processes in the proof, it does raise concerns regarding the robustness of our findings when this assumption is not met in real-world scenarios. Our simulation results indicate that the Gaussian assumption appears non-essential and our asymptotic theory approximates finite sample behavior even with relatively small sizes, typically a few hundred observations. This observation aligns with similar findings in random matrix theory, where asymptotic properties initially derived for Gaussian ensembles were subsequently

---

<sup>3</sup>While [Donoho and Jin \(2004\)](#) and [Hall and Jin \(2010\)](#) study variable selection in rare and weak signals, we focus on prediction in asymptotic settings where identifying nonzero coefficients is infeasible.

shown to extend to a wider spectrum of random matrices—a phenomenon referred to as the universality. Notably, when  $\Sigma_1 = \mathbb{I}$ , [Dobriban and Wager \(2018\)](#) use random matrix theory to bypass the Gaussian assumption. However, their technique appears only applicable to Ridge. Our objective is to compare Ridge and Lasso under a unified framework, which necessitates the Gaussian assumption.

Next, we specify the assumption regarding  $\varepsilon$ :

**Assumption 2.** *Let  $\varepsilon = \Sigma_\varepsilon^{1/2}z$ , where  $z$  comprises i.i.d. variables with mean zero, variance one and finite fourth moment and  $\Sigma_\varepsilon$  is a positive definite matrix satisfying  $c_\varepsilon \leq \lambda_i(\Sigma_\varepsilon) \leq C_\varepsilon$ ,  $i = 1, 2, \dots, n$ , for some fixed positive constants  $c_\varepsilon$  and  $C_\varepsilon$ .*

This assumption allows for autocorrelations and heteroscedasticity. If  $\Sigma_\varepsilon$  is a diagonal matrix, its spectral norm is evidently bounded under the condition that each entry of  $\varepsilon$  has finite variance. Moreover, if  $\varepsilon$  follows a stationary process characterized by exponentially decaying autocorrelations, it can be shown that the spectral norm of  $\Sigma_\varepsilon$  remains bounded.

Under Assumptions 1 and 2, it follows that  $\|X\beta_0\| \asymp_{\mathbb{P}} \sqrt{n} \|\beta_0\|$  and  $\|\varepsilon\| \asymp_{\mathbb{P}} \sqrt{n}$ . This indicates that the magnitude of each entry in  $X$  and  $\varepsilon$  neither explode nor vanish asymptotically. Consequently, the magnitude of the signal-to-noise ratio (or  $R^2$ ) is entirely dictated by  $\|\beta_0\|$ . Next, we impose an assumption that governs the properties of  $\beta_0$ :

**Assumption 3.** *The vector  $b_0 = \sqrt{p\tau^{-1}}\beta_0$  comprises i.i.d. random variables, each following a prior probability distribution  $F$  belonging to the class  $\mathcal{F}$ . The class  $\mathcal{F}$  is defined such that any included random variable can be represented as  $q^{-1/2}b_1b_2$ , where  $b_1$  and  $b_2$  are independent,  $b_1$  follows a binomial distribution  $B(1, q)$ , and  $b_2$  is a sub-exponential random variable with a mean of zero and a variance denoted as  $\sigma_\beta^2$ .*

Without loss of generality, we use the term  $\sqrt{p\tau^{-1}}$  as the normalization factor, ensuring that  $\|\beta_0\|^2 \asymp_{\mathbb{P}} \tau$ . This normalization facilitates a clearer interpretation of our results. While the i.i.d. assumption may seem strong, it offers greater transparency by simplifying more complex technical assumptions necessary to derive essential probability bounds. In particular, this assumption allows for important classes of models, such as a spike-and-slab prior for  $b_0$ , extensively studied by [Giannone et al. \(2022\)](#) to examine the empirical relevance of sparsity in economic datasets. Each element of  $b_0$  follows a mixed distribution, such as when  $q = 1$ , with  $b_2$  modeled by  $(1 - v)\psi_0 + v\psi_1$ , where  $v$ , a fixed constant within  $[0, 1]$ , modulates the mix between the spike ( $\psi_0$ ) and slab ( $\psi_1$ ) components of the prior. More generally, the formulation  $q^{-1/2}b_1b_2$  accommodates a spike-and-slab model with more extreme sparsity ( $q \rightarrow 0$ ) through the component  $q^{-1/2}b_1$ . This scaling,  $q^{-1/2}$ , ensures the variance of



$q^{-1/2}b_1b_2$  remains finite and non-vanishing. Essentially,  $q$  dictates the sparsity of  $\beta_0$ : when  $P(b_2 = 0) = 0$ ,  $\|\beta_0\|_0 \asymp_P pq$ . In scenarios with strong signals ( $\tau = 1$ ), a DGP with  $q$  nearing zero typically favors Lasso, whereas a  $q$  closer to one suggests a preference for Ridge. Therefore, this framework does not inherently privilege either.

The underlying assumptions that justify Ridge and Lasso are notably distinct, particularly in the context of error-bound analysis. For instance, the analysis of Lasso often requires the approximate sparsity condition and the restricted eigenvalue condition (see, e.g., Belloni et al. (2013b) and Bickel et al. (2009)). On the other hand, the convergence rate of Ridge’s prediction error requires intricate conditions on the eigenvalue structure of the design matrix, as discussed in Tsigler and Bartlett (2023). In contrast, our analysis here compares the asymptotic properties of different estimators within a common framework.

## 2.2 Predictors

We now turn our attention to the discussion of the predictors. In scenarios involving weak signals, characterized by  $\|\beta_0\| \rightarrow 0$ , a straightforward and natural baseline predictor emerges, that is, the naive zero predictor. As an estimator, zero is clearly consistent in terms of the  $\ell_2$ -loss of the estimation error, because it reduces to  $\|\beta_0\|$ , which vanishes in this context.

The zero predictor serves as a passive benchmark for a scenario where no learning occurs. To surpass its performance, any alternative predictor must harness some signals. This indicates the alternative predictor’s capability to successfully identify and leverage weak signals. Therefore, to address the earlier question of whether machines can learn weak signals, we need to compare the machine learning method’s performance with that of the zero predictor. Only if they can do so can they outperform the naive zero predictor.

In our study, we consider Ridge and Lasso as contenders that leverage machine learning techniques. These methods are widely used benchmarks in practice, owing to their simplicity and universality. An in-depth analysis of these predictors provides valuable insights into their specific regularization techniques, which can be extended to more advanced models.

The Ridge estimator, denoted as  $\hat{\beta}_r$ , is the solution to the following optimization problem:

$$\hat{\beta}_r(\lambda_n) := \arg \min_{\beta} \frac{1}{n} \|y - X\beta\|^2 + \frac{p\lambda_n}{n} \|\beta\|^2, \quad (2)$$

where  $\lambda_n$  is its tuning parameter governing the strength of the regularization. In contrast, the Lasso estimator, denoted as  $\hat{\beta}_l$ , is defined as:

$$\hat{\beta}_l(\lambda_n) := \arg \min_{\beta} \frac{1}{n} \|y - X\beta\|^2 + \frac{\lambda_n}{\sqrt{n}} \|\beta\|_1. \quad (3)$$

By convention, and without loss of generality, the terms involving penalties are typically scaled by  $p/n$  in the case of Ridge and by  $1/\sqrt{n}$  in the case of Lasso.

In addition, our theoretical results also encompass the OLS estimator and the Ridgeless estimator, both of which correspond to special cases of Ridge when the tuning parameter  $\lambda_n$  is set to zero. When  $p > n$ , the least squares problem has an infinite number of solutions. Among these solutions, the Ridgeless estimator can be regarded as a minimum-norm interpolating linear predictor, as noted by [Bartlett et al. \(2020\)](#):

$$\hat{\beta}_r(0) = \arg \min_{\beta} \|\beta\|, \quad \text{s.t.} \quad X\beta = y. \quad (4)$$

It is also possible to explore other interpolators, such as the minimum  $\ell_1$ -norm interpolator studied by [Liang and Sur \(2022\)](#). Future research might extend our analysis to other penalized linear estimators, such as Elastic Net, as introduced by [Zou and Hastie \(2005\)](#), or SCAD by [Fan and Li \(2001\)](#).

With  $\beta_0$  estimated by some  $\hat{\beta}$ , it is straightforward to construct corresponding predictor,  $(x^{\text{new}})^\top \hat{\beta}$ , for a new data point  $x^{\text{new}}$ .

### 2.3 Bayes Risk

Now, we proceed to define the metric by which we assess various predictors. For any predictor, our primary interest is its Bayes prediction risk. This risk is related to the expected squared prediction error evaluated at a new, independent data point  $(x^{\text{new}}, y^{\text{new}})$ . In the case of a linear predictor,  $\hat{y}^{\text{new}} = (x^{\text{new}})^\top \hat{\beta}$ , we can write the prediction error explicitly as:

$$\mathbb{E}_F (y^{\text{new}} - \hat{y}^{\text{new}})^2 = \sigma_\varepsilon^2 + \mathbb{E}_F \left[ (x^{\text{new}})^\top (\hat{\beta} - \beta_0) \right]^2 = \sigma_\varepsilon^2 + \mathbb{E}_F \|\Sigma_2^{1/2} (\hat{\beta} - \beta_0)\|^2, \quad (5)$$

where the subscript in the expectation operator  $\mathbb{E}_F(\cdot)$  emphasizes the fact that the expectation is taken with respect to the prior distribution of  $b_0 = \sqrt{p\tau^{-1}}\beta_0$ . Obviously, it is the second term in (5) that governs the relative predictive performance. Barring  $\|\Sigma_2\|$ , the prediction risk is closely tied to the estimation error of  $\hat{\beta}$ , i.e.,  $\|\hat{\beta} - \beta_0\|$ .

Formally, we define the Bayes prediction risk associated with an estimator  $\hat{\beta}$  as

$$\mathcal{R}(\hat{\beta}, F) := \mathbb{E}_F \|\Sigma_2^{1/2} (\hat{\beta} - \beta_0)\|^2.$$

The predictor that minimizes this Bayes risk is termed the Bayes predictor. In our framework, it is straightforward to show (see Chapter 4 of [Berger \(1985\)](#)) that the Bayes predictor corresponds with the predictor that is derived from the posterior mean of  $\beta_0$ , which is represented as  $\mathbb{E}_F(\beta_0|X, y)$ . We denote its Bayes risk as  $\mathcal{R}(F)$ .<sup>4</sup>

Under strong signals, i.e.,  $\tau = 1$ , and suppose that  $\Sigma_1, \Sigma_2$  and  $\Sigma_\varepsilon$  are identity matrices,  $p/n \rightarrow c_0 \in \mathbb{R}^+$ , significant progress has been made in understanding the asymptotic behavior of Bayes risk. For Ridge regression, notable studies by [Dicker \(2016\)](#) and [Dobriban and Wager \(2018\)](#) have derived the asymptotic limit of the Bayes risk.<sup>5</sup> Similarly, in the case of Lasso, several studies, such as those conducted by [Bayati and Montanari \(2012\)](#) and [Thrapoulidis et al. \(2018\)](#), have established its asymptotic Bayes risk limit.<sup>6</sup>

In Figure 2, we present two heatmaps, with the  $y$  and  $x$  axes representing various values of  $p/n$  and  $\|\beta_0\|^2 = \tau\sigma_\beta^2$ . The left heatmap illustrates the ratio of Bayes risk between optimal Ridge and the zero estimator, while the right heatmap represents the ratio of optimal Lasso against the zero estimator. For both Ridge and Lasso, their optimal tuning parameters are selected by minimizing the probability limits of their Bayes risk given by (6) and (7), respectively. A prediction error ratio below 1 suggests that zero is outperformed.

The heatmaps, as anticipated, clearly demonstrate that both optimal Ridge and Lasso surpass the performance of the zero estimator. This superiority is particularly pronounced in scenarios involving strong signals and relatively lower dimensions. Notably, the disparity between these estimators becomes less pronounced as the norm of  $\|\beta_0\|$  approaches zero and the ratio  $p/n$  increases, indicating a shift towards scenarios characterized by weaker signals. The existing result on precise error analysis is primarily built upon the assumption of strong signals, where  $\tau = 1$ . To discern the performance of various estimators under weak signal

---

<sup>4</sup>There exists an extensive body of literature focused on empirical Bayes methods, which explores feasible approaches for implementing  $\mathbb{E}_F(\beta_0|X, y)$ , when  $F$  is unknown, see, e.g., [Robbins \(1964\)](#) and [Efron \(2012\)](#).

<sup>5</sup>The exact form of the limit is given by

$$\lim_{n \rightarrow \infty} \mathcal{R}(\hat{\beta}_r(\lambda_n), F) = c_0 m(-\lambda, c_0) + \lambda(\lambda\sigma_\beta^2 - c_0)m'(-\lambda, c_0), \quad (6)$$

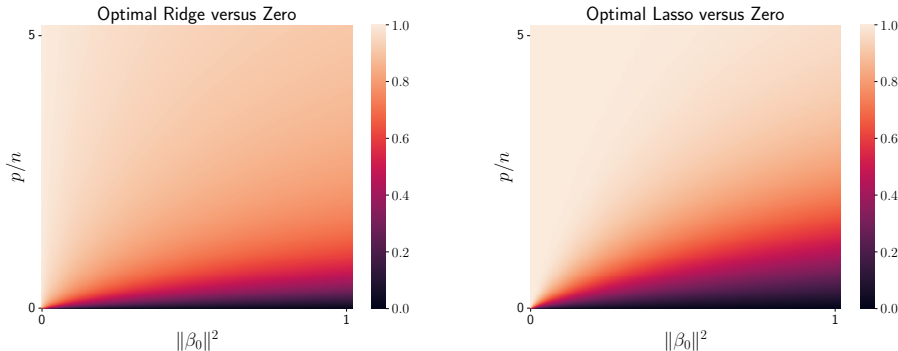
where  $\lambda = \lim_{n \rightarrow \infty} c_0 \lambda_n$  and  $m(-\lambda, c_0) = (-(1 - c_0 + \lambda) + \sqrt{(1 - c_0 + \lambda)^2 + 4c_0\lambda})/2c_0\lambda$ .

<sup>6</sup>The limit in the case of Lasso can be explicitly written as follows:  $\lim_{n \rightarrow \infty} \mathcal{R}(\hat{\beta}_l(\lambda_n), F) = (\alpha^*)^2$ , where

$$\alpha^* = \arg \min_{\alpha \geq 0} \left\{ \inf_{\tau_g > 0} \sup_{\substack{\beta \geq 0 \\ \tau_h > 0}} \frac{\beta\tau_g}{2} + \frac{1}{c_0} L\left(\alpha, \frac{\tau_g}{\beta}\right) - \frac{\alpha\tau_h}{2} - \frac{\alpha\beta^2}{2\tau_h} + \lambda G\left(\frac{\alpha\beta}{\tau_h}, \frac{\alpha\lambda}{\tau_h}\right) \right\}. \quad (7)$$

Here  $\lambda = \lim_{n \rightarrow \infty} \lambda_n/\sqrt{c_0}$ ,  $L(c, \tau) := \mathbb{E}[e_{x^2}(cZ + \varepsilon, \tau) - \varepsilon^2]$ , and  $G(c, \tau) := \mathbb{E}[e_{|x|}(cZ + X_b, \tau) - |X_b|]$ , with  $e_f(y, \tau) := \min_v (y - v)^2/2\tau + f(v)$ , where random variables  $Z, \varepsilon$ , and  $X_b$  follow a standard normal distribution ( $Z \sim \mathcal{N}(0, 1)$ ), the distribution of the noise, and the distribution  $F$ , respectively.

Figure 2: Comparison of Prediction Errors: Optimal Ridge and Lasso vs. the Zero Estimator



Note: The left panel illustrates the ratio of prediction error between the optimal Ridge and the baseline zero estimator. Conversely, the right panel presents a similar comparison for the optimal Lasso estimator against the zero estimator. Both axes,  $y$  and  $x$ , depict a range of  $p/n$  ratios and  $\|\beta_0\|^2$  values, corresponding to data generated in accordance with the model described in (1), with  $\Sigma_1$ ,  $\Sigma_2$ , and  $\Sigma_\varepsilon$  in Assumptions 1 and 2 set as identity matrices. We set  $b_0$  as a Dirac-spike and a Gaussian slab with  $q = 1/5$ . In this context of strong signals, the prediction errors for both optimal Ridge and Lasso are calculated using tuning parameters that are optimally selected to minimize the expected prediction errors' probability limits, as given by (6) and (7).

conditions, a more intricate analysis in the limiting case ( $\tau \rightarrow 0$  and  $p/n \not\rightarrow 0$ ) is necessary.

## 2.4 Zero's Optimality and Relative Prediction Error

Figure 2 also indicates that our attention should be directed towards a regime where the zero estimator exhibits meaningful competitiveness. Otherwise, we may question the appropriateness of our definition of “weak” signals if some machine learning approaches can obviously outperform it by a wide margin.

One might be tempted to define weak signals as instances where the signal strength falls below a certain “detection boundary,” thereby becoming indiscernible through hypothesis testing. Our focus is on prediction, rather than signal detection. This distinction is key because, even when signals are undetectable by hypothesis testing, their collective contribution to prediction can still outperform the zero predictor. The zero predictor serves as a natural benchmark for demonstrating the capacity of machine learning to utilize weak signals.

To motivate our concept of weak signals, we analyze a regime where the zero predictor achieves certain notion of optimality, indicated by its Bayes risk being identical to that of the Bayes predictor. This scenario is delineated more precisely by the assumption below:

**Assumption 4.**  $n^{-1}p \rightarrow c_0 \in (0, \infty]$ ,  $n^{-1}\tau p(\log p)^4 \rightarrow 0$ ,  $n\tau p^{-2/3}(\log p)^{-4} \rightarrow \infty$ , and

$$n^{-1}pq\tau^{-1}(\log p)^{-4} \rightarrow \infty.$$

Assumption 4 covers a wide spectrum of signal strengths and counts, while simultaneously imposing constraints to prevent an excessively large  $p/n$  ratio and overly rapid vanishing of  $\tau$ . The first two constraints imply  $\tau \rightarrow 0$ , while the third imposes a lower bound on  $\tau$ . Together, these constraints require that  $\tau$  is bounded below by  $n^{-1/3}$ .

The final constraint addresses cases of extreme sparsity in  $\beta_0$ . It becomes redundant when  $q$  does not vanish, as it is already implied by the first two constraints. Collectively, these conditions imply that  $(pq) \log p/n$  is bounded below by  $\tau$ , and consequently by  $n^{-1/3}$  (up to a logarithmic factor). Importantly, these constraints do not exclude the sparsity assumptions commonly adopted in the literature for Lasso:  $\|\beta_0\|_0 \log p/n \rightarrow 0$ , where  $\|\beta_0\|_0 \asymp_{\mathbb{P}} pq$ .

These constraints serve to exclude edge cases where the relative performance of different estimators cannot be conclusively determined using our proof technique. In Section 2.9, we investigate scenarios outside the scope of these constraints (e.g., lower bound of  $n^{-1/3}$ ), such as cases of extreme sparsity with only one true signal in the DGP. By employing an alternative proof method that leverages the closed-form solution of Lasso in a special case, we demonstrate that these constraints are not necessary for arriving at our conclusions.

To facilitate the discussion of the optimal estimator in our context, we refer to the definition provided by Robbins (1964).

**Definition 1.** We say  $\hat{\beta}$  is asymptotically optimal relative to  $F$ , if it satisfies

$$\lim_{n \rightarrow \infty} \frac{\mathcal{R}(\hat{\beta}, F)}{\mathcal{R}(F)} = 1.^7$$

**Theorem 1.** Suppose that Assumptions 1–4 hold. Furthermore, assume that the error term  $\varepsilon$  in Assumption 2 follows a Gaussian distribution.<sup>8</sup> Under these conditions, the zero estimator is asymptotically optimal relative to any distribution  $F \in \mathcal{F}$ .

This theorem demonstrates that, unlike the Bayes predictor, which relies on prior knowledge  $F$  and is thus impractical, the zero predictor achieves the same asymptotic optimal prediction risk without requiring such information, making it both feasible and optimal.

The zero estimator can be considered as a special case of both Ridge and Lasso when a sufficiently large tuning parameter is chosen. Given this perspective, and as implied from

---

<sup>7</sup>The definition of asymptotic optimality is provided in terms of a ratio to accommodate more general scenarios where  $\mathcal{R}(F)$  varies with the sample size  $n$ .

<sup>8</sup>The Gaussian assumption on  $\varepsilon$  is only used to facilitate considerations of optimality, which is a standard assumption in the empirical Bayes literature, e.g., Jiang and Zhang (2009). While this assumption motivates our characterization of weak signal regimes, it is not utilized in follow-up analysis of the estimators.

Theorem 1, the relative Bayes risk of the optimal Ridge and Lasso compared to the zero estimator asymptotically approaches one. This suggests that under weak signal conditions, merely comparing their Bayes risk ratios is insufficient to distinguish among these estimators.

As such, we shift our attention to the relative prediction error between any estimator  $\hat{\beta}$  and the zero estimator, defined as follows, in absolute difference rather than their ratio, in the spirit of Bayesian regret. To ensure a meaningful scale in the limit, we multiply the relative error by  $pn^{-1}\tau^{-2}$ , and adopt the following metric for comparison:

$$\Delta(\hat{\beta}) = pn^{-1}\tau^{-2}(\|\Sigma_2^{1/2}(\hat{\beta} - \beta_0)\|^2 - \|\Sigma_2^{1/2}\beta_0\|^2). \quad (8)$$

Based on this definition, if  $\Delta(\hat{\beta}) > 0$  holds with probability approaching one, it indicates that the estimator  $\hat{\beta}$  exhibits inferior prediction performance compared to zero. Conversely, if  $\Delta(\hat{\beta}) < 0$  holds with probability approaching one,  $\hat{\beta}$  outperforms zero.

Before we proceed to present our main results in the following section, we need to provide technical conditions governing the limiting behavior of  $\Sigma_1$ ,  $\Sigma_2$ , and  $\Sigma_\varepsilon$ :

**Assumption 5.** *For matrices  $\Sigma_1$ ,  $\Sigma_2$  and  $\Sigma_\varepsilon$ :*

$$\frac{1}{n} \text{Tr}(\Sigma_1) = 1 + O(n^{-1/2}), \quad \frac{1}{p} \text{Tr}(\Sigma_2) = \sigma_x^2 + O(p^{-1/2}), \quad \frac{1}{n} \text{Tr}(\Sigma_\varepsilon) = \sigma_\varepsilon^2 + O(n^{-1/2}).$$

*Additionally, there exist constants  $\theta_1$  to  $\theta_4$  such that*

$$\begin{aligned} \frac{1}{n} \text{Tr}(\Sigma_\varepsilon \Sigma_1) &= \sigma_\varepsilon^2 \theta_1 + o(n\tau/p), \\ \frac{1}{p} \text{Tr}(\Sigma_2^2) &= \sigma_x^4 \theta_2 + o(1), \quad \frac{1}{n} \text{Tr}(\Sigma_\varepsilon \Sigma_1^2) = \sigma_\varepsilon^2 \theta_3 + o(n/p), \quad \frac{1}{n} \text{Tr}(\Sigma_1^2) = \theta_4 + o(n/p). \end{aligned}$$

As  $\Sigma_1$ ,  $\Sigma_2$ , and  $\Sigma_\varepsilon$  are positive definite, all of these constants  $\theta_i$ ,  $i = 1, 2, 3$ , and 4, are positive. The condition concerning  $\Sigma_2$  can be verified through a more primitive condition: the existence of the limit of  $\Sigma_2$ 's empirical spectral distribution (see [Dobriban and Wager \(2018\)](#)). In situations where all three matrices reduce to identity matrices, a common setting in the literature on precise error analysis, Assumption 5 holds trivially.

## 2.5 Analysis of the Ridge Predictor

In this section, we present the results of Ridge in the context of weak signals. We begin by presenting the relative error of Ridge for any tuning parameter value:

**Theorem 2.** *Assuming that Assumptions 1–5 hold, and setting  $\lambda_n = \tau^{-1}\lambda$ , we establish:*

$$\Delta(\hat{\beta}_r(\lambda_n)) \xrightarrow{\text{P}} \alpha^* := 2\theta_2\sigma_x^4 \left( \frac{\sigma_\varepsilon^2\theta_1}{2\lambda^2} - \frac{\sigma_\beta^2}{\lambda} \right).$$

This theorem yields several important findings. First, by minimizing  $\alpha^*$  with respect to  $\lambda$ , we can determine the optimal tuning parameter value:  $\lambda_n^{\text{opt}} = \tau^{-1}\sigma_\varepsilon^2\theta_1/\sigma_\beta^2$ . Furthermore, with the optimal tuning parameter in place,  $\alpha^*$  is negative, indicating that Ridge can effectively learn weak signals when the tuning parameter is chosen appropriately.

Second, when we set  $\lambda \rightarrow \infty$  (equivalently,  $\tau\lambda_n \rightarrow \infty$ ), the value of  $\alpha^*$  converges to zero. This outcome is expected, as the use of a large tuning parameter makes the predictor’s performance increasingly resemble that of the zero predictor. Nevertheless, it is noteworthy that  $\alpha^* \rightarrow 0^-$ ; in other words, as  $\lambda$  increases, Ridge consistently outperforms the zero predictor until it gradually becomes indistinguishable from it in the limit.

Third, as  $\lambda \rightarrow 0$ , in which case  $\lambda_n = o(\tau^{-1})$ ,  $\alpha^*$  approaches positive infinity. This indicates that Ridge’s performance deteriorates to the point where the Ridgeless predictor (corresponding to  $\lambda = 0$ ) is surpassed by the zero predictor. This is a significant departure from the strong signal setup in which Ridgeless can still outperform the zero predictor, as shown by [Hastie et al. \(2022\)](#). Notably, Ridgeless, defined by  $\lambda = 0$ , is not devoid of regularization; it applies implicit regularization by minimizing the  $\ell_2$ -norm of the interpolator. This form of regularization allows Ridgeless to effectively control variance when the number of predictors  $p$  exceeds the sample size  $n$ , ensuring desirable performance in strong signal scenarios. However, under weak signals this implicit regularization fails to adequately control variance, causing  $\|\hat{\beta}\|$  to be much larger than  $\|\beta_0\|$ , leading to poor performance of  $\hat{\beta}$ .

Furthermore, given that Ridgeless is the interpolator that minimizes  $\|\hat{\beta}\|$ , all linear interpolators, including, for instance, the one that minimizes the  $\ell_1$ -norm, result in even larger  $\|\hat{\beta}\|$ . Therefore, these interpolators also fail to outperform zero in contexts with weak signals.

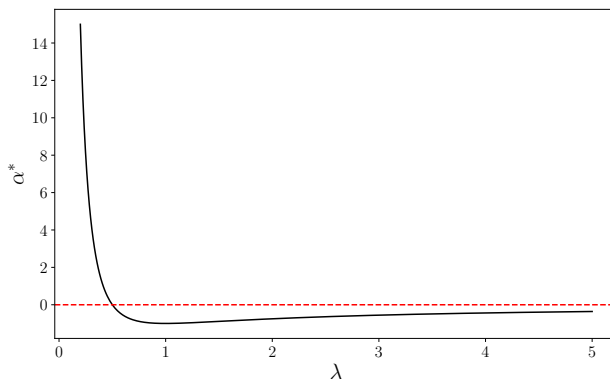
Figure 3 provides an illustrative example of the relationship between the relative error of Ridge and the tuning parameter  $\lambda$ , showcasing the theoretical insights we have discussed. Corollary 1 below summarizes the result on the Ridgeless estimator:

**Corollary 1.** *Under the same assumptions as in Theorem 2, the Ridgeless estimator, defined by (4), satisfies:*

$$\Delta(\hat{\beta}_r(0)) \xrightarrow{\text{P}} \infty.$$

Given Ridge’s ability to effectively learn weak signals with an appropriately tuned parameter, the data-dependent selection of this parameter becomes crucial. A paradigmatic

Figure 3: Ridge vs. Zero Predictor’s Relative Precise Error



Note: In this plot, the black curve represents the probability limit of  $\Delta(\hat{\beta}_r(\lambda_n))$ , denoted as  $\alpha^*$ , as a function of the tuning parameter  $\lambda$ , defined in Theorem 2, in the context of weak signals. To create this plot, we set all parameters  $(c_0, \theta_1, \theta_2, \sigma_x, \sigma_\beta, \sigma_\varepsilon)$  to one for simplicity.

approach for this purpose is  $K$ -fold CV.

In cases where the signals are strong, [Hastie et al. \(2022\)](#) demonstrate the effectiveness of CV for Ridge. Specifically, the cross-validated tuning parameter converges in probability to the optimal value within a pre-specified interval. The fact that this optimal value lies in some known interval simplifies the derivation of the theoretical properties of CV. In scenarios with weak signals, however, the optimal tuning parameter tends to diverge as the sample size increases. The rate of divergence depends on the unknown strength of the weak signal,  $\tau$ . As we show next, CV remains a valid and useful tool in this case. To narrow our focus to the matter of weak signals without delving into a complicated CV procedure, we consider the case where both  $\Sigma_\varepsilon$  and  $\Sigma_1$  are identity matrices. This assumption of no temporal dependence in the data facilitates a more straightforward CV procedure for i.i.d. data.

To determine the optimal tuning parameter using  $K$ -fold CV, denoted as  $\hat{\lambda}^{K-CV}$ , the rows of the design matrix  $X$  are partitioned into  $K$  distinct subsets, labeled as  $X_{(1)}, \dots, X_{(K)}$ . For each  $i \in \{1, \dots, K\}$ , the submatrix  $X_{(-i)}$  is formed by excluding the rows corresponding to  $X_{(i)}$ . Similarly, the associated subvectors are  $y_{(i)}$ ,  $\varepsilon_{(i)}$ ,  $y_{(-i)}$ , and  $\varepsilon_{(-i)}$ . The Ridge estimator for each fold,  $\hat{\beta}_r^i(\lambda_n)$ , is defined as the solution to the optimization problem for  $i = 1, \dots, K$ :

$$\hat{\beta}_r^i(\lambda_n) = \arg \min_{\beta} \left\{ \frac{1}{n} \|y_{(-i)} - X_{(-i)}\beta\|^2 + \frac{p\lambda_n}{n} \|\beta\|^2 \right\}.$$

Consequently, the tuning parameter selected by  $K$ -fold CV is given by



$$\hat{\lambda}_n^{K-CV} = \arg \min_{\lambda_n \in [\epsilon, \infty)} \frac{1}{n} \sum_{i=1}^K \|y_{(i)} - X_{(i)} \hat{\beta}_r^i(\lambda_n)\|^2,$$

where  $\epsilon > 0$  is an arbitrary small constant. The following theorem provides its justification:

**Theorem 3.** *Under the same assumptions as in Theorem 2, if we also assume that  $\Sigma_1 = \mathbb{I}$ ,  $\Sigma_\epsilon = \sigma_\epsilon^2 \mathbb{I}$ ,  $\epsilon$  follows a sub-exponential distribution, and that  $q^{-1} \tau^{-1} n^{-1/2} \log(p)$ ,  $q^{-1/2} \tau^{-3/2} n^{-1/2} \log(p) \rightarrow 0$ , then we can establish that:*

$$\tau \hat{\lambda}_n^{K-CV} \xrightarrow{P} \lambda^{opt} = \sigma_\epsilon^2 / \sigma_\beta^2 \quad \text{and} \quad \Delta(\hat{\beta}_r(\hat{\lambda}_n^{K-CV})) - \Delta(\hat{\beta}_r(\lambda_n^{opt})) \xrightarrow{P} 0.$$

This theorem justifies the use of  $\hat{\lambda}_n^{K-CV}$  as an approximation for the optimal tuning parameter  $\lambda_n^{opt} = \lambda^{opt} / \tau$  ( $\theta_1 = 1$  in this case) for Ridge. Importantly, this result does not require prior knowledge of  $\tau$ , making the CV approach directly applicable in practical scenarios. The additional constraints on  $q$  become relevant only when  $q$  vanishes; otherwise, they naturally follow from Assumption 4. These conditions ensure uniform convergence across the spectrum of tuning parameter values, a prerequisite for the results of Theorem 3.

With our analysis of Ridge concluded, we will now turn our attention to Lasso.

## 2.6 Analysis of the Lasso Predictor

Unlike Ridge, the analysis of Lasso is more intricate, primarily because it lacks a closed-form formula. In the special case where  $\Sigma_1$ ,  $\Sigma_2$ , and  $\Sigma_\epsilon$  are identity matrices, several studies, including [Bayati and Montanari \(2012\)](#) and [Thrampoulidis et al. \(2018\)](#), have established Lasso's precise error given by (7). Additionally, based on (7), [Wang et al. \(2020\)](#) conducted a small-signal Taylor expansion of  $\alpha^*$  with respect to  $\sigma_\beta^2$ , which affects  $\alpha^*$  through the prior distribution  $F$ . They concluded that the optimal Lasso estimator fails to outperform optimal Ridge.<sup>9</sup> In the general case we consider, pinpointing the exact precise error appears a daunting task. Instead, we seek probability bounds of the limit. This turns out sufficient for us to conclude that Lasso cannot outperform zero for all values of its tuning parameter in the context of weak signals. The next theorem summarizes our main findings:

**Theorem 4.** *Assume that Assumptions 1–5 are satisfied and the tuning parameter  $\lambda_n$  is chosen such that the following equation holds for some  $C_\lambda > 0$ :*

---

<sup>9</sup>Their analysis does not address the scenario of Lasso with an arbitrary tuning parameter, nor does it elucidate its relative performance compared to zero.

$$pn^{-2}\tau^{-2}\mathbb{E}_{U\sim\mathcal{N}(0,\Sigma_2)}\left\|\left(2\sigma_\varepsilon\sqrt{\theta_1}|U|-\lambda_n\right)_+\right\|^2=C_\lambda.^{10} \quad (9)$$

Then, with probability approaching one, we have  $c_\alpha \leq \Delta(\hat{\beta}_l(\lambda_n)) \leq C_\alpha$ , where  $c_\alpha$  and  $C_\alpha$  are the solutions to the following equation in terms of  $x$ :

$$x - \sqrt{\frac{2C_\lambda}{c_2}}x = -\frac{C_\lambda}{100C_2}, \quad (10)$$

where  $c_2$  and  $C_2$  are constants defined in Assumption 1.

Equation (9) implicitly determines the rate at which  $\lambda_n$  diverges to infinity. For any fixed  $C_\lambda > 0$ , we can solve for the tuning parameter  $\lambda_n$  from (9), and derive the upper and lower bounds,  $C_\alpha$  and  $c_\alpha$ , from equation (10). Furthermore, equation (10) directly implies that  $C_\alpha$  and  $c_\alpha$  are non-negative, indicating that Lasso does not outperform the zero predictor for any given tuning parameter value in the context of weak signals.

Moreover, as  $\lambda_n$  approaches zero,  $C_\lambda$  diverges to infinity, leading to a simultaneous divergence of both  $c_\alpha$  and  $C_\alpha$ . This suggests that Lasso behaves increasingly worse compared to the zero predictor. Conversely, a larger tuning parameter  $\lambda_n$  causes  $C_\lambda$  to converge to zero from the positive side. As a result, both  $c_\alpha$  and  $C_\alpha$  converge to zero while remaining non-negative. This implies that Lasso’s performance improves but remains inferior to the zero predictor, until they become equivalent in the limit. Figure 4 illustrates upper and lower bounds for the relative error of Lasso versus the zero predictor across different  $\lambda_n$  values.

Lasso struggles with weak signals because it has difficulty distinguishing true signals from spurious ones. Its inability to select true weak signals does not significantly impact its performance compared to the zero predictor, which ignores these signals entirely. The main issue with Lasso is its ineffectiveness in penalizing irrelevant signals. While increasing the tuning parameter might help, our theory indicates that Lasso only reaches the zero predictor’s effectiveness when the penalty is large enough to essentially disregard all signals.

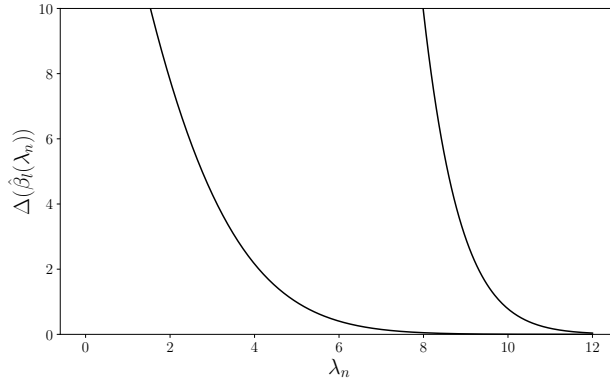
Given Lasso’s limitations with weak signals, the elastic net—merging the  $\ell_1$  and  $\ell_2$  norms in its penalty—also underperforms compared to Ridge, in settings involving weak signals.<sup>11</sup>

In predictive regressions with many covariates, Lasso is often viewed as the successor to OLS. However, our analysis reveals a critical caveat: Lasso, regardless of tuning parameter choice, underperforms a zero predictor in weak-signal settings. This has important implica-

<sup>10</sup>When applied to a vector,  $|\cdot|$  and  $(\cdot)_+$  represent element-wise operations.

<sup>11</sup>Although a formal justification for this observation can be provided in the setting of  $\Sigma_2 = \mathbb{I}$ , we omit it here due to space constraints.

Figure 4: Lasso vs. Zero Predictor: Relative Precise Error Bounds



Note: In this plot, the black curves represent the lower and upper probability bounds on  $\Delta(\hat{\beta}_l(\lambda_n))$ , i.e.,  $c_\alpha$  and  $C_\alpha$ , as a function of the tuning parameter  $\lambda_n$  in the context of weak signals. In this setup, we fix  $n = p = 2,000$ . We assign the elements of  $b_0$  to follow a standard Gaussian distribution and set  $\Sigma_2 = \mathbb{I}$ . As in Figure 3, we set all parameters ( $c_0, \theta_1, \theta_2, \sigma_x, \sigma_\beta, \sigma_\varepsilon$ ) to one. Finally, we select  $\tau = 0.001$ , which results in a population  $R^2$  around 0.1%.

tions for economics and finance, where large-scale regressions with low signal-to-noise ratios are increasingly common. Instead, our results favor Ridge in such scenarios, emphasizing the need to assess signal strength when selecting regularization techniques.

## 2.7 Assessing Signal-to-Noise Ratio

In line with this perspective, we delve into the assessment of the signal-to-noise ratio, a measure that can offer insights into the viability of different machine learning techniques. Our preceding analyses provide indirect guidance in this regard. Specifically, if Lasso underperforms the zero predictor, it implies a potential issue with the data’s signal strength.

A more conventional and direct approach to evaluating the signal-to-noise ratio is through  $R^2$ . However, in-sample  $R^2$  is prone to overfitting, and as such, out-of-sample  $R^2$  is commonly used in machine learning. This metric essentially involves the comparison of mean-squared errors between two predictors. For our specific application, we have chosen to use zero as the reference predictor and define this metric for a given estimator  $\hat{\beta}$  as follows:

$$R_{\text{OOS}}^2(\hat{\beta}) = 1 - \frac{\sum_{i \in \text{OOS}} (y_i - X_i \hat{\beta})^2}{\sum_{i \in \text{OOS}} y_i^2}, \quad (11)$$

where “OOS” represents the out-of-sample data.

Since a model’s predictive performance hinges on the signal-to-noise ratio, this metric is a natural choice for evaluating the signal strength inherent in the DGP. In strong-signals scenarios, the out-of-sample  $R^2$  consistently estimates the population  $R^2$ , irrespective of the specific estimator  $\hat{\beta}$ , provided it is consistent with respect to  $\beta_0$ , i.e.,  $\|\hat{\beta} - \beta_0\| = o_P(1)$ . However, in weak-signal settings, the outcome critically depends on the estimator chosen, beyond the signal-to-noise ratio itself. As an illustration, while both Lasso and Ridge are consistent as their prediction errors asymptotically diminish, the out-of-sample  $R^2$  for Lasso can become non-positive, indicating either no improvement or underperformance relative to the zero predictor, as demonstrated in Theorem 4.

In the context of weak signals, the following proposition provides theoretical support for the relevance of optimal Ridge’s  $R_{\text{os}}^2$  in assessing the signal-to-noise ratio in the data.

**Proposition 1.** *Under the same assumptions as Theorem 3, and assuming that the out-of-sample data follows the same DGP as the in-sample data, if  $n_{\text{os}}p^{-2}n^2\tau^2 \rightarrow \infty$ , where  $n_{\text{os}}$  is the size of the out-of-sample data, then for the optimal Ridge predictor, it holds that*

$$R_{\text{os}}^2(\hat{\beta}_r(\lambda_n^{\text{opt}})) = p^{-1}n\theta_2(R^2)^2(1 + o_P(1)),$$

where  $R^2$  denotes the population  $R$ -squared, given by  $\tau\sigma_x^2\sigma_\beta^2/(\tau\sigma_x^2\sigma_\beta^2 + \sigma_\varepsilon^2)$  in this context.

Interestingly, when  $n_{\text{os}}$  is sufficiently large, to the extent that the estimation error in  $R_{\text{os}}^2$  does not mask the performance differential between the optimal Ridge and the zero predictor,  $R_{\text{os}}^2$  is approximately proportional to  $(R^2)^2$ . While it does not exactly mirror  $R^2$ ,  $R_{\text{os}}^2$  still serves as an indicator of signal strength. The reason for their discrepancy is that in the weak signal case, the numerator of the  $R_{\text{os}}^2$ —which reflects the relative prediction error between the two predictors—decreases more rapidly than the numerator of  $R^2$ . Therefore, the numerator of  $R_{\text{os}}^2$  only provides a higher-order characterization of signal strength.

## 2.8 Mixed Signal Strengths and Alternative Benchmarks

In the preceding sections, our analysis primarily focuses on scenarios where all signals are weak, leading us to consider the zero predictor as our natural benchmark. This section, however, expands our analysis to include models where potentially strong signals serve as benchmarks. Consider another DGP:

$$y = W\gamma_0 + X\beta_0 + \varepsilon, \tag{12}$$

where  $W \in \mathbb{R}^{n \times d}$  represents a predefined set of covariates. These covariates include potentially strong signals and form the basis of the benchmark model. We allow the dimension  $d$  to increase to  $\infty$ , however, it does so at a slower rate compared to  $n$ , ensuring that OLS of  $y$  against  $W$  remains a viable method for estimation.

In many cases,  $W$  could simply be a vector of ones, allowing us to remain agnostic about the magnitude of the regression's intercept. In our empirical analysis,  $W$  can be motivated from economic theory, whose impact on the response variable is of central interest. Alternatively,  $W$  can encompass lagged values of  $y$ , thereby facilitating the inclusion of temporal dependence in the benchmark model. This setup is particularly relevant when using an autoregressive model as a benchmark for forecasting economic variables. Exploring the possibility of a data-driven selection of  $W$  is an intriguing direction for future research.

Building on these considerations, our focus now shifts to evaluating and comparing the performance against the OLS benchmark with covariates in  $W$ . In this context, the OLS benchmark predictor,  $\hat{y}_b^{new}$ , for a new observation  $(w^{new}, x^{new})$  is defined as follows:

$$\hat{y}_b^{new} = (w^{new})^\top \hat{\gamma}, \quad \text{where} \quad \hat{\gamma} = (W^\top W)^{-1} W^\top y. \quad (13)$$

The inclusion of  $W$  leads us to explore the following Ridge estimator with regularization only imposed on coefficients of  $X$ :

$$\begin{aligned} \hat{\beta}(\lambda_n) &:= \arg \min_{\beta} \left\{ \min_{\gamma} \left( \frac{1}{n} \|y - W\gamma - X\beta\|^2 + \frac{p\lambda_n}{n} \|\beta\|^2 \right) \right\} \\ &= \arg \min_{\beta} \left\{ \frac{1}{n} \|\mathcal{M}_W y - \mathcal{M}_W X\beta\|^2 + \frac{p\lambda_n}{n} \|\beta\|^2 \right\}, \end{aligned} \quad (14)$$

where  $\mathcal{M}_W = \mathbb{I} - W(W^\top W)^{-1}W^\top$ . Consequently, the estimator for  $\gamma$  is thus given by

$$\hat{\gamma}(\lambda_n) = (W^\top W)^{-1} W^\top (y - X\hat{\beta}(\lambda_n)). \quad (15)$$

The construction for Lasso is similar. Therefore, utilizing the estimated parameters  $(\hat{\beta}(\lambda_n), \hat{\gamma}(\lambda_n))$  we are able to formulate a predictor for  $y$  as

$$\hat{y}^{new} = (w^{new})^\top \hat{\gamma}(\lambda_n) + (x^{new})^\top \hat{\beta}(\lambda_n) = \hat{y}_b^{new} + (\hat{x}^{new})^\top \hat{\beta}(\lambda_n), \quad (16)$$

where  $\hat{x}^{new} = x^{new} - X^\top W(W^\top W)^{-1}w^{new}$ .

Notably, equation (16) illuminates the role of the second term,  $(\hat{x}^{new})^\top \hat{\beta}(\lambda_n)$ , in cap-

turing the contribution of weak signals relative to the OLS benchmark,  $\hat{y}_b^{new}$ . Moreover, a comparison with the previously analyzed zero-benchmark scenario reveals a key distinction: the modification of the response and covariates in equation (14). Here, the approach involves regressing  $\mathcal{M}_W y$  on  $\mathcal{M}_W X$ , effectively predicting the residuals of the benchmark model using covariates adjusted to eliminate dependence on  $W$ . While our earlier conclusions are likely still valid, the inclusion of generated covariates introduces an additional layer of statistical error that warrants careful examination. The forthcoming theorem will demonstrate that this extra error does not compromise our prior conclusions. Given this context and our prior comparative analysis, we focus on the optimal Ridge in this scenario for reason of space.

**Theorem 5.** *Suppose that  $X = W\eta_0 + U$ . Assume that the triplet  $(U, \beta_0, \varepsilon)$  follows the same distribution as  $(X, \beta_0, \varepsilon)$  in Theorem 2. Additionally, the matrix  $W$  is independent from  $U$ ,  $\beta_0$ , and  $\varepsilon$ . Each covariate within  $W$  is assumed to have a finite second moment. Furthermore, we assume that  $d = o(n^2 p^{-1} \tau)$ , and the eigenvalues of  $n^{-1} W^\top W$  are lower bounded by some positive constant. Given these assumptions, the predictor,  $\hat{y}^{new}$ , as defined in (16) and based on the Ridge estimator from (14) with  $\lambda_n = \tau^{-1} \lambda$ , and the benchmark predictor  $\hat{y}_b^{new}$  from (13), satisfy the following:*

$$pn^{-1} \tau^2 \left( \mathbb{E}_F [(\hat{y}^{new} - y^{new})^2 | \mathcal{I}] - \mathbb{E}_F [(\hat{y}_b^{new} - y^{new})^2 | \mathcal{I}] \right) \xrightarrow{P} \alpha^*, \quad (17)$$

where  $\mathcal{I}$  denotes the information set generated by  $(W, X, y, \gamma_0, \beta_0)$ ,  $\alpha^*$  is defined in Theorem 2, and the tuple  $(y^{new}, w^{new}, x^{new})$  satisfies (12).

This result shows that a Ridge-augmented model outperforms the benchmark model alone, aligning with our earlier conclusion that Ridge can effectively leverage weak signals.

## 2.9 Ridge vs. Lasso in Extremely Sparse Settings

This section clarifies that our asymptotic restrictions are mainly technical. Even in extreme sparsity beyond our main analysis, Lasso fails to outperform zero when signals are sufficiently weak, while Ridge retains a non-negligible probability of learning them. However, the setting is stylized, as the proof relies on a distinct approach using Lasso's closed-form solution.

Consider the Gaussian sequence model where  $y = \beta_0 + \varepsilon$  with  $\varepsilon \sim \mathcal{N}(0, \mathbb{I})$ . For this model,  $X = \mathbb{I}$  and  $n = p$ . We let  $\beta_0 = (\sqrt{n\tau}, 0, \dots, 0)^\top$ , so that  $s = \|\beta_0\|_0 = 1$  and  $R^2 = \frac{\|X\beta_0\|^2}{\|X\beta_0 + \varepsilon\|^2} \asymp \tau \rightarrow 0$ . This represents the most extreme form of sparsity, with only one true signal. Moreover, we relax the  $n^{-1/3}$  lower bound restriction on  $\tau$ .

**Proposition 2.** Assume  $\tau \leq n^{-1} \log(n)/100$ . There exists  $n_0 > 0$  such that  $n > n_0$ , Lasso satisfies

$$P\left(\|\hat{\beta}_l(\lambda) - \beta_0\|^2 - \|\beta_0\|^2 \geq 0, \forall \lambda \geq 0\right) \geq 1 - n^{-1/2}, \quad (18)$$

while for Ridge,

$$P\left(\exists \lambda > 1 \text{ s.t. } \|\hat{\beta}_r(\lambda) - \beta_0\|^2 - \|\beta_0\|^2 < 0\right) \geq 0.5. \quad (19)$$

This result highlights that in extreme sparsity, Ridge can still learn weak signals, whereas Lasso remains ineffective regardless of tuning.

### 3 Monte Carlo Simulations

In this section, we conduct simulation experiments to assess the finite sample relevance of our asymptotic theory. We begin by examining Ridge and Lasso in a linear model setup.

#### 3.1 Ridge and Lasso for Linear Models

We now detail the DGP specified by (1) for the first simulation exercise. We set  $(\Sigma_1)_{ij} = 2^{-|i-j|}$  for  $1 \leq i, j \leq n$ . We construct  $\Sigma_\varepsilon$  as a diagonal matrix with i.i.d. entries sampled from the uniform distribution  $U(0.5, 1.5)$ . The eigenvalues of  $\Sigma_2$  are also simulated from  $U(0.5, 1.5)$ , with corresponding eigenvectors from a randomly generated orthogonal matrix. These components form  $\Sigma_2$ , which, along with other matrices, remains fixed throughout the simulations. By direct calculations, we have  $\theta_1 = 1$ ,  $\theta_2 = 13/12$ , and  $\theta_3 = \theta_4 = 5/3$ .

We experiment with  $n = 500$  and  $p = 300$ , dimensions that align closely with the first microeconomic example studied below. For each simulated sample, we construct  $\beta_0$  as  $\sqrt{p^{-1}\tau}b_0$ , with  $b_0$  drawn from a spike-and-slab distribution:  $(1 - q)\delta_0 + q\mathcal{N}(0, q^{-1}\sigma_\beta^2)$ . Here,  $\delta_0$  represents the Dirac delta function, and we set  $\sigma_\beta^2 = 1$ . The error term  $\varepsilon$  is sampled from  $\mathcal{N}(0, 1)$ , while the design matrix  $X$  is drawn from  $\mathcal{N}(0, 1)$ , then transformed via pre-multiplication by  $\Sigma_1^{1/2}$  and post-multiplication by  $\Sigma_2^{1/2}$ . We consider two cases for the sparsity parameter,  $q = 0.2$  and  $q = 0.8$ . To represent weak and strong signal scenarios, we calibrate two values of  $\tau$  to achieve  $R^2 = 5\%$  and  $50\%$ , respectively. The parameters  $q$  and  $\tau$  (through  $R^2$ ) are varied as they are critical to the asymptotic performance, as highlighted in Assumption 4. A total of 1,000 Monte Carlo repetitions are conducted.

For each simulated sample, we compute the relative prediction error,  $\Delta(\hat{\beta}(\hat{\lambda}_n^{K-CV}))$ , as

defined in equation (8), with the tuning parameter for each method selected via 10-fold CV. The histograms of these errors are presented in Figure 5. Additionally, Table 1 provides summary statistics, including quantiles and the proportion of values classified as “zeros” (where “zero” is defined as being within machine precision).

The histograms reveal notable differences in the performance of Ridge and Lasso when  $R^2$  is low, while showing that both methods outperform the zero predictor when  $R^2$  is high. Ridge displays a clear probability mass on the negative side of the error distribution, even in weak signal settings, and its performance is largely unaffected by changes in sparsity levels.

In contrast, Lasso struggles to capture weak signals, as evidenced by a substantial probability mass on the positive side of the y-axis. While increasing sparsity (lowering  $q$ ) leads to modest improvement in Lasso’s performance, it still lags behind Ridge overall. This pattern is supported by Table 1, which presents quantiles of these distributions. In finite samples with high sparsity, some probability mass for Lasso can fall on the negative side. Nevertheless, Lasso also shows a heavier probability mass at zero compared to Ridge, consistent with the theoretical prediction that optimal Lasso collapses to zero in weak-signal settings.

The above results use tuning parameters selected via 10-fold CV. To validate our theoretical findings independently of tuning parameter selection, Appendix A.1 presents experiments with fixed tuning parameters. Additionally, Further simulations in Appendix A.2 support our theoretical predictions regarding the  $R_{\text{os}}^2$  of the optimal Ridge. Appendix A.3 presents evidence indicating that Type I error primarily influences Lasso’s performance relative to the zero estimator. As  $\lambda$  increases, Type I error diminishes, leading to an improvement in Lasso’s performance, which ultimately becomes identical to that of the zero estimator, while Type II error persists at a high level.

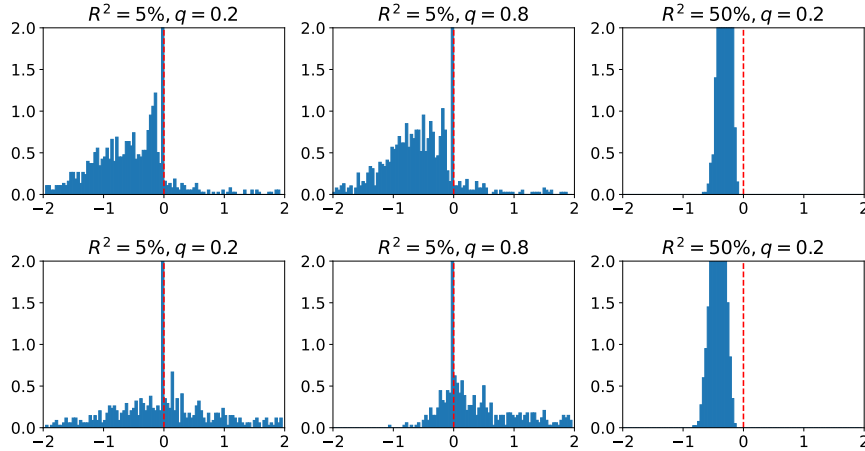
To examine the robustness of our theory in cases of extreme sparsity, Appendix A.4 examines the effect of further reducing Lasso’s sparsity level to  $q = 0.1, 0.05,$  and  $0.02$ . The results reveal that for each sparsity level, as  $R^2$  decreases, Ridge’s performance improves, whereas Lasso’s performance deteriorates. This suggests that even under extreme sparsity conditions, the relative performance is dictated by the strength of the signal. Ridge continues to outperform both the zero predictor and Lasso when the signal is sufficiently weak.

### 3.2 Advanced Machine Learning Methods for Nonlinear Models

In this section, we extend our investigation to nonlinear machine learning methodologies, including RF, GBRT, and NNs, through simulation experiments. While providing a precise theoretical analysis of errors for these algorithms remains challenging—and this part there-



Figure 5: Simulation Results for Ridge and Lasso in Linear DGPs



Note: The histograms depict the relative prediction error  $\Delta(\hat{\beta}_r(\hat{\lambda}_n^{K-CV}))$  (top) and  $\Delta(\hat{\beta}_l(\hat{\lambda}_n^{K-CV}))$  (bottom) following equation (8) across 1,000 Monte Carlo samples. We analyze three setups of  $(R^2, q)$ .

Table 1: Summary Statistics for Ridge and Lasso in Linear DGPs

$q$	$R^2$ (%)	Lasso				Ridge			
		Q1	Q2	Q3	#Zero	Q1	Q2	Q3	#Zero
0.2	5%	-0.127	0.000	0.521	360	-0.992	-0.501	-0.129	97
0.8	5%	0.000	0.000	0.975	400	-0.918	-0.541	-0.169	88
0.2	50%	-0.508	-0.414	-0.331	0	-0.375	-0.300	-0.233	0

Note: The table illustrate the summary statistics (quantiles and the percentage of zeros) of relative prediction error  $\Delta(\hat{\beta}(\hat{\lambda}_n^{K-CV}))$ , for Ridge and Lasso, based on 1,000 Monte Carlo samples.

fore involves some degree of speculation—we draw on insights from linear models to interpret and contextualize our simulation findings.

We simulate the following DGP, expressed explicitly in element-wise form:

$$y_i = \sum_{j=1}^p \beta_{0,j} f(Z_{ij}) + \varepsilon_i, \quad i = 1, \dots, n, \quad (20)$$

where  $y_i$  denotes the  $i$ th observation of the response variable,  $\beta_{0,j}$  represents the coefficient associated with a function  $f(\cdot)$  of the predictor variable  $Z_{ij}$ . We adopt the following procedure for simulating this model:  $Z_{ij}$  is generated by applying an inverse transform to  $X_{ij}$ , previously simulated in Section 3.1. Specifically,  $Z_{ij} = f^{-1}(X_{ij})$ , where  $X$  is constructed using the same DGP as before.<sup>12</sup> Additionally, both the coefficients  $\beta_0$  and the error term

<sup>12</sup>We present results for  $f(x) = \tan(x)$ . The results are nearly identical to those obtained with cubic or

$\varepsilon_i$  are drawn from the same baseline DGP. This approach guarantees the replication of the exact simulation results observed when regressing  $y$  on  $X$ . However, the focus now shifts to predicting  $y$  using nonlinear models of  $Z$  without prior knowledge of  $f(\cdot)$ . The effective signal-to-noise ratio decreases due to the added complexity of learning an unknown function.

### 3.2.1 Simulations with Tree Algorithms

Tree algorithms are essential in machine learning for handling complex DGPs with discrete variables, nonlinearities, and intricate interactions. However, single tree models often underperform, prompting the use of ensemble methods to enhance predictions.

Two popular ensemble techniques are RF and GBRT. RF uses bagging, where multiple trees are trained independently on bootstrap samples, and their predictions are averaged to improve performance. In contrast, GBRT employs boosting, iteratively fitting residuals from prior trees to build a strong ensemble from weak learners.

Since trees are invariant to monotonic transformations, it suffices to report their prediction results for the linear DGP, as these are identical to those for the nonlinear DGP under consideration. However, tree methods may underperform Ridge and Lasso, partly due to an additional approximation error from using piecewise constant functions to approximate the linear DGP. Therefore, the primary focus here should not be on comparing tree methods with linear models but rather on assessing the effectiveness of different ensemble techniques in capturing weak signals and comparing their performance to the zero predictor.

We repeat the experiments from Section 3.1, this time generating an additional set of  $n_{oos}$  out-of-sample observations to evaluate predictive performance.<sup>13</sup> As illustrated in Figure 6, RF demonstrates its ability to learn weak signals when  $R^2 = 5\%$ , with more than half of the probability mass located to the left of the y-axis. In contrast, GBRT struggles at this signal strength level. Sparsity does not appear to significantly impact either method. Nonetheless, both methods markedly outperform the zero predictor as  $R^2$  increases to 50%.

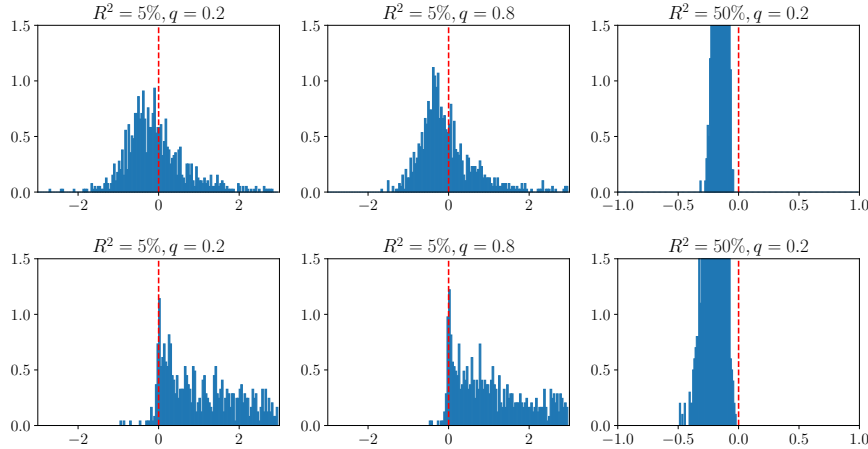
A possible explanation for GBRT’s performance pattern could be an inherent  $\ell_1$ -like regularization in its boosting approach. This conjecture draws support from [Efron et al. \(2004\)](#), which demonstrates a parallel between boosting and the Lasso path in linear regressions. To substantiate this conjecture, we examine the number of active variables (those with a non-zero importance score) from both tree methods under the benchmark case ( $q = 0.2$ ,  $R^2 = 5\%$ ). For RF, the average count of active variables is 300. This mirrors that of

---

hyperbolic sine functions, and are therefore omitted for brevity.

<sup>13</sup>We set  $n_{oos} \approx \lceil \tau^{-3} \rceil$  to meet the assumption in Proposition 1.

Figure 6: Simulation Results for RF and GBRT in Linear DGPs



Note: The histograms depict the relative prediction error,  $pn^{-1}\tau^{-2}n_{oos}^{-1}\sum_{i\in OOS}((y_i - \hat{y}_i)^2 - y_i^2)$ , from 1,000 Monte Carlo samples. We consider RF and GBRT in settings of  $n = 500$ ,  $q = 300$ ,  $n_{oos} = 10,000$ , across three  $(R^2, q)$  configurations. The red dashed line marks the y axis for reference.

Ridge. Conversely, GBRT demonstrates a significantly lower count of active variables, with an average of 27.9, aligning more with the variable selection feature of Lasso. Based on our theoretical findings that Ridge outperforms Lasso in weak signal scenarios, we can infer that RF is more adept than GBRT in settings characterized by low signal strengths.

### 3.2.2 Simulations with Neural Networks

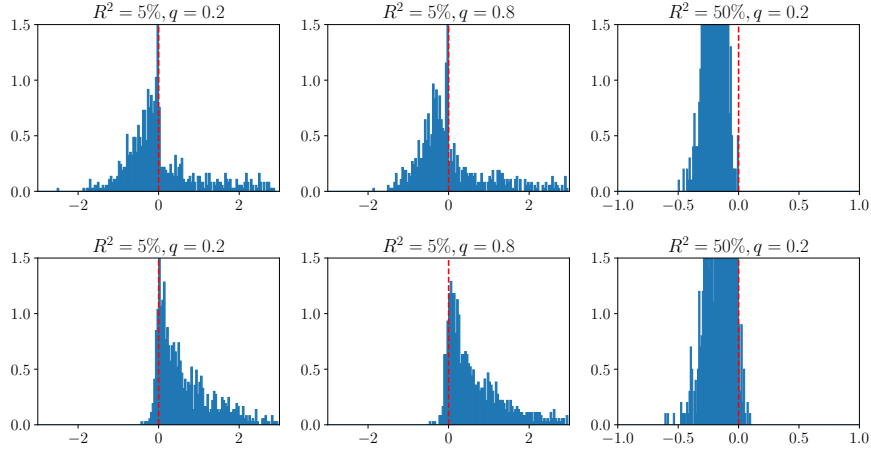
Next, we study fully-connected feed-forward NNs. Still, we fix  $n = 500$  and  $p = 300$ . Such parameters lead to an input layer in the NN configured with 300 neurons. Our chosen architecture features a single hidden layer, which includes 16 neurons.

The training process of these NNs often incorporates a sophisticated mix of optimization and regularization techniques crucial for enhancing performance.<sup>14</sup> To isolate and assess the impact of  $\ell_1$  and  $\ell_2$  regularization, we use plain stochastic gradient descent (SGD), deliberately avoiding other optimization techniques to minimize interference, albeit at the

<sup>14</sup>Key methods include stochastic gradient descent (SGD) with Adam (Kingma and Ba (2014)), which expedites the optimization process through an adaptive learning rate. Early stopping, as discussed in Goodfellow et al. (2016), is employed to prevent overfitting by halting training when validation performance starts to decline. Dropout (Srivastava et al. (2014)) is utilized for better generalization, achieved by randomly deactivating neurons. Batch normalization (Ioffe and Szegedy (2015)) aids in stabilizing the training process. Moreover, ensembling over various random seeds is implemented to reduce the variances in model outputs. Furthermore, the integration of  $\ell_1$  and  $\ell_2$  penalties with these techniques helps regulate the NN parameters.

expense of not fully exploiting the NN’s potential.<sup>15</sup>

Figure 7: Simulation Results for NNs in Nonlinear DGPs



Note: The histograms show the relative prediction error,  $pn^{-1}\tau^{-2}n_{oos}^{-1}\sum_{i\in OOS}((y_i - \hat{y}_i)^2 - y_i^2)$ , across 1,000 Monte Carlo samples, using  $\ell_2$  (top) and  $\ell_1$  (bottom) penalties. The settings include  $n = 500$ ,  $q = 300$ ,  $n_{oos} = 10,000$  for three  $(R^2, q)$  configurations. The red dashed line marks the y axis for reference.

The histograms in Figure 7 display the relative prediction errors of NNs. We observe that  $\ell_2$  regularization effectively leverages weak signals, as indicated by most of the probability mass in their histograms being on the negative side of the y-axis when  $R^2 = 5\%$ . In contrast,  $\ell_1$  regularization shows a notable decline in performance, which is consistent with our theoretical insights from linear models.

## 4 Empirical Analysis of Six Economic Datasets

In this section, we demonstrate the practical relevance of our theoretical insights by applying seven machine learning methods—Ridge, Lasso, OLS/Ridgeless, RF, GBRT, NNs with both  $\ell_1$  and  $\ell_2$  penalties—across six datasets from finance, macroeconomics, and microeconomics, with two datasets per field. We use five datasets similar to those in Giannone et al. (2022), updated with the latest data when feasible, and introduce an updated dataset from Gu et al. (2020) for our second finance example. Our empirical strategy differs notably from Giannone et al. (2022), who estimate a parametric model using a Spike-and-Slab prior within

<sup>15</sup>We avoid early stopping, as our simulations (not reported due to space constraints) reveal it functions similarly to  $\ell_2$ -regularization by shrinking parameter values towards their initial, smaller magnitudes.

a Bayesian framework. In contrast, our study, more aligned with [Gu et al. \(2020\)](#), focuses on a comparative analysis of these methods.

At the outset of each empirical exercise, we face a variety of decisions regarding our implementation strategy. These include setting the in-sample and out-of-sample periods, choosing between a rolling window or an expanding window approach, selecting a CV procedure, and deciding on covariate normalization.<sup>16</sup> We follow the frameworks established by [Giannone et al. \(2022\)](#) and [Gu et al. \(2020\)](#) to limit degrees of freedom, thereby improving the robustness, comparability, and reproducibility of our findings. In applying each machine learning method, choosing the right grid for tuning parameters is crucial, balancing performance optimization with computational efficiency. Finer and wider grids can potentially improve performance but also increase computational demands. Details on our model configuration and tuning parameter selection are provided in Online Appendix B.

Below we summarize the empirical findings from six distinct datasets, each analyzed and reported separately. The primary summary statistics,  $R_{\text{oos}}^2$ s, are collected in Table 2. Additionally, we include variable importance plots in Figure 8 as supplementary evidence to decode the performance of different methods. The notion of variable importance is not universally established and varying across different contexts. Our approach diverges from the well-known method associated with RF, originally presented in [Breiman \(2001\)](#). In our analysis, variable importance is quantified as the reduction in  $R_{\text{oos}}^2$  resulting from setting each variable, one at a time, to zero (its mean value post-normalization), with this metric normalized across all variables. For each method, the most significant variable, as per this definition, is assigned a value of one, and a color gradient is employed to visually represent the relative importance of each variable.

#### 4.1 Finance 1: Market Equity Premium

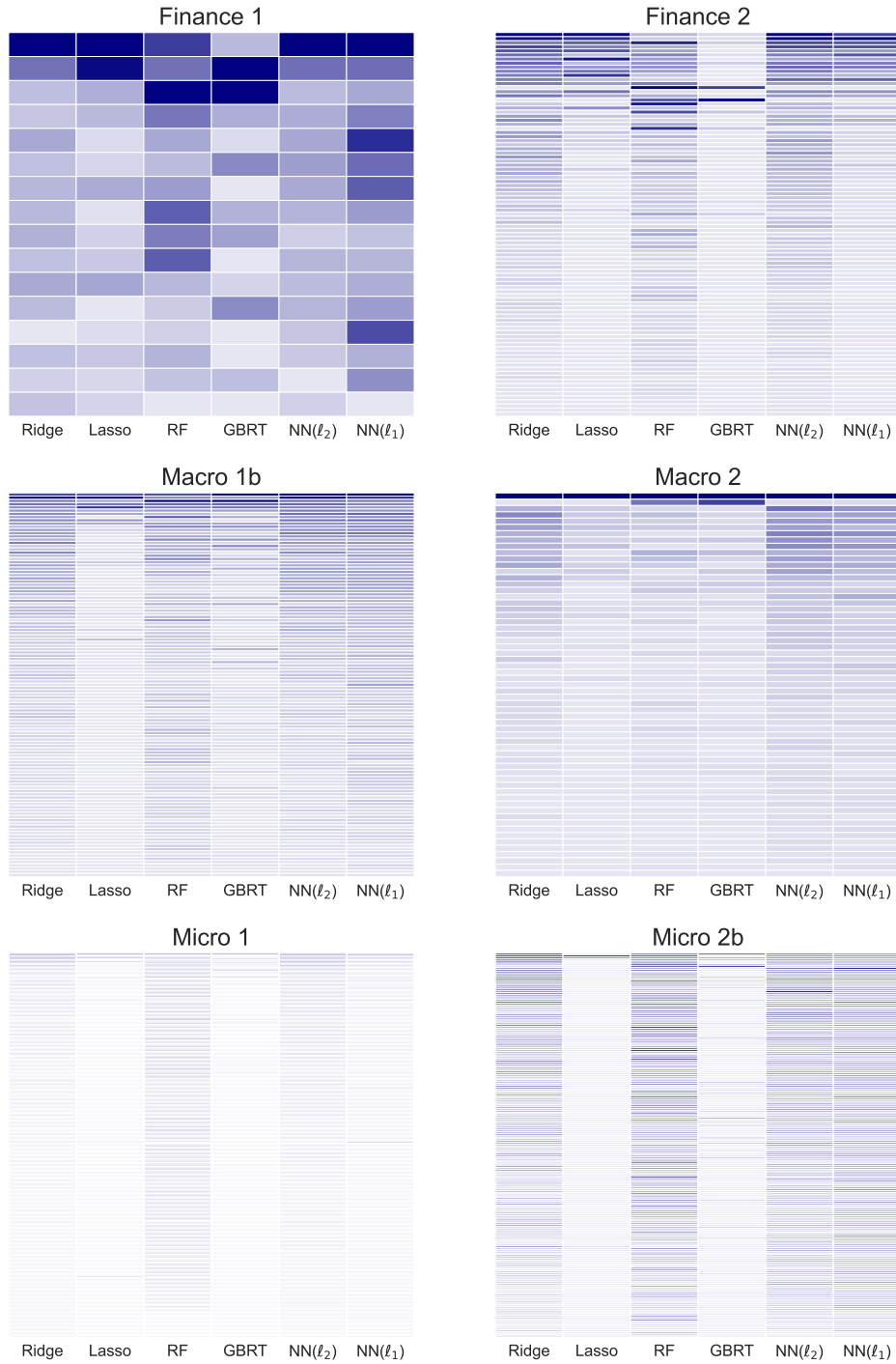
In the first analysis, we focus on predicting market equity returns using a dataset of financial and macroeconomic indicators compiled by [Welch and Goyal \(2007\)](#).<sup>17</sup> This dataset comprises 16 predictors and includes 74 annual observations, covering a period from 1948 to 2021. Despite [Welch and Goyal \(2007\)](#) reporting a consistently negative  $R_{\text{oos}}^2$  for this dataset, many other studies, such as those by [Campbell and Thompson \(2007\)](#), [Ferreira and Santa-Clara \(2011\)](#), [Rapach et al. \(2010\)](#), [Kelly and Pruitt \(2013\)](#), and [Kelly et al. \(2023\)](#),

---

<sup>16</sup>Normalizing covariates is essential before using machine learning methods, as it standardizes their scales, aids in regularization, and improves the convergence of optimization algorithms. To prevent forward-looking bias, normalization is performed using each covariate’s in-sample mean and standard deviation.

<sup>17</sup>The data, sourced from Amit Goyal’s website, was processed using [Giannone et al. \(2022\)](#)’s methodology.

Figure 8: Variable Importance Plots



Note: This figure illustrates the variable importance across six empirical studies, using color gradients to show the relative reductions in  $R^2_{\text{00s}}$  by each covariate. For the first example in Macroeconomics and the second example in Microeconomics, we only present the cases with a more complex benchmark model.

have developed forecasting strategies resulting in economically meaningful  $R_{\text{OOS}}^2$  values, which translate into significant economic gains through simple market timing strategies.

Table 2: Out-of-sample R-squared Values in Empirical Studies

	Ridge	Lasso	OLS/Ridgeless	RF	GBRT	NN( $\ell_2$ )	NN( $\ell_1$ )
Finance 1	0.80	-12.19	-81.08	1.30	-14.21	1.41	-10.31
Finance 2	0.19	0.10	-1.25	0.10	-0.30	0.26	0.14
Macro 1a	15.29	15.40	-1375	25.04	16.44	16.94	19.09
Macro 1b	3.49	3.69	-2939	9.08	1.11	7.09	5.39
Macro 2	6.58 (4.83)	-14.58 (43.74)	-837 (854)	9.03 (9.92)	1.28 (14.04)	4.00 (18.42)	1.92 (13.36)
Micro 1	0.48 (0.84)	-1.01 (2.01)	-13198 (12479)	-1.75 (2.70)	-5.07 (6.60)	0.49 (0.27)	-6.77 (17.87)
Micro 2a	26.27 (7.50)	20.37 (6.41)	-12729 (9213)	27.69 (6.15)	16.44 (3.40)	23.87 (10.07)	23.37 (10.09)
Micro 2b	1.89 (3.09)	-3.43 (5.25)	-14724 (10506)	2.86 (3.67)	-6.45 (6.83)	1.11 (2.20)	-1.73 (5.09)

Note: This table reports  $R_{\text{OOS}}^2$  values, presented in percentages, for Ridge, Lasso, OLS/Ridgeless, RF, GBRT, and NNs with respective  $\ell_1$  and  $\ell_2$  penalties, across six empirical studies spanning Finance, Macroeconomics, Microeconomics. For the first example in Macroeconomics and the second example in Microeconomics, two benchmark models are considered. Standard deviations are provided in parentheses when applicable.

We revisit this exercise, by following the empirical framework of [Giannone et al. \(2022\)](#), the initial training set spans from 1948 to 1964, with the model tested on the 1965 data. Subsequently, data from 1965 is added to the training set, and the process is repeated, testing the next model on the 1966 data. This procedure is conducted 57 times, progressively incorporating one additional year into the training set and shifting the test sample forward by one year each time. We evaluate  $R_{\text{OOS}}^2$  based on 57 different predictions using annually expanding windows.<sup>18</sup> The  $R_{\text{OOS}}^2$  results are summarized in Table 2 and align with our theoretical predictions. Specifically, Ridge records an  $R_{\text{OOS}}^2$  of 0.80%, significantly outperforming Lasso’s -12.19%. With the smallest sample size being 17—just sufficient to run OLS with 16 predictors—OLS produces a highly negative  $R_{\text{OOS}}^2$  of -81.08%. Our RF attains an  $R_{\text{OOS}}^2$  of 1.30% while GBRT show performance similar to Lasso, with an  $R_{\text{OOS}}^2$  of -14.21%. The NN with an  $\ell_2$  penalty, NN( $\ell_2$ ), achieves the highest  $R_{\text{OOS}}^2$  of 1.41%, whereas NN( $\ell_1$ ) has -10.31%. As indicated in Figure 8, Ridge and NN( $\ell_2$ ) appear to assign similar weights to these covariates. The leading covariate, eqis—the equity issuing activity ratio—is closely followed in importance by the dividend-price ratio, d/p.

<sup>18</sup>In this example,  $R_{\text{OOS}}^2 = 1 - \frac{\sum_{t=1965}^{2021} (y_t - \hat{y}_t)^2}{\sum_{t=1965}^{2021} (y_t - \bar{y}_t)^2}$ , where  $\bar{y}_t = \frac{\sum_{s=1948}^{t-1} y_s}{(t - 1948)}$ .

## 4.2 Finance 2: Cross-Section of Expected Returns

In our second analysis, we build upon the predictors utilized by [Gu et al. \(2020\)](#) for predicting monthly individual equity returns. Our dataset spans from March 1957 to December 2021, incorporating a total of 920 covariates and averaging over 6,200 stocks per month. Following [Gu et al. \(2020\)](#), our initial training phase utilizes data from 1957 to 1986, followed by performance evaluation using 1987 data. We employ an expanding-window procedure annually, with each iteration expanding the training sample by an additional year and shifting the evaluation period forward by one year.

Table 2 compares the model performance in terms of  $R_{\text{OOS}}^2$ , where the zero predictor serves as the benchmark.<sup>19</sup> In this case,  $\text{NN}(\ell_2)$  emerges as the leading model, closely followed by Ridge, achieving  $R_{\text{OOS}}^2$  of 0.26% and 0.19%, respectively. These models dominate  $\text{NN}(\ell_1)$  and Lasso, which records an  $R_{\text{OOS}}^2$  of 0.14% and 0.10%. The performance of the OLS continues to be underwhelming in this exercise. For the tree-based models, RF outperforms GBRT, achieving an  $R_{\text{OOS}}^2$  of 0.10%, while GBRT yields a negative  $R_{\text{OOS}}^2$  of -0.30%. Figure 8 reveals an intriguing pattern: there appears to be a relationship between the relatively stronger performance among these pairs—Ridge vs Lasso, RF vs GBRT, and  $\text{NN}(\ell_2)$  vs  $\text{NN}(\ell_1)$ —and their respective patterns of sparsity in variable importance plots. Models with denser variable weights outperform those with sparser ones.

To illustrate the economic significance of these relatively low  $R_{\text{OOS}}^2$ s, we adopt a stock selection portfolio strategy. This strategy involves going long on the top 10% and shorting the bottom 10% of stocks, sorted based on their predicted returns for the upcoming month, with equal weighting applied to each stock every month. Over 35 out-of-sample years,  $\text{NN}(\ell_2)$  achieves the highest Sharpe ratio at 2.13, followed closely by Ridge at 1.64.  $\text{NN}(\ell_1)$  also performs well, with a Sharpe ratio of 1.55. By contrast, GBRT exhibits the least impressive performance, with the lowest Sharpe ratio of 0.80.

## 4.3 Macro 1: Macroeconomic Forecasting

The prediction of US macroeconomic activity using a wide range of predictors has been a topic of significant interest since its exploration by [Stock and Watson \(2002\)](#). In our current study, we utilize the [FRED-MD](#) dataset, compiled by [McCracken and Ng \(2016\)](#), to forecast the monthly growth rate of US industrial production (IP). This dataset includes 119 potential predictors and extends from February 1960 to December 2019. Our evaluation

---

<sup>19</sup>In this pooled regression setting, we define  $R_{\text{OOS}}^2 = 1 - \sum_{i,t \in \text{OOS}} (y_{i,t} - \hat{y}_{i,t})^2 / \sum_{i,t \in \text{OOS}} y_{i,t}^2$ .



methodology aligns with the prediction procedure outlined by [Giannone et al. \(2022\)](#). We begin by training all models using data from February 1960 to December 1974 and then evaluate their performance on data from the subsequent year. This process is repeated 45 times, using a similar expanding-window approach on an annual basis. We follow the guidelines outlined by [McCracken and Ng \(2016\)](#) for covariate transformation and data quality management, including outlier removal and imputation of missing data.

We start with a benchmark predictor that only uses an intercept term. In this scenario, all machine learning methods significantly outperform this benchmark, achieving  $R_{\text{oos}}^2$  values ranging from 15.29% to 25.04%.<sup>20</sup> In contrast, OLS overfits the data, resulting in a negative  $R_{\text{oos}}^2$  of -14. This result indicates the existence of strong signals within the covariates. The benchmark model’s lack of competitiveness aligns with our expectations, given the temporal dependence prevalent in macroeconomic time series.

We thereby propose an alternative benchmark that incorporates lagged values of IP growth. Within each training sample, we fit an AR model to the IP growth, selecting the order based on the AIC. The residuals from this model then serve as our prediction target. As discussed in [Section 2.8](#), this approach combines the predictions from the AR model with those from our machine learning methods, yielding a hybrid out-of-sample forecast.<sup>21</sup> Consequently, the new benchmark becomes the direct use of AR model predictions. In this alternative setup, the comparison of  $R_{\text{oos}}^2$  values reveals a pattern somewhat associated with scenarios of weak signal strength: NN( $\ell_2$ ) and RF emerge as the top performers, achieving  $R_{\text{oos}}^2$  values of 7.09% and 9.08%, respectively. Following closely are NN( $\ell_1$ ) at 5.39%, while GBRT lags with a considerably lower  $R_{\text{oos}}^2$  of 1.11%. This disparity in performance appears associated with the findings in [Figure 8](#), which illustrates GBRT’s tendency towards sparser models in comparison to their counterparts. In this example, linear models, specifically Ridge and Lasso, demonstrate comparable performance, achieving  $R_{\text{oos}}^2$  values of 3.49% and 3.69%, respectively. This suggests that, for this particular case, the challenges associated with signal weakness are primarily evident in nonlinear features.

---

<sup>20</sup>In this case, the definition of  $R_{\text{oos}}^2$  is similar to how it is defined in the Finance 1 case.

<sup>21</sup>In implementing advanced machine learning methods alongside a linear component  $W\gamma$ , we adopt a methodology that parallels the one used in [Eq. \(14\)](#). This approach is based on a DGP assumption that  $\mathcal{M}_W y$  is a general function of  $\mathcal{M}_W X$ . This assumption plays a critical role in streamlining the implementation of these machine learning methods, ensuring that the results are directly comparable to those obtained in linear settings. However, it is generally not equivalent to the assumption that  $y - W\gamma$  is a function of  $X$ .

#### 4.4 Macro 2: Economic Growth Across Countries

Next, we explore a dataset originally compiled by [Barro and Lee \(1994\)](#), which includes 60 socio-economic, institutional, and geographical covariates across 90 countries. This dataset is utilized for predicting long-term economic growth, specifically measured by the growth rate of GDP per capita from 1960 to 1985. A pivotal aspect of this analysis involves testing a key prediction of the classical Solow-Swan-Ramsey growth model, which concerns the effect of an initial (lagged) GDP per capita level on subsequent growth rates. By incorporating the logarithm of each country’s GDP per capita in 1960 alongside a constant term, our prediction model includes a total of 62 potential covariates.

[Belloni et al. \(2013b\)](#) implement the Square-root-Lasso technique in their regression, anticipating sparsity among the control variables. This methodology results in a remarkable sparse model, characterized by the inclusion of a singular control variable: the log of the black market premium, a measure of trade openness. In contrast, [Giannone et al. \(2022\)](#) employ a Bayesian approach with a spike-and-slab prior, concluding that a dense model, which includes all covariates, yields the best log-predictive score.

In our predictive analysis, we adopt the same empirical strategy outlined by [Giannone et al. \(2022\)](#). We begin by randomly selecting half of the data samples for model estimation and then proceed to assess the performance of these models using the remaining samples. This process is repeated 100 times. The average out-of-sample  $R_{\text{OOS}}^2$  from these 100 repetitions, in comparison to a benchmark model that includes only the intercept, is presented in [Table 2](#), accompanied by their standard deviations, provided in parentheses.<sup>2223</sup>

Our empirical findings align with [Giannone et al. \(2022\)](#), indicating that dense models, specifically Ridge, RF, and  $\text{NN}(\ell_2)$ , exhibit superior performance compared to their sparse counterparts, such as Lasso, GBRT, and  $\text{NN}(\ell_1)$ . The limited sample size appears to disadvantage complex NN models, rendering them less effective than the simpler Ridge regression. RF demonstrates strong performance, achieving an  $R_{\text{OOS}}^2$  of 9.03%, although it concurrently introduces a twofold increase in the variability of  $R_{\text{OOS}}^2$  values compared to those based on Ridge. The variance of Lasso is pronounced, driven by a handful of extreme values; excluding these, its  $R_{\text{OOS}}^2$  improves but remains notably low at -0.56%. Across all evaluated models, the black market premium consistently emerges as the most influential variable in [Figure 8](#), aligning with the sole variable selected by [Belloni et al. \(2013b\)](#).

---

<sup>22</sup>Here  $R_{\text{OOS}}^2 = 1 - \sum_{i \in \text{OOS}} (y_i - \hat{y}_i)^2 / \sum_{i \in \text{OOS}} (y_i - \bar{y})^2$ , where  $\bar{y}$  is in-sample average of  $y_i$ .

<sup>23</sup>We may also consider a benchmark model with GDP per capita in 1960 included, as predicted by theory. Interestingly, mandating this variable’s inclusion reduces predictive performance in all models, resulting in a negative  $R_{\text{OOS}}^2$  compared to a model with just an intercept.

## 4.5 Micro 1: Crime Rates across US States

Our first microeconomic case revisits the study by [Donohue and Levitt \(2001\)](#), which investigates the impact of abortion legalization post-Roe vs. Wade on the decline in crime rates. They analyze the change in log per-capita murder rates from 1986 to 1997 across 48 states, totaling 576 observations. [Belloni et al. \(2013a\)](#) expand this analysis by including a broader set of 284 control variables, such as interactions and higher-order terms, to mitigate potential confounders. They apply Lasso to select control variables, finding none selected.<sup>24</sup> Similarly, [Giannone et al. \(2022\)](#) observe from their Bayesian analysis of this regression that the posterior density is concentrated on very low probability values of the slab component, which suggests that the regression model is sparse with high likelihood.

We employ the same benchmark model and sample splitting strategy outlined by [Giannone et al. \(2022\)](#). For the initial estimation, we use data spanning from 1986 to 1989, covering all states. Additionally, we incorporate data from a randomly selected 50% of the states for the period from 1990 to 1997. The remaining 50% of the states from 1990 to 1997 are set aside for evaluating the model. This procedure is iterated 8 times, with each iteration expanding the training sample to include one additional year of data, starting from 1990, while correspondingly adjusting the evaluation sample.<sup>25</sup> The entire sequence is carried out 13 times in total, yielding 104 distinct training and evaluation samples. We report the mean and standard deviation of  $R_{\text{oos}}^2$ s in [Table 2](#).

Our findings reveal that  $\text{NN}(\ell_2)$  exhibits a slight edge over Ridge, attaining a  $R_{\text{oos}}^2$  of 0.49%, compared to Ridge’s 0.48%. Apart from these two models, all other models tested demonstrate negative  $R_{\text{oos}}^2$  values. Together, these results indicate an absence of strong signals in the data. We interpret the scarcity of significant signals reported in the literature as a result of their inherent weakness. Empirical evidence does not definitively categorize the underlying DGP as either dense or sparse—rather, it may simply be that no individual signals are particularly strong. Although a null model might initially seem appropriate, our findings indicate that while individual signals may be weak, their combined predictive power should not be overlooked. This is supported by [Figure 8](#), which shows that both Ridge and  $\text{NN}(\ell_2)$  assign small weights to nearly all covariates.

---

<sup>24</sup>[Belloni et al. \(2013a\)](#) proposes a double-Lasso estimator to make inference on the effect of abortion on murder rate. Part of their procedure involves a Lasso regression of murder rate on control variables. Notably, they use differences as the dependent variable, but observe no substantial changes when using levels instead.

<sup>25</sup>Here and after, we calculate  $R_{\text{oos}}^2$  in the same way as Finance 2 case.

## 4.6 Micro 2: Eminent Domain and Economic Outcomes

In our final study, we consider a causal analysis pertinent to eminent domain. Previous research by [Chen and Yeh \(2012\)](#), and subsequently [Belloni et al. \(2012\)](#), employ instrumental variable regressions to understand the impact of eminent domain decisions on economic outcomes. Differing from their broader focus, we closely follow [Giannone et al. \(2022\)](#), focusing on the first stage of predicting pro-plaintiff decisions in takings law cases based on judicial panel characteristics, using a dataset of 312 observations with 138 covariates. Following their strategy, we train the model using data from 1979 to 1984 for all circuits, supplemented with a random 50% of circuits' data from 1985 to 2004. Performance is evaluated on 1985 data from circuits excluded from training. This process is repeated 20 times, sequentially adding one year's data to the training set and updating the evaluation set. The procedure is independently repeated five times, yielding  $20 \times 5 = 100$  unique training and evaluation datasets.

We initially consider a benchmark model with only an intercept. In this setting, all machine learning methods successfully identify predictive signals and outperform this benchmark, as evidenced by large  $R^2_{\text{oos}}$ s. RF performs best with an  $R^2_{\text{oos}}$  of 27.69%, while other models also perform well, except for GBRT. However, the scenario shifts markedly when the benchmark is expanded to include a dummy variable for the absence of cases in a given circuit-year and the number of takings appellate decisions, for a total of three covariates. Against this expanded benchmark, the incremental predictive power of the remaining covariates drops sharply. Ridge's  $R^2_{\text{oos}}$  falls to 1.89%, RF to 2.86%, and  $\text{NN}(\ell_2)$  to 1.11%, while all other methods yield negative  $R^2_{\text{oos}}$  values. Intriguingly, as highlighted in [Figure 8](#), these results seem to associate with the distinct approaches these methods take in weighting covariates. Ridge, RF, and  $\text{NN}(\ell_2)$  distribute small weights across covariates. Meanwhile,  $\text{NN}(\ell_1)$  also opts for a model with a considerable number of coefficients, resulting in a performance that slightly surpasses both Lasso and GBRT, which favor more sparse models in this case.

## 5 Conclusion

In this paper, we scrutinize the performance of machine learning techniques in contexts characterized by low signal-to-noise ratios, a situation frequently observed in economics and finance. Our theoretical analysis indicates that while Lasso is often considered a modern alternative to traditional ordinary least squares, its application in these areas should be

approached cautiously, primarily due to its lessened effectiveness with weak signals.

Our research complements and expands upon the arguments made by [Giannone et al. \(2022\)](#), who cast doubt on the prevalence of sparsity in economic datasets. We take this debate further by showing that it is signal weakness, not necessarily the absence of sparsity, that more significantly contributes to the observed limitations of Lasso in economic applications. Furthermore, the lack of significant variables in empirical studies may be attributed more to signal weakness than to the sparse nature of the underlying DGP.

Our analysis also reveals a marked difference in the performance of Ridge regression. Notably, Ridge demonstrates superior resilience and effectiveness in these environments. Our theoretical findings are further substantiated by simulation studies encompassing a range of advanced machine learning techniques, including trees and neural networks. These experiments consistently reveal that algorithms designed to exploit sparsity tend to underperform in environments where signals are inherently weak. Broadly, our findings emphasize the importance of a nuanced, context-sensitive application of machine learning techniques, adapting to the distinctive data characteristics encountered across various domains.

## References

- Barro, R. J. and J.-W. Lee (1994). Sources of economic growth. *Carnegie-Rochester Conference Series on Public Policy* 40, 1–46.
- Bartlett, P. L., P. M. Long, G. Lugosi, and A. Tsigler (2020). Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences* 117(48), 30063–30070.
- Bayati, M. and A. Montanari (2012). The lasso risk for gaussian matrices. *IEEE Transactions on Information Theory* 58(4), 1997–2017.
- Belloni, A., D. Chen, V. Chernozhukov, and C. Hansen (2012). Sparse models and methods for optimal instruments with an application to eminent domain. *Econometrica* 80(6), 2369–2429.
- Belloni, A., V. Chernozhukov, and C. Hansen (2013a). Inference on Treatment Effects after Selection among High-Dimensional Controls. *The Review of Economic Studies* 81(2), 608–650.
- Belloni, A., V. Chernozhukov, and C. B. Hansen (2013b). *Inference for High-Dimensional*

- Sparse Econometric Models*, Volume 3 of *Econometric Society Monographs*, pp. 245–295. Cambridge University Press.
- Berger, J. (1985). *Statistical Decision Theory and Bayesian Analysis*. Springer Series in Statistics. Springer.
- Bickel, P. J., Y. Ritov, and A. B. Tsybakov (2009). Simultaneous analysis of Lasso and Dantzig selector. *The Annals of Statistics* 37(4), 1705 – 1732.
- Breiman, L. (2001). Random forests. In *Machine Learning*, pp. 5–32.
- Campbell, J. Y. and S. B. Thompson (2007). Predicting Excess Stock Returns Out of Sample: Can Anything Beat the Historical Average? *The Review of Financial Studies* 21(4), 1509–1531.
- Chen, D. L. and S. Yeh (2012). Growth under the shadow of expropriation? the economic impacts of eminent domain. Mimeo, Toulouse School of Economics.
- Dicker, L. H. (2016). Ridge regression and asymptotic minimax estimation over spheres of growing dimension. *Bernoulli* 22(1), 1 – 37.
- Dobriban, E. and S. Wager (2018). High-dimensional asymptotics of prediction: Ridge regression and classification. *The Annals of Statistics* 46(1), 247 – 279.
- Donoho, D. and J. Jin (2004). Higher criticism for detecting sparse heterogeneous mixtures. *The Annals of Statistics* 32(3), 962 – 994.
- Donohue, John J., I. and S. D. Levitt (2001). The Impact of Legalized Abortion on Crime. *The Quarterly Journal of Economics* 116(2), 379–420.
- Efron, B. (2012). *Large-Scale Inference: Empirical Bayes Methods for Estimation, Testing, and Prediction*. Institute of Mathematical Statistics Monographs. Cambridge University Press.
- Efron, B., T. Hastie, I. Johnstone, and R. Tibshirani (2004). Least angle regression. *The Annals of Statistics* 32(2), 407 – 499.
- Fan, J. and R. Li (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association* 96(456), 1348–1360.

- Ferreira, M. A. and P. Santa-Clara (2011). Forecasting stock market returns: The sum of the parts is more than the whole. *Journal of Financial Economics* 100(3), 514–537.
- Giannone, D., M. Lenza, and G. E. Primiceri (2022). Economic predictions with big data: The illusion of sparsity. *Econometrica* 89(5), 2409–2437.
- Goodfellow, I., Y. Bengio, and A. Courville (2016). *Deep Learning*. MIT Press.
- Gordon, Y. (1988). On milman’s inequality and random subspaces which escape through a mesh in  $\mathbb{R}_n$ . In *Geometric Aspects of Functional Analysis*, Berlin, Heidelberg, pp. 84–106. Springer Berlin Heidelberg.
- Gu, S., B. Kelly, and D. Xiu (2020). Empirical Asset Pricing via Machine Learning. *The Review of Financial Studies* 33(5), 2223–2273.
- Hall, P. and J. Jin (2010). Innovated higher criticism for detecting sparse signals in correlated noise. *The Annals of Statistics* 38(3), 1686 – 1732.
- Hastie, T., A. Montanari, S. Rosset, and R. J. Tibshirani (2022). Surprises in high-dimensional ridgeless least squares interpolation. *The Annals of Statistics* 50(2), 949 – 986.
- Hoerl, A. E. and R. W. Kennard (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics* 12(1), 55–67.
- Ioffe, S. and C. Szegedy (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37*, ICML’15, pp. 448–456. JMLR.org.
- Jiang, W. and C.-H. Zhang (2009). General maximum likelihood empirical Bayes estimation of normal means. *The Annals of Statistics* 37(4), 1647 – 1684.
- Kelly, B. and S. Pruitt (2013). Market expectations in the cross-section of present values. *The Journal of Finance* 68(5), 1721–1756.
- Kelly, B. T., S. Malamud, and K. Zhou (2023). The virtue of complexity in return prediction. *The Journal of Finance*. Forthcoming.

- Kingma, D. and J. Ba (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Kolesár, M., U. K. Müller, and S. T. Roelsgaard (2024). The fragility of sparsity.
- Liang, T. and P. Sur (2022). A precise high-dimensional asymptotic theory for boosting and minimum- $\ell_1$ -norm interpolated classifiers. *The Annals of Statistics* 50(3), 1669 – 1695.
- McCracken, M. W. and S. Ng (2016). Fred-md: A monthly database for macroeconomic research. *Journal of Business & Economic Statistics* 34(4), 574–589.
- Miolane, L. and A. Montanari (2021). The distribution of the Lasso: Uniform control over sparse balls and adaptive parameter tuning. *The Annals of Statistics* 49(4), 2313 – 2335.
- Newey, W. K. (1991). Uniform convergence in probability and stochastic equicontinuity. *Econometrica* 59(4), 1161–1167.
- Rapach, D. E., J. K. Strauss, and G. Zhou (2010). Out-of-sample equity premium prediction: Combination forecasts and links to the real economy. *The Review of Financial Studies* 23(2), 821–862.
- Robbins, H. (1964). The Empirical Bayes Approach to Statistical Decision Problems. *The Annals of Mathematical Statistics* 35(1), 1 – 20.
- Srivastava, N., G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov (2014). Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research* 15(56), 1929–1958.
- Stock, J. H. and M. W. Watson (2002). Forecasting using principal components from a large number of predictors. *Journal of the American Statistical Association* 97(460), 1167–1179.
- Su, W., M. Bogdan, and E. Candès (2017). False discoveries occur early on the lasso path. *The Annals of Statistics* 45(5), 2133–2150.
- Thrampoulidis, C., E. Abbasi, and B. Hassibi (2018). Precise error analysis of regularized  $m$ -estimators in high dimensions. *IEEE Transactions on Information Theory* 64(8), 5592–5628.
- Thrampoulidis, C., S. Oymak, and B. Hassibi (2015). Regularized linear regression: A precise analysis of the estimation error. In *Proceedings of The 28th Conference on Learning*



*Theory*, Volume 40 of *Proceedings of Machine Learning Research*, Paris, France, pp. 1683–1709. PMLR.

Tsigler, A. and P. L. Bartlett (2023). Benign overfitting in ridge regression. *Journal of Machine Learning Research* 24(123), 1–76.

Wainwright, M. (2019). *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press.

Wang, S., H. Weng, and A. Maleki (2020). Which bridge estimator is the best for variable selection? *The Annals of Statistics* 48(5), 2791 – 2823.

Welch, I. and A. Goyal (2007). A Comprehensive Look at The Empirical Performance of Equity Premium Prediction. *The Review of Financial Studies* 21(4), 1455–1508.

Zou, H. and T. Hastie (2005). Regularization and Variable Selection Via the Elastic Net. *Journal of the Royal Statistical Society Series B: Statistical Methodology* 67(2), 301–320.

## A Mathematical Proofs

### A.1 Proof of Theorem 1

*Proof.* Throughout the proof, we employ the shorthand notation “w.a.p.1” to denote “with probability approaching one.” For two random variables  $X$  and  $Y$ , we write  $X \perp Y$  when they are independent and  $X \stackrel{d}{=} Y$  when they have the same distribution. For convenience, we omit the subscript  $F$  from the expectation operator  $\mathbb{E}_F(\cdot)$ .

Our objective is to demonstrate that  $\mathbb{E}\|\Sigma_2^{1/2}(\mathbb{E}(\beta_0|X, y) - \beta_0)\|^2 / \mathbb{E}\|\Sigma_2^{1/2}\beta_0\|^2 \rightarrow 1$ , as  $n \rightarrow \infty$ . This can be shown by proving  $\mathbb{E}\|\Sigma_2^{1/2}\mathbb{E}(\beta_0|X, y)\|^2 = o(\mathbb{E}\|\Sigma_2^{1/2}\beta_0\|^2)$ . Given that the eigenvalues of  $\Sigma_2$  are bounded away from zero and positive infinity, it suffices to establish that  $\mathbb{E}\|\mathbb{E}(\beta_0|X, y)\|^2 = o(\tau)$ . Therefore, we need to prove for all  $1 \leq i \leq p$ ,  $\mathbb{E}(\mathbb{E}(\beta_{0,i}|X, y))^2 = o(p^{-1}\tau)$ , or, equivalently,  $\mathbb{E}(\mathbb{E}(b_{0,i}|X, y))^2 = o(1)$ .

By the inequality  $\mathbb{E}(\mathbb{E}(A|\mathcal{F})^2) \leq \mathbb{E}(\mathbb{E}(A|\mathcal{G})^2)$  for  $\mathcal{F} \subset \mathcal{G}$ , and that  $\beta_0$  is i.i.d., we have

$$\mathbb{E}(\mathbb{E}(b_{0,i}|X, y))^2 \leq \mathbb{E}(\mathbb{E}(b_{0,i}|X, y, \beta_{0,-i}))^2 = \mathbb{E}(\mathbb{E}(b_{0,i}|X_{\cdot,i}, \beta_{0,i}\Sigma_\varepsilon^{-1/2}X_{\cdot,i} + z))^2,$$

where  $z$  is defined in Assumption 2,  $X_{\cdot,i}$  represents the  $i$ -th column of  $X$ , and  $\beta_{0,-i}$  denotes the subvector of  $\beta$  without the  $i$ th entry. Denote the information set generated by

$\{X_{\cdot,i}, \beta_{0,i} \Sigma_\varepsilon^{-1/2} X_{\cdot,i} + z\}$  as  $\mathcal{G}_i$ . By Assumption 3,  $b_{0,i}$  can be written as  $q^{-1/2} b_{1i} b_{2i}$  where  $b_{1i} \sim B(1, q)$  and  $b_{2i}$  is a sub-exponential random variable with mean zero and variance  $\sigma_\beta^2$ , whose distribution function is denoted by  $F_{b_2}$ . For any  $M_1 < 0$ , find  $M_2$  (a function of  $M_1$ ) such that  $\mathbb{E} b_{0,i} \mathbf{1}_{q^{1/2} b_{0,i} \in [M_1, M_2]} = q^{1/2} \mathbb{E} b_{2i} \mathbf{1}_{b_{2i} \in [M_1, M_2]} = 0$ . This is always feasible because  $\mathbb{E} b_{2i} = 0$ . By Cauchy-Schwarz inequality, we have

$$\begin{aligned} \mathbb{E}(\mathbb{E}(b_{0,i} | \mathcal{G}_i))^2 &\leq 3\mathbb{E}(\mathbb{E}(b_{0,i} \mathbf{1}_{q^{1/2} b_{0,i} \in [M_1, M_2]} | \mathcal{G}_i))^2 + 3\mathbb{E}(\mathbb{E}(b_{0,i} \mathbf{1}_{q^{1/2} b_{0,i} > M_2} | \mathcal{G}_i))^2 \\ &\quad + 3\mathbb{E}(\mathbb{E}(b_{0,i} \mathbf{1}_{q^{1/2} b_{0,i} < M_1} | \mathcal{G}_i))^2 =: 3S_{1n} + 3S_{2n} + 3S_{3n}. \end{aligned} \quad (21)$$

Lemma 13 proves that for any given  $M_1$ ,  $\lim_{n \rightarrow \infty} S_{1n} = 0$ . Observe that  $S_{2n} = \mathbb{E}(\mathbb{E}(b_{0,i} \mathbf{1}_{q^{1/2} b_{0,i} > M_2} | \mathcal{G}_i))^2 \leq \mathbb{E} b_{0,i}^2 \mathbf{1}_{q^{1/2} b_{0,i} > M_2}$ . As a result, (21) implies

$$\lim_{n \rightarrow \infty} \mathbb{E}(\mathbb{E}(b_{0,i} | \mathcal{G}_i))^2 \leq 3\mathbb{E} b_{0,i}^2 \mathbf{1}_{q^{1/2} b_{0,i} > M_2} + 3\mathbb{E} b_{0,i}^2 \mathbf{1}_{q^{1/2} b_{0,i} < M_1} = 3\mathbb{E} b_{2i}^2 \mathbf{1}_{b_{2i} > M_2} + 3\mathbb{E} b_{2i}^2 \mathbf{1}_{b_{2i} < M_1}.$$

Since  $b_{2i}$  has finite variance, the right-hand-side of the above inequality can be arbitrarily small by letting  $M_1 \rightarrow -\infty$ , which completes the proof.  $\square$

## A.2 Proof of Theorem 2

*Proof.* For ease of notation, we let  $\hat{\beta} := \hat{\beta}_r(\lambda_n)$  and  $c_n := p/n$ . Additionally, define  $\delta_1^* := 2\sqrt{\sigma_\varepsilon^2 \theta_1}$ ,  $\delta_2^* := (2\lambda \sigma_x^2 \sigma_\beta^2 \theta_4 - 4\sigma_\varepsilon^2 \sigma_x^2 \theta_3) / \delta_1^* \lambda$ ,  $\mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2^*) := (\delta_1^* \delta_2^* - 2\sigma_x^2 \sigma_\beta^2 \theta_4) / 4\sigma_\varepsilon^2 \theta_3$ , and  $C_n^\phi := c_n \tau^{-1} \sigma_x^2 \sigma_\beta^2 - c_n \tau^{-1} \sigma_x^2 (\delta_1^*)^2 (4\lambda)^{-1} + c_n \sigma_\varepsilon^2 \sigma_x^4 \theta_3 \lambda^{-2} - c_n \sigma_x^4 \sigma_\beta^2 \theta_4 \lambda^{-1}$ .

We first show that it is sufficient to establish that

$$c_n \tau^{-3/2} (\|\Sigma_2^{1/2}(\hat{\beta} - \beta_0)\| - \|\Sigma_2^{1/2} \beta_0\|) \xrightarrow{P} \alpha_2^* := \theta_2 \sigma_x^3 \left( \frac{\sigma_\varepsilon^2 \theta_1}{2\lambda^2 \sigma_\beta} - \frac{\sigma_\beta}{\lambda} \right). \quad (22)$$

This is because Eq. (22) implies that  $\|\Sigma_2^{1/2}(\hat{\beta} - \beta_0)\|^2 = \|\Sigma_2^{1/2} \beta_0\|^2 + 2c_n^{-1} \tau^{3/2} \alpha_2^* \|\Sigma_2^{1/2} \beta_0\| + o_P(c_n^{-1} \tau^2)$ . Additionally, using Lemma 2 and  $q^{-1/2} p^{-1/2} = o(1)$  by Assumption 4, we have

$$\|\Sigma_2^{1/2} \beta_0\| = \tau^{1/2} \sigma_x \sigma_\beta + O_P(q^{-1/2} p^{-1/2} \tau^{1/2}). \quad (23)$$

The above two equations together yield the desired result of the theorem. To prove Eq. (22), by incorporating (23) and  $c_n q^{-1/2} p^{-1/2} \tau^{-1} = o(1)$  by Assumption 4, it reduces to

$$c_n \tau^{-1} (\tau^{-1/2} \|\Sigma_2^{1/2}(\hat{\beta} - \beta_0)\| - \sigma_x \sigma_\beta) \xrightarrow{P} \alpha_2^*. \quad (24)$$

Set  $w = \tau^{-3/2}\Sigma_2^{1/2}(\beta - \beta_0)$  and  $\hat{w} = \tau^{-3/2}\Sigma_2^{1/2}(\hat{\beta} - \beta_0)$ . By Eq. (2), we have

$$\hat{w} = \arg \min_w \frac{c_n}{n} \left\| \tau^{1/2}\Sigma_1^{1/2}Zw - \tau^{-1}\varepsilon \right\|^2 + c_n^2\lambda \left\| \Sigma_2^{-1/2}w + \tau^{-3/2}\beta_0 \right\|^2 - \frac{c_n\tau^{-2}}{n} \|\varepsilon\|^2 - C_n^\phi, \quad (25)$$

where subtracting  $c_n\tau^{-2}\|\varepsilon\|^2/n$  and  $C_n^\phi$  from the objective function does not alter the solution. Using the definition of  $\hat{w}$ , proving (24) is equivalent to proving  $c_n\|\hat{w}\| - c_n\tau^{-1}\sigma_x\sigma_\beta \xrightarrow{P} \alpha_2^*$ . Equivalently, we need to prove for all  $\epsilon > 0$ , w.p.a.1,

$$\alpha_2^* - \epsilon \leq c_n\|\hat{w}\| - c_n\tau^{-1}\sigma_x\sigma_\beta \leq \alpha_2^* + \epsilon. \quad (26)$$

Note that Eq. (25) is a high-dimensional optimization problem. By employing CGMT, Lemma 14 connects it to a scalar optimization problem below:

$$\begin{aligned} & \min_{c_n|\alpha - \tau^{-1}\sigma_x\sigma_\beta| \leq K_\alpha} \max_{\substack{\gamma > 0 \\ 0 \leq \delta \leq 4\tau^{-1}\sqrt{c_1c_\varepsilon}}} -\frac{c_n\delta^2}{4}\mu_n(\alpha, \delta) + \frac{c_n}{n}(\tau^{1/2}\alpha g - \tau^{-1}\Sigma_1^{-1/2}\varepsilon)^\top (\Sigma_1^{-1} - \mu_n(\alpha, \delta)\mathbb{I})^{-1} \\ & \times (\tau^{1/2}\alpha g - \tau^{-1}\Sigma_1^{-1/2}\varepsilon) - \frac{c_n\gamma}{2} + \frac{c_n^2\lambda^2}{4} \left( \tau^{-3/2}\Sigma_2^{1/2}\beta_0 + \frac{\alpha^2\delta\tau^{1/2}}{\sqrt{n}\gamma}h \right)^\top \left( \frac{\lambda}{4}\Sigma_2 + \frac{c_n\alpha^2\lambda^2}{2\gamma}\mathbb{I} \right)^{-1} \\ & \times \left( \tau^{-3/2}\Sigma_2^{1/2}\beta_0 + \frac{\alpha^2\delta\tau^{1/2}}{\sqrt{n}\gamma}h \right) - \frac{c_n\tau\alpha^2\delta^2}{2\gamma n} \|h\|^2 - \frac{c_n\tau^{-2}}{n} \|\varepsilon\|^2 - C_n^\phi, \end{aligned} \quad (27)$$

where  $K_\alpha$  is a sufficiently large fixed constant,  $\mu_n$  is defined in Eq. (B9), and  $g \in \mathbb{R}^n$  and  $h \in \mathbb{R}^p$  are standard Gaussian vectors independent of the other random variables. Denote the objective function as  $Q_n(\alpha, \delta, \gamma)$ . Let us define  $\gamma = \tau^{-1}\gamma_1$ ,  $\delta = \tau^{-1}\delta_1^* + \delta_2^* + c_n^{-1/2}\delta_3$  and  $\alpha = \tau^{-1}\sigma_x\sigma_\beta + c_n^{-1}\alpha_2$ . Subsequently, we present the modified objective function  $\tilde{Q}_n(\alpha_2, \delta_3, \gamma_1)$ :

$$\tilde{Q}_n(\alpha_2, \delta_3, \gamma_1) := Q_n(\tau^{-1}\sigma_x\sigma_\beta + c_n^{-1}\alpha_2, \tau^{-1}\delta_1^* + \delta_2^* + c_n^{-1/2}\delta_3, \tau^{-1}\gamma_1). \quad (28)$$

By Lemma 14,  $Q_n(\alpha, \delta, \gamma)$  is convex with respect to  $\alpha$  and jointly concave with respect to  $(\delta, \gamma)$ . Therefore, it is evident that  $\tilde{Q}_n$  remains convex with respect to  $\alpha_2$  and jointly concave with respect to  $(\delta_3, \gamma_1)$ . Lemma 14 implies that if the following claims hold

$$\begin{aligned} & \text{(i) } \forall \text{ compact } A \subset [-K_\alpha, K_\alpha], \min_{\alpha_2 \in A} \max_{\substack{\gamma_1 > 0 \\ \delta_3 \in K_{\delta_3}}} \tilde{Q}_n(\alpha_2, \delta_3, \gamma_1) \xrightarrow{P} \min_{\alpha_2 \in A} \max_{\substack{\gamma_1 > 0 \\ \delta_3 \in \mathbb{R}}} \tilde{Q}_n(\alpha_2, \delta_3, \gamma_1), \\ & \text{(ii) } \min_{\alpha_2 \in [-K_\alpha, K_\alpha]} \max_{\substack{\gamma_1 > 0 \\ \delta_3 \in K_{\delta_3}}} \tilde{Q}_n(\alpha_2, \delta_3, \gamma_1) < \min_{\alpha_2 \in [-K_\alpha, \alpha_2^* - \epsilon] \cup [\alpha_2^* + \epsilon, K_\alpha]} \max_{\substack{\gamma_1 > 0 \\ \delta_3 \in \mathbb{R}}} \tilde{Q}_n(\alpha_2, \delta_3, \gamma_1) \end{aligned} \quad (29)$$

where  $\tilde{Q}(\alpha_2, \delta_3, \gamma_1) := -\frac{\delta_3^2 \theta_1}{4\theta_3} + 2\sigma_x \sigma_\beta \alpha_2 - \frac{\gamma_1^2}{4\sigma_x^2 \sigma_\beta^2 \lambda} \theta_2 - \frac{\gamma_1 \alpha_2}{\sigma_x \sigma_\beta} + \frac{(\delta_1^*)^2 \gamma_1}{8\lambda^2 \sigma_x^2 \sigma_\beta^2} \theta_2$  and  $K_{\delta_3} := [-c_n^{1/2}(\tau^{-1}\delta_1^* + \delta_2^*), 4c_n^{1/2}\tau^{-1}\sqrt{C_1 C_\varepsilon} - c_n^{1/2}(\tau^{-1}\delta_1^* + \delta_2^*)]$ , then Eq. (26) holds. The above two conditions are verified in Lemma 18, thereby completing the proof.  $\square$

### A.3 Proof of Theorem 3

*Proof.* For convenience, we define the shorthand notation  $\hat{\beta}_{\lambda_n}^i := \hat{\beta}_r^i(\lambda_n)$ . Also, we define  $\hat{R}^{K-CV}(\lambda_n) := \frac{1}{n} \sum_{i=1}^K \|y_{(i)} - X_{(i)} \hat{\beta}_r^i(\lambda_n)\|^2$ . By Lemmas 2, 3 and 6, w.p.a.1, we have

$$n^{-1} \|X_{(-i)}^\top X_{(-i)}\| \leq n^{-1} C_2 \|Z_{(-i)}^\top Z_{(-i)}\| \leq C_2 (1 + \sqrt{c_n})^2, \quad i = 1, \dots, K, \quad (30)$$

$$n^{-1} \|\varepsilon\|^2 \leq 2\sigma_\varepsilon^2 \quad \text{and} \quad n^{-1} \|y\|^2 \leq 2\sigma_\varepsilon^2. \quad (31)$$

Based on these inequalities, Lemma 19 proves

$$\inf_{\lambda \in [\varepsilon, \tilde{c}\tau^{-1}]} pn^{-1}\tau^{-2} \left\{ \hat{R}^{K-CV}(\lambda) - \frac{1}{n} \|\varepsilon\|^2 - \|\Sigma_2^{1/2} \beta_0\|^2 \right\} > 0, \quad (32)$$

w.p.a.1 for some constant  $\tilde{c} > 0$ . Additionally, for any fixed  $\lambda > 0$ , Lemma 19 also proves

$$pn^{-1}\tau^{-2} \left\{ \hat{R}^{K-CV}(\tau^{-1}\lambda) - \frac{1}{n} \|\varepsilon\|^2 - \|\Sigma_2^{1/2} \beta_0\|^2 \right\} \xrightarrow{P} \frac{2(K-1)}{K} \theta_2 \sigma_x^4 \left( \frac{\sigma_\varepsilon^2}{2\lambda^2} - \frac{\sigma_\beta^2}{\lambda} \right). \quad (33)$$

Using (33) and  $\lambda^{opt}$ 's definition,  $pn^{-1}\tau^{-2} \left\{ \hat{R}^{K-CV}(\tau^{-1}\lambda^{opt}) - \frac{1}{n} \|\varepsilon\|^2 - \|\Sigma_2^{1/2} \beta_0\|^2 \right\} < 0$ . Together with Eq. (32), the minimizer of  $\hat{R}^{K-CV}(\lambda)$  must satisfy  $\hat{\lambda}_n^{K-CV} \geq \tilde{c}\tau^{-1}$ , w.a.p.1, so that,  $\hat{\lambda}_n^{K-CV} = \arg \min_{\lambda_n \in [\tilde{c}\tau^{-1}, \infty)} \hat{R}^{K-CV}(\lambda_n)$ . Moreover, it also implies that  $\tau^{-1}\lambda^{opt} \notin [\varepsilon, \tilde{c}\tau^{-1}]$ , that is,  $\lambda^{opt} \geq \tilde{c}$ .

Next, we re-parametrize the above optimization problem:  $\tilde{\mu} = \arg \min_{\mu \in [0, \tilde{c}^{-1}]} \tilde{R}(\mu)$ , where  $\tilde{R}(\mu) := \hat{R}^{K-CV}(\tau^{-1}\mu^{-1})$ , and we extend the domain of  $\tilde{R}(\cdot)$  to include 0:  $\tilde{R}(0) := \lim_{\mu \rightarrow 0} \tilde{R}(\mu) = \|y\|^2/n$ . Lemma 20 implies that  $pn^{-1}\tau^{-2}\tilde{R}(\mu)$  satisfies stochastic equicontinuity. Using this fact and Theorem 1 of Newey (1991), the convergence of  $pn^{-1}\tau^{-2} \left\{ \tilde{R}(\mu) - \frac{1}{n} \|\varepsilon\|^2 - \|\Sigma_2^{1/2} \beta_0\|^2 \right\} \xrightarrow{P} \frac{2(K-1)}{K} \theta_2 \sigma_x^4 \left( \frac{\sigma_\varepsilon^2 \mu^2}{2} - \sigma_\beta^2 \mu \right)$  holds uniformly over the interval  $[0, \tilde{c}^{-1}]$ . Since  $(\lambda^{opt})^{-1}$  is a unique minimizer of the right-hand-side and is nonzero, it follows that  $\tilde{\mu} \xrightarrow{P} (\lambda^{opt})^{-1}$  and  $\tilde{\mu} = \tau^{-1}(\hat{\lambda}_n^{K-CV})^{-1}$ , w.a.p.1, which implies that  $\hat{\lambda}_n^{K-CV}/\lambda_n^{opt} \xrightarrow{P} 1$ .

Now we prove  $\Delta(\hat{\beta}_r(\hat{\lambda}_n^{K-CV})) - \Delta(\hat{\beta}_r(\lambda_n^{opt})) \xrightarrow{P} 0$ . For notational simplicity, we write  $\hat{\beta}_r(\hat{\lambda}_n^{K-CV})$  as  $\hat{\beta}_{cv}$  and  $\hat{\beta}_r(\lambda_n^{opt})$  as  $\hat{\beta}_{opt}$ . We need to prove  $c_n \tau^{-2} (\|\Sigma_2^{1/2}(\hat{\beta}_{cv} - \beta_0)\|^2 - \|\Sigma_2^{1/2}(\hat{\beta}_{opt} -$

$\beta_0\|^2) = o_P(1)$ . By direct calculation, we have

$$\begin{aligned}
& c_n \tau^{-2} (\|\Sigma_2^{1/2}(\hat{\beta}_{cv} - \beta_0)\|^2 - \|\Sigma_2^{1/2}(\hat{\beta}_{opt} - \beta_0)\|^2) \\
&= c_n \tau^{-2} (\hat{\beta}_{cv}^\top \Sigma_2 \hat{\beta}_{cv} - \hat{\beta}_{opt}^\top \Sigma_2 \hat{\beta}_{opt} + 2(\hat{\beta}_{opt} - \hat{\beta}_{cv})^\top \Sigma_2 \beta_0) \\
&= c_n \tau^{-2} ((\hat{\beta}_{cv} - \hat{\beta}_{opt})^\top \Sigma_2 \hat{\beta}_{cv} + \hat{\beta}_{opt}^\top \Sigma_2 (\hat{\beta}_{cv} - \hat{\beta}_{opt}) + 2(\hat{\beta}_{opt} - \hat{\beta}_{cv})^\top \Sigma_2 \beta_0) =: A_1 + A_2 + A_3.
\end{aligned} \tag{34}$$

Lemma 21 proves that these terms converge to zero w.p.a.1, which completes the proof.  $\square$

#### A.4 Proof of Theorem 4

*Proof.* For ease of notation, we let  $\hat{\beta} := \hat{\beta}_l(\lambda_n)$ . We adopt the same notation  $\delta_1^*$  and  $\mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2)$  as used in the proof of Theorem 2. Moreover, we define  $\delta_2^* = 2\sigma_x^2 \sigma_\beta^2 \theta_4 / \delta_1^*$  and  $C_n^\phi = c_n \tau^{-1} \sigma_x^2 \sigma_\beta^2$ , respectively. Similar to Theorem 2, it is essential to establish that the following inequality holds w.a.p.1 for any sufficiently small  $\epsilon > 0$ :

$$\frac{c_\alpha}{2\sigma_\beta} + \epsilon \leq c_n \tau^{-1} (\tau^{-1/2} \|\Sigma_2^{1/2}(\hat{\beta} - \beta_0)\| - \sigma_x \sigma_\beta) \leq \frac{C_\alpha}{2\sigma_\beta} - \epsilon. \tag{35}$$

Define  $w = \tau^{-3/2} \Sigma_2^{1/2}(\beta - \beta_0)$ . Using this, we can rewrite (3) as the following problem:

$$\hat{w} = \arg \min_w \frac{c_n}{n} \|\tau^{1/2} \Sigma_1^{1/2} Z w - \tau^{-1} \epsilon\|^2 + \frac{c_n \tau^{-1/2} \lambda_n}{\sqrt{n}} \|\Sigma_2^{-1/2} w + \tau^{-3/2} \beta_0\|_1 - \frac{c_n \tau^{-2}}{n} \|\epsilon\|^2 - C_n^\phi.$$

By CGMT, Lemma 22 establishes its connection with the following optimization problem:

$$\begin{aligned}
& \min_{\alpha \in K_\alpha} \max_{\substack{\gamma > 0 \\ 0 \leq \delta \leq 4\tau^{-1} \sqrt{C_1 C_\epsilon}}} -\frac{c_n \delta^2}{4} \mu_n(\alpha, \delta) + \frac{c_n}{n} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \epsilon)^\top (\Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I})^{-1} \\
& \times (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \epsilon) - \frac{c_n \gamma}{2} + \frac{c_n \gamma \tau^{-3}}{2\alpha^2} \beta_0 \Sigma_2 \beta_0 + \frac{c_n \tau^{-1} \delta}{\sqrt{n}} h^\top \Sigma_2^{1/2} \beta_0 \\
& - \min_{\|v\|_\infty \leq 1} \frac{c_n \alpha^2}{2\gamma} \left\| n^{-1/2} \tau^{-1/2} \lambda_n \Sigma_2^{-1/2} v - n^{-1/2} \tau^{1/2} \delta h - \frac{\gamma}{\alpha^2} \tau^{-3/2} \Sigma_2^{1/2} \beta_0 \right\|^2 - \frac{c_n \tau^{-2}}{n} \|\epsilon\|^2 - C_n^\phi,
\end{aligned} \tag{36}$$

where  $K_\alpha := \{\alpha | c_n \alpha - c_n \tau^{-1} \sigma_x \sigma_\beta \in [c_\alpha / 4\sigma_\beta, C_\alpha / \sigma_\beta]\}$ ,  $\mu_n$  is defined in Eq. (B9), and  $g \in \mathbb{R}^n$  and  $h \in \mathbb{R}^p$  are standard Gaussian vectors independent of the other random variables. Denote the objective function as  $Q_n(\alpha, \delta, \gamma)$ . Similar to Theorem 2, we define  $\gamma = \tau^{-1} \gamma_1$ ,  $\delta = \tau^{-1} \delta_1^* + \delta_2^* + c_n^{-1/2} \delta_3$ , and  $\alpha = \tau^{-1} \sigma_x \sigma_\beta + c_n^{-1} \alpha_2$ , obtaining the modified objective function:

$$\tilde{Q}_n(\alpha_2, \delta_3, \gamma_1) := Q_n(\tau^{-1} \sigma_x \sigma_\beta + c_n^{-1} \alpha_2, \tau^{-1} \delta_1^* + \delta_2^* + c_n^{-1/2} \delta_3, \tau^{-1} \gamma_1). \tag{37}$$

Note that  $\delta_3 \in K_{\delta_3} := [-c_n^{1/2}(\tau^{-1}\delta_1^* + \delta_2^*), 4c_n^{1/2}\tau^{-1}\sqrt{C_1C_\varepsilon} - c_n^{1/2}(\tau^{-1}\delta_1^* + \delta_2^*)]$ . Lemma 22 implies that if the following inequalities

$$\begin{aligned} \min_{\alpha_2 \in [\frac{c_\alpha}{4\sigma_\beta}, \frac{c_\alpha}{\sigma_\beta}]} \max_{\substack{\gamma_1 > 0 \\ \delta_3 \in K_{\delta_3}}} \tilde{Q}_n(\alpha_2, \delta_3, \gamma_1) &< -\frac{C_\lambda}{8C_2} + \eta, \\ \min_{\alpha_2 \in [\frac{c_\alpha}{4\sigma_\beta}, \frac{c_\alpha}{2\sigma_\beta} + \epsilon] \cup [\frac{c_\alpha}{2\sigma_\beta} - \epsilon, \frac{c_\alpha}{\sigma_\beta}]} \max_{\substack{\gamma_1 > 0 \\ \delta_3 \in K_{\delta_3}}} \tilde{Q}_n(\alpha_2, \delta_3, \gamma_1) &> -\frac{C_\lambda}{100C_2} - \eta, \end{aligned} \quad (38)$$

hold for sufficiently small  $\epsilon > 0$  and  $\eta > 0$ , then Eq. (35) holds. Finally, Lemma 24 verifies Eq. (38), thereby completing the proof.  $\square$

### A.5 Proof of Proposition 1

*Proof.* Since the out-of-sample data are mutually independent, Lemmas 2 and 3 lead to  $\sum_{i \in \text{OOS}} y_i^2 = n_{\text{OOS}}(\sigma_\varepsilon^2 + \tau\sigma_x^2\sigma_\beta^2) + o_P(n_{\text{OOS}}\tau)$  and  $\sum_{i \in \text{OOS}} y_i^2 - (y_i - X_i\hat{\beta}_r(\lambda_n^{\text{opt}}))^2 = n_{\text{OOS}}p^{-1}n\tau^2\theta_2\sigma_x^4\sigma_\beta^4\sigma_\varepsilon^{-2}(1 + o_P(1))$ , where we use  $\Delta(\hat{\beta}_r(\lambda_n^{\text{opt}})) = -\theta_2\sigma_x^4\sigma_\beta^4\sigma_\varepsilon^{-2} + o_P(1)$  by Theorem 2 and  $n_{\text{OOS}}p^{-2}n^2\tau^2 \rightarrow \infty$ . The estimates above offer the key components for deriving the limit of  $R_{\text{OOS}}^2$ :

$$R_{\text{OOS}}^2 = \left( \sum_{i \in \text{OOS}} y_i^2 - (y_i - X_i\hat{\beta}_r(\lambda_n^{\text{opt}}))^2 \right) \left( \sum_{i \in \text{OOS}} y_i^2 \right)^{-1} = p^{-1}n\theta_2(R^2)^2(1 + o_P(1)). \quad \square$$

### A.6 Proof of Theorem 5

*Proof.* For convenience, let  $\hat{\beta} := \hat{\beta}_r(\lambda_n)$ . We write the prediction error of the benchmark as:

$$\begin{aligned} y^{\text{new}} - \hat{y}_b^{\text{new}} &= (w^{\text{new}})^\top \gamma_0 + (x^{\text{new}})^\top \beta_0 + \varepsilon^{\text{new}} - (w^{\text{new}})^\top (W^\top W)^{-1} W^\top (W\gamma_0 + X\beta_0 + \varepsilon) \\ &= ((w^{\text{new}})^\top - (w^{\text{new}})^\top (W^\top W)^{-1} W^\top U) \beta_0 + (\varepsilon^{\text{new}} - (w^{\text{new}})^\top (W^\top W)^{-1} W^\top \varepsilon). \end{aligned}$$

Similarly, for the Ridge estimator, we have

$$y^{\text{new}} - \hat{y}^{\text{new}} = ((w^{\text{new}})^\top - (w^{\text{new}})^\top (W^\top W)^{-1} W^\top U) (\beta_0 - \hat{\beta}) + (\varepsilon^{\text{new}} - (w^{\text{new}})^\top (W^\top W)^{-1} W^\top \varepsilon).$$

As a result, with simple algebra we can rewrite  $\mathbb{E}[(y^{\text{new}} - \hat{y}^{\text{new}})^2 | \mathcal{I}] - \mathbb{E}[(y^{\text{new}} - \hat{y}_b^{\text{new}})^2 | \mathcal{I}]$  as

$$\mathbb{E} \left[ \left( (w^{\text{new}})^\top - (w^{\text{new}})^\top (W^\top W)^{-1} W^\top U \right) (\beta_0 - \hat{\beta}) \right]^2 | \mathcal{I} \right]$$

$$\begin{aligned}
& - \mathbb{E} \left[ \left( (w^{new})^\top - (w^{new})^\top (W^\top W)^{-1} W^\top U \right) \beta_0 \right]^2 | \mathcal{I} ] \\
& - 2 \mathbb{E} \left[ \left( (w^{new})^\top - (w^{new})^\top (W^\top W)^{-1} W^\top U \right) \hat{\beta} (\varepsilon^{new} - (w^{new})^\top (W^\top W)^{-1} W^\top \varepsilon) | \mathcal{I} \right] \\
& := S_1 - S_2 - S_3.
\end{aligned}$$

Below we analyze  $S_1$  to  $S_3$  one by one. With respect to  $S_1$ , we note that  $\hat{\beta} = n^{-1}(n^{-1}X^\top \mathcal{M}_W X + n^{-1}p\tau^{-1}\lambda \mathbb{I})^{-1}X^\top \mathcal{M}_W(U\beta_0 + \varepsilon)$ . Define  $R_X = (n^{-1}X^\top \mathcal{M}_W X + n^{-1}p\tau^{-1}\lambda \mathbb{I})^{-1}$ . By direct calculations, we have

$$\begin{aligned}
& \mathbb{E} \left[ \left( (w^{new})^\top (W^\top W)^{-1} W^\top U \hat{\beta} \right)^2 | \mathcal{I} \right] \asymp \|(W^\top W)^{-1} W^\top U \hat{\beta}\|^2 \\
& \leq 2n^{-2} \|(W^\top W)^{-1} W^\top U R_X X^\top \mathcal{M}_W U \beta_0\|^2 + 2n^{-2} \|(W^\top W)^{-1} W^\top U R_X X^\top \mathcal{M}_W \varepsilon\|^2. \quad (39)
\end{aligned}$$

For the second term in (39), we first note that for any constant  $\lambda > 0$ ,

$$\|R_X X^\top \mathcal{M}_W X R_X\| = \|R_U U^\top \mathcal{M}_W U R_U\| = \frac{n\lambda_1(n^{-1}U^\top \mathcal{M}_W U)}{(\lambda_1(n^{-1}U^\top \mathcal{M}_W U) + n^{-1}p\tau^{-1}\lambda)^2} \asymp_P p^{-1}n^2\tau^2,$$

since  $\|n^{-1}U^\top \mathcal{M}_W U\| \lesssim_P n^{-1}\|U\|^2 \lesssim_P 1 + c_n = o(n^{-1}p\tau^{-1})$  by Lemma 6. Therefore, by Lemma 2 and using inequality  $\text{Tr}(AB) \leq \|A\| \text{Tr}(B)$ , for any  $A = A^\top$  and  $B \geq 0$ , we have

$$\begin{aligned}
& n^{-2} \|(W^\top W)^{-1} W^\top U R_X X^\top \mathcal{M}_W \varepsilon\|^2 \\
& \asymp_P n^{-2} \text{Tr}((W^\top W)^{-1} W^\top U R_X X^\top \mathcal{M}_W X R_X U^\top W (W^\top W)^{-1}) \\
& = n^{-2} \text{Tr}(R_X X^\top \mathcal{M}_W X R_X U^\top W (W^\top W)^{-2} W^\top U) \\
& \leq n^{-2} \|R_X X^\top \mathcal{M}_W X R_X\| \text{Tr}(U^\top W (W^\top W)^{-2} W^\top U) \\
& \lesssim_P p^{-1}\tau^2 \text{Tr}(U^\top W (W^\top W)^{-2} W^\top U) \leq p^{-1}\tau^2 \|U^\top U\| \text{Tr}((W^\top W)^{-1}) = o_P(p^{-1}n\tau^3).
\end{aligned}$$

Similarly, we can prove that the first term in (39) is of order  $o_P(p^{-1}n\tau^3)$ . Therefore, we have

$$\mathbb{E} \left[ \left( (w^{new})^\top (W^\top W)^{-1} W^\top U \hat{\beta} \right)^2 | \mathcal{I} \right] = o_P(p^{-1}n\tau^3). \quad (40)$$

In addition, using the independence of  $w^{new}$  with  $W$ ,  $U$ ,  $\mathcal{I}$ , and  $\beta_0$ , the facts that  $w^{new}$  has bounded variance,  $\|(W^\top W)^{-1} W^\top \Sigma_1^{1/2} x\|^2 \asymp_P \|(W^\top W)^{-1} W^\top \Sigma_1^{1/2}\|_F^2$  by Lemma 2,  $\text{Tr}(W^\top W)^{-1} = o_P(p^{-1}n\tau)$ , and that  $\|\Sigma_2^{1/2} \beta_0\|^2 \asymp_P \tau$ ,

$$\mathbb{E} \left[ \left( (w^{new})^\top (W^\top W)^{-1} W^\top U \beta_0 \right)^2 | \mathcal{I} \right] \asymp \|(W^\top W)^{-1} W^\top U \beta_0\|^2 = \|(W^\top W)^{-1} W^\top \Sigma_1^{1/2} Z \Sigma_2^{1/2} \beta_0\|^2$$

$$\asymp_{\mathbb{P}} \|\Sigma_2^{1/2}\beta_0\|^2 \|(W^\top W)^{-1}W^\top \Sigma_1^{1/2}\|_{\mathbb{F}}^2 \asymp_{\mathbb{P}} o_{\mathbb{P}}(p^{-1}n\tau^2). \quad (41)$$

With (40) and (41),  $\mathbb{E} \left[ ((w^{new})^\top (W^\top W)^{-1}W^\top U(\beta_0 - \hat{\beta}))^2 | \mathcal{I} \right] = o_{\mathbb{P}}(p^{-1}n\tau^2)$ . In addition, since  $u^{new}$  is independent of  $\mathcal{I}$  and  $w^{new}$ , we have  $\mathbb{E}[(u^{new})^\top ((\beta_0 - \hat{\beta}))((w^{new})^\top (W^\top W)^{-1}W^\top U(\beta_0 - \hat{\beta})) | \mathcal{I}] = 0$ . Therefore, we conclude

$$S_1 = \mathbb{E} \left[ ((u^{new})^\top (\beta_0 - \hat{\beta}))^2 | \mathcal{I} \right] + o_{\mathbb{P}}(p^{-1}n\tau^2) = \|\Sigma_2^{1/2}(\beta_0 - \hat{\beta})\|^2 + o_{\mathbb{P}}(p^{-1}n\tau^2).$$

By applying a similar argument, with the use of the independence of  $w^{new}$  with  $W$ ,  $U$ ,  $\mathcal{I}$ , and  $\beta_0$ , as well as the fact that  $w^{new}$  has bounded variance, it can be shown that

$$S_2 = \mathbb{E} \left[ ((u^{new})^\top \beta_0)^2 | \mathcal{I} \right] + o_{\mathbb{P}}(p^{-1}n\tau^2) = \|\Sigma_2^{1/2}\beta_0\|^2 + o_{\mathbb{P}}(p^{-1}n\tau^2).$$

Finally we bound  $S_3$ . Since  $u^{new}$  and  $\varepsilon^{new}$  are mean zero, mutually independent, and independent of  $\mathcal{I}$ , along with Eq. (40), Lemma 2 and Cauchy-Schwartz inequality, we have

$$\begin{aligned} |S_3| &= \left| 2\mathbb{E} \left[ (w^{new})^\top (W^\top W)^{-1}W^\top U\hat{\beta} (w^{new})^\top (W^\top W)^{-1}W^\top \varepsilon | \mathcal{I} \right] \right| \\ &\leq 2 \left( \mathbb{E} \left[ ((w^{new})^\top (W^\top W)^{-1}W^\top U\hat{\beta})^2 | \mathcal{I} \right] \right)^{1/2} \left( \mathbb{E} \left[ ((w^{new})^\top (W^\top W)^{-1}W^\top \varepsilon)^2 | \mathcal{I} \right] \right)^{1/2} \\ &\asymp_{\mathbb{P}} o_{\mathbb{P}}(p^{-1/2}n^{1/2}\tau^{3/2}) (\text{Tr}(W(W^\top W)^{-2}W^\top))^{1/2} = o_{\mathbb{P}}(p^{-1}n\tau^2). \end{aligned}$$

In sum, we have  $S_1 - S_2 - S_3 = \|\Sigma_2^{1/2}(\beta_0 - \hat{\beta})\|^2 - \|\Sigma_2^{1/2}\beta_0\|^2 + o_{\mathbb{P}}(p^{-1}n\tau^2)$ . Next we prove,

$$pn^{-1}\tau^{-2}(\|\Sigma_2^{1/2}(\hat{\beta} - \beta_0)\|^2 - \|\Sigma_2^{1/2}\beta_0\|^2) \xrightarrow{\mathbb{P}} \alpha^*. \quad (42)$$

By Theorem 2,  $\tilde{\beta} := n^{-1}\tilde{R}_U U^\top (U\beta_0 + \varepsilon)$  satisfies (42) with  $\hat{\beta}$  being replaced by  $\tilde{\beta}$ , where  $\tilde{R}_U := (n^{-1}U^\top U + n^{-1}p\tau^{-1}\lambda\mathbb{I})^{-1}$ . Given that  $\|\Sigma_2^{1/2}(\beta_0 - \hat{\beta})\|^2 = \|\Sigma_2^{1/2}(\beta_0 - \tilde{\beta}) + \Sigma_2^{1/2}(\hat{\beta} - \tilde{\beta})\|^2$  and that  $\|\Sigma_2^{1/2}(\beta_0 - \tilde{\beta})\| \asymp_{\mathbb{P}} \|\Sigma_2^{1/2}\beta_0\| \asymp_{\mathbb{P}} \tau^{1/2}$ , it is easy to verify that (42) follows from

$$\|\hat{\beta} - \tilde{\beta}\|^2 = o(n^2 p^{-2} \tau^3), \quad (43)$$

which is given by Lemma 26. □

## A.7 Proof of Proposition 2

*Proof.* Let  $\tau_0 = \sqrt{n\tau} \leq \sqrt{\log n}/10$  and  $\lambda_0 = \sqrt{n\lambda}/2$ . Note that when  $X = \mathbb{I}$ , Lasso has a closed form solution  $\hat{\beta}_l(\lambda) = (|\varepsilon_1 + \tau_0| - \lambda_0)_+ \text{sgn}(\varepsilon_1 + \tau_0), (|\varepsilon_2| - \lambda_0)_+ \text{sgn}(\varepsilon_2), \dots, (|\varepsilon_n| -$



$\lambda_0)_+ \text{sgn}(\varepsilon_n))^\top$ . Therefore,  $\|\hat{\beta}_l(\lambda) - \beta_0\|^2 - \|\beta_0\|^2 = ((|\varepsilon_1 + \tau_0| - \lambda_0)_+ - \tau_0)^2 - \tau_0^2 + \sum_{i=2}^n ((|\varepsilon_i| - \lambda_0)_+)^2$ . Assume for now that  $\max_{1 \leq i \leq n} |\varepsilon_i| \geq |\varepsilon_1| + 2\tau_0$ . Let  $i_0 = \arg \max |\varepsilon_i|$ , then we have

$$\|\hat{\beta}_l(\lambda) - \beta_0\|^2 - \|\beta_0\|^2 \geq ((|\varepsilon_1 + \tau_0| - \lambda_0)_+ - \tau_0)^2 - \tau_0^2 + ((|\varepsilon_{i_0}| - \lambda_0)_+)^2 := Q$$

Consider the following three cases for  $\varepsilon_1$ : (i)  $\varepsilon_1 \in [-\lambda_0 - \tau_0, \lambda_0 - \tau_0]$ :  $Q = 0 + ((|\varepsilon_{i_0}| - \lambda_0)_+)^2 \geq 0$ ; (ii)  $\varepsilon_1 \in (-\infty, -\lambda_0 - \tau_0)$ :  $Q \geq -\tau_0^2 + ((|\varepsilon_1| + 2\tau_0 - \lambda_0)_+)^2 \geq -\tau_0^2 + 9\tau_0^2 \geq 0$ ; (iii)  $\varepsilon_1 \in (\lambda_0 - \tau_0, \infty)$ :  $Q \geq -\tau_0^2 + ((|\varepsilon_1| + 2\tau_0 - \lambda_0)_+)^2 \geq -\tau_0^2 + \tau_0^2 = 0$ . Therefore, under the event that  $\max_{1 \leq i \leq n} |\varepsilon_i| \geq |\varepsilon_1| + 2\tau_0$ , it holds that  $\|\hat{\beta}_l(\lambda) - \beta_0\|^2 - \|\beta_0\|^2 \geq 0$ . Now we evaluate the probability of this event. Note that

$$\mathbb{P}(\max_{1 \leq i \leq n} |\varepsilon_i| \leq u) = (\mathbb{P}(|\varepsilon_i| \leq u))^n = \left( \text{erf} \left( \frac{u}{\sqrt{2}} \right) \right)^n \leq \exp \left\{ -\frac{n}{2} \exp \left\{ -\frac{2}{\pi} u^2 \right\} \right\},$$

where  $\text{erf}(\cdot)$  represents the Gauss error function. The last inequality uses the fact that  $(\text{erf}(x))^2 \leq 1 - \exp(-4x^2/\pi)$  and  $1 + x \leq e^x$ . Reparametrizing  $u$  in terms of  $\delta$  by solving  $\delta = \exp \left\{ -\frac{n}{2} \exp \left\{ -\frac{2}{\pi} u^2 \right\} \right\}$ , we obtain that, with probability at least  $1 - \delta$ :

$$\max_{1 \leq i \leq n} |\varepsilon_i| \geq \sqrt{\frac{\pi}{2} \log \frac{n}{2} - \frac{\pi}{2} \log \log \frac{1}{\delta}}. \quad (44)$$

Choosing  $\delta = n^{-1}$ , the event  $\mathcal{C} = \left\{ \max_{1 \leq i \leq n} |\varepsilon_i| \geq \sqrt{\frac{\pi}{2} \log \frac{n}{2} - \frac{\pi}{2} \log \log n} \right\}$  happens with probability at least  $1 - n^{-1}$ . Setting  $u = \sqrt{\frac{\pi}{2} \log \frac{n}{2} - \frac{\pi}{2} \log \log n} - 2\tau_0$ . There exists  $n_0 \in \mathbb{N}$ , when  $n \geq n_0$ ,  $u \geq \sqrt{1.7 \log n} - 0.2\sqrt{\log n} \geq \sqrt{\log n}$ . By Mills' inequalities, when  $n \geq n_0$ ,

$$\begin{aligned} \mathbb{P} \left( \max_{1 \leq i \leq n} |\varepsilon_i| - 2\tau_0 \geq |\varepsilon_1| \middle| \mathcal{C} \right) &\geq \mathbb{P} \left( |\varepsilon_1| \leq \sqrt{\frac{\pi}{2} \log \frac{n}{2} - \frac{\pi}{2} \log \log n} - 2\tau_0 \right) \\ &= \mathbb{P} (|\varepsilon_1| \leq u) \geq 1 - \sqrt{\frac{2}{\pi}} \frac{\exp(-u^2/2)}{u} \geq 1 - \sqrt{\frac{2}{\pi \log n}} n^{-1/2}. \end{aligned}$$

There exists  $n_1 \geq 1$ , when  $n \geq n_1$ ,  $(1 - \sqrt{\frac{2}{\pi \log n}} n^{-1/2})(1 - \frac{1}{n}) \geq 1 - n^{-1/2}$ . Therefore,

$$\mathbb{P} \left( \max_{1 \leq i \leq n} |\varepsilon_i| - 2\tau_0 \geq |\varepsilon_1| \right) \geq \mathbb{P}(\mathcal{C}) \cdot \mathbb{P} \left( \max_{1 \leq i \leq n} |\varepsilon_i| - 2\tau_0 \geq |\varepsilon_1| \middle| \mathcal{C} \right) \geq 1 - n^{-1/2},$$

for  $n \geq \max(n_0, n_1)$ , which implies Eq. (18). For Ridge, since  $\hat{\beta}_r(\lambda) = (1 + n\lambda)^{-1}y$ , we have

$$\|\hat{\beta}_r(\lambda) - \beta_0\|^2 - \|\beta_0\|^2 = \frac{-2n\lambda\varepsilon_1\tau_0 - (1 + 2n\lambda)\tau_0^2 + \sum_{i=1}^n \varepsilon_i^2}{(1 + n\lambda)^2}.$$

Observe that with probability equal to 0.5,  $\varepsilon_1\tau_0 \geq 0$ . Under this event,  $\|\hat{\beta}_r(\lambda) - \beta_0\|^2 - \|\beta_0\|^2 \leq \frac{-(1+2n\lambda)\tau_0^2 + \sum_{i=1}^n \varepsilon_i^2}{(1+n\lambda)^2}$ . Therefore,  $\|\hat{\beta}_r(\lambda) - \beta_0\|^2 - \|\beta_0\|^2 < 0$  as long as  $\lambda > (2n)^{-1}(-1 + \tau_0^{-2} \sum_{i=1}^n \varepsilon_i^2)$ .  $\square$

# Online Appendix of Can Machines Learn Weak Signals?

Zhouyu Shen\*

Dacheng Xiu†

Booth School of Business

Booth School of Business

University of Chicago

University of Chicago and NBER

January 21, 2025

## Abstract

Appendix [A](#) presents additional results from Monte Carlo simulations. Appendix [B](#) discusses the selection of tuning parameters. Appendix [C](#) is devoted to the exposition of technical lemmas along with their corresponding proofs.

## A Supplemental Simulation Results

### A.1 Additional Simulations with Fixed Tunings

In the simulation for the main paper, cross-validation is applied for Ridge and Lasso. In this section, we verify our theories with manually selected  $\lambda_n$ . We also experiment with two sample sizes,  $n = 500$  and  $n = 2,500$ , while maintaining  $p/n = 3/5$ . We fix  $q = 0.2$  and  $R^2 = 5\%$ . In the case of Ridge regression, we set  $\lambda$  as 0.5, 1 and 2, where  $\lambda = 1$  corresponding to the optimal tuning. The histograms of relative prediction error are presented in Figure [A1](#).

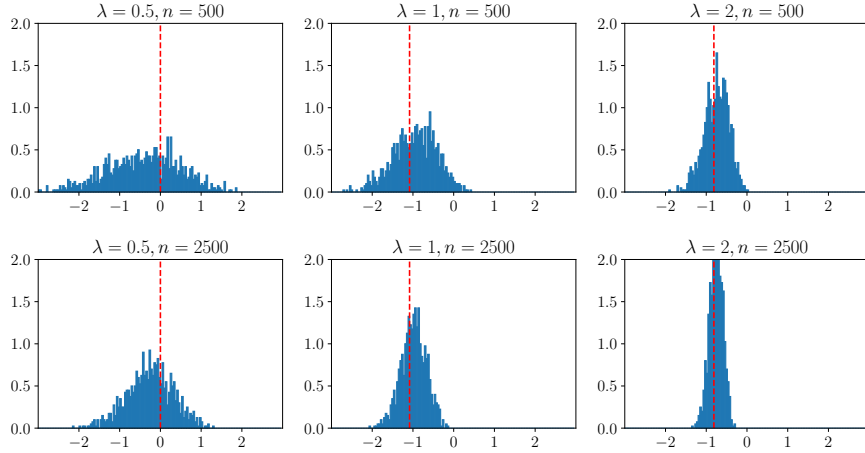
Several noteworthy observations can be made from these histograms. First, across all plots, the probability mass is concentrated around the red vertical line. As the sample size increases from 500 to 2,500 (and dimension increases from 300 to 1,500), the histograms become increasingly concentrated. This aligns with our theory, which predicts that the

---

\*Address: 5807 S Woodlawn Avenue, Chicago, IL 60637 USA. Email: [zshen10@chicagobooth.edu](mailto:zshen10@chicagobooth.edu).

†Address: 5807 S Woodlawn Avenue, Chicago, IL 60637, USA. Email: [dacheng.xiu@chicagobooth.edu](mailto:dacheng.xiu@chicagobooth.edu).

Figure A1: Simulation Results for Ridge with Fixed Tuning Parameters



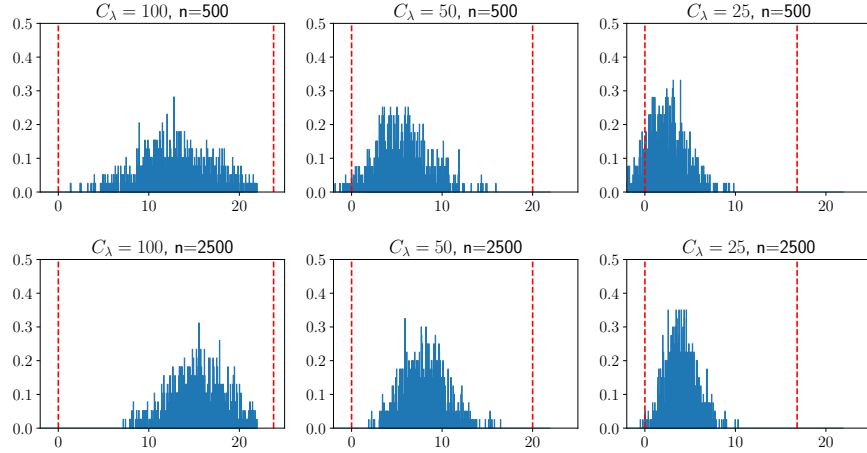
Note: The histograms depict the relative prediction error  $\Delta(\hat{\beta}_r(\lambda_n))$  following equation (8) across 1,000 Monte Carlo samples. We consider two different sample sizes ( $n = 500$  and  $n = 2,500$ ) and examine three different values of  $\lambda$ , where  $\lambda = 0.5, 1$ , and  $2$ . Notably,  $\lambda = 1$  represents the optimal tuning parameter. The red dashed line indicates the values of  $\alpha^*$ .

relative prediction error converges in probability to the limit  $\alpha^*$  as the sample size grows. Second, the value of  $\alpha^*$  corresponding to the optimal tuning parameter  $\lambda = 1$  is the smallest. This is because the optimal Ridge estimator achieves the smallest prediction error. Moreover, almost all the probability mass corresponding to the optimal Ridge estimator is situated on the negative side of the x-axis, indicating that this estimator outperforms the zero estimator with high probability. Third, when  $\lambda = 0.5$ , it results in the worst performance, with a large portion of the probability mass on the positive side of zero. In contrast, for  $\lambda = 2$ ,  $\alpha^*$  gets closer to zero, and the variance of the relative prediction error decreases. This behavior is due to the increasing amount of penalization, which ultimately drives the estimator towards zero, and in turn,  $\alpha^*$  towards zero as well.

In contrast to the results obtained for Ridge regression, our theoretical framework does not provide a precise error limit for Lasso. Instead, Theorem 4 offers high probability bounds on relative prediction errors. Figure A2 displays histograms of these errors for various tuning parameters and sample sizes, accompanied by two red vertical lines in each plot representing the lower and upper bounds,  $c_\alpha$  and  $C_\alpha$ .

These plots yield several interesting findings. First, as the sample size increases, we observe that the probability mass becomes more concentrated and largely falls within the intervals defined by the bounds. Second, regardless of the tuning parameter values, Lasso

Figure A2: Simulation Results for Lasso with Fixed Tuning Parameters



Note: The histograms depict the relative prediction error  $\Delta(\hat{\beta}_l(\lambda_n))$  following equation (8) across 1,000 Monte Carlo samples. We consider two different sample sizes ( $n = 500$  and  $n = 2,500$ ) and examine three different values of  $C_\lambda$ . The two dashed lines in each figure indicate the values of  $c_\alpha$  and  $C_\alpha$  that are solutions to (10).

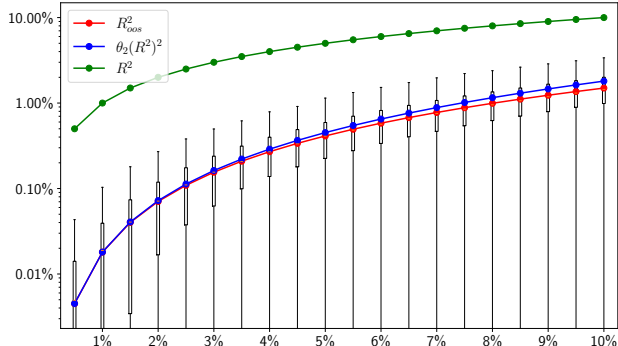
consistently underperforms the zero estimator in almost all samples when the sample size is large. Third, as the tuning parameter increases (indicated by a decrease in  $C_\lambda$ ), both the lower and upper bounds approach zero. This behavior is a consequence of the increased regularization, which, in turn, steers the estimator closer to zero. In the end, Lasso becomes identical to the zero estimator.

## A.2 Out-of-sample $R^2$

Continuing our investigation in the main text, we conduct an experiment to analyze  $R_{\text{OOS}}^2$  based on the optimal Ridge. Proposition 1 describes the expected asymptotic behavior of  $R_{\text{OOS}}^2$ . To empirically test this, we implement the optimal Ridge, setting  $\lambda = 1$ , on a training dataset comprising  $n = 500$  observations. We then calculate  $R_{\text{OOS}}^2$  based on predictions for a separate test dataset of size  $n_{\text{OOS}} = 10,000$ . The comparative analysis between the population  $R^2$ , the empirically estimated  $R_{\text{OOS}}^2$ , and the theoretically derived limit of  $R_{\text{OOS}}^2$  is illustrated in Figure A3. For a clearer visual presentation, we apply a logarithmic transformation to the y-axis. We vary  $\tau$  to compare against a range of population  $R^2$  values from 0.5% to 10% on the x-axis. The red line represents the average  $R_{\text{OOS}}^2$  over 1,000 Monte Carlo simulations. Additionally, we draw boxplots to describe the distributions of  $R_{\text{OOS}}^2$  across these simulations. The theoretical limit, expressed as  $p^{-1}n\theta_2(R^2)^2$ , is traced by the blue

line, and the green line illustrates the population  $R^2$ , which would align with a 45-degree line on a standard scale. Notably, in this weak signal setting, the population  $R^2$  significantly surpasses the empirically achievable  $R_{\text{OOS}}^2$ . Furthermore, the close alignment between the red and blue lines, particularly for scenarios with small  $R^2$  values, substantiates our theoretical predictions.

Figure A3: Out-of-Sample  $R^2$  for Optimal Ridge in Linear DGPs



Note: The figure presents boxplots showing the distributions of  $R_{\text{OOS}}^2$  for optimal Ridge regression ( $\lambda = 1$ ) over 1,000 Monte Carlo repetitions, with  $n = 500$ ,  $p = 300$ ,  $q = 0.2$ , and  $n_{\text{OOS}} = 10,000$ . We explore a range of population  $R^2$  values, from 0.5% to 10% in increments of 0.5% by adjusting  $\tau$ . The plot features red, blue, and green lines to represent the average  $R_{\text{OOS}}^2$  over Monte Carlo samples, the theoretical limit as given by Proposition 1, and the population  $R^2$ . In this plot, we employ a logarithmic scale for the y-axis. Without the logarithmic transformation, the green line would align with a 45-degree line. Additionally, the lower boundaries of the boxplots surpass the axis limits in instances where the  $R_{\text{OOS}}^2$  values are negative.

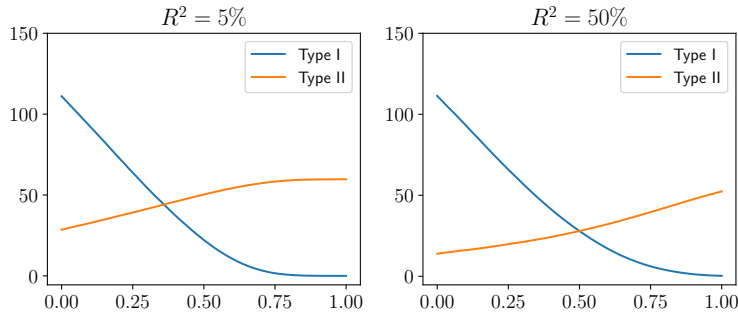
### A.3 Why Lasso Fails?

A plausible explanation for the Lasso’s suboptimal performance with weak signals is its difficulty in distinguishing between genuine and spurious signals. The failure to identify genuine weak signals has a minor impact on Lasso’s performance relative to the zero estimator, which does not utilize any true signals. Hence, the primary challenge for the Lasso lies in its failure to adequately filter out irrelevant signals. This issue could be addressed with a sufficiently large tuning parameter. However, our theory indicates that only when the penalty is so substantial that the Lasso effectively becomes equivalent to the zero estimator does it apply an adequate penalty.

To empirically explore this issue, we quantify Type I and Type II errors in simulations of Lasso’s selection relative to its tuning parameter  $\lambda_n$ . The findings are presented in Figure

A4. Considering our previous discussion, Type I errors represent a significant cost for Lasso. Indeed, a considerable portion of the variables selected by Lasso are incorrectly deemed genuine when  $\lambda_n$  is small. As  $\lambda_n$  increases, Type I errors decrease, enhancing Lasso’s performance. Meanwhile, Type II errors persist and eventually converge to the number of non-zero betas in the DGP.

Figure A4: Lasso’s Type I and Type II Errors



Note: The plots compare average Type I and Type II errors of Lasso using the linear DGP following equation (1) over 1,000 Monte Carlo samples. Two population  $R^2$  are considered:  $R^2 = 5\%$  (left panel) and  $R^2 = 50\%$  (right panel). The horizontal axis of each plot represents the logarithm of  $\lambda_n$ , spanning a range from 0 to 1. The vertical axis measures the count of errors incurred while testing the null hypothesis  $\mathbb{H}_{0,i} : \beta_i = 0$ .

#### A.4 Robustness Check

In our subsequent series of experiments, we intentionally deviate from the assumptions originally established during the development of our theoretical framework. This deviation is aimed at evaluating the robustness and generalizability of our theoretical predictions beyond their premises and initial parameters. To facilitate this evaluation, we introduce specific modifications to the baseline configuration along three key dimensions: First, we adjust  $(R^2, q)$ , exploring more extreme sparsity levels and reducing signal strength accordingly compared to the settings in the main text. Second, we increase the ratio  $p/n$  to 2 by increasing  $p$  while maintaining  $n$ , making it more challenging for both Ridge and Lasso to capture the underlying signals. Third, we modify the distribution of  $Z$  from standard Gaussian to a t-distribution characterized by four degrees of freedom, and with a mean of zero and a variance of one. In addition, we introduce heteroscedasticity into the error distribution, following the configuration outlined by Giannone et al. (2022). The error term’s variance is defined by the function  $\sigma^2 \exp(\alpha X_i^\top \delta / \sqrt{\sum_{i=1}^n (X_i^\top \delta)^2 / n})$  with  $\alpha = 0.5$ . Here,  $X_i$  represents the  $i$ -th row of  $X$ .  $\sigma$  serves as a scaling parameter to standardize the variance and match  $\sigma_\varepsilon^2 = 1$ . The

vector  $\delta$  is a  $p \times 1$  vector with zero elements in the same positions as the zero elements of  $\beta_0$ , while non-zero elements are drawn from a standard Gaussian distribution.

Table A1 compare the summary statistics for various cases under consideration. In Case I, when  $q$  is small, the performance of the Lasso estimator improves relative to the baseline scenario (reproduced from Table 1 for ease of comparison). This improvement is evident at  $R^2 = 5\%$  for all levels of  $q$ , as the Q1 values become negative, indicating that Lasso surpasses zero in predictive accuracy for a larger proportion of Monte Carlo repetitions. However, as  $R^2$  is further reduced to 2%, Lasso once again becomes falls below the performance of zero. In contrast, Ridge’s performance remains largely unaffected by changes in sparsity levels. As expected, its performance deteriorates in finite samples as the signal strength weakens (i.e., as  $R^2$  decreases). Nonetheless, Ridge continues to outperform Lasso, although its relative advantage over the zero estimator diminishes. The theoretical support for these observations is discussed in Section 2.9. In Case II, we observe the increased ratio of  $p/n$  does not affect our conclusion. Case III demonstrates the robustness of our theoretical findings, as it aligns closely with the baseline scenario despite variations in distributional assumptions.

Table A1: Robustness Analysis of Ridge and Lasso in Alternative DGPs

	$q$	$R^2$ (%)	Lasso				Ridge			
			Q1	Q2	Q3	#Zero	Q1	Q2	Q3	#Zero
Case I	0.20	5%	-0.127	0.000	0.521	360	-0.992	-0.501	-0.129	97
	0.10	5%	-0.871	0.000	0.187	327	-0.981	-0.475	-0.077	113
	0.10	2%	0.000	0.000	3.435	493	-0.622	0.000	0.440	237
	0.05	5%	-2.688	-0.305	0.000	255	-1.037	-0.387	0.000	130
	0.05	2%	0.000	0.000	2.948	473	-0.642	0.000	0.426	238
	0.02	5%	-6.542	-2.050	0.000	215	-1.304	-0.230	0.000	149
	0.02	2%	0.000	0.000	1.695	432	-0.605	0.000	0.625	254
Case II	0.20	5%	0.000	0.000	3.228	470	-0.768	-0.416	0.000	183
Case III	0.20	5%	0.000	0.000	0.591	392	-0.848	-0.384	0.000	129

Note: The table illustrate the summary statistics of relative prediction error  $\Delta(\hat{\beta}(\hat{\lambda}_n^{K-CV}))$  for Ridge and Lasso based on 1,000 Monte Carlo samples. We explore several distinct DGPs, each involving the alteration of a specific condition. In Case I, we try a series of different values of  $R^2$  and  $q$ . In Case II, we adjust  $n/p$  to 0.5. In Case III, we introduce t-distributed covariates with heterogeneous variance of  $\varepsilon$ . The benchmark DGP adheres to the following specifications:  $n = 500$ ,  $p = 300$ ,  $p/n = 3/5$ , and complies with Assumptions 1 and 2. 10-fold cross-validation is used throughout these experiments.

## B Choice of Tuning Parameters

In this section, we discuss the selection of tuning parameters for implementing machine learning methods. The selection process aims to balance performance and cost while ensuring



fair method comparisons.

For Ridge and Lasso, each with one tuning parameter, we use the glmnet package, which optimizes it via ten-fold CV. The grid is adaptively selected for efficiency. Our implementation of RF involves three tuning parameters: the depth of each individual tree, the number of randomly selected features used in each tree split, and the proportion of sample data used for bootstrap.<sup>1</sup> <sup>2</sup> For GBRT, we tune tree depth, number of trees, and learning rate.<sup>3</sup> In the case of NNs, we adhere to a uniform architectural choice across our analyses, featuring a single hidden layer. The number of neurons in this hidden layer is approximately equal to the square root of the total number of neurons in the input layer, aligning the architecture with the complexity and dimensions of the dataset. By not tuning the NN architecture extensively, we streamline the model selection process while retaining adequate complexity for effective learning. For the remaining tuning parameters in trees and NNs, we select suitable ranges based on model performance from the cross-validation step. A critical element in selecting our grid is to ensure that the optimal tuning parameters are situated within the median range of the grid. The details regarding model configuration and tuning parameters for empirical studies are provided in Table B2.<sup>4</sup>

## C Technical Lemmas and Their Proofs

For completeness, the following section introduces a collection of lemmas, including proofs for some. We start with the Convex Gaussian Min-max Theorem (CGMT), a pivotal theorem to our proof. For a detailed exposition of its proof, we direct readers to the work of [Thrapoulidis et al. \(2015\)](#). The CGMT pertains to the following optimization problems:

$$\Phi(G) := \min_{w \in \mathcal{S}_w} \max_{u \in \mathcal{S}_u} u^\top G w + \psi(w, u), \text{ and } \phi(g, h) := \min_{w \in \mathcal{S}_w} \max_{u \in \mathcal{S}_u} \|w\| g^\top u - \|u\| h^\top w + \psi(w, u),$$

where  $G \in \mathbb{R}^{m \times n}$ ,  $g \in \mathbb{R}^m$ ,  $h \in \mathbb{R}^n$ ,  $\mathcal{S}_w \subset \mathbb{R}^n$ ,  $\mathcal{S}_u \subset \mathbb{R}^m$ , and  $\psi : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ .

---

<sup>1</sup>This procedure is known as subbagging, which helps address weak signals. In linear regressions, [LeJeune et al. \(2020\)](#) show that the asymptotic risk of subbagging least squares matches that of Ridge regression.

<sup>2</sup>In our simulations, tree depth varies from 5 to 20, selected features from 10 to 300, and bootstrap sample proportion from 0.1 to 0.2. The RF ensemble size is fixed at 5,000 trees, as 10,000 offers no significant improvement.

<sup>3</sup>The learning rate varies from 0.001 to 0.5, tree depth from 1 to 6, and the maximum number of trees is 100, though training usually stops earlier, reflecting GBRT’s preference for shallower and fewer trees.

<sup>4</sup>We follow the same approach for the empirical analysis, except for Finance 2. Due to its scale and computational constraints, we use two-fold CV and a narrower grid:  $\log(\lambda)$  ranges from 6 to 7 for Ridge and from -3.5 to -2.5 for Lasso. We ensure that the optimal tuning parameters fall within the central range of these specified grids.

Table B2: Model Configuration for Machine Learning Methods

	RF	GBRT	NN( $\ell_2$ )	NN( $\ell_1$ )
Finance 1	depth=1~20 #trees=500 #features=1~15 %samples=0.05~1	depth=1~5 #trees=1~10 lr $\in$ {0.01,0.02, 0.05,0.1,0.2,0.5,1}	architecture~{16,4,1} batch size=16 (lr,epochs)={{(0.1,5), (0.01,50), (0.0025,200)}} log( $\lambda$ ) $\in$ [-2, 1]	architecture~{16,4,1} batch size=16 (lr,epochs)={{(0.4,1), (0.08,5), (0.02,20)}} log( $\lambda$ ) $\in$ [-2, 1]
Finance 2	depth=2~12 #trees=500 #features $\in$ {1, 2,3,5} %samples=0.5~1	depth=1~6 #trees=10~400 lr $\in$ {0.0001,0.001, 0.01,0.02,0.05}	architecture~{920,32,1} batch size=10000 (lr,epochs)={{(0.5,2), (0.1,10), (0.067,15)}} log( $\lambda$ ) $\in$ [-4, 0]	architecture~{920,32,1} batch size=10000 (lr,epochs)={{(0.5,2), (0.2,5), (0.067,15), (0.05,20),(0.04,25)}} log( $\lambda$ ) $\in$ [-5, -3]
Macro 1	depth=5~50 #trees=500 #features=1~60 %samples=0.5~1	depth=1~5 #trees=1~600 lr $\in$ {0.005,0.01, 0.02,0.05,0.1, 0.2,0.5}	architecture~{119,8,1} batch size=16 (lr,epochs)={{(0.008,10), (0.004,20), (0.002,40), (0.0008,100), (0.0005,160)}} log( $\lambda$ ) $\in$ [-2, 2]	architecture~{119,8,1} batch size=16 (lr,epochs)={{(0.05,2), (0.02,5),(0.01,10), (0.005,20), (0.002,50)}} log( $\lambda$ ) $\in$ [-10.5, 1.5]
Macro 1b	depth=5~50 #trees=500 #features=5~100 %samples=0.5~1	depth=1~5 #trees=1~200 lr $\in$ {0.01,0.02, 0.05,0.1,0.2,0.5}	architecture~{119,8,1} batch size=16 (lr,epochs)={{(0.024,75), (0.012,150), (0.006,300) (0.003,600)}} log( $\lambda$ ) $\in$ [-1, -0.5]	architecture~{119,8,1} batch size=16 (lr,epochs)={{(0.004,25), (0.002,50),(0.001,100), (0.0005,200)}} log( $\lambda$ ) $\in$ [2, 3]
Macro 2	depth=1~30 #trees=500 #features=1~60 %samples=0.5~1	depth=1~10 #trees=1~500 lr $\in$ {0.01,0.02, 0.05,0.1,0.2,0.5}	architecture~{61,8,1} batch size=16 (lr,epochs)={{(0.02,50), (0.005,200), (0.00125,800)}} log( $\lambda$ ) $\in$ [-3, -3]	architecture~{61,8,1} batch size=16 (lr,epochs)={{(0.05,20), (0.02,50), (0.002,500)}} log( $\lambda$ ) $\in$ [-7, 4]
Micro 1	depth=1~20 #trees=500 #features=1~20 %samples=0.005~0.5	depth=1~5 #trees=1~20 lr $\in$ { $10^{-15}$ , $10^{-14}$ , ...,0.05,0.1,0.2,0.5}	architecture~{297,16,1} batch size=16 (lr,epochs)={{(0.1,1), (0.01,10), (0.001,100), (0.0001,1000), (0.00005,2000)}} log( $\lambda$ ) $\in$ [-11, -7]	architecture~{297,16,1} batch size=16 (lr,epochs)={{(0.4,1), (0.2,2), (0.08,5), (0.04,10), (0.02,20)}} log( $\lambda$ ) $\in$ [-8, 0]
Micro 2	depth=5~30 #trees=500 #features=1~30 %samples=0.5~1.0	depth=1~6 #trees=1~30 lr $\in$ {0.05,0.1, 0.15,...,1}	architecture~{217,16,1} batch size=16 (lr,epochs)={{(0.1,1),(0.02,5), (0.01,10),(0.005,20)}} log( $\lambda$ ) $\in$ [-10, -6]	architecture~{217,16,1} batch size=16 (lr,epochs)={{(0.1,1),(0.01,10), (0.001,100), (0.0001,1000)}} log( $\lambda$ ) $\in$ [-12, -9]
Micro 2b	depth=1~50 #trees=500 #features=1~3 %samples=0.5~1	depth=1~10 #trees=1~50 lr $\in$ { $10^{-10}$ , $10^{-9}$ , ...,0.1,0.2,0.5,1}	architecture~{215,16,1} batch size=16 (lr,epochs)={{(0.04,5), (0.02,10), (0.01,20)}} log( $\lambda$ ) $\in$ [0, 2]	architecture~{215,16,1} batch size=16 (lr,epochs)={{(0.01,50), (0.005,100), (0.0025,200)}} log( $\lambda$ ) $\in$ [0, 2]

Note: The table reports the range of tuning parameters for RF, GBRT, and NNs, as well as the architecture of NNs applied across six datasets. For RF, we fix the number of trees at #trees= 500, and tune three other parameters: the depth of the tree (depth), the number of features (#features), and the ratio of bootstrapped samples (%samples) within a predefined grid. In the case of GBRT, we tune depth and #trees, and the learning rate (lr). For NNs, we adopt a fixed model architecture, denoted by the number of neurons in each layer indicated in brackets. Additionally, we fix the batch size for SGD and focus on jointly tuning the learning rate (lr) and the number of epochs (epochs), as well as the  $\ell_1$ - or  $\ell_2$ -penalty parameter (log( $\lambda$ )).

**Lemma 1** (CGMT). *Suppose that  $\mathcal{S}_w$  and  $\mathcal{S}_u$  are compact sets,  $\psi$  is continuous on  $\mathcal{S}_w \times \mathcal{S}_u$ , and the entries of  $G$ ,  $g$ , and  $h$  are i.i.d. Gaussian. Then we have  $\mathbb{P}(\Phi(G) < c) \leq 2\mathbb{P}(\phi(g, h) \leq c)$ ,  $\forall c \in \mathbb{R}$ . Moreover, if  $\mathcal{S}_w$  and  $\mathcal{S}_u$  are convex sets, and  $\psi$  is convex-concave on  $\mathcal{S}_w \times \mathcal{S}_u$ , then  $\mathbb{P}(\Phi(G) > c) \leq 2\mathbb{P}(\phi(g, h) \geq c)$ ,  $\forall c \in \mathbb{R}$ .*

**Lemma 2** (Lemma B.26 from [Bai and Silverstein \(2009\)](#)). *Let  $x = (x_1, \dots, x_n)^\top$  be a random vector of i.i.d. entries. Assume that  $\mathbb{E}x_i = 0$ ,  $\mathbb{E}x_i^2 = 1$ , and  $\mathbb{E}x_i^4 \leq v_4$ . Then, for any  $A \in \mathbb{R}^{n \times n}$ , it holds that  $x^\top Ax - \text{Tr}(A) = O_{\mathbb{P}}(\sqrt{v_4 \text{Tr}(AA^\top)})$ .*

**Lemma 3.** *Let  $x = (x_1, \dots, x_n)^\top$  and  $y = (y_1, \dots, y_m)^\top$  be two independent random vectors with i.i.d. entries. Assume that each element has a mean of zero and a variance of one. Then, for any  $A \in \mathbb{R}^{n \times m}$ , it holds that  $x^\top Ay = O_{\mathbb{P}}(\sqrt{\text{Tr}(AA^\top)})$ .*

*Proof.* The conclusion follows from the fact that  $\mathbb{E}(x^\top Ay)^2 = \text{Tr}(AA^\top)$ .  $\square$

The lemma below pertains to the Neumann series. See [Meyer \(2000\)](#) for a detailed proof.

**Lemma 4.** *If  $A$  is a square matrix with  $\|A\| < 1$ , then  $\mathbb{I} - A$  is nonsingular and  $(\mathbb{I} - A)^{-1} = \sum_{k=0}^{\infty} A^k$ . As a consequence,  $\|(\mathbb{I} - A)^{-1} - \sum_{k=0}^{\ell} A^k\| \leq \sum_{k=\ell+1}^{\infty} \|A\|^k = \|A\|^{\ell+1}/(1 - \|A\|)$ .*

**Lemma 5.** *Assume  $x = (x_1, \dots, x_n)^\top$  and  $y = (y_1, \dots, y_p)^\top$  are two independent random vectors with i.i.d. sub-exponential random variables with their sub-exponential norm bounded by  $K$ . Then for any  $A \in \mathbb{R}^{n \times n}$  and  $B \in \mathbb{R}^{n \times p}$ , there exists a constant  $c > 0$  such that*

$$\mathbb{P}(|x^\top Ax - \mathbb{E}x^\top Ax| \geq t) \leq 2 \exp\left(-c \min\left\{\frac{t^2}{K^4 \|A\|_{\mathbb{F}}^2}, \frac{t^{1/2}}{K \|A\|^{1/2}}\right\}\right), \quad (\text{C1})$$

$$\mathbb{P}(|x^\top By| \geq t) \leq 2 \exp\left(-c \min\left\{\frac{t^2}{K^4 \|B\|_{\mathbb{F}}^2}, \frac{t^{1/2}}{K \|B\|^{1/2}}\right\}\right). \quad (\text{C2})$$

*Proof.* Inequality (C1) is given by Proposition 1.1 presented in [Götze et al. \(2021\)](#) for the case of symmetric  $A$ . To extend it for the asymmetric case, we use the identity that  $x^\top Ax = x^\top (A + A^\top)x/2$ , which allows us to apply (C1) to  $(A + A^\top)/2$ . Using triangle inequalities, we have  $\|(A + A^\top)/2\|_{\mathbb{F}}^2 \leq \|A\|_{\mathbb{F}}^2$  and  $\|(A + A^\top)/2\|^{1/2} \leq \|A\|^{1/2}$ . Thus, (C1) holds for asymmetric  $A$ . For (C2), let  $z = (x^\top, y^\top)^\top$  and  $C = \begin{pmatrix} 0_{n \times n} & B \\ 0_{p \times n} & 0_{p \times p} \end{pmatrix}$ . By (C1), we obtain

$$\begin{aligned} \mathbb{P}(|x^\top By| \geq t) &= \mathbb{P}(|z^\top Cz| \geq t) \leq 2 \exp\left(-c \min\left\{\frac{t^2}{K^4 \|C\|_{\mathbb{F}}^2}, \frac{t^{1/2}}{K \|C\|^{1/2}}\right\}\right) \\ &= 2 \exp\left(-c \min\left\{\frac{t^2}{K^4 \|B\|_{\mathbb{F}}^2}, \frac{t^{1/2}}{K \|B\|^{1/2}}\right\}\right). \quad \square \end{aligned}$$

The next lemma is established in [Bai and Silverstein \(2009\)](#) and [Chen and Pan \(2012\)](#).

**Lemma 6.** *Suppose  $Z$  is an  $n \times p$  matrix with i.i.d. Gaussian entries. Then for any positive constant  $\epsilon > 0$ , it holds that  $n^{-1}Z^\top Z \leq (1 + \epsilon)(1 + \sqrt{c_n})^2$ , w.p.a.1, for  $c_n = p/n \in [0, \infty]$ .*

**Lemma 7** (Convexity). *Let  $O \subseteq \mathbb{R}^d$  be open and convex and  $D$  be a dense subset of  $O$ . For  $\theta \in O$ , both  $M_n(\theta)$  and  $M(\theta)$  are convex in  $\theta$ . If  $M_n(\theta) \xrightarrow{\text{P}} M(\theta)$ , for any  $\theta \in D$ , then  $\sup_{\theta \in K} |M_n(\theta) - M(\theta)| \xrightarrow{\text{P}} 0$ , for any compact subset  $K \subset O$ .*

This lemma has been shown by Lemma 7.75 of [Liese and Miescke \(2008\)](#). Next, we present a min-convergence theorem for functions defined on an open set  $(0, \infty)$ , as shown by Lemma 10 of [Thrampoulidis et al. \(2018\)](#).

**Lemma 8.** *Consider a sequence of proper, convex stochastic functions  $M_n : \mathbb{R}^+ \rightarrow \mathbb{R}$ , and a deterministic function  $M : \mathbb{R}^+ \rightarrow \mathbb{R}$ , satisfying (a)  $M_n(x) \xrightarrow{\text{P}} M(x)$ ,  $\forall x > 0$ ; (b) there exists  $z > 0$  such that  $M(x) > \inf_{y>0} M(y)$ ,  $\forall x \geq z$ . Then we have  $\inf_{x>0} M_n(x) \xrightarrow{\text{P}} \inf_{x>0} M(x)$ .*

Relatedly, we introduce a lemma for functions on a diverging sequence of closed sets.

**Lemma 9.** *Consider a sequence of closed intervals  $\{[x_n, y_n]\}_{n=1}^\infty$  such that  $\lim_{n \rightarrow \infty} x_n = -\infty$  and  $\lim_{n \rightarrow \infty} y_n = +\infty$ . Additionally, let there be a sequence of proper random and convex functions  $M_n : [x_n, y_n] \rightarrow \mathbb{R}$ , and a convex, continuous, and deterministic function  $M : \mathbb{R} \rightarrow \mathbb{R}$  that satisfy: (a)  $M_n(x) \xrightarrow{\text{P}} M(x)$  for every  $x \in \mathbb{R}$ ; (b) there exists  $z > 0$  such that  $M(x) > \inf_{y \in \mathbb{R}} M(y)$  holds for all  $|x| \geq z$ . Then it holds that  $\inf_{x \in [x_n, y_n]} M_n(x) \xrightarrow{\text{P}} \inf_{x \in \mathbb{R}} M(x)$ .*

*Proof.* For  $n$  sufficiently large,  $z \in [x_n, y_n]$ . Assume  $x^* \in [-z, z]$  minimizes  $M(x)$ . Condition (b) in fact implies that  $x^* \in (-z, z)$  and that  $M(x^*) = \inf_{x \in \mathbb{R}} M(x)$ . Consider the event  $\inf_{|x|>z, x \in [x_n, y_n]} M_n(x) < M_n(x^*)$ . Under this event, there exists  $|z_n| > z$  and  $z_n \in [x_n, y_n]$  such that  $M_n(z_n) < M_n(x^*)$ . The geometry implies that there exists  $\theta_n \in (0, 1)$ , such that either  $z_n\theta_n + x^*(1 - \theta_n) = z$  or  $z_n\theta_n + x^*(1 - \theta_n) = -z$  holds. Using convexity, we have  $\min(M_n(z), M_n(-z)) \leq \theta_n M_n(z_n) + (1 - \theta_n) M_n(x^*) < M_n(x^*)$ . By taking limits on both sides, we have  $\min(M(z), M(-z)) \leq M(x^*)$ , which contradicts condition (b). Therefore, w.p.a.1, we have  $\inf_{|x|>z, x \in [x_n, y_n]} M_n(x) \geq M_n(x^*)$ . Furthermore, by Lemma 7, for all arbitrarily small  $\epsilon > 0$ , w.p.a.1,  $\sup_{|x| \leq z} |M_n(x) - M(x)| < \epsilon$ . In addition, by definition, there exists a sequence of  $z_n$ , such that  $|z_n| \leq z$  and  $\inf_{|x| \leq z} M_n(x) \geq M_n(z_n) - \epsilon$ . Combining these two inequalities with the fact that  $M(x^*)$  minimizes  $M$  on  $\mathbb{R}$  leads to  $\inf_{|x| \leq z} M_n(x) \geq M_n(z_n) - \epsilon \geq M(z_n) - 2\epsilon \geq M(x^*) - 2\epsilon$ , w.p.a.1. On the other hand,  $\inf_{|x| \leq z} M_n(x) \leq M_n(x^*) \xrightarrow{\text{P}} M(x^*)$ . Since  $\epsilon$

is arbitrary, we have  $\inf_{|x| \leq z} M_n(x) \xrightarrow{P} M(x^*)$ . Along with  $\inf_{|x| > z, x \in [x_n, y_n]} M_n(x) \geq M_n(x^*)$  and  $M_n(x^*) \xrightarrow{P} M(x^*)$  by (a), we have

$$\inf_{x \in [x_n, y_n]} M_n(x) = \min \left( \inf_{|x| \leq z} M_n(x), \inf_{|x| > z, x \in [x_n, y_n]} M_n(x) \right) \xrightarrow{P} M(x^*). \quad \square$$

**Lemma 10.** *Suppose  $X$  is a standard Gaussian random variable, then for  $x > 0$ ,*

$$\frac{\sqrt{2}}{\sqrt{\pi}} \exp\left(-\frac{x^2}{2}\right) (2x^{-3} - 12x^{-5} - 15x^{-7}) \leq \mathbb{E}(|X| - x)_+^2 \leq \frac{\sqrt{2}}{\sqrt{\pi}} \exp\left(-\frac{x^2}{2}\right) (2x^{-3} + 3x^{-5}).$$

*Proof.* With integration by parts, we find

$$\mathbb{E}(|X| - x)_+^2 = \sqrt{\frac{2}{\pi}} \int_x^\infty (t - x)^2 \exp\left(-\frac{t^2}{2}\right) dt = \sqrt{\frac{2}{\pi}} \left( -x \exp\left(-\frac{x^2}{2}\right) + (x^2 + 1) f_G(x) \right),$$

where  $f_G(x) = \int_x^\infty \exp(-t^2/2) dt$ . Lemma 10 then follows from the tail inequality:

$$\exp\left(-\frac{x^2}{2}\right) \left( \frac{1}{x} - \frac{1}{x^3} + \frac{3}{x^5} - \frac{15}{x^7} \right) \leq f_G(x) \leq \exp\left(-\frac{x^2}{2}\right) \left( \frac{1}{x} - \frac{1}{x^3} + \frac{3}{x^5} \right). \quad \square$$

**Lemma 11.** *Given that  $X$  is a standard Gaussian random variable, the following inequalities hold when  $x > 0$  and  $x$  is sufficiently large:*

$$\mathbb{E}|X|(|X| - x)_+^2 \leq 2x\mathbb{E}(|X| - x)_+^2 \quad \text{and} \quad \mathbb{E}X^2(|X| - x)_+^2 \leq 2x^2\mathbb{E}(|X| - x)_+^2.$$

*Proof.* The proof is analogous to that of Lemma 10 and is therefore omitted.  $\square$

**Definition 1.** *A centered random variable  $X$  belongs to the sub-exponential class  $SE(\nu^2, \alpha)$  with  $\nu > 0$  and  $\alpha > 0$ , if  $\mathbb{E}e^{\lambda X} \leq e^{\frac{\lambda^2 \nu^2}{2}}$ , for all  $\lambda$  such that  $|\lambda| < \alpha^{-1}$ .*

**Lemma 12.** *Let  $\{x_k\}_{k=1}^\infty$  be a sequence of diverging positive numbers. Then as  $p \rightarrow \infty$ , we have w.p.a.1,  $\|\Sigma_2 b_0\|_\infty < x_p q^{-1/2} \log(p)$  and  $\|\Sigma_2^{1/2} h\|_\infty < x_p \sqrt{\log(p)}$ , where  $\Sigma_2$  and  $b_0$  are defined in Assumptions 1 and 3, respectively, and  $h \in \mathbb{R}^p$  is a standard Gaussian vector.*

*Proof.* We only present the proof for the first inequality, noting that the proof for the second inequality follows similarly. By definition, there exist  $b_{1i} \sim B(1, q)$  and a sub-exponential random variable  $b_{2i}$  such that  $b_{0,i} = q^{-1/2} b_{1i} b_{2i}$ . Note that  $b_{1i} b_{2i}$  is still sub-exponential. Without loss of generality, assume  $q^{1/2} b_{0,i} = b_{1i} b_{2i} \in SE(1, 1)$ .

Write the  $(i, j)$ -th element of  $\Sigma_2$  as  $\Sigma_{2,ij}$ . By the properties of sub-exponential variables, we have  $(\Sigma_2 q^{1/2} b_0)_i \in SE\left(\sum_{j=1}^p \Sigma_{2,ij}^2, \max_j |\Sigma_{2,ij}|\right)$ . Given that  $\sum_{j=1}^p \Sigma_{2,ij}^2 = (\Sigma_2^2)_{i,i} \leq$

$\lambda_1(\Sigma_2^2) = C_2^2$  and  $\max_j |\Sigma_{2,ij}| \leq C_2$ , we conclude that  $(\Sigma_2 q^{1/2} b_0)_i \in \text{SE}(C_2^2, C_2)$ . The tail bound of sub-exponential variables yields  $\mathbb{P}(|(\Sigma_2 q^{1/2} b_0)_i| > x_p \log(p)) \leq 2 \exp\left(-\frac{x_p \log(p)}{2C_2}\right)$ . The conclusion follows by union bound inequality and  $2p \exp\left(-\frac{x_p \log(p)}{2C_2}\right) \rightarrow 0$ .  $\square$

**Lemma 13.** *For any given  $M_1$ , it holds that  $\lim_{n \rightarrow \infty} S_{1n} = 0$  for  $S_{1n}$  defined in Eq. (21).*

*Proof.* Write  $\tilde{x}_k = (\Sigma_\varepsilon^{-1/2} X_{\cdot,i})_k$  and  $\tilde{y}_k = \beta_{0,i} \tilde{x}_k + z_k$  for  $k = 1, \dots, n$ . By definition, we have

$$\begin{aligned} \mathbb{E}(b_{0,i} \mathbf{1}_{q^{1/2} b_{0,i} \in [M_1, M_2]} | \mathcal{G}_i) &= \frac{\int b \mathbf{1}_{q^{1/2} b \in [M_1, M_2]} \exp\left(-\sum_{k=1}^n (\tilde{y}_k - p^{-1/2} \tau^{1/2} \tilde{x}_k b)^2 / 2\right) dF(b)}{\int \exp\left(-\sum_{k=1}^n (\tilde{y}_k - p^{-1/2} \tau^{1/2} \tilde{x}_k b)^2 / 2\right) dF(b)} \\ &= \frac{\int b \mathbf{1}_{q^{1/2} b \in [M_1, M_2]} \exp\left(-p^{-1} \tau b^2 \sum_{k=1}^n \tilde{x}_k^2 / 2 + b p^{-1/2} \tau^{1/2} \sum_{k=1}^n \tilde{y}_k \tilde{x}_k\right) dF(b)}{\int \exp\left(-p^{-1} \tau b^2 \sum_{k=1}^n \tilde{x}_k^2 / 2 + b p^{-1/2} \tau^{1/2} \sum_{k=1}^n \tilde{y}_k \tilde{x}_k\right) dF(b)} := \frac{Q_{1n}}{Q_{2n}}, \end{aligned}$$

where  $F$  is the distribution function of  $b_{0,i}$ . By the facts that  $\int b \mathbf{1}_{q^{1/2} b \in [M_1, M_2]} dF(b) = 0$  and  $dF(b) = (1-q)\delta_0 + q dF_{b_2}(q^{1/2}b)$ , we have

$$\begin{aligned} |Q_{1n}| &= \left| \int q b \mathbf{1}_{q^{1/2} b \in [M_1, M_2]} \left[ \exp\left(-p^{-1} \tau b^2 \sum_{k=1}^n \tilde{x}_k^2 / 2 + b p^{-1/2} \tau^{1/2} \sum_{k=1}^n \tilde{y}_k \tilde{x}_k\right) - 1 \right] dF_{b_2}(q^{1/2}b) \right| \\ &\leq q^{1/2} \tilde{M} \int \left| \exp\left(-p^{-1} \tau q^{-1} \tilde{b}^2 \sum_{k=1}^n \tilde{x}_k^2 / 2 + \tilde{b} q^{-1/2} p^{-1/2} \tau^{1/2} \sum_{k=1}^n \tilde{y}_k \tilde{x}_k\right) - 1 \right| dF_{b_2}(\tilde{b}), \end{aligned}$$

where  $\tilde{M} := \max(|M_1|, |M_2|)$ . Define the event

$$A_n := \left\{ \left| p^{-1/2} \tau^{1/2} \sum_{k=1}^n \tilde{y}_k \tilde{x}_k \right| \leq \tilde{C} p^{-1/2} \tau^{1/2} n^{1/2} \log^2(p) \text{ and } p^{-1} \tau \sum_{k=1}^n \tilde{x}_k^2 \leq \tilde{C} p^{-1} n \tau \right\}, \quad (\text{C3})$$

where  $\tilde{C} := 5C_1 C_2 c_\varepsilon^{-1}$ . Under this event, we observe that

$$\begin{aligned} &\left| \exp\left(-p^{-1} \tau q^{-1} \tilde{b}^2 \sum_{k=1}^n \tilde{x}_k^2 / 2 + \tilde{b} q^{-1/2} p^{-1/2} \tau^{1/2} \sum_{k=1}^n \tilde{y}_k \tilde{x}_k\right) - 1 \right| \\ &\leq \exp\left(\tilde{C} |\tilde{b}| p^{-1/2} \tau^{1/2} n^{1/2} q^{-1/2} \log^2(p)\right) - \exp\left(-\tilde{C} \tilde{b}^2 p^{-1} n \tau q^{-1} - \tilde{C} |\tilde{b}| p^{-1/2} \tau^{1/2} n^{1/2} q^{-1/2} \log^2(p)\right). \end{aligned}$$

Since  $p^{-1/2}\tau^{1/2}n^{1/2}q^{-1/2}\log^2(p) \rightarrow 0$  and  $p^{-1}n\tau q^{-1} \rightarrow 0$  by Assumption 4, and given that  $F_{b_2}$  follows a sub-exponential distribution, the integral of both terms on the right-hand-side converges to zero as  $n \rightarrow \infty$ . Therefore, for any  $\epsilon > 0$ , there exists  $n_0$  such that for all  $n > n_0$ , we have  $|Q_{1n}| \leq \epsilon$  under the event  $A_n$ . Similarly, it can be proven that there exists  $n_1$  such that for all  $n > n_1$ , we obtain  $Q_{2n} \geq 1/2$  under the event  $A_n$ .

Next we analyze the event  $A_n$ . We start with the second inequality in  $A_n$ . Note that

$$\begin{aligned} \sum_{k=1}^n \tilde{x}_k^2 &= X_{\cdot,i}^\top \Sigma_\epsilon^{-1} X_{\cdot,i} \leq c_\epsilon^{-1} X_{\cdot,i}^\top X_{\cdot,i} = c_\epsilon^{-1} e_i^\top \Sigma_2^{1/2} Z^\top \Sigma_1 Z \Sigma_2^{1/2} e_i \\ &\leq c_\epsilon^{-1} C_1 \|Z \Sigma_2^{1/2} e_i\|^2 \stackrel{d}{=} c_\epsilon^{-1} C_1 \|\Sigma_2^{1/2} e_i\|^2 \chi^2(n) \leq c_\epsilon^{-1} C_1 C_2 \chi^2(n), \end{aligned}$$

where  $e_i$  is the  $i$ -th standard basis vector. By Lemma 5, with probability at least  $1 - 2 \exp(-cp)$  for some fixed constant  $c > 0$ ,  $\chi^2(n) \leq 5n$ , which implies the second inequality.

For the first inequality in  $A_n$ , using the second inequality, we observe that,

$$\begin{aligned} p^{-1/2}\tau^{1/2} \left| \sum_{k=1}^n \tilde{x}_k \tilde{y}_k \right| &= p^{-1/2}\tau^{1/2} \left| \beta_{0,i} \sum_{k=1}^n \tilde{x}_k^2 + \sum_{k=1}^n \tilde{x}_k z_k \right| \leq \tilde{C} n p^{-1/2}\tau^{1/2} |\beta_{0,i}| + p^{-1/2}\tau^{1/2} \left| \sum_{k=1}^n \tilde{x}_k z_k \right| \\ &= \tilde{C} p^{-1}\tau n |q^{-1/2} b_{1i} b_{2i}| + p^{-1/2}\tau^{1/2} \left| \sum_{k=1}^n \tilde{x}_k z_k \right| \leq \tilde{C} p^{-1} q^{-1/2} \tau n |b_{2i}| + p^{-1/2}\tau^{1/2} \left| \sum_{k=1}^n \tilde{x}_k z_k \right|. \end{aligned}$$

Using the property of a sub-exponential random variable, for some constant  $c > 0$ , with probability at least  $1 - 2 \exp(-c \log^2(p))$ , we have  $|b_2| \leq \log^2(p)/2$ , which implies  $\tilde{C} p^{-1} q^{-1/2} \tau n |b_{2i}| \leq \tilde{C} p^{-1} q^{-1/2} \tau n \log^2(p)/2 = o(\tilde{C} p^{-1/2} \tau^{1/2} n^{1/2} \log^2(p)/2)$  by Assumption 4. In addition, by Lemma 5, with probability at least  $1 - 2 \exp(-c \log^2(p))$ , we have  $\left| \sum_{k=1}^n \tilde{x}_k z_k \right| \leq \tilde{C} n^{1/2} \log^2(p)/2$ , which implies  $p^{-1/2}\tau^{1/2} \left| \sum_{k=1}^n \tilde{x}_k z_k \right| \leq \tilde{C} p^{-1/2} \tau^{1/2} n^{1/2} \log^2(p)/2$ .

In sum, using the facts that  $\max(\exp(-cp), \exp(-c \log^2(p))) = o(p^{-1})$ , we conclude that with probability at least  $1 - p^{-1}$  as  $p \rightarrow \infty$ ,  $A_n$  holds. Hence we have

$$\begin{aligned} \lim_{n \rightarrow \infty} S_{1n} &= \lim_{n \rightarrow \infty} \mathbb{E} \left( \mathbb{E}(b_{0,i} \mathbf{1}_{q^{1/2} b_{0,i} \in [M_1, M_2]} | \mathcal{G}_i) \right)^2 \mathbf{1}_{A_n} + \lim_{n \rightarrow \infty} \mathbb{E} \left( \mathbb{E}(b_{0,i} \mathbf{1}_{q^{1/2} b_{0,i} \in [M_1, M_2]} | \mathcal{G}_i) \right)^2 \mathbf{1}_{A_n^c} \\ &\leq 4\epsilon^2 + \lim_{n \rightarrow \infty} q^{-1} \tilde{M}^2 \mathbb{P}(A_n^c) \leq 4\epsilon^2 + \lim_{n \rightarrow \infty} p^{-1} q^{-1} \tilde{M}^2 = 4\epsilon^2. \end{aligned}$$

The conclusion then follows from the arbitrariness of  $\epsilon$ . □

**Lemma 14.** *The objective function in Eq. (27) is convex with respect to  $\alpha$  and jointly concave with respect to  $(\delta, \gamma)$ . Additionally, as long as Eq. (29) holds, we have Eq. (26).*

*Proof.* By Lemma 15, it suffices to prove Eq. (26) holds for  $\hat{w}^B$ , which equals

$$\arg \min_{w \in S_w^n} \frac{c_n}{n} \left\| \tau^{1/2} \Sigma_1^{1/2} Z w - \tau^{-1} \varepsilon \right\|^2 + c_n^2 \lambda \left\| \Sigma_2^{-1/2} w + \tau^{-3/2} \beta_0 \right\|^2 - \frac{c_n \tau^{-2}}{n} \|\varepsilon\|^2 - C_n^\phi, \quad (\text{C4})$$

and  $S_w^n = \{w \mid c_n \tau^{-1} \sigma_x \sigma_\beta - K_\alpha \leq c_n \|w\| \leq c_n \tau^{-1} \sigma_x \sigma_\beta + K_\alpha\}$  for some sufficiently large  $K_\alpha$ . With slight abuse of notation, we refer to the optimal solution  $\hat{w}$  instead of using  $\hat{w}^B$ .

Note that for any vector  $x$ ,  $\|x\|^2 = \max_u \sqrt{n} u^\top x - n \|u\|^2/4$ , where its argmax is  $2x/\sqrt{n}$ , and similarly  $\|x\|^2 = \max_v v^\top x - \|v\|^2/4$ . Applying these equalities to  $\|\tau^{1/2} \Sigma_1^{1/2} Z w - \tau^{-1} \varepsilon\|^2$  and  $\|\Sigma_2^{-1/2} w + \tau^{-3/2} \beta_0\|^2$ , setting  $\tilde{u} = \Sigma_1^{1/2} u$ , and  $\tilde{v} = \Sigma_2^{-1/2} v$ , we can rewrite (C4) as

$$\begin{aligned} \min_{w \in S_w^n} \max_{\tilde{u}, \tilde{v}} & \frac{c_n \tau^{1/2}}{\sqrt{n}} \tilde{u}^\top Z w - \frac{c_n \tau^{-1}}{\sqrt{n}} \tilde{u}^\top \Sigma_1^{-1/2} \varepsilon - \frac{c_n \|\Sigma_1^{-1/2} \tilde{u}\|^2}{4} + c_n^2 \lambda \tilde{v}^\top w + c_n^2 \lambda \tau^{-3/2} \tilde{v}^\top \Sigma_2^{1/2} \beta_0 \\ & - \frac{c_n^2 \lambda \|\Sigma_2^{1/2} \tilde{v}\|^2}{4} - \frac{c_n \tau^{-2}}{n} \|\varepsilon\|^2 - C_n^\phi. \end{aligned} \quad (\text{C5})$$

To simplify notation and without ambiguity, we continue using  $u$  and  $v$  in place of  $\tilde{u}$  and  $\tilde{v}$ .

For a given  $w$ , the argmax of Eq. (C5), denoted by  $\hat{u}$ , is equal to  $\frac{2}{\sqrt{n}}(\tau^{1/2} \Sigma_1 Z w - \tau^{-1} \Sigma_1^{1/2} \varepsilon)$ . Given the definition of  $S_w^n$  and Assumptions 1 and 2, we have  $\|w\| \leq \tau^{-1} \sigma_x \sigma_\beta + c_n^{-1} K_\alpha$ ,  $\|\Sigma_1\| \leq C_1$ ,  $\|\Sigma_\varepsilon\| \leq C_\varepsilon$ . Furthermore, w.p.a. 1,  $\|z\| \leq \sqrt{2n}$  by the law of large numbers, which implies  $\|\varepsilon\| \leq \sqrt{2C_\varepsilon n}$ . Together with Lemma 6 and that  $\tau c_n \rightarrow 0$  by Assumption 4, we have the following upper bound for  $\|\hat{u}\|$  as  $n$  is large enough:  $\|\hat{u}\| \leq \frac{2\tau^{1/2}}{\sqrt{n}} \|\Sigma_1 Z w\| + \frac{2}{\sqrt{n}} \|\tau^{-1} \Sigma_1^{1/2} \varepsilon\| \leq 4\tau^{-1} \sqrt{C_1 C_\varepsilon}$ . Let  $S_u^n = \{u \mid \|u\| \leq 4\tau^{-1} \sqrt{C_1 C_\varepsilon}\}$ . Based on the above result, w.a.p.1, the following optimization problem is equivalent to (C5):

$$\begin{aligned} \min_{w \in S_w^n} \max_{u \in S_u^n, v} & \frac{c_n \tau^{1/2}}{\sqrt{n}} u^\top Z w - \frac{c_n \tau^{-1}}{\sqrt{n}} u^\top \Sigma_1^{-1/2} \varepsilon - \frac{c_n \|\Sigma_1^{-1/2} u\|^2}{4} + c_n^2 \lambda v^\top w + c_n^2 \lambda \tau^{-3/2} v^\top \Sigma_2^{1/2} \beta_0 \\ & - \frac{c_n^2 \lambda \|\Sigma_2^{1/2} v\|^2}{4} - \frac{c_n \tau^{-2}}{n} \|\varepsilon\|^2 - C_n^\phi. \end{aligned} \quad (\text{C6})$$

Next, we need introduce an auxiliary problem for the purpose of applying CGMT:



$$\begin{aligned}
\phi(g, h) &= \max_{0 \leq \delta \leq 4\tau^{-1}\sqrt{c_1 c_\varepsilon}} \min_{w \in S_w^n} \max_{\|u\|=\delta} \mathcal{R}_n(w, v, u), \quad \text{where} \\
\mathcal{R}_n(w, v, u) &= \frac{c_n \tau^{1/2}}{\sqrt{n}} \|w\| g^\top u - \frac{c_n \tau^{1/2}}{\sqrt{n}} \delta h^\top w - \frac{c_n \tau^{-1}}{\sqrt{n}} u^\top \Sigma_1^{-1/2} \varepsilon - \frac{c_n \|\Sigma_1^{-1/2} u\|^2}{4} \\
&\quad + c_n^2 \lambda v^\top w + c_n^2 \lambda \tau^{-3/2} v^\top \Sigma_2^{1/2} \beta_0 - \frac{c_n^2 \lambda \|\Sigma_2^{1/2} v\|^2}{4} - \frac{c_n \tau^{-2}}{n} \|\varepsilon\|^2 - C_n^\phi,
\end{aligned} \tag{C7}$$

and  $g \in \mathbb{R}^n$  and  $h \in \mathbb{R}^p$  are standard Gaussian vectors, independent of the other random variables. Similarly, let  $\tilde{\mathcal{S}}_n := \{w \mid |c_n \|w\| - c_n \tau^{-1} \sigma_x \sigma_\beta - \alpha_2^*| < \epsilon\}$ , define  $\phi_{\tilde{\mathcal{S}}_n^c}(g, h)$  as the optimal value of an analogous optimization problem to (C7), with  $w$  restricted to  $S_w^n \cap \tilde{\mathcal{S}}_n^c$ .

Lemma 16 characterizes the limiting behavior of the optimal solution to (C6),  $\hat{w}$ , and in turn, proves the desired (26), under conditions pertaining to the optimization problem (C7). Therefore, we only need show that conditions outlined in Lemma 16 hold as long as (29) holds. That is, under (29), we need to prove the existence of the constants  $\bar{\phi} < \bar{\phi}_{\tilde{\mathcal{S}}_n^c}$  such that for all  $\eta > 0$ , w.p.a.1 in the limit of  $n \rightarrow \infty$ ,  $\phi(g, h) < \bar{\phi} + \eta$  and  $\phi_{\tilde{\mathcal{S}}_n^c}(g, h) > \bar{\phi}_{\tilde{\mathcal{S}}_n^c} - \eta$ .

Let  $\bar{u} = u/\delta$ , maximizing part of  $\mathcal{R}_n(w, v, u)$  over  $u$  simplifies to the following problem:

$$\begin{aligned}
&\max_{\|u\|=\delta} \frac{c_n \tau^{1/2}}{\sqrt{n}} \|w\| g^\top u - \frac{c_n \tau^{-1}}{\sqrt{n}} u^\top \Sigma_1^{-1/2} \varepsilon - \frac{c_n \|\Sigma_1^{-1/2} u\|^2}{4} \\
&= \max_{\|\bar{u}\|=1} \frac{c_n \delta}{\sqrt{n}} (\tau^{1/2} \|w\| g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon)^\top \bar{u} - \frac{c_n \delta^2}{4} \bar{u}^\top \Sigma_1^{-1} \bar{u}.
\end{aligned}$$

The latter is a quadratic programming problem, which has been studied in Gander et al. (1989). The optimal value associated with this problem is given by:

$$-\frac{c_n \delta^2}{4} \mu_n(\alpha, \delta) + \frac{c_n}{n} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon)^\top (\Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I})^{-1} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon) \tag{C8}$$

where  $\alpha := \|w\|$  and  $\mu_n(\alpha, \delta)$  is the solution to

$$\frac{1}{n} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon)^\top (\Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I})^{-2} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon) - \frac{\delta^2}{4} = 0, \tag{C9}$$

under the condition  $\Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I}$  is positive semidefinite. With this, Eq. (C7) equals:

$$\begin{aligned}
&\max_{0 \leq \delta \leq 4\tau^{-1}\sqrt{c_1 c_\varepsilon}} \min_{w \in S_w^n} -\frac{c_n \delta^2}{4} \mu_n(\alpha, \delta) + \frac{c_n}{n} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon)^\top (\Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I})^{-1} \\
&\quad \times (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon) - \frac{c_n \tau^{1/2}}{\sqrt{n}} \delta h^\top w + c_n^2 \lambda v^\top w
\end{aligned}$$

$$+ c_n^2 \lambda \tau^{-3/2} v^\top \Sigma_2^{1/2} \beta_0 - \frac{c_n^2 \lambda \|\Sigma_2^{1/2} v\|^2}{4} - \frac{c_n \tau^{-2}}{n} \|\varepsilon\|^2 - C_n^\phi.$$

Solving the inside minimization problem with respect to  $w/\alpha$  while fixing  $\alpha$  leads to

$$\begin{aligned} \max_{0 \leq \delta \leq 4\tau^{-1} \sqrt{C_1 C_\varepsilon}} \min_{|c_n \alpha - c_n \tau^{-1} \sigma_x \sigma_\beta| \leq K_\alpha} & - \frac{c_n \delta^2}{4} \mu_n(\alpha, \delta) + \frac{c_n}{n} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon)^\top (\Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I})^{-1} \\ & \times (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon) - c_n \|c_n \lambda v - n^{-1/2} \tau^{1/2} \delta h\| \alpha \quad (\text{C10}) \\ & + c_n^2 \lambda \tau^{-3/2} v^\top \Sigma_2^{1/2} \beta_0 - \frac{c_n^2 \lambda \|\Sigma_2^{1/2} v\|^2}{4} - \frac{c_n \tau^{-2}}{n} \|\varepsilon\|^2 - C_n^\phi. \end{aligned}$$

By Lemma 17, the objective function of the above optimization is convex in  $\alpha$  and jointly concave in  $(\delta, v)$ . As a result, we can switch the order of min and max by Corollary 3.3 in Sion (1958). Also, note that for any vector  $x$ ,  $\|x\| = \min_{\gamma > 0} \frac{1}{2\gamma} \|x\|^2 + \frac{\gamma}{2}$ . Applying this equation to  $\|c_n \lambda v - n^{-1/2} \tau^{1/2} \delta h\| \alpha$ , Eq. (C10) becomes

$$\begin{aligned} \min_{c_n |\alpha - \tau^{-1} \sigma_x \sigma_\beta| \leq K_\alpha} \max_{\gamma > 0} \max_v & - \frac{c_n \delta^2}{4} \mu_n(\alpha, \delta) + \frac{c_n}{n} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon)^\top \\ & \times (\Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I})^{-1} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon) - \frac{c_n \gamma}{2} - \frac{c_n \alpha^2}{2\gamma} \|c_n \lambda v - n^{-1/2} \tau^{1/2} \delta h\|^2 \\ & + c_n^2 \lambda \tau^{-3/2} v^\top \Sigma_2^{1/2} \beta_0 - \frac{c_n^2 \lambda \|\Sigma_2^{1/2} v\|^2}{4} - \frac{c_n \tau^{-2}}{n} \|\varepsilon\|^2 - C_n^\phi. \end{aligned}$$

Note that the objective function above is jointly concave in  $(\delta, \gamma, v)$ . To see why this is true, it is sufficient to prove that  $-\frac{\alpha^2}{2\gamma} \|c_n \lambda v - n^{-1/2} \tau^{1/2} \delta h\|^2$  is jointly concave in  $(\delta, \gamma, v)$ , which follows by Lemma 13 in Thrampoulidis et al. (2018). Consequently, after solving the first maximization problem over  $v$ , the resulting function remains jointly concave in  $(\delta, \gamma)$ . Maximizing over  $v$  is again a standard quadratic programming problem, which leads to Eq. (27). Thus, we conclude that (27) is convex with respect to  $\alpha$  and jointly concave with respect to  $(\delta, \gamma)$ .

For any compact set  $A$ , define  $\phi_A(g, h) := \min_{\alpha_2 \in A} \max_{\substack{\gamma_1 > 0 \\ \delta_3 \in K_{\delta_3}}} \tilde{Q}_n(\alpha_2, \delta_3, \gamma_1)$ . Based on the above argument and condition (i) in (29), we can deduce

$$\begin{aligned} \phi(g, h) &= \phi_{[-K_\alpha, K_\alpha]}(g, h) \xrightarrow{\text{P}} \min_{\alpha_2 \in [-K_\alpha, K_\alpha]} \max_{\substack{\gamma_1 > 0 \\ \delta_3 \in \mathbb{R}}} \tilde{Q}(\alpha_2, \delta_3, \gamma_1), \\ \phi_{\tilde{S}_n^c}(g, h) &= \min\{\phi_{[-K_\alpha, \alpha_2^* - \epsilon]}(g, h), \phi_{[\alpha_2^* + \epsilon, K_\alpha]}(g, h)\} \xrightarrow{\text{P}} \min_{\alpha_2 \in [-K_\alpha, \alpha_2^* - \epsilon] \cup [\alpha_2^* + \epsilon, K_\alpha]} \max_{\substack{\gamma_1 > 0 \\ \delta_3 \in \mathbb{R}}} \tilde{Q}(\alpha_2, \delta_3, \gamma_1). \end{aligned}$$

Together with condition (ii), the conditions in Lemma 16 are satisfied by letting  $\bar{\phi} = \min_{\alpha_2 \in [-K_\alpha, K_\alpha]} \max_{\substack{\gamma_1 > 0 \\ \delta_3 \in \mathbb{R}}} \tilde{Q}(\alpha_2, \delta_3, \gamma_1)$  and  $\bar{\phi}_{\tilde{S}_n^c} = \min_{\alpha_2 \in [-K_\alpha, \alpha_2^* - \epsilon] \cup [\alpha_2^* + \epsilon, K_\alpha]} \max_{\substack{\gamma_1 > 0 \\ \delta_3 \in \mathbb{R}}} \tilde{Q}(\alpha_2, \delta_3, \gamma_1)$ .  $\square$

We now introduce two lemmas whose proofs follow the same reasoning as those of Lemma 5 and Lemma 7 in [Thrampoulidis et al. \(2018\)](#), and are therefore omitted here.

**Lemma 15.** *Under the conditions of Theorem 2, define  $S_w^n := \{w \mid c_n \tau^{-1} \sigma_x \sigma_\beta - K_\alpha \leq c_n \|w\| \leq c_n \tau^{-1} \sigma_x \sigma_\beta + K_\alpha\}$  for some  $K_\alpha$  such that  $|\alpha_2^*| < K_\alpha$ . If the solution  $\hat{w}^B$  to*

$$\arg \min_{w \in S_w^n} \frac{c_n}{n} \left\| \tau^{1/2} \Sigma_1^{1/2} Z w - \tau^{-1} \varepsilon \right\|^2 + c_n^2 \lambda \left\| \Sigma_2^{-1/2} w + \tau^{-3/2} \beta_0 \right\|^2 - \frac{c_n \tau^{-2}}{n} \|\varepsilon\|^2 - C_n^\phi$$

satisfies  $c_n \|\hat{w}^B\| - c_n \tau^{-1} \sigma_x \sigma_\beta \rightarrow \alpha_2^*$ , then the same holds true for  $\hat{w}$  of Eq. (25).

**Lemma 16.** *Let  $\hat{w}$  denote an optimal solution of Eq. (C6). Regarding  $\phi(g, h)$  and  $\phi_{\tilde{S}_n^c}(g, h)$ , as introduced and discussed in relation to Eq. (C7), suppose there are constants  $\bar{\phi}$  and  $\bar{\phi}_{\tilde{S}_n^c}$  with  $\bar{\phi} < \bar{\phi}_{\tilde{S}_n^c}$ , such that for all  $\eta > 0$ , the following hold w.a.p.1 as  $n \rightarrow \infty$ : (a)  $\phi(g, h) < \bar{\phi} + \eta$ , (b)  $\phi_{\tilde{S}_n^c}(g, h) > \bar{\phi}_{\tilde{S}_n^c} - \eta$ . Under these conditions, we have  $\hat{w} \in \tilde{S}_n$  w.p.a.1.*

**Lemma 17.** *The objective function of Eq. (C10) is convex in  $\alpha$  and jointly concave in  $(\delta, v)$ .*

*Proof.* First, we prove the objective function is convex in  $\alpha = \|w\|$ . Revisiting the term  $f(\alpha, u) := \frac{c_n \tau^{1/2}}{\sqrt{n}} \alpha g^\top u - \frac{c_n \tau^{-1}}{\sqrt{n}} u^\top \Sigma_1^{-1/2} \varepsilon - \frac{c_n \|\Sigma_1^{-1/2} u\|^2}{4}$  in Eq. (C7), we observe that it is convex in  $\alpha$ . After maximizing over the direction of  $u$ , the term remains convex in  $\alpha$  since  $\max_{\|u\|=\delta} f(\theta \alpha_1 + (1-\theta) \alpha_2, u) \leq \max_{\|u\|=\delta} \{\theta f(\alpha_1, u) + (1-\theta) f(\alpha_2, u)\} \leq \theta \max_{\|u\|=\delta} f(\alpha_1, u) + (1-\theta) \max_{\|u\|=\delta} f(\alpha_2, u)$ , for  $\theta \in (0, 1)$ . Note that from Eq. (C8),  $\max_{\|u\|=\delta} f(\alpha, u) = -\frac{c_n \delta^2}{4} \mu_n(\alpha, \delta) + \frac{c_n}{n} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon)^\top (\Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I})^{-1} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon)$ , which yield the first two terms in Eq. (C10). Meanwhile, the term  $-\|c_n \lambda v - n^{-1/2} \tau^{1/2} \delta h\| \alpha$  is also convex in  $\alpha$ . Consequently, we deduce that the objective function of Eq. (C10) is convex in  $\alpha$ .

Next, we demonstrate that this function is jointly concave in  $(\delta, v)$ . It is easy to verify that  $-\|c_n \lambda v - n^{-1/2} \tau^{1/2} \delta h\| \alpha$  is jointly concave in  $(\delta, v)$ , since  $\alpha \geq 0$ . Moreover,  $\lambda \tau^{-3/2} v^\top \Sigma_2^{1/2} \beta_0 - \lambda \|\Sigma_2^{1/2} v\|^2 / 4$  is concave in  $v$ . Therefore, it suffices to prove

$$-\frac{\delta^2}{4} \mu_n(\alpha, \delta) + \frac{1}{n} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon)^\top (\Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I})^{-1} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon) \quad (\text{C11})$$

is concave in  $\delta$ . Let the eigenvalues and normalized eigenvectors of  $\Sigma_1$  be  $\{(\lambda_i, v_i)\}_{i=1}^n$ , and let  $w_i = (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon)^\top v_i$ , for  $i = 1, 2, \dots, n$ . Then (C11) equals  $-\frac{\delta^2}{4} \mu_n(\alpha, \delta) +$

$\frac{1}{n} \sum_{i=1}^n \frac{1}{1/\lambda_i - \mu_n(\alpha, \delta)} w_i^2$ . The first order derivative of this equation with respect to  $\delta$  is

$$-\frac{\delta}{2} \mu_n(\alpha, \delta) - \frac{\delta^2}{4} \partial_\delta \mu_n(\alpha, \delta) + \frac{\partial_\delta \mu_n(\alpha, \delta)}{n} \sum_{i=1}^n \frac{1}{(1/\lambda_i - \mu_n(\alpha, \delta))^2} w_i^2 = -\frac{\delta}{2} \mu_n(\alpha, \delta), \quad (\text{C12})$$

where the last equation follows from the definition of the function  $\mu_n(\alpha, \delta)$ :

$$\frac{1}{n} \sum_{i=1}^n \frac{1}{(1/\lambda_i - \mu_n(\alpha, \delta))^2} w_i^2 = \frac{\delta^2}{4}. \quad (\text{C13})$$

Further, the second-order derivative with respect to  $\delta$  can be calculated as:  $-\frac{1}{2} \mu_n(\alpha, \delta) - \frac{\delta}{2} \partial_\delta \mu_n(\alpha, \delta)$ . By the chain rule of differentiation,  $\partial_\delta \mu_n(\alpha, \delta)$  is the reciprocal of  $\partial \mu_n \delta$ . The latter can be calculated directly using the definition of  $\mu_n$  via Eq. (C13):  $\partial_{\mu_n} \delta = \left( \frac{1}{n} \sum_{i=1}^n \frac{1}{(1/\lambda_i - \mu_n)^2} w_i^2 \right)^{-1/2} \cdot \frac{2}{n} \sum_{i=1}^n \frac{1}{(1/\lambda_i - \mu_n)^3} w_i^2$ . With this, we can write the second-order derivative as follows:

$$-\frac{1}{2} \mu_n(\alpha, \delta) - \frac{\delta}{2} \partial_\delta \mu_n(\alpha, \delta) = -\frac{1}{2} \mu_n - \frac{1}{2} \cdot \frac{\sum_{i=1}^n \frac{1}{(1/\lambda_i - \mu_n)^2} w_i^2}{\sum_{i=1}^n \frac{1}{(1/\lambda_i - \mu_n)^3} w_i^2} = -\frac{1}{2} \cdot \frac{\sum_{i=1}^n \frac{1/\lambda_i}{(1/\lambda_i - \mu_n)^3} w_i^2}{\sum_{i=1}^n \frac{1}{(1/\lambda_i - \mu_n)^3} w_i^2}.$$

Since  $\Sigma_1^{-1} - \mu_n \mathbb{I}$  is positive semidefinite, the right-hand-side is no larger than zero, which concludes the proof.  $\square$

**Lemma 18.** For  $\tilde{Q}_n = \tilde{Q}_n(\alpha_2, \delta_3, \gamma_1)$  in Eq. (28), Eq. (29) holds.

*Proof.* The notation below is defined in the proof of Theorem 2. Let  $\delta_2 = \delta_2^* + c_n^{-1/2} \delta_3$ . First, we demonstrate that  $c_n \tau^{-1} \mu_n(\alpha, \delta) - c_n \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) \xrightarrow{P} 0$ . Let  $f(x) := \frac{1}{n} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon)^\top (\Sigma_1^{-1} - x \mathbb{I})^{-2} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon)$ . Recall that  $\mu_n(\alpha, \delta)$  is the solution to  $f(x) = \delta^2/4$ . Note that  $f(x)$  exhibits a monotonic increase in  $x$  when  $x \leq 1/C_1$ . Therefore, it suffices to show that, given any arbitrarily small  $\epsilon > 0$ , w.p.a.1, the following inequalities hold:  $c_n \tau f(\tau \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) + \tau c_n^{-1} \epsilon) - c_n \delta^2 \tau / 4 > c_+ > 0$  and  $c_n \tau f(\tau \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) - \tau c_n^{-1} \epsilon) - c_n \delta^2 \tau / 4 < c_- < 0$ , for some constants  $c_+$  and  $c_-$ .

By Lemmas 2 and 3, we can deduce the following equations:

$$\begin{aligned} & \frac{c_n \tau^{-1}}{n} \varepsilon^\top \Sigma_1^{-1/2} (\Sigma_1^{-1} - \tau \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) \mathbb{I} - \tau c_n^{-1} \epsilon \mathbb{I})^{-2} \Sigma_1^{-1/2} \varepsilon \\ & - \frac{c_n \tau^{-1}}{n} \text{Tr} \left[ \Sigma_\varepsilon^{1/2} \Sigma_1^{-1/2} (\Sigma_1^{-1} - \tau \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) \mathbb{I} - \tau c_n^{-1} \epsilon \mathbb{I})^{-2} \Sigma_1^{-1/2} \Sigma_\varepsilon^{1/2} \right] = O_P(c_n \tau^{-1} n^{-1/2}), \\ & \frac{c_n \tau^2}{n} \alpha^2 g^\top (\Sigma_1^{-1} - \tau \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) \mathbb{I} - \tau c_n^{-1} \epsilon \mathbb{I})^{-2} g \end{aligned}$$

$$\begin{aligned}
& -\frac{c_n \alpha^2 \tau^2}{n} \text{Tr} \left[ \left( \Sigma_1^{-1} - \tau \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) \mathbb{I} - \tau c_n^{-1} \epsilon \mathbb{I} \right)^{-2} \right] = O_{\text{P}}(c_n n^{-1/2}), \\
& \frac{c_n \tau^{1/2} \alpha}{n} \varepsilon \Sigma_1^{-1/2} \left( \Sigma_1^{-1} - \tau \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) \mathbb{I} - \tau c_n^{-1} \epsilon \mathbb{I} \right)^{-2} g = O_{\text{P}}(c_n \tau^{-1/2} n^{-1/2}).
\end{aligned}$$

Therefore, using the definition of  $f(\cdot)$  we can deduce that:

$$\begin{aligned}
& c_n \tau f(\tau \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) + \tau c_n^{-1} \epsilon) \\
& - \frac{c_n \tau^{-1}}{n} \text{Tr} \left[ \Sigma_\varepsilon^{1/2} \Sigma_1^{-1/2} \left( \Sigma_1^{-1} - \tau \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) \mathbb{I} - \tau c_n^{-1} \epsilon \mathbb{I} \right)^{-2} \Sigma_1^{-1/2} \Sigma_\varepsilon^{1/2} \right] \\
& - \frac{c_n \alpha^2 \tau^2}{n} \text{Tr} \left[ \left( \Sigma_1^{-1} - \tau \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) \mathbb{I} - \tau c_n^{-1} \epsilon \mathbb{I} \right)^{-2} \right] = O_{\text{P}}(c_n \tau^{-1} n^{-1/2}) = o_{\text{P}}(1). \quad (\text{C14})
\end{aligned}$$

Note that for sufficiently small  $x$  such that  $x \|\Sigma_1\| < 1$ ,

$$\begin{aligned}
& \frac{\tau^{-1}}{n} \text{Tr} \left[ \Sigma_\varepsilon^{1/2} \Sigma_1^{-1/2} \left( \Sigma_1^{-1} - x \mathbb{I} \right)^{-2} \Sigma_1^{-1/2} \Sigma_\varepsilon^{1/2} - \Sigma_\varepsilon^{1/2} \Sigma_1^{1/2} \left( \mathbb{I} + 2x \Sigma_1 \right) \Sigma_1^{1/2} \Sigma_\varepsilon^{1/2} \right] \\
& = \frac{\tau^{-1}}{n} \text{Tr} \left[ \Sigma_\varepsilon^{1/2} \Sigma_1^{1/2} \left( \mathbb{I} - x \Sigma_1 \right)^{-2} \Sigma_1^{1/2} \Sigma_\varepsilon^{1/2} - \Sigma_\varepsilon^{1/2} \Sigma_1^{1/2} \left( \mathbb{I} + 2x \Sigma_1 \right) \Sigma_1^{1/2} \Sigma_\varepsilon^{1/2} \right] \\
& \leq \tau^{-1} C_1 C_\varepsilon \left\| \left( \mathbb{I} - x \Sigma_1 \right)^{-2} - \left( \mathbb{I} + 2x \Sigma_1 \right) \right\| \lesssim \tau^{-1} x^2,
\end{aligned}$$

where we apply Lemma 4 in the last inequality. As a consequence, we have:

$$\begin{aligned}
& \frac{\tau^{-1}}{n} \text{Tr} \left[ \Sigma_\varepsilon^{1/2} \Sigma_1^{-1/2} \left( \Sigma_1^{-1} - \tau \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) \mathbb{I} - \tau c_n^{-1} \epsilon \mathbb{I} \right)^{-2} \Sigma_1^{-1/2} \Sigma_\varepsilon^{1/2} \right] \\
& = \frac{1}{n} \text{Tr} \left[ \Sigma_\varepsilon^{1/2} \Sigma_1^{-1/2} \left( \tau^{-1} \Sigma_1^2 + 2 \left( \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) + c_n^{-1} \epsilon \right) \Sigma_1 \right) \Sigma_1^{-1/2} \Sigma_\varepsilon^{1/2} \right] + O(\tau) \\
& = \tau^{-1} \sigma_\varepsilon^2 \theta_1 + 2 \left( \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) + c_n^{-1} \epsilon \right) \sigma_\varepsilon^2 \theta_3 + O(\tau) + o(c_n^{-1}), \quad (\text{C15})
\end{aligned}$$

where the last equation follows by Assumption 5. By the same argument, it follows that:

$$\frac{\alpha^2 \tau^2}{n} \text{Tr} \left[ \left( \Sigma_1^{-1} - \tau \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) \mathbb{I} - \tau c_n^{-1} \epsilon \mathbb{I} \right)^{-2} \right] = \sigma_x^2 \sigma_\beta^2 \theta_4 + O(\tau) + o(c_n^{-1}).$$

Applying the above estimates to the left-hand-side of (C14), we can deduce that:

$$\begin{aligned}
& c_n \tau f(\tau \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) + \tau c_n^{-1} \epsilon) - \frac{c_n \delta^2 \tau}{4} \\
& = 2c_n \left( \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) + c_n^{-1} \epsilon \right) \sigma_\varepsilon^2 \theta_3 + c_n \sigma_x^2 \sigma_\beta^2 \theta_4 - \frac{c_n \delta_1^* \delta_2}{2} + o_{\text{P}}(1).
\end{aligned}$$

By the definition of  $\mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2)$ , the right-hand side of the above equation is positive

w.p.a.1. The proof of the other inequality is similar. Hence, we have proved

$$c_n \tau^{-1} \mu_n(\alpha, \delta) - c_n \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) \xrightarrow{P} 0. \quad (\text{C16})$$

Next, to analyze  $\tilde{Q}_n$ , we first investigate the limiting behavior of:

$$-\frac{\delta^2}{4} \mu_n(\alpha, \delta) + \frac{1}{n} \left( \tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon \right)^\top \left( \Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I} \right)^{-1} \left( \tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon \right).$$

By (C16), we have  $\|\mu_n(\alpha, \delta) \Sigma_1\| = O_P(\tau)$ . Applying Lemma 4 again, we deduce:

$$\begin{aligned} & \left\| \left( \Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I} \right)^{-2} - \Sigma_1^{-2} - 2\mu_n(\alpha, \delta) \Sigma_1^{-3} - 3\mu_n^2(\alpha, \delta) \Sigma_1^{-4} \right\| \lesssim_P \tau^3 \\ & \left\| \left( \Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I} \right)^{-1} - \Sigma_1^{-1} - \mu_n(\alpha, \delta) \Sigma_1^{-2} - \mu_n^2(\alpha, \delta) \Sigma_1^{-3} \right\| \lesssim_P \tau^3. \end{aligned}$$

Furthermore, by the fact that  $\|\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon\| = O_P(n\tau^{-1})$  and Eq. (C9), we have

$$\frac{\delta^2}{4} \mu_n(\alpha, \delta) = \frac{1}{n} \left( \tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon \right)^\top \left( \mu_n(\alpha, \delta) \Sigma_1^2 + 2\mu_n^2(\alpha, \delta) \Sigma_1^3 \right) \left( \tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon \right) + O_P(\tau).$$

With a similar approach, we have

$$\begin{aligned} & \frac{1}{n} \left( \tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon \right)^\top \left( \Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I} \right)^{-1} \left( \tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon \right) \\ &= \frac{1}{n} \left( \tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon \right)^\top \left( \Sigma_1 + \mu_n(\alpha, \delta) \Sigma_1^2 + \mu_n^2(\alpha, \delta) \Sigma_1^3 \right) \left( \tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon \right) + O_P(\tau). \end{aligned}$$

As a consequence, based on Lemmas 2 and 3, as well as the definition of  $\alpha_2$  and the fact that  $c_n \tau^{-1} \mu_n(\alpha, \delta) - c_n \mu(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) \xrightarrow{P} 0$ , we have:

$$\begin{aligned} & -\frac{c_n \delta^2}{4} \mu_n(\alpha, \delta) + \frac{c_n}{n} \left( \tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon \right)^\top \left( \Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I} \right)^{-1} \left( \tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon \right) \\ &= \frac{c_n}{n} \left( \tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon \right)^\top \left( \Sigma_1 - \mu_n^2(\alpha, \delta) \Sigma_1^3 \right) \left( \tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon \right) + O_P(c_n \tau) \\ &= c_n \tau^{-1} \sigma_x^2 \sigma_\beta^2 - c_n \sigma_\varepsilon^2 \theta_3 \mu^2(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) + 2\sigma_x \sigma_\beta \alpha_2 + \frac{c_n \tau^{-2}}{n} \|\varepsilon\|^2 + o_P(1). \quad (\text{C17}) \end{aligned}$$

Finally, we examine the remainder term that contributes to  $\tilde{Q}_n$ :

$$\frac{c_n^2 \lambda^2}{4} \left( \tau^{-3/2} \Sigma_2^{1/2} \beta_0 + \frac{\alpha^2 \delta \tau^{1/2}}{\sqrt{n} \gamma} h \right)^\top \left( \frac{\lambda}{4} \Sigma_2 + \frac{c_n \alpha^2 \lambda^2}{2\gamma} \mathbb{I} \right)^{-1} \left( \tau^{-3/2} \Sigma_2^{1/2} \beta_0 + \frac{\alpha^2 \delta \tau^{1/2}}{\sqrt{n} \gamma} h \right) - \frac{c_n \tau \alpha^2 \delta^2}{2\gamma n} \|h\|^2.$$

Using Lemmas 2-3,  $p^{1/2}\tau^{-1}n^{-1}q^{-1/2} = o(1)$  by Assumption 4, and the assumptions on  $\Sigma_2$ , this term converges in probability to:

$$\begin{aligned} & \frac{c_n^2 \lambda^2 \tau^{-2} \sigma_\beta^2}{4p} \text{Tr} \left[ \Sigma_2^{1/2} \left( \frac{\lambda}{4} \Sigma_2 + \frac{c_n \alpha^2 \lambda^2}{2\gamma} \mathbb{I} \right)^{-1} \Sigma_2^{1/2} \right] + c_n \tau \text{Tr} \left[ \frac{c_n \lambda^2 \alpha^4 \delta^2}{4n \gamma_h^2} \left( \frac{\lambda}{4} \Sigma_2 + \frac{c_n \alpha^2 \lambda^2}{2\gamma} \mathbb{I} \right)^{-1} - \frac{\alpha^2 \delta^2}{2\gamma n} \mathbb{I} \right] \\ &= \frac{c_n \gamma n_1}{2} \tau^{-1} - \frac{\gamma_1^2}{4\sigma_\beta^2 \lambda} \theta_2 - \frac{\gamma_1 \alpha_2}{\sigma_x \sigma_\beta} - \tau^{-1} \frac{(\delta_1^*)^2 \sigma_x^2 c_n}{4\lambda} - \frac{c_n \sigma_x^2 \delta_1^* \delta_2}{2\lambda} + \frac{(\delta_1^*)^2 \gamma_1 \sigma_x^2}{8\lambda^2 \sigma_\beta^2} \theta_2 + o_n(1), \end{aligned}$$

where we apply Lemma 4 and the same argument in proving Eq. (C15). Combining this estimate with (C17) we conclude that

$$\tilde{Q}_n = -\frac{\delta_3^2 \theta_1}{4\theta_3} + 2\sigma_x \sigma_\beta \alpha_2 - \frac{\gamma_1^2}{4\sigma_\beta^2 \lambda} \theta_2 - \frac{\gamma_1 \alpha_2}{\sigma_x \sigma_\beta} + \frac{(\delta_1^*)^2 \gamma_1 \sigma_x^2}{8\lambda^2 \sigma_\beta^2} \theta_2 + o_P(1).$$

We now proceed to establish Claims (i) to (ii) in Eq. (29). Fix  $\alpha_2 \in A$  and  $\gamma_1 > 0$ , and observe that  $\lim_{\delta_3 \rightarrow \pm\infty} \tilde{Q}(\alpha_2, \delta_3, \gamma_1) \rightarrow -\infty$ . By the concave version of Lemma 9, we conclude that  $\max_{\delta_3 \in K_{\delta_3}} \tilde{Q}_n(\alpha_2, \delta_3, \gamma_1) \xrightarrow{P} \max_{\delta_3 \in \mathbb{R}} \tilde{Q}(\alpha_2, \delta_3, \gamma_1)$ . Since  $\tilde{Q}_n$  is jointly concave in  $(\delta_3, \gamma_1)$ , after maximizing with respect to  $\delta_3$ , the function should remain concave in  $\gamma_1$ . Moreover, consider the following equation:  $\max_{\delta_3 \in \mathbb{R}} \tilde{Q}(\alpha_2, \delta_3, \gamma_1) = 2\sigma_x \sigma_\beta \alpha_2 - \frac{\gamma_1^2}{4\sigma_\beta^2 \lambda} \theta_2 - \frac{\gamma_1 \alpha_2}{\sigma_x \sigma_\beta} + \frac{(\delta_1^*)^2 \gamma_1 \sigma_x^2}{8\lambda^2 \sigma_\beta^2} \theta_2$ . As a result,  $\lim_{\gamma_1 \rightarrow \infty} \max_{\delta_3 \in \mathbb{R}} \tilde{Q}(\alpha_2, \delta_3, \gamma_1) \rightarrow -\infty$ . By Lemma 8, we conclude that  $\max_{\delta_3 \in K_{\delta_3}} \tilde{Q}_n(\alpha_2, \delta_3, \gamma_1) \xrightarrow{P} \max_{\delta_3 \in \mathbb{R}} \tilde{Q}(\alpha_2, \delta_3, \gamma_1)$ . Since  $\tilde{Q}_n(\alpha_2, \delta_3, \gamma_1)$  is convex in  $\alpha_2$ , it should retain its convexity in  $\alpha_2$  after being maximized with respect to  $\delta_2$  and  $\gamma_1$ . Since the above equation holds for any  $\alpha_2 \in A$ , by Lemma 7, we conclude that Claim (i) holds.

The first-order condition implies a unique solution:  $\alpha_2^* := \arg \min_{\alpha_2} \max_{\delta_3 \in \mathbb{R}} \tilde{Q}(\alpha_2, \delta_3, \gamma_1)$ , which is given by  $\theta_2 \sigma_x^3 \left( \frac{\sigma_\beta^2 \theta_1}{2\lambda^2 \sigma_\beta} - \frac{\sigma_\beta}{\lambda} \right)$ . Thus, Claim (ii) holds true, concluding the proof.  $\square$

**Lemma 19.** *Under the conditions of Theorem 3, there exists a constant  $\tilde{c} > 0$  that depends solely on fixed constants, such that w.p.a.1, inequality (32) holds. In addition, as  $n \rightarrow \infty$ , for any given fixed  $\lambda > 0$ , Eq. (33) holds.*

*Proof.* Under the condition that  $\lambda_n \geq \epsilon$ , using (30) and (31), we have, w.p.a.1,

$$\begin{aligned} \|\hat{\beta}_{\lambda_n}^i\|^2 &= \left\| \frac{1}{n} \left( \frac{1}{n} X_{(-i)}^\top X_{(-i)} + c_n \lambda_n \mathbb{I} \right)^{-1} X_{(-i)} y_{(-i)} \right\|^2 \leq \frac{\|X_{(-i)} y_{(-i)}\|^2}{c_n^2 \epsilon^2 n^2} \leq \frac{2C_2 \sigma_\epsilon^2 (1 + \sqrt{c_n})^2}{c_n^2 \epsilon^2}, \\ \|\hat{\beta}_{\lambda_1}^i - \hat{\beta}_{\lambda_2}^i\|^2 &= c_n^2 (\lambda_1 - \lambda_2)^2 \left\| \frac{1}{n} \left( \frac{1}{n} X_{(-i)}^\top X_{(-i)} + c_n \lambda_2 \right)^{-1} \left( \frac{1}{n} X_{(-i)}^\top X_{(-i)} + c_n \lambda_1 \right)^{-1} X_{(-i)} y_{(-i)} \right\|^2 \end{aligned}$$

$$\leq \frac{(\lambda_1 - \lambda_2)^2}{n^2 c_n^2 \lambda_1^2 \lambda_2^2} \|X_{(-i)} y_{(-i)}\|^2 \leq \frac{2C_2 \sigma_\varepsilon^2 (1 + \sqrt{c_n})^2 (\lambda_1 - \lambda_2)^2}{c_n^2 \varepsilon^4}. \quad (\text{C18})$$

With the inequalities above and triangle inequalities, we obtain, w.a.p.1,

$$\begin{aligned} |\hat{R}^{K-CV}(\lambda_1) - \hat{R}^{K-CV}(\lambda_2)| &= \frac{1}{n} \left| \sum_{i=1}^K \left( \|y_{(i)} - X_{(i)} \hat{\beta}_{\lambda_1}^i\|^2 - \|y_{(i)} - X_{(i)} \hat{\beta}_{\lambda_2}^i\|^2 \right) \right| \\ &\leq \frac{2}{n} \sum_{i=1}^K \|X_{(i)}\| \|\hat{\beta}_{\lambda_1}^i - \hat{\beta}_{\lambda_2}^i\| \left( \|y_{(i)}\| + \frac{2C_2^{1/2} \sigma_\varepsilon (1 + \sqrt{c_n})}{c_n \varepsilon} \|X_{(i)}\| \right) \leq \tilde{C} |\lambda_1 - \lambda_2|, \end{aligned} \quad (\text{C19})$$

where  $\tilde{C}$  is some fixed constant.

Let us fix a constant  $\tilde{c}$  such that the inequality  $\frac{c_2^2 \sigma_\varepsilon^2}{2K \tilde{c}^2} - 4C_2^2 \frac{(1 + \sqrt{c_n})^2}{c_n \tilde{c}} > 100$  remains true as  $n, p \rightarrow \infty$ . This is possible because  $(1 + \sqrt{c_n})^2 / c_n$  is bounded as  $n, p \rightarrow \infty$ . In addition, let  $S := \{\lambda_j = \epsilon + p^{-9}(j-1) : 1 \leq j \leq 1 + [p^9(\tilde{c}\tau^{-1} - \epsilon)]\}$ . Given  $\tau^{-1} = o(p)$ , the cardinality of the set satisfies  $|S| \leq p^{10}$ . By definition, for any  $\lambda \in [\epsilon, \tilde{c}\tau^{-1}]$ , there exists a  $\lambda_{j^*} \in S$  such that  $|\lambda - \lambda_{j^*}| \leq p^{-9}$ . By Eq. (C19), we have  $|\hat{R}^{K-CV}(\lambda) - \hat{R}^{K-CV}(\lambda_{j^*})| \leq \tilde{C} |\lambda - \lambda_{j^*}| \leq \tilde{C} p^{-9}$ . Therefore, if we show that

$$\inf_{\lambda_j \in S} \left\{ \hat{R}^{K-CV}(\lambda_j) - \frac{1}{n} \|\varepsilon\|^2 - \|\Sigma_2^{1/2} \beta_0\|^2 \right\} > np^{-1} \tau^2 \quad (\text{C20})$$

holds w.p.a.1, we have

$$\inf_{\lambda \in [\epsilon, \tilde{c}\tau^{-1}]} \left\{ \hat{R}^{K-CV}(\lambda) - \frac{1}{n} \|\varepsilon\|^2 - \|\Sigma_2^{1/2} \beta_0\|^2 \right\} > np^{-1} \tau^2 - \tilde{C} p^{-9} > \frac{np^{-1} \tau^2}{2}, \quad (\text{C21})$$

which implies Eq. (32). By the definition of  $\hat{R}^{K-CV}$ , it is easy to verify that we need prove:

$$\inf_{\lambda_j \in S} \left\{ n^{-1} K \|Z_{(i)} \Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\|^2 - 2n^{-1} K \varepsilon_{(i)} Z_{(i)} \Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0) - \|\Sigma_2^{1/2} \beta_0\|^2 \right\} > np^{-1} \tau^2$$

holds w.p.a.1 for all  $i = 1, \dots, K$ . By the independence of  $Z_{(i)}$  and  $\hat{\beta}_{\lambda_j}^i$ , the first term on the left-hand-side is distributed as:  $n^{-1} K \|Z_{(i)} \Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\|^2 \stackrel{d}{=} n^{-1} K \chi^2(K^{-1}n) \|\Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\|^2$ , where  $\chi^2(K^{-1}n)$  denotes a Chi-squared random variable with  $K^{-1}n$  degrees of freedom. Consequently, we can deduce that:

$$\mathbb{P} \left( \left| n^{-1} K \|Z_{(i)} \Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\|^2 - \|\Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\|^2 \right| \geq \frac{\log(p)}{\sqrt{n}} \|\Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\|^2 \right)$$



$$= \mathbb{P} \left( \left| n^{-1} K \chi^2 (K^{-1} n) - 1 \right| \geq \frac{\log(p)}{\sqrt{n}} \right) \leq 2 \exp(-\tilde{c}_1 \log^2(p)),$$

where the last step uses Lemma 5, and  $\tilde{c}_1$  is a fixed positive constant. Analogously, we have:

$$\mathbb{P} \left( \left| n^{-1} K \varepsilon_{(i)} Z_{(i)} \Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0) \right| \geq \frac{\log(p)}{\sqrt{n}} \|\Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\| \right) \leq 2 \exp(-\tilde{c}_2 \log^2(p)),$$

with  $\tilde{c}_2$  being another fixed positive constant. For simplicity, we consolidate the constants  $\tilde{c}_1$  and  $\tilde{c}_2$  into a unified constant denoted as  $\tilde{c}_1$ . By the union bound inequality, we have that with probability exceeding  $1 - 4p^{10} \exp(-\tilde{c}_1 \log^2(p))$ , the following relation holds:

$$\begin{aligned} & \inf_{\lambda_j \in S} \left\{ n^{-1} K \|\Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\|^2 - 2n^{-1} K \varepsilon_{(i)} Z_{(i)} \Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0) - \|\Sigma_2^{1/2} \beta_0\|^2 \right\} \\ & \geq \inf_{\lambda_j \in S} \left\{ \left( 1 - \frac{\log(p)}{\sqrt{n}} \right) \|\Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\|^2 - \frac{\log(p)}{\sqrt{n}} \|\Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\| - \|\Sigma_2^{1/2} \beta_0\|^2 \right\}. \end{aligned}$$

Assume for now that  $\|\Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\|^2 - \|\Sigma_2^{1/2} \beta_0\|^2 \geq 50np^{-1}\tau^2$  holds. In this scenario,  $\left( 1 - \frac{\log(p)}{\sqrt{n}} \right) \|\Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\|^2 - \frac{\log(p)}{\sqrt{n}} \|\Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\|$  is monotonically increasing in  $\|\Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\|$  since  $\log(p)/\sqrt{n} = o(\tau)$ , hence it achieves its minimum when  $\|\Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\|^2 - \|\Sigma_2^{1/2} \beta_0\|^2 = 50np^{-1}\tau^2$ . As a result, it can be shown that  $\left( 1 - \frac{\log(p)}{\sqrt{n}} \right) \|\Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\|^2 - \frac{\log(p)}{\sqrt{n}} \|\Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\| - \|\Sigma_2^{1/2} \beta_0\|^2 \geq np^{-1}\tau^2$ . Therefore, we only need to prove  $\inf_{\lambda_j \in S} \{ \|\Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\|^2 - \|\Sigma_2^{1/2} \beta_0\|^2 \} \geq 50np^{-1}\tau^2$  holds w.p.a.1.

We now establish a uniform lower bound for  $\|\Sigma_2^{1/2} (\hat{\beta}_{\lambda_j}^i - \beta_0)\|^2 - \|\Sigma_2^{1/2} \beta_0\|^2$ , which can be written as:  $\|\Sigma_2^{1/2} \hat{\beta}_{\lambda_j}^i\|^2 - 2\beta_0^\top \Sigma_2 \hat{\beta}_{\lambda_j}^i$ . By direct calculation, we have for each  $i$ ,

$$\begin{aligned} \|\Sigma_2^{1/2} \hat{\beta}_{\lambda_j}^i\|^2 & \geq c_2 \|\hat{\beta}_{\lambda_j}^i\|^2 = c_2 \left\| \frac{1}{n} \left( \frac{1}{n} X_{(-i)}^\top X_{(-i)} + c_n \lambda_j \mathbb{I} \right)^{-1} X_{(-i)} y_{(-i)} \right\|^2 \\ & \geq \frac{c_2}{n^2} \left\| \frac{1}{n} X_{(-i)}^\top X_{(-i)} + c_n \lambda_j \mathbb{I} \right\|^{-2} \|X_{(-i)} y_{(-i)}\|^2 \geq \frac{c_2}{n^2} (C_2(1 + \sqrt{c_n})^2 + c_n \lambda_j)^{-2} \|X_{(-i)}^\top y_{(-i)}\|^2. \end{aligned}$$

Further, by Lemmas 2 and 3, we have

$$\|X_{(-i)}^\top y_{(-i)}\|^2 = \sigma_\varepsilon^2 \text{Tr}(X_{(-i)} X_{(-i)}^\top) + p^{-1} \tau \sigma_\beta^2 \text{Tr}(X_{(-i)}^\top X_{(-i)} X_{(-i)}^\top X_{(-i)}) + o_{\mathbb{P}}(n^{-1/2}).$$

By the fact that  $\lambda_{\min}(A) \text{Tr}(B) \leq \text{Tr}(AB) \leq \lambda_{\max}(A) \text{Tr}(B)$  when  $A, B$  are positive semidefinite, we have  $c_2 \text{Tr}(Z_{(-i)} Z_{(-i)}^\top) \leq \text{Tr}(X_{(-i)} X_{(-i)}^\top) \leq C_2 \text{Tr}(Z_{(-i)} Z_{(-i)}^\top)$ , which, along with

$(np)^{-1} \text{Tr}(Z_{(-i)}Z_{(-i)}^\top) \xrightarrow{P} (K-1)/K$  and Eq. (30), imply that

$$p^{-1}\tau \text{Tr}(X_{(-i)}^\top X_{(-i)} X_{(-i)}^\top X_{(-i)}) \leq p^{-1}\tau \|X_{(-i)}^\top X_{(-i)}\| \text{Tr}(X_{(-i)}^\top X_{(-i)}) \lesssim_P \tau pn = o_P(np).$$

Therefore, w.p.a.1, we obtain

$$\frac{c_2\sigma_\varepsilon^2 pn}{2K} \leq \|X_{(-i)}^\top y_{(-i)}\|^2 \leq 2C_2\sigma_\varepsilon^2 pn. \quad (\text{C22})$$

Consequently, uniformly over  $\lambda_j \in S$ , we deduce:

$$\|\Sigma_2^{1/2} \hat{\beta}_{\lambda_j}^i\|^2 \geq \frac{c_2^2\sigma_\varepsilon^2 p}{2nK} (C_2(1 + \sqrt{c_n})^2 + c_n\lambda_j)^{-2}. \quad (\text{C23})$$

On the other hand, we have

$$\begin{aligned} |\beta_0^\top \Sigma_2 \hat{\beta}_{\lambda_j}^i| &\leq \frac{1}{n} \left| \beta_0^\top \Sigma_2 \left( \frac{1}{n} X_{(-i)}^\top X_{(-i)} + c_n \lambda_j \mathbb{I} \right)^{-1} X_{(-i)}^\top X_{(-i)} \beta_0 \right| \\ &\quad + \frac{1}{n} \left| \varepsilon_{(-i)}^\top X_{(-i)} \left( \frac{1}{n} X_{(-i)}^\top X_{(-i)} + c_n \lambda_j \mathbb{I} \right)^{-1} \Sigma_2 \beta_0 \right| =: L_1 + L_2. \end{aligned} \quad (\text{C24})$$

To bound  $L_1$ , note that  $\text{Tr}(AB) \leq \|AB\| \text{rank}(AB) \leq \|A\| \|B\| \text{rank}(B)$ , which implies

$$\begin{aligned} &\frac{1}{n} \text{Tr} \left( \Sigma_2 \left( \frac{1}{n} X_{(-i)}^\top X_{(-i)} + c_n \lambda_j \mathbb{I} \right)^{-1} X_{(-i)}^\top X_{(-i)} \right) \\ &\leq \|\Sigma_2\| \left\| \left( \frac{1}{n} X_{(-i)}^\top X_{(-i)} + c_n \lambda_j \mathbb{I} \right)^{-1} \frac{1}{n} X_{(-i)}^\top X_{(-i)} \right\| \text{rank}(X_{(-i)}^\top X_{(-i)}) \leq \frac{nC_2^2(1 + \sqrt{c_n})^2}{C_2(1 + \sqrt{c_n})^2 + c_n\lambda_j}, \end{aligned}$$

where the last inequality uses  $\lambda_1((A + \mathbb{I})^{-1}A) = (\lambda_1(A) + 1)^{-1}\lambda_1(A)$  and  $n^{-1}\|X_{(-i)}^\top X_{(-i)}\| \leq C_2(1 + \sqrt{c_n})^2$ . In addition, by Lemma 5 and the fact that the sub-exponential norm of  $b_{0,i}$  is of order  $O(q^{-1/2})$ , we have, with probability exceeding  $1 - 2p^{10} \exp(-\tilde{c}_1 \log^2(p))$ ,

$$\sup_{\lambda_j \in S} \left| L_1 - \frac{p^{-1}\tau}{n} \text{Tr} \left( \Sigma_2 \left( \frac{1}{n} X_{(-i)}^\top X_{(-i)} + c_n \lambda_j \mathbb{I} \right)^{-1} X_{(-i)}^\top X_{(-i)} \right) \right| \leq q^{-1} p^{-1} \tau n^{1/2} \log(p).$$

Combining the above two inequalities, we have

$$L_1 \leq np^{-1}\tau C_2^2 \frac{(1 + \sqrt{c_n})^2}{C_2(1 + \sqrt{c_n})^2 + c_n\lambda_j} + q^{-1} p^{-1} \tau n^{1/2} \log(p), \quad \forall \lambda_j \in S.$$

To bound  $L_2$  in (C24), by definition, we have  $L_2 = |n^{-1}p^{-1/2}\tau^{1/2}q^{-1/2}z_{(-i)}^\top X_{(-i)}(\frac{1}{n}X_{(-i)}^\top X_{(-i)} + c_n\lambda_j\mathbb{I})^{-1}\Sigma_2(\sqrt{qb_0})|$ . Using the facts that  $\lambda_{\min}(n^{-1}X_{(-i)}^\top X_{(-i)} + c_n\lambda_j\mathbb{I}) \geq c_n\lambda_j \geq c_n\epsilon$ ,  $\|A\|_F \leq \sqrt{\text{rank}(A)}\|A\|$ , and Eq. (30), we have

$$\begin{aligned} \|X_{(-i)}(\frac{1}{n}X_{(-i)}^\top X_{(-i)} + c_n\lambda_j\mathbb{I})^{-1}\Sigma_2\| &\leq C_2c_n^{-1}\epsilon^{-1}\|X_{(-i)}\| \lesssim np^{-1/2}, \\ \|X_{(-i)}(\frac{1}{n}X_{(-i)}^\top X_{(-i)} + c_n\lambda_j\mathbb{I})^{-1}\Sigma_2\|_F^2 &\lesssim \text{rank}(X_{(-i)})n^2p^{-1} \lesssim n^3p^{-1}. \end{aligned}$$

Therefore, by Lemma 5 and the fact that  $\sqrt{qb_{0,i}}$  has bounded sub-exponential norm, it holds that, for some constant  $\tilde{c}_1$ ,  $\mathbb{P}(L_2 > q^{-1/2}n^{1/2}\tau^{1/2}p^{-1}\log(p)) \leq 2\exp(-\tilde{c}_1\log^2(p))$ . As a consequence, with probability at least  $1 - 2p^{10}\exp(-\tilde{c}_1\log^2(p))$ , we have  $\sup_{\lambda_j \in S} L_2 \leq q^{-1/2}n^{1/2}\tau^{1/2}p^{-1}\log(p)$ . Therefore, taking bounds for  $L_1$  and  $L_2$  together, we have, w.p.a.1,

$$\begin{aligned} |\beta_0^\top \Sigma_2 \hat{\beta}_{\lambda_j}^i| &\leq np^{-1}\tau C_2^2 \frac{(1 + \sqrt{c_n})^2}{C_2(1 + \sqrt{c_n})^2 + c_n\lambda_j} + p^{-1}q^{-1}\tau n^{1/2}\log(p) + q^{-1/2}n^{1/2}\tau^{1/2}p^{-1}\log(p) \\ &\leq 2np^{-1}\tau C_2^2 \frac{(1 + \sqrt{c_n})^2}{C_2(1 + \sqrt{c_n})^2 + c_n\lambda_j}, \end{aligned} \quad (\text{C25})$$

for each  $\lambda_j \in S$ . In the last inequality, we use  $p^{-1}q^{-1}\tau n^{1/2}\log(p)$ ,  $q^{-1/2}n^{1/2}\tau^{1/2}p^{-1}\log(p) = o(np^{-1}\tau^2)$  by the assumptions of Theorem 3. With (C23) and (C25), we have

$$\begin{aligned} \|\Sigma_2^{1/2}(\hat{\beta}_{\lambda_j}^i - \beta_0)\|^2 - \|\Sigma_2^{1/2}\beta_0\|^2 &= \|\Sigma_2^{1/2}\hat{\beta}_{\lambda_j}^i\|^2 - 2\beta_0^\top \Sigma_2 \hat{\beta}_{\lambda_j}^i \\ &\geq \frac{c_n^2\sigma_\epsilon^2 p}{2nK} (C_2(1 + \sqrt{c_n})^2 + c_n\lambda_j)^{-2} - 4np^{-1}\tau C_2^2 \frac{(1 + \sqrt{c_n})^2}{C_2(1 + \sqrt{c_n})^2 + c_n\lambda_j}. \end{aligned}$$

This inequality holds w.p.a. 1 as  $n, p \rightarrow \infty$  uniformly for all  $\lambda_j \in S$ . Given our initial choice for  $\tilde{c}$ , it is easy to check that the right-hand side exceeds  $50np^{-1}\tau^2$ , which implies Eq. (32).

To prove Eq. (33), note that

$$\frac{1}{n}\|y_{(i)} - X_{(i)}\hat{\beta}_{\tau^{-1}\lambda}^i\|^2 - \frac{1}{n}\|\varepsilon_{(i)}\|^2 = \frac{1}{n}\|Z_{(i)}\Sigma_2^{1/2}(\hat{\beta}_{\tau^{-1}\lambda}^i - \beta_0)\|^2 + \frac{2}{n}\varepsilon_{(i)}^\top Z_{(i)}\Sigma_2^{1/2}(\hat{\beta}_{\tau^{-1}\lambda}^i - \beta_0).$$

By the facts  $Z_{(i)} \perp \hat{\beta}_{\tau^{-1}\lambda}^i - \beta_0$  and  $n^{-1}\chi^2(K^{-1}n) = K^{-1} + O_p(n^{-1/2})$ , we have

$$\begin{aligned} \frac{1}{n}\|Z_{(i)}\Sigma_2^{1/2}(\hat{\beta}_{\tau^{-1}\lambda}^i - \beta_0)\|^2 &\stackrel{d}{=} \frac{1}{n}\chi^2(K^{-1}n)\|\Sigma_2^{1/2}(\hat{\beta}_{\tau^{-1}\lambda}^i - \beta_0)\|^2 \\ &= \frac{1}{K}\|\Sigma_2^{1/2}(\hat{\beta}_{\tau^{-1}\lambda}^i - \beta_0)\|^2 + O_P\left(\frac{1}{\sqrt{n}}\|\Sigma_2^{1/2}(\hat{\beta}_{\tau^{-1}\lambda}^i - \beta_0)\|^2\right). \end{aligned}$$

Additionally, by Theorem 2, we deduce:

$$\|\Sigma_2^{1/2}(\hat{\beta}_{\tau^{-1}\lambda}^i - \beta_0)\|^2 - \|\Sigma_2^{1/2}\beta_0\|^2 = \frac{2(K-1)}{K}np^{-1}\tau^2\theta_2\sigma_x^4\left(\frac{\sigma_\varepsilon^2}{2\lambda^2} - \frac{\sigma_\beta^2}{\lambda}\right) + o_{\mathbb{P}}(\tau^2np^{-1}).$$

Hence, using the fact that  $\|\Sigma_2^{1/2}\beta_0\|^2 \asymp_{\mathbb{P}} \tau$ , we derive the following equation:

$$\frac{1}{n}\sum_{i=1}^K\|Z_{(i)}\Sigma_2^{1/2}(\hat{\beta}_{\tau^{-1}\lambda}^i - \beta_0)\|^2 - \|\Sigma_2^{1/2}\beta_0\|^2 = \frac{2(K-1)}{K}np^{-1}\tau^2\theta_2\sigma_x^4\left(\frac{\sigma_\varepsilon^2}{2\lambda^2} - \frac{\sigma_\beta^2}{\lambda}\right) + o_{\mathbb{P}}(\tau^2np^{-1}).$$

Thus, to prove Eq. (33), it remains to show that:  $\frac{2}{n}\varepsilon_{(i)}Z_{(i)}\Sigma_2^{1/2}(\hat{\beta}_{\tau^{-1}\lambda}^i - \beta_0) = o_{\mathbb{P}}(\tau^2np^{-1})$ . Given that  $\varepsilon_{(i)} \perp Z_{(i)}\Sigma_2^{1/2}(\hat{\beta}_{\tau^{-1}\lambda}^i - \beta_0)$  and  $n^{-3/2}\tau^{-3/2}p \rightarrow 0$  by Assumption 4, we have

$$\frac{2}{n}\varepsilon_{(i)}Z_{(i)}\Sigma_2^{1/2}(\hat{\beta}_{\tau^{-1}\lambda}^i - \beta_0) \stackrel{d}{=} \frac{2}{n}\|\Sigma_2^{1/2}(\hat{\beta}_{\tau^{-1}\lambda}^i - \beta_0)\|\varepsilon_{(i)}^\top x = O_{\mathbb{P}}(n^{-1/2}\tau^{1/2}) = o_{\mathbb{P}}(\tau^2np^{-1}),$$

where  $x$  is a standard Gaussian vector independent of  $\varepsilon_{(i)}$ . This concludes the proof.  $\square$

**Lemma 20.** *There exists a constant  $\tilde{C}_1$  such that, w.p.a.1, uniformly for  $\mu_1, \mu_2 \in [0, \tilde{c}^{-1}]$ ,*

$$pn^{-1}\tau^{-2}|\tilde{R}^{K-CV}(\mu_1) - \tilde{R}^{K-CV}(\mu_2)| \leq \tilde{C}_1|\mu_1 - \mu_2| + o_{\mathbb{P}}(pn^{-1}\tau^{-2}).$$

*Proof.* By the Woodbury identity, we deduce that

$$\left(\frac{1}{n}X_{(-i)}^\top X_{(-i)} + c_n\tau^{-1}\mu^{-1}\mathbb{I}\right)^{-1} - c_n^{-1}\tau\mu\mathbb{I} = -\frac{c_n^{-2}\tau^2\mu^2}{n}X_{(-i)}^\top \left(\mathbb{I} + \frac{c_n^{-1}\tau\mu}{n}X_{(-i)}X_{(-i)}^\top\right)^{-1}X_{(-i)}.$$

Hence, we arrive at:

$$\begin{aligned} & \sup_{\substack{1 \leq i \leq K \\ \mu \in [0, \tilde{c}^{-1}]}} c_n\tau^{-3}\log^{-1}(p)\left\|\left(\frac{1}{n}X_{(-i)}^\top X_{(-i)} + c_n\tau^{-1}\mu^{-1}\mathbb{I}\right)^{-1} - c_n^{-1}\tau\mu\mathbb{I} + \frac{c_n^{-2}\tau^2\mu^2}{n}X_{(-i)}^\top X_{(-i)}\right\| \\ &= \sup_{\substack{1 \leq i \leq K \\ \mu \in [0, \tilde{c}^{-1}]}} c_n^{-1}\mu^2\tau^{-1}\log^{-1}(p)\left\|\frac{1}{n}X_{(-i)}^\top \left[\left(\mathbb{I} + \frac{c_n^{-1}\tau\mu}{n}X_{(-i)}X_{(-i)}^\top\right)^{-1} - \mathbb{I}\right]X_{(-i)}\right\| \\ &\leq \sup_{\substack{1 \leq i \leq K \\ \mu \in [0, \tilde{c}^{-1}]}} \mu^3c_n^{-2}\log^{-1}(p)\left\|\frac{1}{n}X_{(-i)}^\top X_{(-i)}\right\|^2 \xrightarrow{\mathbb{P}} 0. \end{aligned} \tag{C26}$$

The last inequality is a consequence of Eq. (30) and the fact that

$$\begin{aligned} \left\| \left( \mathbb{I} + \frac{c_n^{-1}\tau\mu}{n} X_{(-i)} X_{(-i)}^\top \right)^{-1} - \mathbb{I} \right\| &\leq \left\| \left( \mathbb{I} + \frac{c_n^{-1}\tau\mu}{n} X_{(-i)} X_{(-i)}^\top \right)^{-1} \right\| \cdot c_n^{-1}\tau\mu \left\| \frac{1}{n} X_{(-i)}^\top X_{(-i)} \right\| \\ &\leq c_n^{-1}\tau\mu \left\| \frac{1}{n} X_{(-i)}^\top X_{(-i)} \right\|. \end{aligned}$$

On the other hand, by direct calculation we have that  $\widetilde{R}^{K-CV}(\mu_1) - \widetilde{R}^{K-CV}(\mu_2)$  equals:

$$\begin{aligned} &\sum_{i=1}^K \left( \frac{1}{n} \|X_{(i)} \hat{\beta}_{\tau^{-1}\mu_1^{-1}}^i\|^2 - \frac{1}{n} \|X_{(i)} \hat{\beta}_{\tau^{-1}\mu_2^{-1}}^i\|^2 \right) - \frac{2}{n} y_{(i)}^\top X_{(i)} (\hat{\beta}_{\tau^{-1}\mu_1^{-1}}^i - \hat{\beta}_{\tau^{-1}\mu_2^{-1}}^i) \\ &=: \sum_{i=1}^K W_{1i}(\mu_1, \mu_2) - W_{2i}(\mu_1, \mu_2). \end{aligned}$$

We next investigate  $W_{1i}(\mu_1, \mu_2)$  and  $W_{2i}(\mu_1, \mu_2)$  separately. For  $W_{1i}(\mu_1, \mu_2)$ , we have

$$\begin{aligned} W_{1i}(\mu_1, \mu_2) &= \frac{1}{n} (\hat{\beta}_{\tau^{-1}\mu_1^{-1}}^i - \hat{\beta}_{\tau^{-1}\mu_2^{-1}}^i)^\top X_{(i)}^\top X_{(i)} \hat{\beta}_{\tau^{-1}\mu_1^{-1}}^i + \frac{1}{n} \hat{\beta}_{\tau^{-1}\mu_2^{-1}}^i X_{(i)}^\top X_{(i)} (\hat{\beta}_{\tau^{-1}\mu_1^{-1}}^i - \hat{\beta}_{\tau^{-1}\mu_2^{-1}}^i) \\ &\leq \frac{1}{n} \left\| X_{(i)} (\hat{\beta}_{\tau^{-1}\mu_1^{-1}}^i - \hat{\beta}_{\tau^{-1}\mu_2^{-1}}^i) \right\| \cdot \|X_{(i)} \hat{\beta}_{\tau^{-1}\mu_1^{-1}}^i\| + \frac{1}{n} \|X_{(i)} \hat{\beta}_{\tau^{-1}\mu_2^{-1}}^i\| \cdot \left\| X_{(i)} (\hat{\beta}_{\tau^{-1}\mu_1^{-1}}^i - \hat{\beta}_{\tau^{-1}\mu_2^{-1}}^i) \right\|. \end{aligned}$$

Define  $\tilde{\beta}_{\tau^{-1}\mu_1^{-1}}^i = \frac{1}{n} \left[ c_n^{-1}\tau\mu_1 \mathbb{I} - \frac{c_n^{-2}\tau^2\mu_1^2}{n} X_{(-i)}^\top X_{(-i)} \right] X_{(-i)}^\top y_{(-i)}$ . Observe that

$$\begin{aligned} &\sup_{\mu_1 \in [0, \tilde{c}^{-1}]} \frac{1}{\sqrt{n}} \left\| X_{(i)} \hat{\beta}_{\tau^{-1}\mu_1^{-1}}^i - X_{(i)} \tilde{\beta}_{\tau^{-1}\mu_1^{-1}}^i \right\| \\ &\leq \sup_{\mu_1 \in [0, \tilde{c}^{-1}]} \frac{1}{n^{3/2}} \|X_{(i)}\| \left\| \left( \frac{1}{n} X_{(-i)}^\top X_{(-i)} + c_n \tau^{-1} \mu_1^{-1} \right)^{-1} - c_n^{-1} \tau \mu_1 \mathbb{I} + \frac{c_n^{-2} \tau^2 \mu_1^2}{n} X_{(-i)}^\top X_{(-i)} \right\| \\ &\quad \times \|X_{(-i)}^\top y_{(-i)}\| = O_P(\tau^3 \log(p)) = o_P(c_n^{-1/2} \tau), \end{aligned} \tag{C27}$$

where we use Eq. (C26), Eq. (30), and Eq. (C22). Additionally, it is easy to verify that

$$\begin{aligned} \sup_{\mu_1 \in [0, \tilde{c}^{-1}]} \frac{1}{\sqrt{n}} \left\| \frac{1}{n} X_{(i)} \tilde{\beta}_{\tau^{-1}\mu_1^{-1}}^i \right\| &= \sup_{\mu_1 \in [0, \tilde{c}^{-1}]} \frac{1}{\sqrt{n}} \left\| \frac{1}{n} X_{(i)} \left[ c_n^{-1} \tau \mu_1 \mathbb{I} - \frac{c_n^{-2} \tau^2 \mu_1^2}{n} X_{(-i)}^\top X_{(-i)} \right] X_{(-i)}^\top y_{(-i)} \right\| \\ &\leq \frac{c_n^{-1} \tau \tilde{c}^{-1}}{n \sqrt{n}} \|X_{(i)} X_{(-i)}^\top y_{(-i)}\| + \frac{c_n^{-2} \tau^2 \tilde{c}^{-2}}{n^2 \sqrt{n}} \|X_{(i)} X_{(-i)}^\top X_{(-i)} X_{(-i)}^\top y_{(-i)}\|. \end{aligned}$$

For the first term, by Eq. (C22), it equals

$$\frac{c_n^{-1} \tau \tilde{c}^{-1}}{n \sqrt{n}} \left\| Z_{(i)} \Sigma_2^{1/2} X_{(-i)}^\top y_{(-i)} \right\| \stackrel{d}{=} \frac{c_n^{-1} \tau \tilde{c}^{-1}}{n \sqrt{n}} \sqrt{\chi^2(n/K)} \left\| \Sigma_2^{1/2} X_{(-i)}^\top y_{(-i)} \right\| \leq \frac{\tilde{C}_1}{2} c_n^{-1/2} \tau,$$

w.p.a.1 for some constant  $\tilde{C}_1$  that only depends on fixed constants. The second term can be bounded in the same way. Therefore, we have

$$\sup_{\mu_1 \in [0, \tilde{c}^{-1}]} \frac{1}{\sqrt{n}} \|X_{(i)} \hat{\beta}_{\tau^{-1}\mu_1^{-1}}^i\| \leq \tilde{C}_1 c_n^{-1/2} \tau + o_{\mathbb{P}}(c_n^{-1/2} \tau). \quad (\text{C28})$$

Analogously, we can prove that  $\frac{1}{\sqrt{n}} \|X_{(i)} (\hat{\beta}_{\tau^{-1}\mu_1^{-1}}^i - \hat{\beta}_{\tau^{-1}\mu_2^{-1}}^i)\| \leq \tilde{C}_1 |\mu_1 - \mu_2| c_n^{-1/2} \tau + o_{\mathbb{P}}(c_n^{-1/2} \tau)$  holds uniformly for  $\mu_1, \mu_2 \in [0, \tilde{c}^{-1}]$ , where  $\tilde{C}_1$  is a fixed constant that may vary from line to line. In light of this, we deduce that:  $\sup_{1 \leq i \leq K} W_{1i}(\mu_1, \mu_2) \leq \tilde{C}_1^2 c_n^{-1} \tau^2 |\mu_1 - \mu_2| + o_{\mathbb{P}}(c_n^{-1} \tau^2)$  holds w.p.a.1 uniformly for  $\mu_1, \mu_2 \in [0, \tilde{c}^{-1}]$ .

To bound  $W_{2i}(\mu_1, \mu_2)$ , we first define  $\tilde{W}_{2i}(\mu_1, \mu_2) = \frac{2}{n} y_{(i)}^\top X_{(i)} (\tilde{\beta}_{\tau^{-1}\mu_1^{-1}}^i - \tilde{\beta}_{\tau^{-1}\mu_2^{-1}}^i)$ . By Eq. (C27), it holds that

$$\begin{aligned} \sup_{\mu_1, \mu_2 \in [0, \tilde{c}^{-1}]} |\tilde{W}_{2i}(\mu_1, \mu_2) - W_{2i}(\mu_1, \mu_2)| &\leq \frac{2}{n} \|y_{(i)}^\top\| \|X_{(i)} (\tilde{\beta}_{\tau^{-1}\mu_1^{-1}}^i - \hat{\beta}_{\tau^{-1}\mu_1^{-1}}^i)\| \\ &\quad + \frac{2}{n} \|y_{(i)}^\top\| \|X_{(i)} (\tilde{\beta}_{\tau^{-1}\mu_1^{-1}}^i - \hat{\beta}_{\tau^{-1}\mu_2^{-1}}^i)\| = O_{\mathbb{P}}(\tau^3 \log(p)) = o_{\mathbb{P}}(c_n^{-1} \tau^2). \end{aligned}$$

Moreover, employing a similar argument to that used in proving Eq. (C28), we have

$$\begin{aligned} &\sup_{\mu_1, \mu_2 \in [0, \tilde{c}^{-1}]} |\tilde{W}_{2i}(\mu_1, \mu_2)| \\ &= \sup_{\mu_1, \mu_2 \in [0, \tilde{c}^{-1}]} \left| \frac{2}{n} y_{(i)}^\top X_{(i)} \frac{1}{n} \left[ c_n^{-1} \tau (\mu_1 - \mu_2) \mathbb{I} - \frac{c_n^{-2} \tau^2 (\mu_1^2 - \mu_2^2)}{n} X_{(-i)}^\top X_{(-i)} \right] X_{(-i)}^\top y_{(-i)} \right| \\ &\lesssim |\mu_1 - \mu_2| \frac{c_n^{-1} \tau}{n^2} |y_{(i)}^\top X_{(i)} X_{(-i)}^\top y_{(-i)}| + |\mu_1 - \mu_2| \frac{c_n^{-2} \tau^2}{n^3} |y_{(i)}^\top X_{(i)} X_{(-i)}^\top X_{(-i)} X_{(-i)}^\top y_{(-i)}|. \end{aligned}$$

For the first term, by Lemmas 2 and 3, it is easy to verify that

$$\begin{aligned} \frac{c_n^{-1} \tau}{n^2} |y_{(i)}^\top X_{(i)} X_{(-i)}^\top y_{(-i)}| &\leq \frac{c_n^{-1} \tau}{n^2} |\varepsilon_{(i)}^\top X_{(i)} X_{(-i)}^\top \varepsilon_{(-i)}| + \frac{c_n^{-1} \tau}{n^2} |\varepsilon_{(i)}^\top X_{(i)} X_{(-i)}^\top X_{(-i)} \beta_0| \\ &\quad + \frac{c_n^{-1} \tau}{n^2} |\beta_0^\top X_{(i)}^\top X_{(i)} X_{(-i)}^\top \varepsilon_{(-i)}| + \frac{c_n^{-1} \tau}{n^2} |\beta_0^\top X_{(i)}^\top X_{(i)} X_{(-i)}^\top X_{(-i)} \beta_0| \leq \tilde{C}_1 c_n^{-1} \tau^2, \end{aligned}$$

for some constant  $\tilde{C}_1$  w.p.a.1. The second term can be shown analogously. As a result, we have  $\sup_{1 \leq i \leq K} W_{2i}(\mu_1, \mu_2) \leq \tilde{C}_1 c_n^{-1} \tau^2 |\mu_1 - \mu_2| + o_{\mathbb{P}}(c_n^{-1} \tau^2)$  w.p.a.1, uniformly for  $\mu_1, \mu_2 \in [0, \tilde{c}^{-1}]$ . Combining the bounds for  $W_{1i}(\mu_1, \mu_2)$  and  $W_{2i}(\mu_1, \mu_2)$  concludes the proof.  $\square$

**Lemma 21.**  $A_1$  to  $A_3$  defined in (34) converge to zero w.p.a. 1.

*Proof.* For  $A_1$ , by using the same argument as Eq. (C18) and noting that  $\lambda_n^{opt} \asymp \tau^{-1}$ ,

$\hat{\lambda}_n^{K-CV} \asymp_P \tau^{-1}$  and  $\tau(\lambda_n^{opt} - \hat{\lambda}_n^{K-CV}) = o_P(1)$ , we have

$$|A_1| \lesssim c_n \tau^{-2} \|\hat{\beta}_{cv} - \hat{\beta}_{opt}\| \|\hat{\beta}_{cv}\| \lesssim_P c_n \tau^{-2} \cdot c_n^{-1/2} |\lambda_n^{opt} - \hat{\lambda}_n^{K-CV}| \tau^2 \cdot c_n^{-1/2} \tau = o_P(1).$$

Similarly,  $A_2 = o_P(1)$ . To prove  $A_3 = o_P(1)$ , define  $\tilde{\beta}(\lambda_n) = \frac{1}{n} \left[ c_n^{-1} \lambda_n^{-1} \mathbb{I} - \frac{c_n^{-2} \lambda_n^{-2}}{n} X^\top X \right] X^\top y$  and write  $\tilde{\beta}(\hat{\lambda}_n^{K-CV})$  as  $\tilde{\beta}_{cv}$  and  $\tilde{\beta}(\lambda_n^{opt})$  as  $\tilde{\beta}_{opt}$  for simplicity. By using a similar result as Eq. (C26) as well as the fact  $n^{-1} \|X^\top y\| \lesssim n^{-1} \|X\| \|y\| \lesssim_P c_n^{1/2}$  according to Lemma 6, we have

$$\begin{aligned} c_n \tau^{-2} |(\hat{\beta}_{opt} - \tilde{\beta}_{opt})^\top \Sigma_2 \beta_0| &\lesssim c_n \tau^{-2} \|\hat{\beta}_{opt} - \tilde{\beta}_{opt}\| \|\beta_0\| \\ &\lesssim c_n \tau^{-2} \left\| \left( \frac{1}{n} X^\top X + c_n \lambda_n^{opt} \mathbb{I} \right)^{-1} - c_n^{-1} (\lambda_n^{opt})^{-1} \mathbb{I} + \frac{c_n^{-2} (\lambda_n^{opt})^{-2}}{n} X^\top X \right\| \left\| \frac{1}{n} X^\top y \right\| \cdot \|\beta_0\| \\ &= o_P(c_n \tau^{-2} \cdot c_n^{-1} \tau^3 \log(p) \cdot c_n^{1/2} \cdot \tau^{1/2}) = o_P(c_n^{1/2} \tau^{3/2} \log(p)) = o_P(1), \end{aligned}$$

where the last equation holds by Assumption 4. Similarly,  $c_n \tau^{-2} |(\hat{\beta}_{cv} - \tilde{\beta}_{cv})^\top \Sigma_2 \beta_0| = o_P(1)$ . Therefore,  $A_3 = o_P(1)$  follows by  $c_n \tau^{-2} (\tilde{\beta}_{opt} - \tilde{\beta}_{cv})^\top \Sigma_2 \beta_0 = o_P(1)$ . Note that

$$\begin{aligned} c_n \tau^{-2} (\tilde{\beta}_{opt} - \tilde{\beta}_{cv})^\top \Sigma_2 \beta_0 &= n^{-1} \tau^{-2} \left( (\hat{\lambda}_n^{K-CV})^{-1} - (\lambda_n^{opt})^{-1} \right) \beta_0^\top \Sigma_2 X^\top y \\ &\quad - n^{-2} c_n^{-1} \tau^{-2} \left( (\hat{\lambda}_n^{K-CV})^{-2} - (\lambda_n^{opt})^{-2} \right) \beta_0^\top \Sigma_2 X^\top X X^\top y =: B_1 + B_2. \end{aligned}$$

By Lemma 2 and Lemma 3, we have

$$\begin{aligned} B_1 &= n^{-1} \tau^{-2} \left( (\hat{\lambda}_n^{K-CV})^{-1} - (\lambda_n^{opt})^{-1} \right) \beta_0^\top \Sigma_2 X^\top X \beta_0 + n^{-1} \tau^{-2} \beta_0^\top \Sigma_2 \left( (\hat{\lambda}_n^{K-CV})^{-1} - (\lambda_n^{opt})^{-1} \right) X^\top \varepsilon \\ &\asymp_P n^{-1} \tau^{-1} p^{-1} \left( (\hat{\lambda}_n^{K-CV})^{-1} - (\lambda_n^{opt})^{-1} \right) \text{Tr}(\Sigma_2 X^\top X) = o_P(\tau^{-1} \left( (\hat{\lambda}_n^{K-CV})^{-1} - (\lambda_n^{opt})^{-1} \right)) = o_P(1). \end{aligned}$$

The same argument proves  $B_2 = o_P(1)$ , leading to  $A_3 = o_P(1)$ , which concludes the proof.  $\square$

**Lemma 22.** *The objective function in (36) is convex with respect to  $\alpha$  and jointly concave with respect to  $(\delta, \gamma)$ . Additionally, as long as Eq. (38) holds, we have Eq. (35).*

*Proof.* Define  $S_w^n = \{w \mid c_n \tau^{-1} \sigma_x \sigma_\beta + c_\alpha / 4 \sigma_\beta \leq c_n \|w\| \leq c_n \tau^{-1} \sigma_x \sigma_\beta + C_\alpha / \sigma_\beta\}$ . Analogous to the result proved by Lemma 15, if the solution  $\hat{w}^B$  to the following problem

$$\min_{w \in S_w^n} \frac{c_n}{n} \|\tau^{1/2} \Sigma_1^{1/2} Z \tilde{w} - \tau^{-1} \varepsilon\|^2 + \frac{c_n \tau^{-1/2} \lambda_n}{\sqrt{n}} \|\Sigma_2^{-1/2} w + \tau^{-3/2} \beta_0\|_1 - \frac{c_n \tau^{-2}}{n} \|\varepsilon\|^2 - C_n^\phi \quad (\text{C29})$$

satisfies  $c_n \|\hat{w}^B\| - c_n \tau^{-1} \sigma_x \sigma_\beta \in [c_\alpha / 2 \sigma_\beta + \epsilon, C_\alpha / 2 \sigma_\beta - \epsilon]$  w.a.p.1, then the same holds true for  $\hat{w}$ , which leads to the desired result, (35). In light of this, without ambiguity we now

directly focus on (C29), and refer to  $\hat{w}^B$  as  $\hat{w}$  for ease of notation.

Note that for any vector  $x$ , it holds that  $\|x\|^2 = \max_u \sqrt{n} u^\top x - n\|u\|^2/4$ , and  $\|x\|_1 = \max_{\|v\|_\infty \leq 1} v^\top x$ . By applying these equations to  $\|\tau^{1/2} \Sigma_1^{1/2} Z \tilde{w} - \tau^{-1} \varepsilon\|^2$  and  $\|\Sigma_2^{-1/2} w + \tau^{-3/2} \beta_0\|_1$ , and letting  $\tilde{u} := \Sigma_1^{1/2} u$ , the problem (C29) can be reformulated as:

$$\begin{aligned} \min_{w \in S_w^n} \max_{\substack{\tilde{u} \\ \|v\|_\infty \leq 1}} & \frac{c_n \tau^{1/2}}{\sqrt{n}} \tilde{u}^\top Z w - \frac{c_n \tau^{-1}}{\sqrt{n}} \tilde{u}^\top \Sigma_1^{-1/2} \varepsilon - \frac{c_n \|\Sigma_1^{-1/2} \tilde{u}\|^2}{4} + \frac{c_n \tau^{-2} \lambda_n}{\sqrt{n}} v^\top \beta_0 \\ & + \frac{c_n \tau^{-1/2} \lambda_n}{\sqrt{n}} v^\top \Sigma_2^{-1/2} w - \frac{c_n \tau^{-2}}{n} \|\varepsilon\|^2 - C_n^\phi. \end{aligned} \quad (\text{C30})$$

For convenience, we shall continue to employ  $u$  in place of  $\tilde{u}$  throughout the remainder of the proof. Let  $S_u^n = \{u \mid \|u\| \leq 4\tau^{-1} \sqrt{C_1 C_\varepsilon}\}$ . Similar to the proof of Lemma 14, w.a.p.1,

$$\begin{aligned} \min_{w \in S_w^n} \max_{\substack{u \in S_u^n \\ \|v\|_\infty \leq 1}} & \frac{c_n \tau^{1/2}}{\sqrt{n}} u^\top Z w - \frac{c_n \tau^{-1}}{\sqrt{n}} u^\top \Sigma_1^{-1/2} \varepsilon - \frac{c_n \|\Sigma_1^{-1/2} u\|^2}{4} + \frac{c_n \tau^{-2} \lambda_n}{\sqrt{n}} v^\top \beta_0 \\ & + \frac{c_n \tau^{-1/2} \lambda_n}{\sqrt{n}} v^\top \Sigma_2^{-1/2} w - \frac{c_n \tau^{-2}}{n} \|\varepsilon\|^2 - C_n^\phi \end{aligned} \quad (\text{C31})$$

is equivalent to Eq. (C30). Next, we construct an auxiliary optimization problem:

$$\begin{aligned} \phi(g, h) &= \max_{\substack{0 \leq \delta \leq 4\tau^{-1} \sqrt{C_1 C_\varepsilon} \\ \|v\|_\infty \leq 1}} \min_{w \in S_w^n} \max_{\|u\|=\delta} \mathcal{R}_n(w, v, u), \quad \text{where} \\ \mathcal{R}_n(w, v, u) &= \frac{c_n \tau^{1/2}}{\sqrt{n}} \|w\| g^\top u - \frac{c_n \tau^{1/2}}{\sqrt{n}} \|u\| h^\top w - \frac{c_n \tau^{-1}}{\sqrt{n}} u^\top \Sigma_1^{-1/2} \varepsilon - \frac{c_n \|\Sigma_1^{-1/2} u\|^2}{4} \\ &+ \frac{c_n \tau^{-2} \lambda_n}{\sqrt{n}} v^\top \beta_0 + \frac{c_n \tau^{-1/2} \lambda_n}{\sqrt{n}} v^\top \Sigma_2^{-1/2} w - \frac{c_n \tau^{-2}}{n} \|\varepsilon\|^2 - C_n^\phi, \end{aligned} \quad (\text{C32})$$

and both  $g \in \mathbb{R}^n$  and  $h \in \mathbb{R}^p$  are standard Gaussian vectors, independent of all other random variables. Moreover, let  $\tilde{\mathcal{S}}_n := \{w \mid c_\alpha/2\sigma_\beta + \epsilon < c_n \|w\| - c_n \tau^{-1} \sigma_x \sigma_\beta < C_\alpha/2\sigma_\beta - \epsilon\}$ , define  $\phi_{\tilde{\mathcal{S}}_n^c}(g, h)$  as the optimal value of the optimization problem (C32), with  $w \in S_w^n \cap \tilde{\mathcal{S}}_n^c$ .

Lemma 23 characterizes the limiting behavior of the optimal solution to (C31),  $\hat{w}$ , and in turn, proves the desired (35), under conditions pertaining to the optimization problem (C32). Therefore, we only need show that conditions outlined in Lemma 23 hold as long as (38) holds. That is, under (38), we need to prove the existence of the constants  $\bar{\phi} < \bar{\phi}_{\tilde{\mathcal{S}}_n^c}$  such that for all  $\eta > 0$ , w.a.p.1,  $\phi(g, h) < \bar{\phi} + \eta$  and  $\phi_{\tilde{\mathcal{S}}_n^c}(g, h) > \bar{\phi}_{\tilde{\mathcal{S}}_n^c} - \eta$ .

Following the same argument as in the proof of Lemma 14, after maximizing over the direction of  $u$  and minimizing over the direction of  $w$ , Eq. (C32) becomes equivalent to:



$$\begin{aligned}
& \max_{\substack{0 \leq \delta \leq 4\tau^{-1}\sqrt{C_1 C_\varepsilon} \\ \|v\|_\infty \leq 1}} \min_{\alpha \in K_\alpha} - \frac{c_n \delta^2}{4} \mu_n(\alpha, \delta) + \frac{c_n}{n} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon)^\top (\Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I})^{-1} \\
& \quad \times (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon) - c_n \left\| n^{-1/2} \tau^{1/2} \delta h - n^{-1/2} \tau^{-1/2} \lambda_n \Sigma_2^{-1/2} v \right\| \alpha \\
& \quad + \frac{c_n \tau^{-2} \lambda_n}{\sqrt{n}} v^\top \beta_0 - \frac{c_n \tau^{-2}}{n} \|\varepsilon\|^2 - C_n^\phi,
\end{aligned}$$

where  $K_\alpha = \{\alpha | c_n \alpha - c_n \tau^{-1} \sigma_x \sigma_\beta \in [c_\alpha/4\sigma_\beta, C_\alpha/\sigma_\beta]\}$ . By Lemma 17, the objective function of the above optimization problem is convex in  $\alpha$  and jointly concave in  $(\delta, v)$ . Consequently, we can interchange the order of min and max by applying Corollary 3.3 in Sion (1958). Applying  $\|x\| = \min_{\gamma > 0} \frac{1}{2\gamma} \|x\|^2 + \frac{\gamma}{2}$  to  $\left\| n^{-1/2} \tau^{1/2} \delta h^\top - n^{-1/2} \tau^{-1/2} \lambda_n v^\top \Sigma_2^{-1/2} \right\| \alpha$ , and By completing the square for terms associated with  $v$ , we can rewrite this problem as (36). As a consequence, we conclude that (36) is convex with respect to  $\alpha$  and jointly concave with respect to  $(\delta, \gamma)$ .

Finally, based on the above argument, we conclude that under (38), for all  $\eta > 0$ , w.a.p.1,  $\phi(g, h) < \bar{\phi} + \eta$  and  $\phi_{\tilde{\mathcal{S}}_n^c}(g, h) > \bar{\phi}_{\tilde{\mathcal{S}}_n^c} - \eta$  by choosing  $\bar{\phi} = -\frac{C_\lambda}{8C_2}$  and  $\bar{\phi}_{\tilde{\mathcal{S}}_n^c} = -\frac{C_\lambda}{100C_2}$ , thereby verifying conditions outlined in Lemma 23.  $\square$

Next, we introduce a lemma that resembles Lemma 16.

**Lemma 23.** *Let  $\hat{w}$  denote an optimal solution of Eq. (C31). Regarding  $\phi(g, h)$  and  $\phi_{\tilde{\mathcal{S}}_n^c}(g, h)$ , as introduced and discussed in relation to Eq. (C32), suppose there are constants  $\bar{\phi}$  and  $\bar{\phi}_{\tilde{\mathcal{S}}_n^c}$  with  $\bar{\phi} < \bar{\phi}_{\tilde{\mathcal{S}}_n^c}$ , such that for all  $\eta > 0$ , the following hold w.a.p.1 as  $n \rightarrow \infty$ : (a)  $\phi(g, h) < \bar{\phi} + \eta$ , (b)  $\phi_{\tilde{\mathcal{S}}_n^c}(g, h) > \bar{\phi}_{\tilde{\mathcal{S}}_n^c} - \eta$ . Under these conditions, we have  $\hat{w} \in \tilde{\mathcal{S}}_n$  w.p.a.1.*

**Lemma 24.** *There exists some sufficiently small  $\epsilon > 0$ , such that for any  $\eta > 0$ , w.p.a.1, the inequalities in (38) hold.*

*Proof.* By Eq. (C17) in Lemma 18, we have the following result:

$$\begin{aligned}
& - \frac{c_n \delta^2}{4} \mu_n(\alpha, \delta) + \frac{c_n}{n} \left( \tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon \right)^\top (\Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I})^{-1} \left( \tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon \right) \\
& = c_n \tau^{-1} \sigma_x^2 \sigma_\beta^2 - c_n \sigma_\varepsilon^2 \theta_3 \mu^2(\sigma_x \sigma_\beta, \delta_1^*, \delta_2) + 2\sigma_x \sigma_\beta \alpha_2 + \frac{c_n \tau^{-2}}{n} \|\varepsilon\|^2 + o_P(1).
\end{aligned}$$

Additionally, by Lemmas 2-3 and  $p^{1/2} \tau^{-1} n^{-1} q^{-1/2} = o(1)$  by Assumption 4, we deduce that  $-\frac{c_n \gamma}{2} + \frac{c_n \gamma \tau^{-3}}{2\alpha^2} \beta_0 \Sigma_2 \beta_0 + \frac{c_n \tau^{-1} \delta}{\sqrt{n}} h^\top \Sigma_2^{1/2} \beta_0 \xrightarrow{P} -\frac{\gamma_1 \alpha_2}{\sigma_x \sigma_\beta}$ . In the sequel, we examine the asymptotic behavior of the remaining term in  $\tilde{Q}_n(\alpha_2, \delta_3, \gamma_1)$ :

$$\min_{\|v\|_\infty \leq 1} \left\{ \frac{c_n \alpha^2}{2\gamma} \left\| n^{-1/2} \tau^{-1/2} \lambda_n \Sigma_2^{-1/2} v - n^{-1/2} \tau^{1/2} \delta h - \frac{\gamma}{\alpha^2} \tau^{-3/2} \Sigma_2^{1/2} \beta_0 \right\|^2 \right\}. \quad (\text{C33})$$

By using  $\|\Sigma_2^{-1}\| \leq c_2^{-1}$ , we see (C33) is upper bounded by

$$\begin{aligned} & \frac{c_n \alpha^2}{2\gamma c_2} \min_{\|v\|_\infty \leq 1} \left\{ \left\| n^{-1/2} \tau^{-1/2} \lambda_n v - n^{-1/2} \tau^{1/2} \Sigma_2^{1/2} \delta h - \frac{\gamma}{\alpha^2} \tau^{-3/2} \Sigma_2 \beta_0 \right\|^2 \right\} \\ &= \frac{c_n \alpha^2}{2\gamma c_2} \left\| \left( \left| n^{-1/2} \tau^{1/2} \Sigma_2^{1/2} \delta h + \frac{\gamma}{\alpha^2} \tau^{-3/2} \Sigma_2 \beta_0 \right| - n^{-1/2} \tau^{-1/2} \lambda_n \right)_+ \right\|^2. \end{aligned}$$

Similarly, with  $\lambda_{\min} \Sigma_2^{-1} \geq C_2^{-1}$ , (C33) is lower bounded by  $\frac{c_n \alpha^2}{2\gamma C_2} \left\| \left( \left| n^{-1/2} \tau^{1/2} \Sigma_2^{1/2} \delta h + \frac{\gamma}{\alpha^2} \tau^{-3/2} \Sigma_2 \beta_0 \right| - n^{-1/2} \tau^{-1/2} \lambda_n \right)_+ \right\|^2$ . Together with Lemma 25, we deduce that, w.p.a.1, (C33) lies in  $\left[ \frac{\sigma_x^2 \sigma_\beta^2}{4\gamma_1 C_2} C_\lambda, \frac{\sigma_x^2 \sigma_\beta^2}{\gamma_1 c_2} C_\lambda \right]$ .

Recall that  $\tilde{Q}_n(\alpha_2, \delta_3, \gamma_1)$  is defined in (37). We introduce  $\tilde{Q}_n^{\text{upper}}(\alpha_2, \delta_3, \gamma_1)$ , defined as:

$$\begin{aligned} & -\frac{c_n \delta^2}{4} \mu_n(\alpha, \delta) + \frac{c_n}{n} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon)^\top (\Sigma_1^{-1} - \mu_n(\alpha, \delta) \mathbb{I})^{-1} (\tau^{1/2} \alpha g - \tau^{-1} \Sigma_1^{-1/2} \varepsilon) \\ & -\frac{c_n \gamma}{2} - \frac{\sigma_x^2 \sigma_\beta^2}{4\gamma_1 C_2} C_\lambda + \frac{c_n \gamma \tau^{-3}}{2\alpha^2} \beta_0 \Sigma_2 \beta_0 + \frac{c_n \tau^{-1} \delta}{\sqrt{n}} h^\top \Sigma_2^{1/2} \beta_0 - \frac{c_n \tau^{-2}}{n} \|\varepsilon\|^2 - C_n^\phi. \end{aligned}$$

Similarly, we define  $\tilde{Q}_n^{\text{lower}}$  with the term  $-\frac{\sigma_x^2 \sigma_\beta^2}{4\gamma_1 C_2} C_\lambda$  in  $\tilde{Q}_n^{\text{upper}}$  being replaced by  $-\frac{\sigma_x^2 \sigma_\beta^2}{\gamma_1 c_2} C_\lambda$ . Consequently,  $\tilde{Q}_n^{\text{lower}} \leq \tilde{Q}_n \leq \tilde{Q}_n^{\text{upper}}$ . Note also that  $\tilde{Q}_n^{\text{lower}}(\alpha_2, \delta_2, \gamma_1)$  and  $\tilde{Q}_n^{\text{upper}}(\alpha_2, \delta_3, \gamma_1)$  maintain their convexity in  $\alpha_2$  and joint concavity in  $(\delta_3, \gamma_1)$ . By employing a similar line of reasoning as presented in Lemma 18, alongside the definitions of  $c_\alpha$  and  $C_\alpha$ , it becomes evident that there exists a sufficiently small  $\epsilon > 0$  such that

$$\begin{aligned} & \min_{\alpha_2 \in [\frac{c_\alpha}{4\sigma_\beta}, \frac{C_\alpha}{\sigma_\beta}]} \max_{\substack{\gamma_1 > 0 \\ \delta_3 \in K_{\delta_3}}} \tilde{Q}_n^{\text{upper}} \xrightarrow{\text{P}} \min_{\alpha_2 \in [\frac{c_\alpha}{4\sigma_\beta}, \frac{C_\alpha}{\sigma_\beta}]} \max_{\substack{\gamma_1 > 0 \\ \delta_3 \in K_{\delta_3}}} -\frac{\delta_3^2 \theta_1}{4\theta_3} + 2\sigma_x \sigma_\beta \alpha_2 - \frac{\gamma_1 \alpha_2}{\sigma_x \sigma_\beta} - \frac{\sigma_x^2 \sigma_\beta^2}{4\gamma_1 C_2} C_\lambda = -\frac{C_\lambda}{8C_2}, \\ & \text{and} \quad \min_{\alpha_2 \in [\frac{c_\alpha}{4\sigma_\beta}, \frac{c_\alpha}{2\sigma_\beta} + \epsilon] \cup [\frac{C_\alpha}{2\sigma_\beta} - \epsilon, \frac{C_\alpha}{\sigma_\beta}]} \max_{\substack{\gamma_1 > 0 \\ \delta_3 \in K_{\delta_3}}} \tilde{Q}_n^{\text{lower}} \xrightarrow{\text{P}} -\frac{C_\lambda}{100C_2}. \end{aligned}$$

These results immediately yield the desired inequalities.  $\square$

**Lemma 25.** For any  $\alpha_2, \delta_3 \in \mathbb{R}$  and  $\gamma_1 > 0$ , w.p.a.1, we have

$$\frac{C_\lambda}{2} \leq c_n \tau^{-1} \left\| \left( \left| n^{-1/2} \tau^{1/2} \Sigma_2^{1/2} \delta h + \frac{\gamma}{\alpha^2} \tau^{-3/2} \Sigma_2 \beta_0 \right| - n^{-1/2} \tau^{-1/2} \lambda_n \right)_+ \right\|^2 \leq 2C_\lambda. \quad (\text{C34})$$

*Proof.* We first establish the following:

$$c_n n^{-1} \tau^{-2} \left\{ \left\| \left( \left| \Sigma_2^{1/2} \delta_1^* h \right| - \lambda_n \right)_+ \right\|^2 - \mathbb{E} \left\| \left( \left| \Sigma_2^{1/2} \delta_1^* h \right| - \lambda_n \right)_+ \right\|^2 \right\} \xrightarrow{P} 0. \quad (\text{C35})$$

Let  $\tilde{h} := \Sigma_2^{1/2} h \sim \mathcal{N}(0, \Sigma_2)$ . Let us denote the  $(i, j)$ -th element of  $\Sigma_2$  as  $\Sigma_{2,ij}$ , thus we have  $\tilde{h}_j | \tilde{h}_i \stackrel{d}{=} \Sigma_{2,ij} \Sigma_{2,ii}^{-1} \tilde{h}_i + \sqrt{\Sigma_{2,jj} - \Sigma_{2,ii}^{-1} \Sigma_{2,ij}^2} g_1$ , where  $g_1$  is a standard Gaussian random variable independent of  $\tilde{h}_i$ . Consequently,  $\text{Cov} \left( (|\delta_1^* \tilde{h}_i| - \lambda_n)_+, (|\delta_1^* \tilde{h}_j| - \lambda_n)_+ \right)$  equals

$$\begin{aligned} & \mathbb{E} \left\{ (|\delta_1^* \tilde{h}_i| - \lambda_n)_+ \mathbb{E} \left[ (|\delta_1^* \tilde{h}_j| - \lambda_n)_+ - \mathbb{E} (|\delta_1^* \tilde{h}_j| - \lambda_n)_+ \mid \tilde{h}_i \right] \right\} \\ &= \mathbb{E} \left\{ (|\delta_1^* \tilde{h}_i| - \lambda_n)_+ \mathbb{E} \left[ \eta(\tilde{h}_i)^2 - \eta(g_2)^2 \mid \tilde{h}_i \right] \right\}, \end{aligned}$$

where  $\eta(x) := \left( \left| \delta_1^* \left( \Sigma_{2,ij} \Sigma_{2,ii}^{-1} x + \sqrt{\Sigma_{2,jj} - \Sigma_{2,ii}^{-1} \Sigma_{2,ij}^2} g_1 \right) \right| - \lambda_n \right)_+$  and  $g_2 \sim \mathcal{N}(0, \Sigma_{2,ii})$  is independent of both  $g_1$  and  $\tilde{h}_i$ . In addition, note that  $\left| \eta(\tilde{h}_i)^2 - \eta(g_2)^2 \right| \leq |\delta_1^* \Sigma_{2,ij} \Sigma_{2,ii}^{-1} (\tilde{h}_i - g_2)| \cdot \left| \eta(\tilde{h}_i) + \eta(g_2) \right|$ . Applying the Cauchy-Schwarz inequality to the above inequality yields

$$\begin{aligned} \mathbb{E} \left[ \eta(\tilde{h}_i)^2 - \eta(g_2)^2 \mid \tilde{h}_i \right] &\lesssim \left( \mathbb{E} \left( \left| \Sigma_{2,ij} \Sigma_{2,ii}^{-1} (\tilde{h}_i - g_2) \right|^2 \mid \tilde{h}_i \right) \right)^{1/2} \left( \mathbb{E} \left( \eta(g_2)^2 + \eta(\tilde{h}_i)^2 \mid \tilde{h}_i \right) \right)^{1/2} \\ &\lesssim |\Sigma_{2,ij} \Sigma_{2,ii}^{-1}| \sqrt{\Sigma_{2,ii} + \tilde{h}_i^2} \sqrt{\Sigma_{2,ij}^2 \Sigma_{2,ii}^{-2} \tilde{h}_i^2 + \mathbb{E}_{Y \sim \mathcal{N}(0,1)} (|\delta_1^* \Sigma_{2,jj} Y| - \lambda_n)_+^2} \\ &\lesssim |\Sigma_{2,ij}| (1 + |\tilde{h}_i|) \left( |\Sigma_{2,ij}| |\tilde{h}_i| + \sqrt{\mathbb{E}_{Y \sim \mathcal{N}(0,1)} (|\delta_1^* \Sigma_{2,jj} Y| - \lambda_n)_+^2} \right), \end{aligned}$$

where the last step is due to  $c_2 \leq \Sigma_{2,ii} \leq C_2$ . Therefore, by Lemma 11, we have

$$\begin{aligned} & \text{Cov} \left( (|\delta_1^* \tilde{h}_i| - \lambda_n)_+, (|\delta_1^* \tilde{h}_j| - \lambda_n)_+ \right) \\ &\lesssim |\Sigma_{2,ij}| (\lambda_n + 1) \sqrt{\mathbb{E}_{Y \sim \mathcal{N}(0,1)} (|\delta_1^* \Sigma_{2,jj} Y| - \lambda_n)_+^2} \mathbb{E} (|\delta_1^* \tilde{h}_i| - \lambda_n)_+ + \Sigma_{2,ij}^2 (\lambda_n^2 + \lambda_n) \mathbb{E} (|\delta_1^* \tilde{h}_i| - \lambda_n)_+^2. \end{aligned}$$

Further, by Lemma 10 and Eq. (9), we have  $\lambda_n = o(\log(p))$ . The above inequality leads to:

$$\begin{aligned} & \text{Var} \left( c_n n^{-1} \tau^{-2} \left\| \left( \left| \Sigma_2^{1/2} \delta_1^* h \right| - \lambda_n \right)_+ \right\|^2 \right) = \sum_{i,j=1}^p c_n^2 n^{-2} \tau^{-4} \text{Cov} \left( (|\delta_1^* \tilde{h}_i| - \lambda_n)_+, (|\delta_1^* \tilde{h}_j| - \lambda_n)_+ \right) \\ &\lesssim c_n^2 n^{-2} \tau^{-4} \log(p) \sum_{i=1}^p \mathbb{E} (|\delta_1^* \tilde{h}_i| - \lambda_n)_+^2 \left( \sum_{j=1}^p |\Sigma_{2,ij}| \sqrt{\mathbb{E}_{Y \sim \mathcal{N}(0,1)} (|\delta_1^* \Sigma_{2,jj} Y| - \lambda_n)_+^2} \right) \\ &+ c_n^2 n^{-2} \tau^{-4} (\log(p))^2 \sum_{i=1}^p \mathbb{E} (|\delta_1^* \tilde{h}_i| - \lambda_n)_+^2 \sum_{j=1}^p \Sigma_{2,ij}^2 \end{aligned}$$

$$\begin{aligned} &\leq c_n^2 n^{-2} \tau^{-4} \left\{ \log(p) C_2 \sum_{i=1}^p \mathbb{E}(|\delta_1^* \tilde{h}_i| - \lambda_n)_+^2 \left( \sum_{j=1}^p \mathbb{E}_{Y \sim \mathcal{N}(0,1)}(|\delta_1^* \Sigma_{2,jj} Y| - \lambda_n)_+^2 \right)^{1/2} \right. \\ &\quad \left. + (\log(p))^2 C_2^2 \sum_{i=1}^p \mathbb{E}(|\delta_1^* \tilde{h}_i| - \lambda_n)_+^2 \right\} = O(c_n^{1/2} n^{-1/2} \tau^{-1} \log(p) + c_n n^{-1} \tau^{-2} \log^2(p)) = o_n(1), \end{aligned}$$

where we use  $\sum_{j=1}^p \Sigma_{2,jj}^2 \leq C_2^2$  and Cauchy–Schwartz inequality in the second step. This leads to Eq. (C35). Using the same approach, we can prove

$$\begin{aligned} &\left\| \left( \left| \Sigma_2^{1/2} \delta_1^* h \right| - \lambda_n + \log^{-1}(p) \right)_+ \right\|^2 - \mathbb{E} \left\| \left( \left| \Sigma_2^{1/2} \delta_1^* h \right| - \lambda_n + \log^{-1}(p) \right)_+ \right\|^2 = o_{\mathbb{P}}(c_n^{-1} n \tau^2), \\ &\left\| \left( \left| \Sigma_2^{1/2} \delta_1^* h \right| - \lambda_n - \log^{-1}(p) \right)_+ \right\|^2 - \mathbb{E} \left\| \left( \left| \Sigma_2^{1/2} \delta_1^* h \right| - \lambda_n - \log^{-1}(p) \right)_+ \right\|^2 = o_{\mathbb{P}}(c_n^{-1} n \tau^2). \end{aligned}$$

Now we are ready to establish Eq. (C34). Note that w.p.a.1, we have

$$\begin{aligned} &c_n \tau^{-1} \left\| \left( \left| n^{-1/2} \tau^{1/2} \Sigma_2^{1/2} \delta h + \frac{\gamma}{\alpha^2} \tau^{-3/2} \Sigma_2 \beta_0 \right| - n^{-1/2} \tau^{-1/2} \lambda_n \right)_+ \right\|^2 \\ &\leq 2c_n n^{-1} \tau^{-2} \left\| \left( \left| \Sigma_2^{1/2} \delta_1^* h \right| - \lambda_n + \log^{-1}(p) \right)_+ \right\|^2, \end{aligned}$$

where the inequality is given by Lemma 12 and the facts that  $\tau \log^4(p) = o(1)$  and  $n^{1/2} \tau^{1/2} p^{-1/2} q^{-1/2} \log^2(p) = o(1)$  by Assumption 4. Similarly, the left-hand-side is lower bounded by  $\frac{1}{2} c_n n^{-1} \tau^{-2} \mathbb{E} \left\| \left( \left| \Sigma_2^{1/2} \delta_1^* h \right| - \lambda_n - \log^{-1}(p) \right)_+ \right\|^2$  w.p.a.1. Finally, by Lemma 10 and the fact that  $\lambda_n = o(\log(p))$ , it is not hard to verify that  $\mathbb{E} \left\| \left( \left| \Sigma_2^{1/2} \delta_1^* h \right| - \lambda_n - \log^{-1}(p) \right)_+ \right\|^2 = \mathbb{E} \left\| \left( \left| \Sigma_2^{1/2} \delta_1^* h \right| - \lambda_n \right)_+ \right\|^2 (1 + o(1))$  and  $\mathbb{E} \left\| \left( \left| \Sigma_2^{1/2} \delta_1^* h \right| - \lambda_n + \log^{-1}(p) \right)_+ \right\|^2 = \mathbb{E} \left\| \left( \left| \Sigma_2^{1/2} \delta_1^* h \right| - \lambda_n \right)_+ \right\|^2 (1 + o(1))$ . Together with the definition that  $C_\lambda = \lim_{n \rightarrow \infty} p n^{-2} \tau^{-2} \mathbb{E} \left\| \left( \left| \Sigma_2^{1/2} \delta_1^* h \right| - \lambda_n \right)_+ \right\|^2$  given by Eq. (9), we conclude the proof.  $\square$

**Lemma 26.** *Eq. (43) defined in the proof of Theorem 5 holds as  $n \rightarrow \infty$ .*

*Proof.* Note that  $\|\hat{\beta} - \tilde{\beta}\|^2$  is upper bounded by

$$\begin{aligned} &\frac{2}{n^2} \|R_U U^\top \mathcal{M}_W (U \beta_0 + \epsilon) - R_U U^\top (U \beta_0 + \epsilon)\|^2 + \frac{2}{n^2} \|(R_U - \tilde{R}_U) U^\top (U \beta_0 + \epsilon)\|^2 \\ &\leq \frac{4}{n^2} \|R_U U^\top \mathcal{M}_W \epsilon - R_U U^\top \epsilon\|^2 + \frac{4}{n^2} \|R_U U^\top \mathcal{M}_W U \beta_0 - R_U U^\top U \beta_0\|^2 \\ &\quad + \frac{4}{n^2} \|(R_U - \tilde{R}_U) U^\top \epsilon\|^2 + \frac{4}{n^2} \|(R_U - \tilde{R}_U) U^\top U \beta_0\|^2. \end{aligned}$$

For the first term, using  $\|R_U\| \lesssim_P np^{-1}\tau$ ,  $\|U^\top U\| \lesssim_P p$ , and Lemma 2, we have

$$\begin{aligned} & \frac{4}{n^2} \|R_U U^\top \mathcal{M}_W \varepsilon - R_U U^\top \varepsilon\|^2 \asymp_P \frac{1}{n^2} \text{Tr}((\mathcal{M}_W U - U) R_U^2 (U^\top \mathcal{M}_W - U^\top)) \\ &= \frac{1}{n^2} \text{Tr}(U^\top W (W^\top W)^{-1} W^\top U R_U^2) \leq \frac{1}{n^2} \|R_U^2\| \|U^\top U\| \text{Tr}(W (W^\top W)^{-1} W^\top) \lesssim_P p^{-1} \tau^2 \text{rank}(W). \end{aligned}$$

Similarly, it can be shown that the second term is of order  $O_P(p^{-1} \tau^2 \text{rank}(W))$ . In addition, by Lemma 2, and using the fact that  $\text{Tr}(AB) \leq \|A\| \text{Tr}(B)$  and  $\|\tilde{R}_U\| \lesssim_P np^{-1}\tau$ , we have

$$\begin{aligned} & \frac{4}{n^2} \|(R_U - \tilde{R}_U) U^\top \varepsilon\|^2 \asymp_P \frac{1}{n^2} \text{Tr}(U (R_U - \tilde{R}_U)^2 U^\top) \leq \frac{1}{n^2} \|U^\top U\| \text{Tr}((R_U - \tilde{R}_U)^2) \\ & \lesssim_P \frac{p}{n^2} \text{Tr}((R_U (\tilde{R}_U^{-1} - R_U^{-1}) \tilde{R}_U)^2) = \frac{p}{n^4} \text{Tr}((R_U U^\top W (W^\top W)^{-1} W^\top U \tilde{R}_U)^2) \\ & \leq \frac{p}{n^4} \|R_U\|^2 \|\tilde{R}_U\|^2 \|U^\top U\|^2 \text{Tr}(W (W^\top W)^{-1} W^\top) \lesssim_P p^{-1} \tau^4 \text{rank}(W). \end{aligned}$$

Similarly, the final term is of order  $O_P(p^{-1} \tau^4 \text{rank}(W))$ . To sum up, we have  $\|\hat{\beta} - \tilde{\beta}\|^2 = O(p^{-1} \tau^2 \text{rank}(W)) = o(n^2 p^{-2} \tau^3)$ , since  $\text{rank}(W) = o(n^2 p^{-1} \tau)$ .  $\square$

## References

- Bai, Z. and J. Silverstein (2009). *Spectral Analysis of Large Dimensional Random Matrices*. Springer Series in Statistics. Springer New York.
- Chen, B. and G. Pan (2012). Convergence of the largest eigenvalue of normalized sample covariance matrices when  $p$  and  $n$  both tend to infinity with their ratio converging to zero. *Bernoulli* 18(4), 1405 – 1420.
- Gander, W., G. H. Golub, and U. von Matt (1989). A constrained eigenvalue problem. *Linear Algebra and its Applications* 114-115, 815–839.
- Giannone, D., M. Lenza, and G. E. Primiceri (2022). Economic predictions with big data: The illusion of sparsity. *Econometrica* 89(5), 2409–2437.
- Götze, F., H. Sambale, and A. Sinulis (2021). Concentration inequalities for polynomials in  $\alpha$ -sub-exponential random variables. *Electronic Journal of Probability* 26(none), 1 – 22.
- LeJeune, D., H. Javadi, and R. Baraniuk (2020). The implicit regularization of ordinary least squares ensembles. In *International Conference on Artificial Intelligence and Statistics*, pp. 3525–3535. PMLR.

- Liese, F. and K. Miescke (2008). *Statistical Decision Theory: Estimation, Testing, and Selection*. Springer Series in Statistics. Springer New York.
- Meyer, C. D. (2000). *Matrix analysis and applied linear algebra*. USA: Society for Industrial and Applied Mathematics.
- Sion, M. (1958). On general minimax theorems. *Pacific Journal of Mathematics* 8(1), 171 – 176.
- Thrampoulidis, C., E. Abbasi, and B. Hassibi (2018). Precise error analysis of regularized  $m$ -estimators in high dimensions. *IEEE Transactions on Information Theory* 64(8), 5592–5628.
- Thrampoulidis, C., S. Oymak, and B. Hassibi (2015). Regularized linear regression: A precise analysis of the estimation error. In *Proceedings of The 28th Conference on Learning Theory*, Volume 40 of *Proceedings of Machine Learning Research*, Paris, France, pp. 1683–1709. PMLR.