

---

# On the Optimal Weighted $\ell_2$ Regularization in Overparameterized Linear Regression

---

Denny Wu\*

University of Toronto and Vector Institute  
dennywu@cs.toronto.edu

Ji Xu\*

Columbia University  
jixu@cs.columbia.edu

## Abstract

We consider the linear model  $\mathbf{y} = \mathbf{X}\beta_* + \epsilon$  with  $\mathbf{X} \in \mathbb{R}^{n \times p}$  in the overparameterized regime  $p > n$ . We estimate  $\beta_*$  via generalized (weighted) ridge regression:  $\hat{\beta}_\lambda = (\mathbf{X}^\top \mathbf{X} + \lambda \Sigma_w)^\dagger \mathbf{X}^\top \mathbf{y}$ , where  $\Sigma_w$  is the weighting matrix. Under a random design setting with general data covariance  $\Sigma_x$  and general prior on the true coefficients  $\mathbb{E}\beta_*\beta_*^\top = \Sigma_\beta$ , we provide an exact characterization of the prediction risk  $\mathbb{E}(y - \mathbf{x}^\top \hat{\beta}_\lambda)^2$  in the proportional asymptotic limit  $p/n \rightarrow \gamma \in (1, \infty)$ . Our general setup leads to a number of interesting findings. We outline precise conditions that decide the sign of the optimal choice  $\lambda_{\text{opt}}$  of the ridge parameter  $\lambda$ , based on the *alignment* between  $\Sigma_x$  and  $\Sigma_\beta$ ; this rigorously justifies the surprising empirical observation that  $\lambda_{\text{opt}}$  can be *negative* in the overparameterized regime. We also discuss the risk monotonicity of optimally tuned ridge regression, and confirm the double descent phenomenon for principal component regression (PCR) under anisotropic  $\mathbf{X}$  and  $\beta_*$ . Finally, we determine the optimal  $\Sigma_w$  for both the ridgeless ( $\lambda \rightarrow 0$ ) and optimally regularized ( $\lambda = \lambda_{\text{opt}}$ ) case, and demonstrate the advantage of the weighted objective over standard ridge regression and PCR.

## 1 Introduction

In this work we consider learning the target signal  $\beta_*$  in the following linear regression model:

$$y_i = \mathbf{x}_i^\top \beta_* + \epsilon_i, \quad i = 1, 2, \dots, n$$

where each feature vector  $\mathbf{x}_i \in \mathbb{R}^p$  and noise  $\epsilon_i \in \mathbb{R}$  are drawn i.i.d. from the two independent random variables  $\tilde{\mathbf{x}}$  and  $\tilde{\epsilon}$  satisfying  $\mathbb{E}\tilde{\epsilon} = 0$ ,  $\mathbb{E}\tilde{\epsilon}^2 = \tilde{\sigma}^2$ ,  $\tilde{\mathbf{x}} = \Sigma_x^{1/2} \mathbf{z} / \sqrt{n}$ , and the components of  $\mathbf{z}$  are i.i.d. random variables with zero mean, unit variance, and bounded 12th absolute central moment. To estimate  $\beta_*$  from  $(\mathbf{x}_i, y_i)$ , we consider the following generalized ridge regression estimator:

$$\hat{\beta}_\lambda = (\mathbf{X}^\top \mathbf{X} + \lambda \Sigma_w)^\dagger \mathbf{X}^\top \mathbf{y}, \quad (1.1)$$

in which  $\mathbf{X} \in \mathbb{R}^{n \times p}$  is the feature matrix,  $\mathbf{y}$  is vector of the observations,  $\Sigma_w$  is a positive definite weighting matrix, and the symbol  $\dagger$  denotes the Moore-Penrose pseudo-inverse. When  $\lambda \geq 0$ ,  $\hat{\beta}_\lambda$  minimizes the squared loss plus a weighted  $\ell_2$  regularization:  $\min_{\beta} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \beta)^2 + \lambda \beta^\top \Sigma_w \beta$ . Note that  $\Sigma_w = \mathbf{I}_d$  reduces the objective to standard ridge regression.

While the standard ridge regression estimator is relatively well-understood in the data-abundant regime ( $n > p$ ), several interesting properties have been recently discovered in high dimensions, especially when  $p > n$ . For instance, the double descent phenomenon suggests that overparameterization may not result in overfitting due to the *implicit regularization* of the least squares estimator [HMRT19, BLLT19]. This implicit regularization also relates to the surprising empirical finding that the optimal ridge parameter  $\lambda$  can be negative in the overparameterized regime [KLS20].

---

\*Equal contribution; alphabetical ordering.

Motivated by the observations above, we analyze the estimator  $\hat{\beta}_\lambda$  in the proportional limit:  $p/n \rightarrow \gamma \in (1, \infty)^2$  as  $n, p \rightarrow \infty$ . We place a general prior on the true parameters (independent of  $\tilde{\mathbf{x}}$  and  $\tilde{\epsilon}$ ):  $\mathbb{E}\beta_*\beta_*^\top = \Sigma_\beta$ , which covers both *random* and *deterministic*  $\beta_*$ . Our goal is to study the prediction risk of  $\hat{\beta}_\lambda$ :  $\mathbb{E}_{\tilde{\mathbf{x}}, \tilde{\epsilon}, \beta_*} (\tilde{y} - \tilde{\mathbf{x}}^\top \hat{\beta}_\lambda)^2$ , where  $\tilde{y} = \tilde{\mathbf{x}}^\top \beta_* + \tilde{\epsilon}$ <sup>3</sup>. Compared to previous high-dimensional analysis of ridge regression [DW18], our setup is generalized in two important aspects:

**Anisotropic  $\Sigma_x$  and  $\Sigma_\beta$ .** Our analysis handles general prior  $\Sigma_\beta$  and data covariance  $\Sigma_x$ , in contrast to previous works which assume either isotropic features or signal (e.g., [DW18, HMRT19]). Note that the isotropic assumption on the signal or features implies that each component is roughly of the same magnitude, which may not hold true in practice. For instance, the optimal ridge penalty is *provably* non-negative when either the signal  $\Sigma_\beta$  [DW18, Theorem 2.1] or the features  $\Sigma_x$  [HMRT19, Theorem 5] is isotropic. On the other hand, it has been empirically demonstrated that the optimal ridge for real-world data can be negative [KLS20]. While this observation cannot be captured by previous works, our less restrictive assumptions lead to a concise description of this phenomenon.

**Weighted  $\ell_2$  Regularization.** We consider generalized ridge regression instead of simple isotropic shrinkage. While the generalized formulation has also been studied (e.g., [HK70, Cas80]), to the best of our knowledge, no existing work computes the exact risk in the overparameterized proportional limit and decides the corresponding optimal  $\Sigma_w$ . Our setting is also inspired by recent observations in deep learning that weighted  $\ell_2$  regularization often achieves better generalization compare to isotropic weight decay [ZWXG18]. Our analysis illustrates the benefit of weighted  $\ell_2$  regularization.

Under the general setup (1.1), the contributions of this work can be summarized as (see Figure 1):

- **Exact Asymptotic Risk.** In Section 4 we derive the prediction risk  $R(\lambda)$  of our estimator (1.1) in its bias-variance decomposition (see Figure 2). We also characterize principal component regression (PCR) and confirm the double descent phenomenon under more general setting than [XH19].
- **“Negative Ridge” Phenomenon.** In Section 5, we analyze the optimal regularization strength  $\lambda_{\text{opt}}$  under different  $\Sigma_w$ , and provide precise conditions under which the optimal  $\lambda_{\text{opt}}$  is negative in the overparameterized regime. In brief, we show that  $\lambda_{\text{opt}}$  is negative when SNR is large and the large directions of  $\Sigma_x$  and  $\beta_*$  are *aligned* (see Figure 4), and vice versa. On the other hand, we show that the optimal ridge penalty is always non-negative in the underparameterized regime ( $p < n$ ); this implies an implicit  $\ell_2$  regularization effect of overparameterization for certain cases.
- **Optimal Weighting Matrix  $\Sigma_w$ .** In Section 6, we decide the optimal  $\Sigma_w$  for both the optimally regularized ridge estimator ( $\lambda = \lambda_{\text{opt}}$ ) and the ridgeless limit ( $\lambda \rightarrow 0$ ). In the ridgeless limit, based on the bias-variance decomposition, we show that the optimal  $\Sigma_w$  should interpolate between  $\Sigma_x$ , which minimizes the variance, and  $\Sigma_\beta^{-1}$ , which minimizes the bias. Whereas for the optimally regularized case, in many settings the optimal  $\Sigma_w$  is simply  $\Sigma_\beta^{-1}$  (Figure 6) (for more general result see Theorem 10). We demonstrate the benefit of weighted regularization over standard ridge regression and PCR, and also propose a heuristic choice of  $\Sigma_w$  when information of  $\beta_*$  is not present.

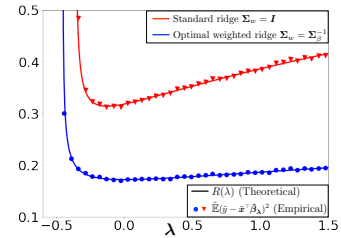


Figure 1: Illustration of the “negative ridge” phenomenon and the advantage of weighted  $\ell_2$  regularization under “aligned”  $\Sigma_x$  and  $\Sigma_\beta$ . We set  $\gamma = 2$ ,  $\tilde{\sigma}^2 = 0$ . Red: standard ridge regression ( $\Sigma_w = \mathbf{I}$ ); note that the lowest prediction risk is achieved when  $\lambda < 0$ . Blue: optimally weighted ridge regression ( $\Sigma_w = \Sigma_\beta^{-1}$ ), which achieves lower risk compared to the isotropic shrinkage.

**Notations:** We denote  $\tilde{\mathbb{E}}$  as taking expectation over  $\beta_*, \tilde{\mathbf{x}}, \tilde{\epsilon}$ . Let  $\mathbf{d}_x, \mathbf{d}_\beta, \mathbf{d}_w$  be the vectors of the eigenvalues of  $\Sigma_x, \Sigma_\beta$  and  $\Sigma_w$  respectively. We use  $\mathbb{I}_S$  as the indicator function of set  $S$ . We write  $\xi = \mathbb{E}(\tilde{\mathbf{x}}^\top \beta_*)^2 / (\gamma \tilde{\sigma}^2)$  as the signal-to-noise ratio (SNR) of the problem.

## 2 Related Works

**Asymptotics of Ridge Regression.** The prediction risk of standard ridge regression ( $\Sigma_w = \mathbf{I}_d$ ) in the proportional asymptotics has been widely studied. When the data is isotropic, precise characterization can be obtained from random matrix theory [Kar13, Dic16, HMRT19], approximate message

<sup>2</sup>Some of our results also apply to the underparameterized case, as we explicitly highlight in the sequel.

<sup>3</sup>When  $\beta_*$  is deterministic,  $\mathbb{E}_{\tilde{\mathbf{x}}, \tilde{\epsilon}, \beta_*} (\tilde{y} - \tilde{\mathbf{x}}^\top \hat{\beta}_\lambda)^2$  reduces to the prediction risk for one fixed  $\beta_*$ .

passing algorithm [DM16], or the convex Gaussian min-max theorem<sup>4</sup>[TAH18]. Under general data covariance, closely related to our work is [DW18], which considered a random effects model with isotropic prior on the target coefficients ( $\Sigma_\beta = \mathbf{I}_d$ ). Our risk characterization is built upon the general random matrix result of [RM11, LP11]. Similar tools have been applied in the analysis of sketching [LD19] and the connection between ridge regression and early stopping [AKT19, Lol20].

**Weighted Regularization.** The formulation (1.1) was first introduced in [HK70], and many choices of weighting matrix have been proposed [Str78, Cas80, MS05, MS18]; yet since these estimators are usually derived in the  $n > p$  setup, their effectiveness in the high-dimensional and overparameterized regime is largely unknown. In semi-supervised linear regression, it is known that weighted matrix estimated from unlabeled data can improve the model performance [RC15, TCG20]. In deep learning, anisotropic Gaussian prior on the parameters enjoyed empirical success [LW17, ZTSG19]. Additionally, decoupled weight decay [LH17] and elastic weight consolidation [KPR<sup>+</sup>17] can both be interpreted as an  $\ell_2$  regularization weighted by an approximate Fisher information matrix [ZWXG18], which relates to the Fisher-Rao norm [LPRS17]. Finally, beyond the  $\ell_2$  penalty, weighted regularization is also effective in LASSO regression [Zou06, CWB08, BVDBS<sup>+</sup>15].

**Benefit of Overparameterization.** Our overparameterized setting is partially motivated by the double descent phenomenon [BHMM18], which can be theoretically explained in linear regression [AS17, HMRT19, BLLT19], random features regression [MM19, dRBK20, DL20], and max-margin classification [MRSY19, DKT19, HMX20], although translation to neural networks can be more nuanced [BES<sup>+</sup>20]. For least squares regression, it has been shown in special cases that overparameterization induces an implicit  $\ell_2$  regularization [KLS20, DLM19], which agrees with the absence of overfitting. This observation also leads to the speculation that the optimal ridge penalty in the overparameterized regime may be negative, to partially cancel out the implicit regularization. While the possibility of negative ridge parameter has been noted in [HG83, BS99], theoretical understanding of its benefit is largely missing, expect for heuristic argument (and empirical evidence) in [KLS20]. We provide a rigorous characterization of this “negative ridge” phenomenon.

**Concurrent Works.** Independent to our work, [RMR20] computed the asymptotic prediction risk under a similar extension of the isotropic assumption on  $\Sigma_\beta$ , but did not consider the sign of  $\lambda_{\text{opt}}$ . We note that their result requires codiagonalizability of the covariances and certain functional relations between the eigenvalues, which is more restrictive than our setting. [TB20] provided a non-asymptotic analysis of ridge regression and constructed a specific spike model<sup>5</sup> in which negative regularization may lead to better generalization bound than interpolation ( $\lambda = 0$ ). In a companion work [ABG<sup>+</sup>20], we connect properties of the ridgeless limit of the generalized ridge regression estimator to the implicit bias of preconditioned gradient descent (e.g., natural gradient descent), which allows us to decide the optimal preconditioner (for generalization) in the interpolation setting.

### 3 Setup and Assumptions

In addition to the prediction risk of the weighted ridge estimator  $\hat{\beta}_\lambda = (\mathbf{X}^\top \mathbf{X} + \lambda \Sigma_w)^\dagger \mathbf{X}^\top \mathbf{y}$ , the setup of which we outlined in Section 1, we also analyze the principal component regression (PCR) estimator: for  $\theta \in [0, 1]$ , the PCR estimator is given as  $\hat{\beta}_\theta = (\mathbf{X}_\theta^\top \mathbf{X}_\theta)^\dagger \mathbf{X}_\theta^\top \mathbf{y}$ , where  $\mathbf{X}_\theta = \mathbf{X} \mathbf{U}_\theta$  and the columns of  $\mathbf{U}_\theta \in \mathbb{R}^{p \times \theta p}$  are the leading  $\theta p$  eigenvectors of  $\Sigma_x$ .

Under the setup on  $(\tilde{x}, \beta_*, \tilde{\epsilon})$  described in Section 1, the prediction risk of (1.1) can be simplified as

$$\begin{aligned} \mathbb{E} \left( \tilde{y} - \tilde{x}^\top \hat{\beta}_\lambda \right)^2 &= \underbrace{\tilde{\sigma}^2 \left( 1 + \frac{1}{n} \text{tr} \left( \Sigma_{x/w} \left( \mathbf{X}_{/w}^\top \mathbf{X}_{/w} + \lambda \mathbf{I} \right)^{-1} - \lambda \Sigma_{x/w} \left( \mathbf{X}_{/w}^\top \mathbf{X}_{/w} + \lambda \mathbf{I} \right)^{-2} \right) \right)}_{\text{Part 1, Variance}} \\ &\quad + \underbrace{\frac{\lambda^2}{n} \text{tr} \left( \Sigma_{x/w} \left( \mathbf{X}_{/w}^\top \mathbf{X}_{/w} + \lambda \mathbf{I} \right)^{-1} \Sigma_{w\beta} \left( \mathbf{X}_{/w}^\top \mathbf{X}_{/w} + \lambda \mathbf{I} \right)^{-1} \right)}_{\text{Part 2, Bias}}, \end{aligned} \quad (3.1)$$

where  $\mathbf{X}_{/w} = \mathbf{X} \Sigma_w^{-1/2}$ ,  $\Sigma_{x/w} = \Sigma_w^{-1/2} \Sigma_x \Sigma_w^{-1/2}$ ,  $\Sigma_{w\beta} = \Sigma_w^{1/2} \Sigma_\beta \Sigma_w^{1/2}$ . Note that the variance term does not depend on the true signal, and the bias is independent of the noise level. Let  $\mathbf{d}_{x/w}$  be

<sup>4</sup>We remark that convergence and uniqueness of AMP and CGMT can be challenging to establish for  $\lambda < 0$ . Also, to our knowledge the current AMP framework does not handle the *joint* relation between  $\Sigma_x$  and  $\Sigma_\beta$ .

<sup>5</sup>Our negative ridge construction relies on the general notion of *alignment*, which subsumes the spike model.

the eigenvalues of  $\Sigma_{x/w}$  and  $\Sigma_{x/w} = U_{x/w} D_{x/w} U_{x/w}^\top$  be the eigendecomposition of  $\Sigma_{x/w}$ , where  $U_{x/w}$  is the eigenvector matrix and  $D_{x/w} = \text{diag}(\mathbf{d}_{x/w})$ . Let  $\mathbf{d}_{w\beta} \triangleq \text{diag}(U_{x/w}^\top \Sigma_{w\beta} U_{x/w})$ . When  $\Sigma_w = \mathbf{I}$ ,  $\mathbf{d}_{w\beta}$  characterizes the strength of the signal  $\beta_*$  along the directions of the eigenvectors of feature covariance  $\Sigma_x$ . To simplify the RHS of (3.1), we make the following assumption:

**Assumption 1.** Let  $d_{x/w,i}$  and  $d_{w\beta,i}$  be the  $i$ th entry of  $\mathbf{d}_{x/w}$  and  $\mathbf{d}_{w\beta}$ , respectively. The empirical distribution of  $(d_{x/w,i}, d_{w\beta,i})$  jointly converges to non-negative random variables  $(h, g)$ . Further,  $\min_i d_{x/w,i} \geq c_l$ ,  $\max_i (d_{x/w,i}, d_{w\beta,i}) \leq c_u$  and  $\|\Sigma_{w\beta}\| \leq c_u$  for some  $c_l, c_u > 0$  independent of  $p$ .

One can check that  $\Sigma_x$  and  $\Sigma_\beta$  studied in [DW18, HMRT19, XH19] (with  $\Sigma_w = \mathbf{I}$ ) are special cases of Assumption 1 with either  $h$  or  $g$  being a point mass. Our Assumption 1 thus covers much more general settings of  $\Sigma_x$  and  $\Sigma_\beta$ , which allows us to precisely analyze the *negative ridge* phenomenon.

## 4 Risk Characterization

With the aforementioned assumptions, we now present our characterization of the prediction risk.

**Theorem 1.** Under Assumption 1, the asymptotic prediction risk is given as

$$\mathbb{E}(\tilde{y} - \tilde{\mathbf{x}}^\top \hat{\beta}_\lambda)^2 \xrightarrow{p} \frac{m'(-\lambda)}{m^2(-\lambda)} \cdot \left( \gamma \mathbb{E} \frac{gh}{(h \cdot m(-\lambda) + 1)^2} + \tilde{\sigma}^2 \right) := R(\lambda), \quad \forall \lambda > -c_0 \quad (4.1)$$

where  $c_0 = (\sqrt{\gamma} - 1)^2 c_l$ , and  $m(z)$  is the Stieltjes transform of the limiting distribution of the eigenvalues of  $\mathbf{X}_{/w} \mathbf{X}_{/w}^\top$ . Additionally,  $m(-\lambda), m'(-\lambda) > 0$  satisfy the self-consistent equations:

$$\lambda = \frac{1}{m(-\lambda)} - \gamma \mathbb{E} \frac{h}{1 + h \cdot m(-\lambda)} \quad (4.2)$$

$$1 = \left( \frac{1}{m^2(-\lambda)} - \gamma \mathbb{E} \frac{h^2}{(h \cdot m(-\lambda) + 1)^2} \right) m'(-\lambda). \quad (4.3)$$

Note that the condition  $\lambda > -c_0$  ensures both  $m(-\lambda)$  and  $m'(-\lambda)$  exist and are positive. Furthermore, it can be shown from prior works [DW18, XH19] that the variance term (part 1) in (3.1), converges to  $\tilde{\sigma}^2 \frac{m'(-\lambda)}{m^2(-\lambda)}$ . Our main contribution is to characterize the bias term, Part 2, under significantly less restrictive assumption on  $(\Sigma_x, \Sigma_\beta, \Sigma_w)$ . In particular, we show that

$$\text{Part 2} \xrightarrow{p} \frac{m'(-\lambda)}{m^2(-\lambda)} \cdot \gamma \mathbb{E} \frac{gh}{(h \cdot m(-\lambda) + 1)^2}, \quad \forall \lambda > -c_0.$$

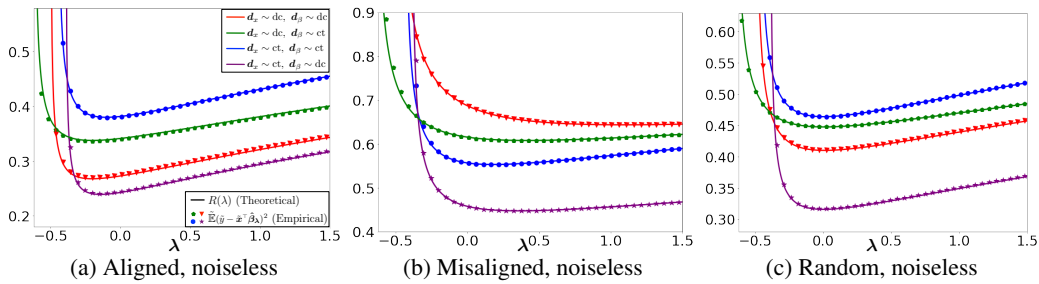


Figure 2: Finite sample prediction risk  $\mathbb{E}(\tilde{y} - \tilde{\mathbf{x}}^\top \hat{\beta}_\lambda)^2$  (experiment) and the asymptotic risk  $R(\lambda)$  (theory) against  $\lambda$  for standard ridge regression ( $\Sigma_w = \mathbf{I}_d$ ). We set  $\gamma = 2$  and  $(n, p) = (300, 600)$ . ‘dc’ and ‘ct’ stand for discrete and continuous distribution, respectively. We write ‘aligned’ if  $\mathbf{d}_x$  and  $\mathbf{d}_\beta$  have the same order, ‘misaligned’ for the reverse, and ‘random’ for random order. Colors indicate different combinations of  $\mathbf{d}_x$  and  $\mathbf{d}_\beta$ . Note that our derived risk  $R(\lambda)$  matches the experimental values, and in the aligned and noiseless case, the optimal risk is achieved when  $\lambda < 0$  (predicted by Theorem 4). The noisy case is presented in Appendix D.

We illustrate the results of Theorem 1 in Figure 2 (noiseless case) and Figure 8 (noisy case) for both discrete and continuous design for  $\mathbf{d}_x$  and  $\mathbf{d}_\beta$  with  $\Sigma_x = \text{diag}(\mathbf{d}_x)$ ,  $\Sigma_\beta = \text{diag}(\mathbf{d}_\beta)$  and  $\Sigma_w = \mathbf{I}$  (see design details in Appendix D). Note that Assumption 1 specifies a joint relation between  $\mathbf{d}_x (= \mathbf{d}_{x/w})$  and  $\mathbf{d}_\beta (= \mathbf{d}_{w\beta})$ . In the following section, we mainly consider the three following relations, which allow us to precisely determine the sign of  $\lambda_{\text{opt}}$ .

**Definition 2.** For two vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^p$ , we say  $\mathbf{a}$  is aligned (misaligned) with  $\mathbf{b}$  if the order of  $\mathbf{a}$  is the same as (reverse of) the order of  $\mathbf{b}$ , i.e.,  $a_i \geq a_j$  iff  $b_i \geq (\leq) b_j$  for all  $i, j$ . Additionally, we say  $\mathbf{a}$  and  $\mathbf{b}$  have random relation if given the order of one vector, the order of the other is uniformly permuted at random.

Intuitively, aligned  $\mathbf{d}_x$  and  $\mathbf{d}_\beta$  implies that when one component in  $\mathbf{d}_x$  has large magnitude, then so does the corresponding component in  $\mathbf{d}_\beta$ , and vice versa (see Figure 3). In Figure 2, we plot the prediction risk of all three joint relations defined above (see Appendix D for details).

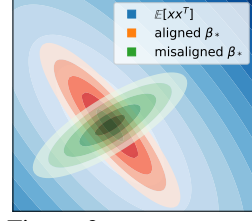


Figure 3: Alignment between  $\mathbf{x}$  and  $\beta_*$  in 2D.

Theorem 1 allows us to compute the risk of the generalized ridge estimator  $\hat{\beta}_\lambda$  and also its ridgeless limit, i.e., the minimum  $\|\hat{\beta}\|_{\Sigma_w}$  norm interpolant (taking  $\Sigma_w = \mathbf{I}$  yields the min  $\ell_2$  norm solution).

**Connection to PCR estimator.** Note that the principal component regression (PCR) estimator is closely related to the ridgeless estimator in the following sense: picking the leading  $\theta p$  eigenvectors of  $\Sigma_x$  (for some  $\theta \in [0, 1]$ ) is equivalent to setting the remaining  $(1 - \theta)p$  eigenvalues of  $\Sigma_w$  to be infinity [HG83]. The following corollary characterizes the prediction risk of the PCR estimator  $\hat{\beta}_\theta$ :

**Corollary 3.** Given Assumption 1 and  $\Sigma_w = \mathbf{I}$ , and  $h$  has continuous and strictly increasing quantile function  $Q_h$ . Then for all  $\theta \in (0, 1]$ , as  $n, p \rightarrow \infty$ ,

$$\tilde{\mathbb{E}}\left(\tilde{y} - \tilde{x}\hat{\beta}_\theta\right)^2 \xrightarrow{p} \begin{cases} \frac{m'_\theta(0)}{m_\theta^2(0)} \cdot \left(\gamma \mathbb{E} \frac{gh}{(h_\theta \cdot m_\theta(0) + 1)^2} + \tilde{\sigma}^2\right), & \theta\gamma > 1 \\ \left(\gamma \mathbb{E}[gh \cdot \mathbb{I}_{h < Q_h(1-\theta)}] + \tilde{\sigma}^2\right) \frac{1}{1 - \theta\gamma}, & \theta\gamma < 1 \end{cases} \quad (4.4)$$

where  $h_\theta = h \cdot \mathbb{I}_{h \geq Q_h(1-\theta)}$  and  $m_\theta(z)$  satisfies  $-z = m_\theta^{-1}(z) - \gamma \mathbb{E} h_\theta \cdot (1 + h_\theta \cdot m_\theta(z))^{-1}$ .

In addition, if  $\mathbb{E}[g|h]$  is a decreasing function of  $h$ , and  $h$  has continuous p.d.f., then the asymptotic prediction risk of  $\hat{\beta}_\theta$  is a decreasing function of  $\theta$  when  $\theta\gamma > 1$ .

Corollary 3 confirms double descent under more general settings of  $(\Sigma_x, \Sigma_\beta)$  than [XH19], i.e. the risk exhibits a spike as  $\theta\gamma \rightarrow 1^-$ , and then decreases as we further overparameterize by increasing  $\theta$ . In Section 6 we compare the PCR estimator  $\hat{\beta}_\theta$  with the minimum  $\|\hat{\beta}\|_{\Sigma_w}$  norm solution.

**Remark.** The PCR estimator [XH19] and the ridgeless regression estimator (considered in [HMRT19]) are fundamentally different in the following way: in ridgeless regression, increasing the model size corresponds to changing  $\gamma$ , which also alters the dimensions of  $\beta_*$ ; in contrast, in PCR, increasing  $\theta$  does not change the data generating process (which is a more natural setting).

In terms of the risk curve, Figure 9(a) shows that the ridgeless regression estimator can exhibit “multiple descent” as  $\gamma$  increases, whereas Corollary 3 and Figure 9(b) demonstrate that in the misaligned case, the PCR risk is monotonically decreasing in the overparameterized regime  $\theta\gamma > 1$ .

## 5 Analysis of Optimal $\lambda_{\text{opt}}$

In this section, we focus on the optimal weighted ridge estimator and determine the sign of the optimal regularization parameter  $\lambda_{\text{opt}}$ . Taking the derivatives of (4.1) yields

$$R'(\lambda) \cdot \frac{m^3(-\lambda)}{2\gamma(m'(-\lambda))^2} = \underbrace{\left(-\tilde{\sigma}^2 \frac{\mathbb{E} \frac{\zeta^2}{(1+\zeta)^3}}{1 - \gamma \mathbb{E} \frac{\zeta^2}{(1+\zeta)^2}}\right)}_{\text{Part 3}} + \underbrace{\left(\mathbb{E} \frac{gh\zeta}{(1+\zeta)^3} - \frac{\gamma \mathbb{E} \frac{\zeta^2}{(1+\zeta)^3} \mathbb{E} \frac{gh}{(1+\zeta)^2}}{1 - \gamma \mathbb{E} \frac{\zeta^2}{(1+\zeta)^2}}\right)}_{\text{Part 4}}, \quad (5.1)$$

where  $\zeta = h \cdot m(-\lambda)$ . For certain special cases, we obtain a closed form solution for  $\lambda_{\text{opt}}$  (see details in Appendix B.1) and recover the result from [HMRT19, DW18]<sup>6</sup> and beyond:

- When  $h \stackrel{\text{a.s.}}{=} c$  (i.e., isotropic features [HMRT19]), the optimal  $\lambda_{\text{opt}}$  is achieved at  $c/\xi$ .
- When  $g \stackrel{\text{a.s.}}{=} c$  (i.e., isotropic signals [DW18]), the optimal  $\lambda_{\text{opt}}$  is achieved at  $\tilde{\sigma}^2/c$ .
- When  $\mathbb{E}[g|h] \stackrel{\text{a.s.}}{=} \mathbb{E}[g]$  (e.g., random order), the optimal  $\lambda_{\text{opt}}$  is achieved at  $\tilde{\sigma}^2/\mathbb{E}[g]$ .

<sup>6</sup>In [HMRT19],  $h \stackrel{\text{a.s.}}{=} 1$ . In [DW18],  $\tilde{\sigma}^2 = 1$  and their signal strength  $\alpha^2$  is equivalent to  $c\gamma$  in our setting.

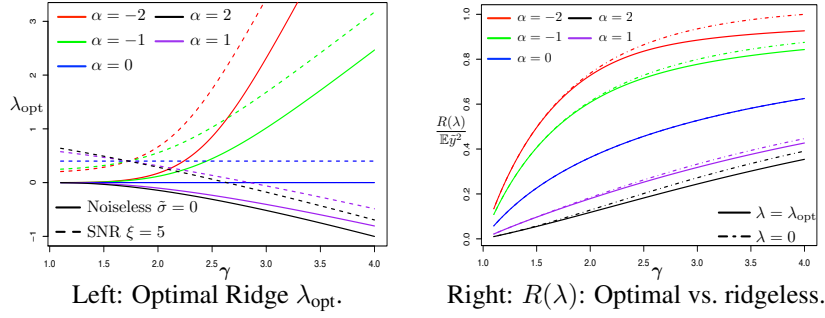


Figure 4: We set  $\Sigma_w = \mathbf{I}$  and  $\Sigma_\beta = \Sigma_x^\alpha$  where  $\mathbf{d}_x$  has two point masses on 1 and 5 with probability 3/4 and 1/4 respectively. **Left:** optimal  $\lambda$ ; solid lines represents the noiseless case  $\tilde{\sigma} = 0$  and dashed lines represents SNR  $\xi = 5$ . **Right:** prediction risk of the ridgeless ( $R(\lambda_{\text{opt}})$ , dashed lines) and optimally regularized ( $R(\lambda_{\text{opt}})$ , solid lines) estimator in the noiseless case. We normalize the prediction risk as  $\mathbb{E}\tilde{y}^2 = \mathbb{E}(\tilde{\mathbf{x}}^\top \tilde{\beta}_*)^2$ .

Although  $\lambda_{\text{opt}}$  may not have a tractable form in general, we may infer the sign of  $\lambda_{\text{opt}}$ . Note that in (5.1), Part 3 is due to the variance term (Part 1) and Part 4 from the bias term (Part 2) in (3.1). We therefore consider the sign of Part 3 and Part 4 separately in the following theorem.

**Theorem 4.** *Under Assumption 1, we have*

- Part 3 (derivative of variance) is negative for all  $\lambda > -c_0$ .
- If  $\mathbb{E}[g|h]$  is an increasing function of  $h$  on its support, then Part 4 (derivative of bias) is positive for all  $\lambda > 0$ . At  $\lambda = 0$ , Part 4 is non-negative and achieves 0 only if  $\mathbb{E}[g|h] \stackrel{\text{a.s.}}{=} \mathbb{E}[g]$ .
- If  $\mathbb{E}[g|h]$  is a decreasing function of  $h$  on its support, then Part 4 is negative for all  $\lambda \in (-c_0, 0)$ . At  $\lambda = 0$ , Part 4 is non-positive and achieves 0 only if  $\mathbb{E}[g|h] \stackrel{\text{a.s.}}{=} \mathbb{E}[g]$ .

The first point in Theorem 4 is consistent with the well-understood variance reduction property of ridge regularization. On the other hand, when the prediction risk is dominated by the bias term (i.e.,  $\tilde{\sigma}^2 = o(1)$ ) and both  $\mathbf{d}_{x/w}$  and  $\mathbf{d}_{w\beta}$  converge to non-trivial distributions, the second and third point of Theorem 4 reveal the following surprising phenomena (see Figure 2 (a) and (b)):

- M1**  $\lambda_{\text{opt}} < 0$  when  $\mathbf{d}_{x/w}$  aligns with  $\mathbf{d}_{w\beta}$ , or in general,  $\mathbb{E}[g|h]$  is a strictly increasing function of  $h$ . In the context of standard ridge regression, it means that shrinkage regularization only *increases the bias* in the overparameterized regime when features are informative, i.e., the projection of the signal is large in the directions where the feature variance is large.
- M2**  $\lambda_{\text{opt}} > 0$  when  $\mathbf{d}_{x/w}$  is misaligned with  $\mathbf{d}_{w\beta}$ , or in general,  $\mathbb{E}[g|h]$  is a strictly decreasing function of  $h$ . This is to say, in standard ridge regression, when features are not informative, i.e., the projection of the signal is small in the directions of large feature variance, shrinkage is beneficial even in the *absence of label noise* (the variance term is zero).

M1 and M2, together with aforementioned special case when  $g$  and  $h$  have random relation, provide a precise characterization of the sign of  $\lambda_{\text{opt}}$ . In particular, M1 confirms the “negative ridge” phenomenon empirically observed in [KLS20] and outlines concise conditions under which it occurs. We emphasize that neither M1 nor M2 would be observed when one of  $\Sigma_{x/w}$  and  $\Sigma_{w\beta}$  is identity. In other words, these observations arise from our more general assumption on  $(\Sigma_{x/w}, \Sigma_{w\beta})$ .

**Implicit regularization of overparameterization.** Taking both the bias and variance into account, Theorem 4 suggests a bias-variance tradeoff between Part 3 and Part 4, and  $\lambda_{\text{opt}}$  will eventually become positive as  $\tilde{\sigma}^2$  increases (i.e., risk is dominated by variance, for which positive  $\lambda$  is beneficial). For certain special cases, we can provide a lower bound for the transition from  $\lambda_{\text{opt}} < 0$  to  $\lambda_{\text{opt}} > 0$ .

**Proposition 5.** *Given Assumption 1, let  $(h, g) = (1, 1)$  with probability  $1 - q$  and  $(h, g) = (h_1, g_1)$  with probability  $q$ , where  $h_1 > 1$  and  $g_1 > 1$ . Denote  $\tilde{\gamma} = \gamma - 1$ . Then  $\lambda_{\text{opt}} < 0$  if*

$$\tilde{\sigma}^2 < (h_1 - 1)(g_1 - 1)h_1 \cdot \max\left(\frac{(\gamma q - 1)^3 \tilde{\gamma}^3 (1 - q)}{(1 - q)\gamma^2 (\tilde{\gamma}^3 q^2 + (\gamma q - 1)^3 h_1^2)}, \frac{\gamma q (1 - q) \tilde{\gamma}^3}{(1 - q)(h_1 + \tilde{\gamma})^3 + q h_1^2 \gamma^3}\right).$$

As  $q$  approaches 0 or 1, the above upper bound goes to 0 because  $\Sigma_x, \Sigma_\beta$  becomes closer to  $\mathbf{I}$ . Otherwise, when  $\gamma q > 1$ , the upper bound suggests  $\tilde{\sigma}^2 = O(g_1 \gamma)$  which implies the SNR  $\xi = \Omega(h_1/\gamma)$ .

Hence, as  $\gamma$  increases,  $\lambda_{\text{opt}}$  remains negative for a lower SNR, which coincides with the intuition that overparameterization has an implicit effect of  $\ell_2$  regularization (Figure 4 Left). Indeed, the following proposition suggests such implicit regularization is only present in the overparameterized regime:

**Proposition 6.** *When  $\gamma < 1$ ,  $\lambda_{\text{opt}}$  on  $(-c_0, \infty)$  is always non-negative under Assumption 1.*

In Figure 4 we confirm our findings in Theorem 4 (for additional results on different distributions see Figure 10). Specifically, we set  $\Sigma_w = \mathbf{I}$ ,  $\Sigma_x = \text{diag}(\mathbf{d}_x)$  and  $\Sigma_\beta = \Sigma_x^\alpha$ . As we increase  $\alpha$  from negative to positive, the relation between  $\mathbf{d}_x$  and  $\mathbf{d}_\beta$  transitions from misaligned to aligned. The left panel shows that the sign of  $\lambda_{\text{opt}}$  is the exact opposite to the sign of  $\alpha$  in the noiseless case (i.e. the variance is 0), which is consistent with M1 and M2. Moreover, when  $\mathbf{d}_x$  aligns with  $\mathbf{d}_\beta$ ,  $\lambda_{\text{opt}}$  decreases as  $\gamma$  becomes larger, which agrees with our observation on the implicit  $\ell_2$  regularization of overparameterization. Last but not least, in Figure 4 (Right) we see that the optimal ridge regression estimator leads to considerable improvement over the ridgeless estimator. We comment that this improvement becomes more significant as  $\gamma$  or condition number of  $\Sigma_x$  and  $\Sigma_\beta$  increases.

**Risk monotonicity of optimal ridge regression.** [DS20, Proposition 6] showed that for isotropic data ( $\Sigma_x = \mathbf{I}$ ), the asymptotic prediction risk of optimally-tuned ridge regression monotonically increases with  $\gamma$ . This is to say, under proper regularization, more training data always helps the test performance [KH92]. Here we extend this result to data with general covariance and isotropic  $\beta_*$ .

**Proposition 7.** *Given  $\mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \Sigma_x$  satisfying Assumption 1 and  $\mathbb{E}[\beta_*\beta_*^\top] = \frac{c}{p}\mathbf{I}^7$ , the asymptotic prediction risk of the optimally-tuned ridge regression estimator (i.e.,  $\Sigma_w = \mathbf{I}$ ) with  $\lambda_{\text{opt}} = \gamma\tilde{\sigma}^2/c$  is an increasing function of  $\gamma \in (0, \infty)$ .*

As shown in Figure 5 (where  $\mathbf{d}_x$  has 3 point masses and  $\mathbf{d}_\beta = 1$ ),  $\ell_2$  regularization can suppress “multiple descent”, and the risk of the optimally-tuned ridge estimator (purple) is monotone w.r.t.  $\gamma$ . We remark that establishing such characterization under general orientation of  $\beta_*$  (anisotropic  $\Sigma_\beta$ ) can be challenging, because the optimal regularization  $\lambda_{\text{opt}}$  may not have a convenient closed-form. We leave the analysis for general  $\Sigma_\beta$  as future work.

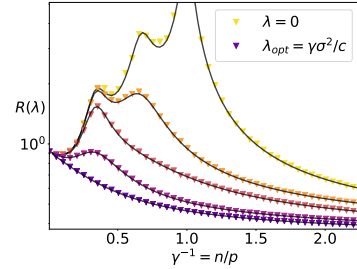


Figure 5: Impact of ridge regularization on the risk curve (SNR=3). Darker color corresponds to larger  $\lambda$ .

## 6 Optimal Weighting Matrix

Having characterized the optimal regularization strength, we now turn to the optimal weighting matrix  $\Sigma_w$ . Toward this goal, we additionally require the following assumptions on  $(\Sigma_x, \Sigma_\beta, \Sigma_w)$ :

**Assumption 2.** *The covariance matrix  $\Sigma_x$  and the weighting matrix  $\Sigma_w$  share the same set of eigenvectors, i.e., we have the following eigendecompositions:  $\Sigma_x = \mathbf{U}\mathbf{D}_x\mathbf{U}^\top$  and  $\Sigma_w = \mathbf{U}\mathbf{D}_w\mathbf{U}^\top$ , where  $\mathbf{U} \in \mathbb{R}^{p \times p}$  is orthogonal, and  $\mathbf{D}_x = \text{diag}(\mathbf{d}_x)$ ,  $\mathbf{D}_w = \text{diag}(\mathbf{d}_w)$ .*

Define  $\bar{\mathbf{d}}_\beta = \text{diag}(\mathbf{U}^\top \Sigma_\beta \mathbf{U})$ . Note that when  $\Sigma_\beta$  also shares the same eigenvector matrix  $\mathbf{U}$ , then  $\bar{\mathbf{d}}_\beta = \mathbf{d}_\beta$ , which is simply the eigenvalues of  $\Sigma_\beta$ .

**Assumption 3.** *Let  $d_{x,i}, \bar{d}_{\beta,i}, d_{w,i}$  be the  $i$ th element of  $\mathbf{d}_x, \bar{\mathbf{d}}_\beta, \mathbf{d}_w$  respectively. We assume that the empirical distribution of  $(d_{x,i}, \bar{d}_{\beta,i}, d_{w,i})$  jointly converges to  $(s, v, s/r)$ , where  $s, v, r$  are non-negative random variables. Further, there exists constants  $c_l, c_u > 0$  independent of  $n$  and  $p$  such that  $\min_i(\min(d_{x,i}, \bar{d}_{\beta,i}, d_{w,i})) \geq c_l$ ,  $\max_i(\max(d_{x,i}, \bar{d}_{\beta,i}, d_{w,i})) \leq c_u$  and  $\|\Sigma_\beta\| \leq c_u$ .*

For notational convenience, we define  $\mathcal{H}_w$  and  $\mathcal{H}_r$  to be the sets of all  $\Sigma_w$  and  $r$ , respectively, that satisfy Assumption 2 and Assumption 3. Additionally, let  $\mathcal{S}_w$  and  $\mathcal{S}_r$  be the subset of  $\mathcal{H}_w$  and  $\mathcal{H}_r$  such that  $r = f(s)$  for some function  $f$  (this represents  $\Sigma_w \in \mathcal{H}_w$  that only depends on  $\Sigma_x$  but not  $\Sigma_\beta$ ). By Assumption 2 and 3, the empirical distribution of  $(d_{x/w,i}, d_{w\beta,i})$  jointly converges to  $(r, sv/r)$  and satisfies the boundedness requirement in Assumption 1. Thus by Theorem 1 we have:

$$R(r, \lambda) \triangleq \frac{m'_r(-\lambda)}{m_r^2(-\lambda)} \cdot \left( \gamma \mathbb{E} \frac{sv}{(r \cdot m_r(-\lambda) + 1)^2} + \tilde{\sigma}^2 \right), \quad (6.1)$$

where  $m_r(-\lambda)$  satisfies the equation  $\lambda = m_r^{-1}(-\lambda) - \gamma \mathbb{E}(1 + r \cdot m_r(-\lambda))^{-1} r$ . It is clear that when  $r \stackrel{\text{a.s.}}{=} s$ , (6.1) reduces to the standard ridge regression with  $\Sigma_w = \mathbf{I}$ , and for  $r \stackrel{\text{a.s.}}{=} 1$ , the equation

<sup>7</sup>Note that the parameter scaling differs from the previous setting by  $\gamma$  to be consistent with that of [DS20].

reduces to the cases of isotropic features ( $\Sigma_w = \Sigma_x$ ). Note that (6.1) indicates that the impact of  $\Sigma_\beta$  on the risk is fully captured by  $\bar{\mathbf{d}}_\beta$ . Hence we define  $\bar{\Sigma}_\beta = \mathbf{U} \text{diag}(\bar{\mathbf{d}}_\beta) \mathbf{U}^\top$ , which corresponds to  $r \stackrel{\text{a.s.}}{=} sv$ , and is equivalent to  $\Sigma_\beta$  when  $\Sigma_\beta$  also shares the same eigenvector matrix  $\mathbf{U}$ . In the following subsections, we discuss the optimal  $\Sigma_w$  for two types of estimator: the minimum  $\|\hat{\beta}\|_{\Sigma_w}$  solution (taking  $\lambda \rightarrow 0$ ), and the optimally weighted ridge estimator ( $\lambda = \lambda_{\text{opt}}$ ). Note that the risk for both estimators is scale-invariant over  $\Sigma_w$  and  $r$ . Hence, when we define a specific choice of  $(\Sigma_w, r)$ , we simultaneously consider all pairs  $(c\Sigma_w, r/c)$  for  $c > 0$ . Finally, we note that the choice of  $r \stackrel{\text{a.s.}}{=} s \cdot \mathbb{E}[v|s] \in \mathcal{S}_r$  plays a key role in our analysis, and its corresponding choice of  $\Sigma_w$  is given as  $\Sigma_w = (f_v(\Sigma_x))^{-1}$ , where  $f_v(s) \triangleq \mathbb{E}[v|s]$  and  $f_v$  applies to the eigenvalues of  $\Sigma_x$ .

## 6.1 Minimum $\|\hat{\beta}\|_{\Sigma_w}$ solution

Taking the ridgeless limit leads to the following bias-variance decomposition of the prediction risk,

$$\text{Bias: } R_b(r) \triangleq \frac{m_r'(0)}{m_r^2(0)} \cdot \gamma \mathbb{E} \frac{sv}{(r \cdot m_r(0) + 1)^2} \quad \text{Variance: } R_v(r) \triangleq \frac{m_r'(0)}{m_r^2(0)} \cdot \tilde{\sigma}^2.$$

In the previous sections we observe a bias-variance tradeoff in choosing the optimal  $\lambda$ . Interesting, the following theorem illustrates a similar bias-variance tradeoff in choosing the optimal  $\Sigma_w$ :

**Theorem 8.** *Given Assumptions 2 and 3,*

- $r \stackrel{\text{a.s.}}{=} sv$  (i.e.,  $\Sigma_w = \bar{\Sigma}_\beta^{-1}$ ) is the optimal choice in  $\mathcal{H}_r$  that minimizes the bias function  $R_b(r)$ . Additionally,  $r \stackrel{\text{a.s.}}{=} \mathbb{E}[v|s] \cdot s$  (i.e.,  $\Sigma_w = (f_v(\Sigma_x))^{-1}$ ) is the optimal in  $\mathcal{S}_r$  that minimizes  $R_b(r)$ .
- $r \stackrel{\text{a.s.}}{=} 1$  (i.e.,  $\Sigma_w = \Sigma_x$ ) is optimal in both  $\mathcal{S}_r$  and  $\mathcal{H}_r$  that minimizes the variance function  $R_v(r)$ .

Theorem 8 implies that the variance is minimized when  $\Sigma_w = \Sigma_x$ . Since the variance term does not depend on  $\beta_*$ , it is not surprising that the optimal  $\Sigma_w$  is also independent of  $\Sigma_\beta$ . Furthermore, this result is consistent with the intuition that to minimize the variance,  $\hat{\beta}_\lambda$  should be penalized more in the higher variance directions of  $\Sigma_x$ , and vice versa. On the other hand, Theorem 8 also implies that the bias is minimized when  $\mathbf{d}_w = 1/\bar{\mathbf{d}}_\beta$  which does not depend on  $\mathbf{d}_x$ . While this characterization may not be intuitive, when  $\bar{\mathbf{d}}_\beta = \mathbf{d}_\beta$  (i.e.,  $\Sigma_\beta$  also shares the same eigenvector matrix  $\mathbf{U}$ ), one analogy is that since the quadratic regularization corresponds to the a Gaussian prior  $\mathcal{N}(\mathbf{0}, \Sigma_w^{-1})$ , it is reasonable to match  $\Sigma_w^{-1}$  with the covariance of  $\beta_*$ , which gives the maximum a posteriori (MAP) estimate. In general, the optimal  $\Sigma_w$  admits a bias-variance tradeoff (i.e., the bias and variance are optimal under different  $\Sigma_w$ ) except for the special case of  $\Sigma_x \Sigma_\beta = \mathbf{I}$ .

Additionally, the following proposition demonstrates the advantage of the minimum  $\|\hat{\beta}\|_{\Sigma_w}$  solution over the PCR estimator in the noiseless case.

**Proposition 9.** *Given Assumption 2 and 3 and  $\tilde{\sigma} = 0$ , suppose  $s$  and  $\mathbb{E}[v|s] \cdot s$  both have continuous and strictly increasing quantile functions. Then the minimum  $\|\hat{\beta}\|_{\Sigma_w}$  solution outperforms the PCR estimator for all  $\theta \in [0, 1)$  when  $\Sigma_w = \bar{\Sigma}_\beta^{-1} \in \mathcal{H}_w$ , or when  $\Sigma_w = (f_v(\Sigma_x))^{-1} \in \mathcal{S}_w$ .*

## 6.2 Optimal weighted ridge estimator

Finally, we consider the optimally-tuned weighted shrinkage and discuss the optimal choice of  $\Sigma_w$ .

**Theorem 10.** *Suppose Assumptions 2 and 3 hold. Then  $r \stackrel{\text{a.s.}}{=} sv$  (i.e.,  $\Sigma_w = \bar{\Sigma}_\beta^{-1}$ ) is the optimal solution in  $\mathcal{H}_r$  that minimizes  $\min_\lambda R(r, \lambda)$ . Additionally,  $r \stackrel{\text{a.s.}}{=} \mathbb{E}[v|s] \cdot s$  (i.e.,  $\Sigma_w = (f_v(\Sigma_x))^{-1}$ ) is the optimal solution in  $\mathcal{S}_r$  that minimizes  $\min_\lambda R(r, \lambda)$ .*

In contrast to the ridgeless setting in Theorem 8, the optimal  $\mathbf{d}_w$  for general  $\lambda_{\text{opt}}$  does not depend on the noise level but only on  $\bar{\mathbf{d}}_\beta$ , the strength of the signal in the directions of the eigenvectors of  $\Sigma_x$ . We conjecture that this is because in the optimally weighted estimator,  $\lambda_{\text{opt}}$  is capable of balancing the bias-variance tradeoff; therefore the weighting matrix may not need to adjust to the label noise and can be chosen solely based on the signal  $\beta_*$ . Indeed, as previously discussed,  $\Sigma_w = \Sigma_\beta^{-1}$  is a preferable choice of prior under the Bayesian perspective when  $\mathbf{d}_\beta = \bar{\mathbf{d}}_\beta$ .



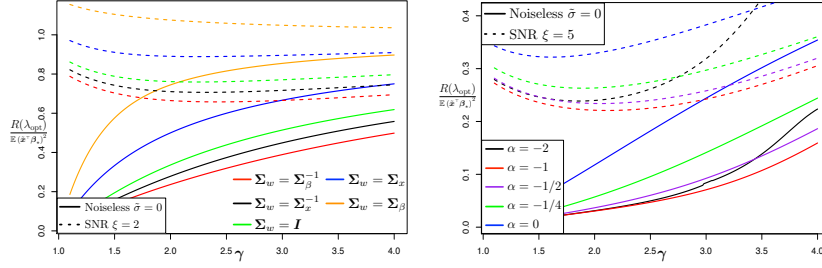


Figure 6:  $R(\lambda_{\text{opt}})/\mathbb{E}(\tilde{\mathbf{x}}^\top \beta_x)^2$  against  $\gamma$  for various weighting matrix  $\Sigma_w$ . Solid lines represent the noiseless case  $\tilde{\sigma} = 0$  and the dashed lines represent the noisy case with fixed SNR  $\xi$ . We set  $\mathbf{d}_x$  to be aligned with  $\mathbf{d}_\beta$  and **Left**:  $\mathbf{d}_x$  to have 4 point masses (1, 2, 3, 4) with equal probabilities and  $\mathbf{d}_\beta$  with 2 point masses on 1 and 5 with probabilities 3/4 and 1/4, respectively; **Right**:  $\mathbf{d}_x$  has 2 point masses on 1 and 5 with probabilities 3/4 and 1/4, respectively, and  $\Sigma_\beta = \Sigma_x^2$ ; we set  $\Sigma_w = \Sigma_\beta^\alpha$ .

Theorem 10 is supported by Figure 6, where we plot the prediction risk of the generalized ridge regression estimator under different  $\Sigma_w$  and optimally tuned  $\lambda_{\text{opt}}$ . We consider a simple discrete construction for aligned  $\mathbf{d}_x$  and  $\mathbf{d}_\beta (= \bar{\mathbf{d}}_\beta)$ . On the left panel, we enumerate a few standard choices of  $\Sigma_w$ :  $\Sigma_x$ ,  $\Sigma_\beta$ ,  $\mathbf{I}$ ,  $\Sigma_x^{-1}$  and the optimal choice  $\Sigma_\beta^{-1}$  (red). On the right, we take  $\Sigma_w$  to be powers of  $\Sigma_\beta$  around the optimal  $\Sigma_\beta^{-1}$ . In both setups, we confirm that  $\Sigma_\beta^{-1}$  achieves the lowest risk uniformly over  $\gamma$ , as predicted by Theorem 10.

Note that our main results require knowledge of  $\Sigma_x$  and  $\bar{\Sigma}_\beta$ . While  $\Sigma_x$  can be obtained in a semi-supervised setting using unlabeled data (e.g., [RC15, TCG20]), it is typically difficult to estimate  $\bar{\Sigma}_\beta$  directly from data. Without prior knowledge on  $\bar{\Sigma}_\beta$ , Theorem 10 suggests that  $r \stackrel{\text{a.s.}}{=} \mathbb{E}[v|s] \cdot s$  is the optimal  $r$  that only depends on  $s$ . That is,  $\Sigma_w = (f_v(\Sigma_x))^{-1}$  is the optimal  $\Sigma_w$  that only depends on  $\Sigma_x$ . In the special case of  $\mathbb{E}[v|s] = \mathbb{E}[v]$ , standard ridge regression ( $\Sigma_w = \mathbf{I}$ ) is optimal in  $\mathcal{S}_w$ . When the exact form of  $f_v(s)$  is also not known, we may use a polynomial or power function of  $s$  to approximate either  $f_v(s)$  or  $1/f_v(s)$ , whose coefficients can be considered as hyper-parameters to be cross-validated. We demonstrate the effectiveness of this heuristic in Figure 7: although our proposed  $\Sigma_w = f_v(\Sigma_x)^{-1}$  (blue) is worse than the actual optimal (red)  $\Sigma_w = \Sigma_\beta^{-1}$  (same as  $\bar{\Sigma}_\beta^{-1}$  due to diagonal design), it is the best choice among weighting matrices that only depend on  $\Sigma_x$ . In addition, we seek the best approximation of  $f_v(s)$  by applying a power transformation on  $\Sigma_x$ , and we observe that certain powers of  $\Sigma_x$  also outperform the standard isotropic regularization.

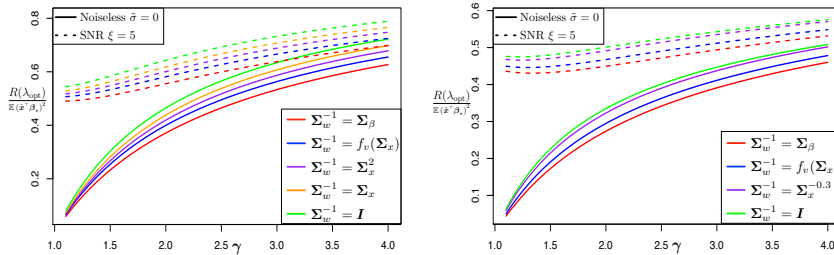


Figure 7:  $R(\lambda_{\text{opt}})/\mathbb{E}(\tilde{\mathbf{x}}^\top \beta_x)^2$  against  $\gamma$  for various weighting matrix  $\Sigma_w$  under noiseless  $\tilde{\sigma} = 0$  (solid lines) and noisy setting with fixed SNR  $\xi$  (dashed lines). **Left**: We set  $f_v(s)$  as an increasing function of  $s$  on its support; **Right**: We set  $f_v(s)$  as a decreasing function of  $s$  on its support. Note that the heuristically chosen weighting matrices often outperform the standard ridge regression estimator (green).

## 7 Conclusion

We provide a precise asymptotic characterization of the prediction risk of generalized ridge regression in the overparameterized regime. Our result greatly generalizes previous high-dimensional analysis of ridge regression, which enables us to discover and theoretically justify various interesting findings, including the negative ridge phenomenon, the implicit regularization of overparameterization, and a concise description of the optimal weighted shrinkage. Future works include extending our analysis to border settings, such as more general eigenvalue conditions [XH19] or the random features regression model [MM19]. Another important direction is to construct weighting matrix  $\Sigma_w$  solely from training data that outperforms isotropic shrinkage in the overparameterized regime.

## 8 Broader Impact

This work does not present any foreseeable direct societal consequence.

## Acknowledgement

The authors would like to thank Murat A. Erdogdu, Daniel Hsu and Taiji Suzuki for comments and suggestions, and also anonymous NeurIPS reviewers 1 and 3 for helpful feedback. DW was partially funded by CIFAR, NSERC and LG Electronics. JX was supported by a Cheung-Kong Graduate School of Business Fellowship.

## References

- [ABG<sup>+</sup>20] Shun-ichi Amari, Jimmy Ba, Roger Grosse, Xuechen Li, Atsushi Nitanda, Taiji Suzuki, Denny Wu, and Ji Xu, *When does preconditioning help or hurt generalization?*, arXiv preprint arXiv:2006.10732 (2020).
- [AKT19] Alnur Ali, J Zico Kolter, and Ryan J Tibshirani, *A continuous-time view of early stopping for least squares*, International Conference on Artificial Intelligence and Statistics, vol. 22, 2019.
- [AS17] Madhu S Advani and Andrew M Saxe, *High-dimensional dynamics of generalization error in neural networks*, arXiv preprint arXiv:1710.03667 (2017).
- [BES<sup>+</sup>20] Jimmy Ba, Murat Erdogdu, Taiji Suzuki, Denny Wu, and Tianzong Zhang, *Generalization of two-layer neural networks: An asymptotic viewpoint*, International Conference on Learning Representations, 2020.
- [BHMM18] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal, *Reconciling modern machine learning and the bias-variance trade-off*, arXiv preprint arXiv:1812.11118 (2018).
- [BLLT19] Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler, *Benign overfitting in linear regression*, arXiv preprint arXiv:1906.11300 (2019).
- [BS99] Anders Björkström and Rolf Sundberg, *A generalized view on continuum regression*, Scandinavian Journal of Statistics **26** (1999), no. 1, 17–30.
- [BVDBS<sup>+</sup>15] Małgorzata Bogdan, Ewout Van Den Berg, Chiara Sabatti, Weijie Su, and Emmanuel J Candès, *Slope—adaptive variable selection via convex optimization*, The annals of applied statistics **9** (2015), no. 3, 1103.
- [Cas80] George Casella, *Minimax ridge regression estimation*, The Annals of Statistics (1980), 1036–1056.
- [CWB08] Emmanuel J Candès, Michael B Wakin, and Stephen P Boyd, *Enhancing sparsity by reweighted  $\ell_1$  minimization*, Journal of Fourier analysis and applications **14** (2008), no. 5-6, 877–905.
- [Dic16] Lee H Dicker, *Ridge regression and asymptotic minimax estimation over spheres of growing dimension*, Bernoulli **22** (2016), no. 1, 1–37.
- [DKT19] Zeyu Deng, Abla Kammoun, and Christos Thrampoulidis, *A model of double descent for high-dimensional binary linear classification*, arXiv preprint arXiv:1911.05822 (2019).
- [DL20] Oussama Dhifallah and Yue M Lu, *A precise performance analysis of learning with random features*, arXiv preprint arXiv:2008.11904 (2020).
- [DLM19] Michał Dereziński, Feynman Liang, and Michael W Mahoney, *Exact expressions for double descent and implicit regularization via surrogate random design*, arXiv preprint arXiv:1912.04533 (2019).
- [DM16] David Donoho and Andrea Montanari, *High dimensional robust  $m$ -estimation: Asymptotic variance via approximate message passing*, Probability Theory and Related Fields **166** (2016), no. 3-4, 935–969.

- [dRBK20] Stéphane d’Ascoli, Maria Refinetti, Giulio Biroli, and Florent Krzakala, *Double trouble in double descent: Bias and variance ( $s$ ) in the lazy regime*, arXiv preprint arXiv:2003.01054 (2020).
- [DS20] Edgar Dobriban and Yue Sheng, *Wonder: Weighted one-shot distributed ridge regression in high dimensions.*, Journal of Machine Learning Research **21** (2020), no. 66, 1–52.
- [DW18] Edgar Dobriban and Stefan Wager, *High-dimensional asymptotics of prediction: Ridge regression and classification*, The Annals of Statistics **46** (2018), no. 1, 247–279.
- [HG83] Tsushung A Hua and Richard F Gunst, *Generalized ridge regression: a note on negative ridge parameters*, Communications in Statistics-Theory and Methods **12** (1983), no. 1, 37–45.
- [HK70] Arthur E Hoerl and Robert W Kennard, *Ridge regression: Biased estimation for nonorthogonal problems*, Technometrics **12** (1970), no. 1, 55–67.
- [HMRT19] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani, *Surprises in high-dimensional ridgeless least squares interpolation*, arXiv preprint arXiv:1903.08560 (2019).
- [HMX20] Daniel Hsu, Vidya Muthukumar, and Ji Xu, *On the proliferation of support vectors in high dimensions*, arXiv preprint arXiv:2009.10670 (2020).
- [Kar13] Noureddine El Karoui, *Asymptotic behavior of unregularized and ridge-regularized high-dimensional robust regression estimators: rigorous results*, arXiv preprint arXiv:1311.2445 (2013).
- [KH92] Anders Krogh and John A Hertz, *A simple weight decay can improve generalization*, Advances in neural information processing systems, 1992, pp. 950–957.
- [KLS20] Dmitry Kobak, Jonathan Lomond, and Benoit Sanchez, *The optimal ridge penalty for real-world high-dimensional data can be zero or negative due to the implicit ridge regularization*, Journal of Machine Learning Research **21** (2020), no. 169, 1–16.
- [KPR<sup>+</sup>17] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al., *Overcoming catastrophic forgetting in neural networks*, Proceedings of the national academy of sciences **114** (2017), no. 13, 3521–3526.
- [LD19] Sifan Liu and Edgar Dobriban, *Ridge regression: Structure, cross-validation, and sketching*, arXiv preprint arXiv:1910.02373 (2019).
- [LH17] Ilya Loshchilov and Frank Hutter, *Decoupled weight decay regularization*, arXiv preprint arXiv:1711.05101 (2017).
- [Lol20] Panagiotis Lolas, *Regularization in high-dimensional regression and classification via random matrix theory*, arXiv preprint arXiv:2003.13723 (2020).
- [LP11] Olivier Ledoit and Sandrine Péché, *Eigenvectors of some large sample covariance matrix ensembles*, Probability Theory and Related Fields **151** (2011), no. 1-2, 233–264.
- [LPRS17] Tengyuan Liang, Tomaso Poggio, Alexander Rakhlin, and James Stokes, *Fisher-rao metric, geometry, and complexity of neural networks*, arXiv preprint arXiv:1711.01530 (2017).
- [LW17] Christos Louizos and Max Welling, *Multiplicative normalizing flows for variational bayesian neural networks*, Proceedings of the 34th International Conference on Machine Learning-Volume 70, JMLR. org, 2017, pp. 2218–2227.
- [MM19] Song Mei and Andrea Montanari, *The generalization error of random features regression: Precise asymptotics and double descent curve*, arXiv preprint arXiv:1908.05355 (2019).
- [MRSY19] Andrea Montanari, Feng Ruan, Youngtak Sohn, and Jun Yan, *The generalization error of max-margin linear classifiers: High-dimensional asymptotics in the over-parametrized regime*, arXiv preprint arXiv:1911.01544 (2019).

- [MS05] Yuzo Maruyama and William E Strawderman, *A new class of generalized bayes minimax ridge regression estimators*, The Annals of Statistics **33** (2005), no. 4, 1753–1770.
- [MS18] Yuichi Mori and Taiji Suzuki, *Generalized ridge estimator and model selection criteria in multivariate linear regression*, Journal of Multivariate Analysis **165** (2018), 243–261.
- [RC15] Kenneth Joseph Ryan and Mark Vere Culp, *On semi-supervised linear regression in covariate shift problems*, The Journal of Machine Learning Research **16** (2015), no. 1, 3183–3217.
- [RM11] Francisco Rubio and Xavier Mestre, *Spectral convergence for a general class of random matrices*, Statistics & probability letters **81** (2011), no. 5, 592–602.
- [RMR20] Dominic Richards, Jaouad Mourtada, and Lorenzo Rosasco, *Asymptotics of ridge (less) regression under general source condition*, arXiv preprint arXiv:2006.06386 (2020).
- [SC95] Jack W Silverstein and Sang-Il Choi, *Analysis of the limiting spectral distribution of large dimensional random matrices*, Journal of Multivariate Analysis **54** (1995), no. 2, 295–309.
- [Str78] William E Strawderman, *Minimax adaptive generalized ridge regression estimators*, Journal of the American Statistical Association **73** (1978), no. 363, 623–627.
- [TAH18] Christos Thrampoulidis, Ehsan Abbasi, and Babak Hassibi, *Precise error analysis of regularized  $m$ -estimators in high dimensions*, IEEE Transactions on Information Theory **64** (2018), no. 8, 5592–5628.
- [TB20] Alexander Tsigler and Peter L Bartlett, *Benign overfitting in ridge regression*, arXiv preprint arXiv:2009.14286 (2020).
- [TCG20] T Tony Cai and Zijian Guo, *Semisupervised inference for explained variance in high dimensional linear regression and its applications*, Journal of the Royal Statistical Society: Series B (Statistical Methodology) (2020).
- [XH19] Ji Xu and Daniel J Hsu, *On the number of variables to use in principal component regression*, Advances in Neural Information Processing Systems, 2019, pp. 5095–5104.
- [Zou06] Hui Zou, *The adaptive lasso and its oracle properties*, Journal of the American statistical association **101** (2006), no. 476, 1418–1429.
- [ZTSG19] Han Zhao, Yao-Hung Hubert Tsai, Russ R Salakhutdinov, and Geoffrey J Gordon, *Learning neural networks with adaptive regularization*, Advances in Neural Information Processing Systems, 2019, pp. 11389–11400.
- [ZWXG18] Guodong Zhang, Chaoqi Wang, Bowen Xu, and Roger Grosse, *Three mechanisms of weight decay regularization*, arXiv preprint arXiv:1810.12281 (2018).