

---

# Near-Interpolators: Fast Norm Growth and Tempered Near-Overfitting

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We study linear regression when the input data population covariance matrix has  
2 eigenvalues  $\lambda_i \sim i^{-\alpha}$  where  $\alpha > 1$ . Under a generic random matrix theory  
3 assumption, we prove that any near-interpolator, i.e.,  $\beta$  whose training error is  
4 below the noise floor, must have its squared  $\ell_2$ -norm growing super-linearly with  
5 the number of samples  $n$ :  $\|\beta\|_2^2 = \Omega(n^\alpha)$ . This implies that existing norm-based  
6 generalization bounds increase as the number of samples increases, matching the  
7 empirical observations from prior work. On the other hand, such near-interpolators  
8 when properly tuned achieve good generalization, where the test errors approach  
9 arbitrarily close to the noise floor. Our work demonstrates that existing norm-based  
10 generalization bounds are vacuous for explaining the generalization capability of  
11 *any* near-interpolators. Moreover, we show that the trade-off between train and test  
12 accuracy is better when the norm growth exponential is smaller.

## 13 1 Introduction

14 Learning algorithms that near-perfectly interpolate the training data such as deep neural networks have  
15 been surprisingly effective in practice despite conventional statistical wisdom suggesting otherwise  
16 [Zhang et al., 2021]. Near-interpolators arise frequently in modern machine learning, e.g., via early  
17 stopping rules [Ji et al., 2021, Kuzborskij and Szepesvári, 2022]. Therefore, understanding the  
18 fundamental trade-offs between near-interpolation and generalization is crucial.

19 Power law spectra assumptions arise commonly in popular settings such as in neural tangent kernels  
20 computed from practical networks. For instance, power law spectra of the neural tangent kernel matrix  
21 has been observed empirically in the MNIST, CIFAR-100 and CALTECH 101 datasets [Velikanov  
22 and Yarotsky, 2021, Wei et al., 2022, Murray et al., 2022]. Power law spectra assumptions provide  
23 a setting amenable to analysis while resembling real datasets, which has been used by previously  
24 Mallinar et al. [2022] to show that *perfect*-interpolators exhibit the so-called *tempered overfitting*  
25 phenomenon.

26 In this work, we analyze *near*-interpolators under power law spectra assumptions. Our result  
27 shows that such near-interpolators have norms increase super-linearly in the number of samples and  
28 exhibit tempered *near*-overfitting. Consequently, current norm-based generalization bounds are not  
29 applicable to explained this tempered near-overfitting behavior, and that tighter bounds are needed in  
30 the power law spectra assumption.

### 31 1.1 Our contributions

32 *Super-linear growth of the squared norm.* We show that when the data population covariance matrix  
33 has a power law spectra  $\lambda_i = i^{-\alpha}$  with exponent  $\alpha > 1$ , *near*-interpolators have squared norm  $\Omega(n^\alpha)$ .

34 In this setting, our work answers the question raised in the the “Discussion” section of Koehler et al.  
 35 [2021] regarding the growth of the norm for near-interpolators.

36 *Tempered near-overfitting of near-interpolators.* Tempered overfitting, coined by Mallinar et al.  
 37 [2022], refers to the situation when estimators perfectly interpolate the training data and achieve  
 38 test error  $c\sigma^2$  for some  $c \in (1, \infty)$ , i.e., proportional to the Bayes optimal error/noise floor  $\sigma^2$ .  
 39 Under the power law spectra  $\lambda_i = i^{-\alpha}$  condition where  $\alpha > 1$ , they show that the proportionality  
 40 constant  $c = \alpha$ . Under this same setting, we show that the *near-interpolators* achieve tempered  
 41 *near-overfitting*. More precisely, properly tuned ridge regression achieve proportionality constant  $c$   
 42 down to the benign regime where  $c = 1$ .

## 43 1.2 Related works

44 The main difference between our work and that of Mallinar et al. [2022] is that our work establishes  
 45 super-linear growth of the squared norm of near-interpolators. Our work is motivated by the empirical  
 46 evidence found by Wei et al. [2022] suggests that norms of kernel ridge regressors grow rapidly  
 47 potentially beyond the purview of norm-based bound. We confirm that bounds similar to the ones in  
 48 Koehler et al. [2021] are indeed vacuous for power-law spectra. Therefore, our work suggests that  
 49 explaining the generalization capability of near-interpolators will likely require new tools.

50 Ghosh and Belkin [2022] provides a lower bound on the *test error* for near-interpolators, demonstrat-  
 51 ing a fundamental trade-off between training and testing error. Our work derives a lower bound on  
 52 the *norm* for near-interpolators. Therefore, our work complements both Mallinar et al. [2022] and  
 53 Ghosh and Belkin [2022].

54 Our result is reminiscent of the result [Belkin et al., 2018, Theorem 1] in *classification*, which  
 55 establishes that the RKHS norm of a “near-interpolating” classifier grows at rate  $\Omega(\exp(n^{1/p}))$ . Note  
 56 that if the number of samples  $n = \Theta(\text{poly}(p))$ , then the lower bound does not grow to infinity and  
 57 thus is only meaningful when  $n = \Omega(\exp(p))$ . In contrast, our result is for *regression*. While our  
 58 results are not directly comparable, our lower bound is meaningful in the more practical  $n \propto p$   
 59 regime.

For more related works, see Appendix Section E.

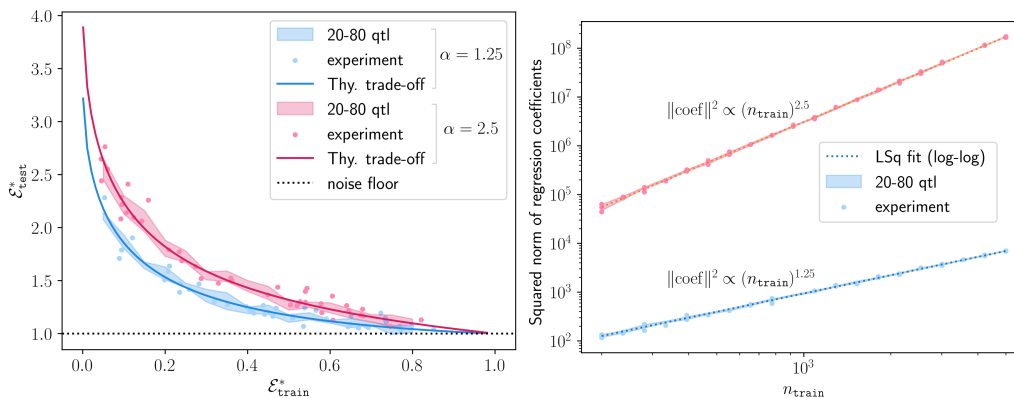


Figure 1: **Left:** Trade-off between the testing and training errors from Proposition 2.9. The solid lines are the parametrized curves  $(x, y) = (\mathcal{E}_{\text{train}}^*, \mathcal{E}_{\text{test}}^*)$  traced out by varying  $k$  (equivalently  $r$ ). The resulting estimators can achieve a continuum regimes of overfitting. The scatter points are empirical results from synthetic experiments on the HDA model (Example 2.4). The value for  $r$  are tuned according to the tuning scheme in Remark A.1 for prescribed training error  $\tau \approx \mathcal{E}_{\text{train}}^*$ . The parameters are  $n_{\text{train}} = n_{\text{test}} = 1000$ ,  $\gamma_* = 0.5$ ,  $\alpha \in \{1.25, 2.5\}$  and  $\sigma^2 = 1$ . See Appendix D for experimental details. **Right:** Synthetic experiments validating the norm lower bound given by Theorem 2.3. See Appendix D for additional experiment details. The squared norms are log-transformed then fitted by least squares to estimate the exponent  $\alpha$ . The estimated exponents matches the true  $\alpha$ 's. Note that the trade-off is better (left) when the corresponding norm growth exponent is smaller (right).

## 61 2 Main results

62 *Assumptions on the data distribution.* Let  $n$  denote the number of samples, treated as the fundamental  
 63 parameter. The feature dimension  $p$  is assumed to depend on  $n$  implicitly. The sample-to-feature  
 64 ratio is denoted  $\gamma := n/p \in (0, 1]$  and the asymptotic sample-to-feature ratio is denoted  $\gamma_* :=$   
 65  $\lim_{n \rightarrow \infty} \gamma \in [0, 1]$ . When  $\gamma_* = 0$ ,  $p$  grows much faster than  $n$ . Denote by  $X = [x_1, \dots, x_n] \in \mathbb{R}^{p \times n}$   
 66 the data matrix and  $y \in \mathbb{R}^n$  the training labels. Suppose that there exists a function  $f : \mathbb{R}^p \rightarrow \mathbb{R}$   
 67 (depending on  $n$ ) such that  $y_i = \varepsilon_i + f(x_i)$  where  $\varepsilon_i \in \mathbb{R}$  denote the noise. For instance, the  
 68 well-specified case corresponds to when  $f(x) = x^\top \beta^*$  for some  $\beta^* \in \mathbb{R}^p$ . Both  $y$  and  $\varepsilon$  are viewed  
 69 as column vectors.

70 *Assumptions on the noise.* Suppose that the noise are independent across samples, has zero mean  
 71  $0 = \mathbb{E}[\varepsilon_1]$  and variance  $\sigma^2 = \mathbb{E}[\varepsilon_1^2] > 0$ . For a positive integer  $p$ , let  $\mathbb{I}_p$  denote the  $p \times p$  identity  
 72 matrix. Thus we have  $\mathbb{E}[\varepsilon\varepsilon^\top] = \sigma^2\mathbb{I}_n$ . Moreover, suppose that  $\varepsilon \perp X$ , i.e., the noise and the data are  
 73 independent.

74 **Definition 2.1.** Ridge regression with regularizer  $\varrho > 0$  is the vector  $\hat{\beta}_\varrho$  defined via the optimization:

$$\hat{\beta}_\varrho := \operatorname{argmin}_{\beta \in \mathbb{R}^p} \frac{1}{n} \|X^\top \beta - y\|_2^2 + \varrho \|\beta\|_2^2. \quad (1)$$

75 Let  $\hat{\Sigma} := n^{-1} X X^\top$  denote the sample covariance matrix,  $\Sigma := \mathbb{E}[\hat{\Sigma}]$  the population covariance and  
 76  $\tilde{G} := n^{-1} X^\top X$  the (scaled) gram matrix.

### 77 2.1 Super-linear growth of the squared norm

78 Our main result is that the expected squared norm of the KRR with  $\varrho := r n^{-\alpha}$  regularizer grows at  
 79 least on the order of  $n^\alpha$  under suitable assumptions which we now introduce:

80 **Assumption 2.2.** Let  $\alpha > 1$ . The *exact eigenvalue decay (EVD)* condition with exponent  $\alpha$  assumes  
 81 that  $\Sigma = \operatorname{diag}(\lambda_1, \dots, \lambda_p)$  where  $\lambda_i = i^{-\alpha}$ .

82 Assumption 2.2 has been analyzed in many different context, most notably recently in being the  
 83 setting for the so-called *tempered overfitting* phenomenon [Mallinar et al., 2022]. See the related  
 84 works section for a detailed discussion.

85 **Theorem 2.3.** Assume that the exact EVD (Assumption 2.2) and certain random matrix-theoretic  
 86 conditions hold. Define regularizers  $\varrho := r n^{-\alpha}$  for the ridge regression (Definition 2.1) where  $r > 0$   
 87 is a positive number. Then, we have  $\mathbb{E}[\|\hat{\beta}_\varrho\|_2^2] = \Omega(n^\alpha)$ .

88 See Figure 1-Left for experimental validation of the lower bound. Below, we will use the term  
 89 “regularizer” to refer to both  $\varrho$  and  $r$  interchangeably.

90 The assumptions made in Theorem 2.3 are satisfied by the so-called HDA model, defined below. This  
 91 is proved later in Lemma B.5.

92 **Example 2.4.** Bai and Silverstein [2010], Dobriban and Wager [2018]. The following is sometimes  
 93 referred to as the *high-dimensional asymptotic (HDA)* model: 1.  $X = \Sigma^{1/2} Z$  where the entries of  
 94  $Z = \{Z_{ij}\} \in \mathbb{R}^{p \times n}$  are i.i.d, have zero mean  $\mathbb{E}[Z_{ij}] = 0$  and unit variance  $\mathbb{E}[Z_{ij}^2] = 1$ . The matrix  
 95  $\Sigma$  is positive semidefinite. 2.  $n/p \rightarrow \gamma_* \in (0, \infty)$ , and 3. Spectral distribution of  $\Sigma$  converges to a  
 96 distribution  $H$  supported on  $\mathbb{R}_{\geq 0}$ .

97 *Remark 2.5.* When the conditions of Theorem 2.3 are met, the expected norm  $\|\hat{\beta}_\varrho\|_2^2 = \Omega(n^\alpha)$ . The  
 98 current state-of-the-art uniform convergence generalization bound [Koehler et al., 2021, Corollary  
 99 1] are of the form  $\|\beta\|_2 / \sqrt{n}$  and are thus vacuous when  $\|\beta\|_2^2 = \Omega(n^\alpha)$  when  $\alpha > 1$ . We note that  
 100 the aforementioned results are for *perfect*-interpolators that achieve zero training error, rather than  
 101 near-interpolators. To our knowledge, no analogous theory for near-interpolators is known. Whether  
 102 the techniques of [Koehler et al., 2021] can be extended to explain near-interpolators is left as future  
 103 work.

104 While stated for the ridge regressor as in Definition 2.1, our lower bound holds for any  $\beta$  that is “as  
 105 good of an interpolator as  $\hat{\beta}_\varrho$ ”, i.e.,  $\beta$  has training error less than that of  $\hat{\beta}_\varrho$ .

106 **Definition 2.6.** Let  $\tau \geq 0$  arbitrary. The *minimum norm  $\tau$ -near-interpolator* is defined as

$$\beta_\tau := \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|\beta\|_2^2 \quad \text{s.t.} \quad \frac{1}{n} \|X^\top \beta - y\|_2^2 \leq \tau. \quad (2)$$

107 A  $\tau$ -near-interpolator is any  $\beta \in \mathbb{R}^p$  that is feasible for Equation (2).

108 **Proposition 2.7.** Let  $\rho > 0$  be arbitrary,  $\hat{\beta}_\rho \in \mathbb{R}^p$  be as in Definition 2.1, and  $\tau := \frac{1}{n} \|X^\top \hat{\beta}_\rho - y\|_2^2$ .  
 109 Consider  $\underline{\beta}_\tau$  as in Definition 2.6. Then  $\|\hat{\beta}_\rho\|_2 = \|\underline{\beta}_\tau\|_2$ . Consequently, if  $\beta \in \mathbb{R}^p$  has less training  
 110 error than  $\|\hat{\beta}_\rho\|_2$ , then  $\|\hat{\beta}_\rho\|_2 \leq \|\beta\|_2$ .

111 For the proof, see Appendix Section F.

## 112 2.2 Near-overfitting: Benign, tempered and everything in between

113 Simon et al. [2022] analyzed certain approximations of the testing and training errors of kernel ridge  
 114 regression. While these approximations, dubbed the *eigenlearning framework*, are non-rigorous  
 115 [Mallinar et al., 2022], they have been shown to be highly predictive in practice [Jacot et al., 2020,  
 116 Bordelon et al., 2020, Canatar et al., 2021].

117 Following Mallinar et al. [2022], we use the eigenlearning framework to calculate the training and  
 118 testing error of the estimators in Theorem 2.3 in terms of the *effective regularizer* [Wei et al., 2022],  
 119 denoted by  $k$ .

120 **Definition 2.8.** Let  $\alpha > 1$  and  $\gamma_* \in [0, \infty)$ . Define functions  $\mathcal{I}(\cdot) \equiv \mathcal{I}_{\alpha, \gamma_*}(\cdot)$  and  $\mathcal{J}(\cdot) \equiv \mathcal{J}_{\alpha, \gamma_*}(\cdot)$  as  
 121  $\mathcal{I}(k) := \int_0^{1/\gamma_*} \frac{dx}{1+kx^\alpha}$ , and  $\mathcal{J}(k) := \int_0^{1/\gamma_*} \frac{dx}{(1+kx^\alpha)^2}$ . When  $\gamma_* = 0$ , we assume that  $1/\gamma_* = +\infty$ .

122 Under Assumption 2.2, these functions from Definition 2.8 can be solved in closed-form given in  
 123 Appendix G. The reason we work with the effective regularizer  $k$  rather than the regularizer  $r$  is  
 124 that it is easier to calculate the approximations  $\mathcal{E}_{\text{test}}^*$ ,  $\mathcal{E}_{\text{train}}^*$  of the testing and training errors in the  
 125 eigenlearning framework:

126 **Proposition 2.9.** In the setting of Section 2, assume further that  $f$  is well-specified, i.e.,  $f(x) = x^\top \beta^*$   
 127 for some  $\beta^*$ . Moreover, suppose that  $\sup_{n=1,2,\dots} \|\beta^*\|_2 < +\infty$ . Assume the exact polynomial EVD  
 128 condition (Assumption 2.2) with exponent  $\alpha > 1$ . For the estimator in Theorem 2.3 we have

$$\mathcal{E}_{\text{test}}^* \equiv \lim_{n \rightarrow \infty} \mathcal{E}_{\text{test}} = \sigma^2 \cdot \frac{1}{1 - \mathcal{J}(k)}, \quad \text{and} \quad \mathcal{E}_{\text{train}}^* \equiv \lim_{n \rightarrow \infty} \mathcal{E}_{\text{train}} = \sigma^2 \cdot \frac{(1 - \mathcal{I}(k))^2}{1 - \mathcal{J}(k)}$$

129 Moreover, there exists  $k_{\text{crit}} \in \mathbb{R}_{>0}$  such that 1. For each  $r > 0$ , there exists a unique  $k \in (k_{\text{crit}}, +\infty)$   
 130 such that  $r = \mathcal{R}(k) := k(1 - \mathcal{I}(k))$ , 2.  $\mathcal{R}$  is monotonically increasing on  $(k_{\text{crit}}, +\infty)$ , 3.  $\mathcal{E}_{\text{test}}^* > \sigma^2$   
 131 for all  $k \in (k_{\text{crit}}, +\infty)$ , 4.  $\lim_{k \rightarrow +\infty} \mathcal{E}_{\text{test}}^* = \sigma^2$ , and 5.  $\frac{d}{d\alpha} \mathcal{E}_{\text{test}}^* > 0$  for any fixed  $k > 0$ .

132 For the proof of Proposition 2.9, see Appendix J. Thus,  $\mathcal{R}$  is a bijection that relates the effective  
 133 regularizer  $k$  and the (ordinary) regularizer  $r$ . Furthermore, note that  $\lim_{k \rightarrow +\infty} \mathcal{E}_{\text{test}}^* = \sigma^2$  precisely  
 134 states that the test error can be made arbitrarily close to the noise floor as  $k$  (equivalently,  $r$ ) goes to  
 135 infinity (See Proposition J.1 and Figure 2-Left).

136 *Remark 2.10* (Trade-off quality vs norm growth exponent). Note that item 5 of Proposition 2.9 makes  
 137 rigorous the observation that in Figure 1-left, the trade-off is better when the corresponding norm  
 138 growth exponent is smaller (see Figure 1-right).

## 139 3 Discussion and limitations

140 *Connection to early stopping.* Typically, early stopping prevents the trained algorithm from perfectly  
 141 interpolating the data. Can early stopped learning theory results, e.g., Ji et al. [2021], Kuzborskij and  
 142 Szepesvári [2022], be applied to analyze near-interpolators?

143 *Near-interpolators and uniform convergence generalization bound.* Is possible to use uniform  
 144 convergence-based approach to give non-vacuous generalization bound under the setting studied in  
 145 this work? This question has already been raised by Dobriban and Wager [2018] in the context of  
 146 classification.

147 **Limitations.** Our work is restricted to analyzing a random matrix model. Understanding the  
 148 phenomenon uncovered in this paper in more general models and real world settings will be needed.  
 149 Moreover, our work does not rule out the existence of uniform convergence generalization bound.

## References

- 150 **References**
- 151 Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining neural  
152 scaling laws. *arXiv preprint arXiv:2102.06701*, 2021.
- 153 Zhidong Bai and Jack W Silverstein. *Spectral analysis of large dimensional random matrices*,  
154 volume 20. Springer, 2010.
- 155 Zhidong Bai and Wang Zhou. Large sample covariance matrices without independence structures in  
156 columns. *Statistica Sinica*, pages 425–442, 2008.
- 157 Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear  
158 regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- 159 Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to understand  
160 kernel learning. In *International Conference on Machine Learning*, pages 541–549. PMLR, 2018.
- 161 Raphaël Berthier, Francis Bach, and Pierre Gaillard. Tight nonparametric convergence rates for  
162 stochastic gradient descent under the noiseless linear model. *Advances in Neural Information  
163 Processing Systems*, 33:2576–2586, 2020.
- 164 Alberto Bietti, Luca Venturi, and Joan Bruna. On the sample complexity of learning with geometric  
165 stability. In *Advances in Neural Information Processing Systems*, 2021.
- 166 Blake Bordelon, Abdulkadir Canatar, and Cengiz Pehlevan. Spectrum dependent learning curves in  
167 kernel regression and wide neural networks. In *International Conference on Machine Learning*,  
168 pages 1024–1034. PMLR, 2020.
- 169 Sébastien Bubeck and Mark Sellke. A universal law of robustness via isoperimetry. *Advances in  
170 Neural Information Processing Systems*, 34:28811–28822, 2021.
- 171 Abdulkadir Canatar, Blake Bordelon, and Cengiz Pehlevan. Spectral bias and task-model align-  
172 ment explain generalization in kernel regression and infinitely wide neural networks. *Nature  
173 communications*, 12(1):1–12, 2021.
- 174 Romain Couillet and Mérouane Debbah. Signal processing in large systems: A new paradigm. *IEEE  
175 Signal Processing Magazine*, 30(1):24–39, 2012.
- 176 Hugo Cui, Bruno Loureiro, Florent Krzakala, and Lenka Zdeborová. Generalization error rates  
177 in kernel regression: The crossover from the noiseless to noisy regime. In *Advances in Neural  
178 Information Processing Systems*, pages 10131–10143, 2021.
- 179 Edgar Dobriban and Stefan Wager. High-dimensional asymptotics of prediction: Ridge regression  
180 and classification. *The Annals of Statistics*, 46(1):247–279, 2018.
- 181 Nikhil Ghosh and Mikhail Belkin. A universal trade-off between the model size, test loss, and training  
182 loss of linear predictors. *arXiv preprint arXiv:2207.11621*, 2022.
- 183 Arthur Jacot, Berfin Simsek, Francesco Spadaro, Clément Hongler, and Franck Gabriel. Kernel  
184 alignment risk estimator: Risk prediction from training data. *Advances in Neural Information  
185 Processing Systems*, 33:15568–15578, 2020.
- 186 Ziwei Ji, Justin Li, and Matus Telgarsky. Early-stopped neural networks are consistent. *Advances in  
187 Neural Information Processing Systems*, 34:1805–1817, 2021.
- 188 Antti Knowles and Jun Yin. Anisotropic local laws for random matrices. *Probability Theory and  
189 Related Fields*, 169(1):257–352, 2017.
- 190 Frederic Koehler, Lijia Zhou, Danica J Sutherland, and Nathan Srebro. Uniform convergence of  
191 interpolators: Gaussian width, norm bounds and benign overfitting. *Advances in Neural Information  
192 Processing Systems*, 34:20657–20668, 2021.
- 193 Ilja Kuzborskij and Csaba Szepesvári. Learning lipschitz functions by gd-trained shallow overparam-  
194 eterized relu neural networks. *arXiv preprint arXiv:2212.13848*, 2022.

- 195 Neil Rohit Mallinar, James B Simon, Amirhesam Abedsoltan, Parthe Pandit, Misha Belkin, and  
196 Preetum Nakkiran. Benign, tempered, or catastrophic: Toward a refined taxonomy of overfitting.  
197 In *Advances in Neural Information Processing Systems*, 2022.
- 198 Michael Murray, Hui Jin, Benjamin Bowman, and Guido Montufar. Characterizing the spectrum of  
199 the NTK via a power series expansion. *arXiv preprint arXiv:2211.07844*, 2022.
- 200 Courtney Paquette, Kiwon Lee, Fabian Pedregosa, and Elliot Paquette. Sgd in the large: Average-case  
201 analysis, asymptotics, and stepsize criticality. In *Conference on Learning Theory*, pages 3548–3626.  
202 PMLR, 2021.
- 203 Courtney Paquette, Bart van Merriënboer, Elliot Paquette, and Fabian Pedregosa. Halting time is  
204 predictable for large models: A universality property and average-case analysis. *Foundations of*  
205 *Computational Mathematics*, pages 1–77, 2022.
- 206 Jack W Silverstein and Sang-Il Choi. Analysis of the limiting spectral distribution of large dimensional  
207 random matrices. *Journal of Multivariate Analysis*, 54(2):295–309, 1995.
- 208 James B Simon, Madeline Dickens, Dhruva Karkada, and Michael R DeWeese. The eigenlearning  
209 framework: A conservation law perspective on kernel regression and wide neural networks. *arXiv*  
210 *preprint arXiv:2110.03922*, 2022.
- 211 Terrence Tao. Intuitive understanding of the Stieltjes transform. MathOverflow, 2011. URL  
212 <https://mathoverflow.net/q/79129>. Version: 2011-10-25.
- 213 Alexander Tsigler and Peter L Bartlett. Benign overfitting in ridge regression. *arXiv preprint*  
214 *arXiv:2009.14286*, 2020.
- 215 Maksim Velikanov and Dmitry Yarotsky. Explicit loss asymptotics in the gradient descent training of  
216 neural networks. *Advances in Neural Information Processing Systems*, 34:2570–2582, 2021.
- 217 Maksim Velikanov and Dmitry Yarotsky. Tight convergence rate bounds for optimization under  
218 power law spectral conditions. *arXiv preprint arXiv:2202.00992*, 2022.
- 219 Alexander Wei, Wei Hu, and Jacob Steinhardt. More than a toy: Random matrix models predict how  
220 real-world neural representations generalize. In *Proceedings of the 39th International Conference*  
221 *on Machine Learning*, pages 23549–23588. PMLR, 2022.
- 222 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep  
223 learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115,  
224 2021.
- 225 Hongyang Zhang, Yihan Wu, and Heng Huang. How many data are needed for robust learning?  
226 *arXiv preprint arXiv:2202.11592*, 2022.

227 **A Additional discussion on implications of Proposition 2.9**

228 Using Proposition 2.9, we illustrate the trade-off between the training error versus the testing error  
 229 in Figure 1-Right using closed-form expression for  $\mathcal{E}_{\text{train}}^*$  and  $\mathcal{E}_{\text{test}}^*$  are presented in Appendix G.  
 230 Figure 2-Left demonstrates that empirical training and test errors from synthetic experiments concentrate  
 231 around the theoretical values  $(\mathcal{E}_{\text{train}}^*, \mathcal{E}_{\text{test}}^*)$  with growing  $n$ .

232 *Remark A.1* (Tuning the regularizer). Proposition 2.9 allows for tuning the ridge parameter  $\varrho := rn^\alpha$   
 233 to achieve a user-specified value of training error  $\tau$  via the following procedure: First, use a binary  
 234 search algorithm to find  $k_\tau$  such that  $\tau = \mathcal{E}_{\text{train}}^*$ . Next, set  $r := \mathcal{R}(k_\tau)$ . Finally, set  $\varrho := rn^\alpha$ .

235 *Remark A.2.* The upshot of Proposition 2.9 is that any trade-off  $(\mathcal{E}_{\text{train}}^*, \mathcal{E}_{\text{test}}^*)$  on along the blue curve  
 236 in Figure 1-Right can be achieved by the tuning algorithm in Remark A.1. For perfect-interpolators,  
 237 Mallinar et al. [2022] shows that estimators with tempered overfitting achieve test error of exactly  
 238  $\alpha\sigma^2$ . In contrast, *near*-overfitting can achieve a continuum of test errors, i.e.,  $c\sigma^2$  where  $c \in (1, c_{\text{max}})$   
 239 belongs to an interval.

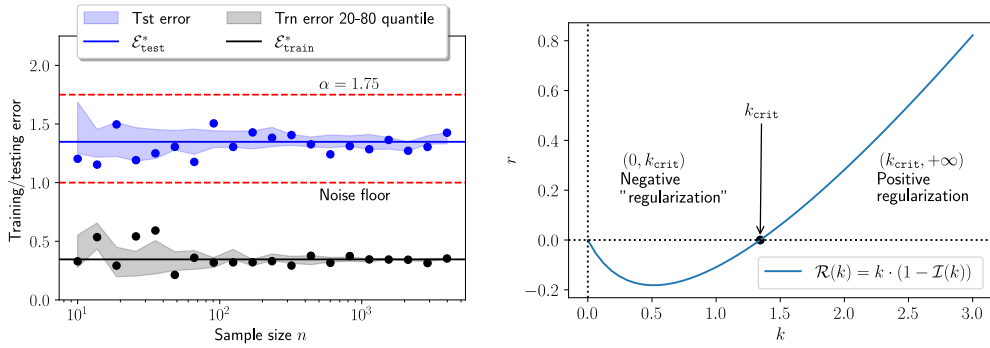


Figure 2: **Left:** Synthetic experiment validating the approximations given by Proposition 2.9 using the same setup as in Figure 1. By setting  $r \approx 3.54$ , we get a test error of  $\approx 1.35$  which is significantly below the tempered overfitting test error of  $\alpha = 1.75$  in [Mallinar et al., 2022, Theorem 3.1]. See Figure 1 and Appendix D for experimental details. **Right:** The  $\mathcal{R}(k)$  function from Proposition 2.9. The  $x$ -axis is the input  $k$ . Note that for  $k < k_{\text{crit}}$  the regularizer  $r$  is negative. Although we are only interested in the  $(k_{\text{crit}}, +\infty)$  portion, negative regularizers have been studied by Tsigler and Bartlett [2020] in the context of benign overfitting.

240 **B Random Matrix Theory and Assumptions**

241 In this section, we review and define the random matrix theory-based assumptions used to establish  
 242 our results. These assumptions, while seemingly restrictive, are common in random matrix theory  
 243 and showing their universality is an ongoing research area. See Remark B.11.

244 For  $c \in \mathbb{R}$ , let  $\delta_c$  denote the *Dirac-delta measure* on  $\mathbb{R}$  at  $c$ . In other words, for a Borel-measurable  
 245 set  $E \subseteq \mathbb{R}$ , we have  $\delta_c(E) = 1$  if  $c \in E$  and  $\delta_c(E) = 0$  otherwise.

246 **Definition B.1** (Empirical spectral measure). Let  $M \in \mathbb{R}^{p \times p}$  be a matrix with real eigenvalues  
 247  $\lambda_1, \dots, \lambda_p$ . The *empirical spectral measure* of  $M$ , denoted by  $\text{esd}(M)$ , is the measure on  $\mathbb{R}$  given  
 248 by  $\text{esd}(M) = \frac{1}{p} \sum_{i=1}^p \delta_{\lambda_i}$ .

249 We now state the weaker eigenvalue decay assumption sufficient for Theorem 2.3:

250 **Assumption B.2** (Asymptotic EVD). Let  $\alpha > 1$ . Suppose that  $\text{esd}(n^\alpha \Sigma)$  converges to a distribution  
 251  $H$  on  $\mathbb{R}_{\geq 0}$ .

252 In Proposition C.7, we show that Assumption B.2 generalizes the earlier Assumption 2.2.

253 Random matrix theory are primarily concerned with analysis of the spectra of large random matrices.  
 254 A key analytic tool is the *Stieltjes transform* of the empirical spectral measures of matrices:

255 **Definition B.3.** Let  $\mu$  be a measure on  $\mathbb{R}$ . The *Stieltjes transform* of  $\mu$  is the (complex-valued)  
 256 function with input  $z \in \mathbb{C}$  given by  $\mathcal{S}_\mu(z) := \int \frac{\mu(t)dt}{t-z}$ .

257 See Bai and Silverstein [2010, Appendix B.2] for reference.

258 For a matrix  $M \in \mathbb{R}^{p \times p}$  with  $p$  real eigenvalues (e.g., when  $M$  is real and symmetric), the following  
259 holds:

$$\mathcal{S}_{\text{esd}(M)}(z) = p^{-1} \text{tr}((M - z\mathbb{I}_p)^{-1}).$$

260 We now state the other assumption made in Theorem 2.3:

261 **Assumption B.4** (Positivity condition). For every  $r > 0$ , suppose that  
262  $\lim_{n \rightarrow \infty} \mathbb{E} \left[ \frac{d}{dr} (r \mathcal{S}_{\text{esd}(n^\alpha \check{G})}(-r)) \right] > 0$ .

263 By leveraging the results of Silverstein and Choi [1995], we prove in Appendix C.1 the following:

264 **Lemma B.5.** Under the HDA model (Example 2.4) and the EVD condition (Assumption B.2), we  
265 have that  $\lim_{n \rightarrow \infty} \mathbb{E} \left[ \frac{d}{dr} (r \mathcal{S}_{\text{esd}(n^\alpha \check{G})}(-r)) \right] > 0$ .

266 Next, we state what is sometimes referred to as the *self-consistent equation* [Tao, 2011]:

267 **Assumption B.6.** For each  $r > 0$ , there exists a unique  $k \equiv k(r) \in \mathbb{R}$  such that the limit exists,

$$\tilde{\mathcal{I}}(k) := \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^p \frac{1}{1 + kn^{-\alpha} \lambda_i^{-1}} \in \mathbb{R},$$

268 and the tuple  $(r, k)$  satisfies the equation

$$1 = \frac{r}{k} + \tilde{\mathcal{I}}(k). \quad (3)$$

269 *Remark B.7.* The functional relationship between  $r$  and  $k$  can be computed explicitly under the exact  
270 eigenvalue decay condition. As we will see in the proof of Proposition J.1, the expression  $\tilde{\mathcal{I}}$  coincide  
271 with  $\mathcal{I}$  from Definition 2.8.

272 Next, we state a version of the classical Marchenko-Pastur law for a random matrix ensemble  $X$ .

273 **Assumption B.8** (Marchenko-Pastur law). In the setting of Assumption B.6, further assume that  
274 almost surely

$$\lim_{n \rightarrow \infty} r \mathcal{S}_{\text{esd}(n^\alpha \check{G})}(-r) = k \mathcal{S}_H(-k)$$

275 and  $\lim_{n \rightarrow \infty} \frac{d}{dr} (r \mathcal{S}_{\text{esd}(n^\alpha \check{G})}(-r)) = \frac{d}{dr} (k \mathcal{S}_H(-k))$ . We note that the  $k$  on the RHS depends on  $r$ .

276 *Remark B.9.* While we assume that the data is generated from the HDA model  $X = \Sigma Z$ , we note  
277 that, when  $\Sigma = \text{diag}(\{\lambda_i = i^{-\alpha}\})$  (Assumption 2.2), the empirical spectral measure of the *scaled*  
278 covariance  $n^\alpha \Sigma$  converges to a limiting distribution. On the other hand, than the unscaled  $\Sigma$  matrix  
279 does not. Thus, the above Assumption B.8 reduces to the standard Marchenko-Pastur law when we  
280 consider the ‘‘fictitious’’ scaled HDA model  $n^{\alpha/2} \sqrt{\Sigma} Z$  which is used in the analysis. The scaling of  
281 the regularizer  $\varrho = rn^{-\alpha}$  in Definition 2.1 is chosen specifically to allow us to shift our analysis to  
282 this ‘‘fictitious’’ scaled HDA model.

283 The following is well-known [Dobriban and Wager, 2018]:

284 **Theorem B.10** (Marchenko-Pastur theorem). Under Example 2.4, both Assumption B.6 and As-  
285 sumption B.8 hold.

286 *Remark B.11.* Many works have demonstrated these so-called universality phenomena for a broad  
287 range of random matrix ensemble beyond the simple HDA model. For instance, the Marchenko-Pastur  
288 law (Assumption B.8) and their variants has been extended to the setting where certain independence  
289 assumptions are dropped [Bai and Zhou, 2008] and when  $\gamma_* = \lim_{n \rightarrow \infty} n/p = 0$  [Knowles and Yin,  
290 2017, Wei et al., 2022]. As such, we expect Assumption B.4 to hold in these broader contexts as well.  
291 We leave this as an important future direction.

292 Having introduced the necessary assumptions, we now turn to proving Theorem 2.3.

## 293 C Norm lower bound in RMT settings

294 The goal of this section is to sketch the proof for Theorem 2.3. Complete proofs of all results are  
295 included in the Appendix. Throughout, we assume the setting of Section 2. The first key technical  
296 step the following:



297 **Proposition C.1.**  $\mathbb{E}\|\hat{\beta}_\varrho\|_2^2 \geq n^{-1}\sigma^2\mathbb{E}[\text{tr}((\hat{\Sigma} + \varrho\mathbb{I}_p)^{-2}\hat{\Sigma})]$ .

298 *Proof sketch of Proposition C.1.* We first simplify  $\|\hat{\beta}_\varrho\|_2^2$  using the well-known formula for ridge  
299 regression:

300 **Lemma C.2.** The closed-form solution for Equation (1) is given by the formula  $\hat{\beta}_\varrho := (\hat{\Sigma} +$   
301  $\varrho\mathbb{I}_p)^{-1}\frac{1}{n}Xy$ .

302 Next, let  $M := (\hat{\Sigma} + \varrho\mathbb{I}_p)^{-1}\frac{1}{n}X$ . Using the independence of  $X$  and  $\varepsilon$ , we get  $\mathbb{E}[\|\hat{\beta}_\varrho\|_2^2] \geq$   
303  $\mathbb{E}[\text{tr}(M^\top M\varepsilon\varepsilon^\top)]$ . Since  $M^\top M$  and  $\varepsilon\varepsilon^\top$  are also independent, we have

$$\mathbb{E}[\text{tr}(M^\top M\varepsilon\varepsilon^\top)] = \sigma^2\mathbb{E}[\text{tr}(M^\top M)].$$

304 By  $M^\top M = \frac{1}{n}(\hat{\Sigma} + \varrho\mathbb{I}_p)^{-1}\hat{\Sigma}(\hat{\Sigma} + \varrho\mathbb{I}_p)^{-1}$  and the cyclic property of trace, we get the desired  
305 inequality.  $\square$

306 The next step towards proving Theorem 2.3 is the following:

307 **Proposition C.3.** Let  $\varrho := rn^{-\alpha}$ . Then we have  $\mathbb{E}\|\hat{\beta}_\varrho\|_2^2 \geq n^\alpha\sigma^2 \cdot \mathbb{E}\left[\frac{d}{dr}(r\mathcal{S}_{\text{esd}(n^\alpha\check{G})}(-r))\right]$ .

308 *Proof sketch of Proposition C.3.* We first relate the quantity  $\text{tr}((\hat{\Sigma} + \varrho\mathbb{I}_p)^{-2}\hat{\Sigma})$  inside the lower bound  
309 in Proposition C.1 to  $\mathcal{S}_{\text{esd}(n^\alpha\hat{\Sigma})}$ , the Stieltjes transform of  $n^\alpha\hat{\Sigma}$ :

310 **Lemma C.4.** Let  $M \in \mathbb{R}^{p \times p}$  be any symmetric matrix and  $z \in \mathbb{R}$ . Then we have

$$\frac{d}{dz}\text{tr}(z(M + z\mathbb{I}_p)^{-1}) = \text{tr}(M(M + z\mathbb{I}_p)^{-2}).$$

311 Next, we use the following well-known result for relating  $\mathcal{S}_{\text{esd}(n^\alpha\hat{\Sigma})}$  and  $\mathcal{S}_{\text{esd}(c\check{G})}$ . For the sake of  
312 completeness, we include the proof in the Appendix.

313 **Lemma C.5** (Gram-to-covariance). Let  $c \in \mathbb{R}$  and  $z \in \mathbb{C}$  be arbitrary, then  $\mathcal{S}_{\text{esd}(c\hat{\Sigma})}(z) = \gamma \cdot$   
314  $\mathcal{S}_{\text{esd}(c\check{G})}(z) - \frac{(1-\gamma)}{z}$ .

315 Using Proposition C.1 and the two preceding Lemmas, the desired inequality follows from algebraic  
316 manipulation.  $\square$

317 Given the lower bound in Proposition C.3, our goal now is to relate the random quantity  $\mathcal{S}_{\text{esd}(n^\alpha\hat{\Sigma})}(\cdot)$   
318 with the deterministic quantity  $\mathcal{S}_{\text{esd}(n^\alpha\Sigma)}(\cdot)$  using random matrix theory. Later, we will see that  
319 a consequence of Proposition C.7 is that  $\mathbb{E}\left[\frac{d}{dr}(r\mathcal{S}_{\text{esd}(n^\alpha\hat{\Sigma})}(-r))\right]$  is positive. This implies that  
320  $\mathbb{E}[\|\hat{\beta}_\varrho\|_2^2] \geq o(n^\alpha)$ . We now conclude with the proof of Theorem 2.3.

321 *Proof of Theorem 2.3.* Let  $L := \lim_{n \rightarrow \infty} \mathbb{E}\left[\frac{d}{dr}(r\mathcal{S}_{\text{esd}(n^\alpha\check{G})}(-r))\right] > 0$  be as in Assumption B.4.  
322 Thus, for all  $n \gg 0$  sufficiently large, we have  $\mathbb{E}\left[\frac{d}{dr}(r\mathcal{S}_{\text{esd}(n^\alpha\check{G})}(-r))\right] > L/2 > 0$ . By Proposi-  
323 tion C.3, we get that  $\mathbb{E}\|\hat{\beta}_\varrho\|_2^2 \geq n^\alpha\sigma^2 \cdot \frac{L}{2}$  for all  $n \gg 0$ , as desired.  $\square$

### 324 C.1 Positivity condition for the HDA model

325 This section will focus on the proof of Lemma B.5. Thus, throughout this section, we assume the  
326 setting of Example 2.4. Using the Marchenko-Pastur law (Assumption B.8) and calculus, we first  
327 show that

$$\lim_{n \rightarrow \infty} \mathbb{E}\left[\frac{d}{dr}(r\mathcal{S}_{\text{esd}(n^\alpha\check{G})}(-r))\right] = \left(\frac{dr}{dk}\right)^{-1} \cdot \frac{d}{dk}(k\mathcal{S}_H(-k))$$

328 where  $r$  and  $k$  are as in Assumption B.6. Thus, we reduce to showing the positivity of  $\frac{dr}{dk}$  and  
329  $\frac{d}{dk}(k\mathcal{S}_H(-k))$ . See Appendix I.  $\square$

330 **C.2 Convergence to limiting distribution**

331 It remains to check that the exact eigenvalue decaying assumption (Assumption 2.2) indeed satisfy  
 332 the condition 3 of Example 2.4.

333 **Definition C.6.** Given a measure  $\mu$  on  $\mathbb{R}$ , we let  $\text{cdf}[\mu]$  denote the cumulative distribution function  
 334 of  $\mu$ .

335 **Proposition C.7.** Under Assumption 2.2, we have the following:

$$\lim_{n \rightarrow \infty} \text{cdf}[\text{esd}(n^\alpha \Sigma)](t) = \begin{cases} 1 - \gamma_* t^{-1/\alpha} & : t \geq \gamma_*^\alpha \\ 0 & \text{otherwise.} \end{cases}$$

336 *Proof of Proposition C.7.* The set of eigenvalues of  $n^\alpha \Sigma$

$$\{(n/i)^\alpha\}_{i=1,\dots,p} = \underbrace{\left\{ \left(\frac{n}{p}\right)^\alpha, \dots, \left(\frac{n}{n+1}\right)^\alpha \right\}}_{=\gamma^\alpha}, \underbrace{\frac{n}{n}}_{=1}, \underbrace{\left\{ \left(\frac{n}{n-1}\right)^\alpha, \dots, \left(\frac{n}{1}\right)^\alpha \right\}}_{=n^\alpha}.$$

337 Thus,  $\text{cdf}[\text{esd}(n^\alpha \Sigma)](t) = 0$  if  $t < \gamma^\alpha$  and  $= 1$  if  $t > n^\alpha$ .

338 Below, let  $t \in [\gamma^\alpha, n^\alpha]$  and  $j(t) \in \{1, \dots, p\}$  be the index such that  $t \approx (n/j(t))^\alpha$  is as close as  
 339 possible. Solving for  $j(t)$ , we have  $j(t) \approx nt^{-1/\alpha}$ . Thus, there are (approximately)  $p - j(t)$  indices,  
 340 denoted by  $i$ , such that  $(n/i)^\alpha < (n/j(t))^\alpha$ . Divide by  $p$ , we get the relative frequency of such  
 341 indices  $i$ , which is  $\text{cdf}[\text{esd}(n^\alpha \Sigma)](t) = 1 - (j(t)/p) \approx 1 - \gamma t^{-1/\alpha}$ . This approximation becomes  
 342 exact as  $n \rightarrow \infty$ .  $\square$

343 **D Experiment**

344 We run experiment with  $\alpha = 1.75$  and  $n/p = \gamma = 0.5$ . We sample  $\beta^* \in \mathbb{R}^p$  such that  $\beta_i^*$  are  
 345 i.i.d Gaussian with zero mean and variance  $= 10/p$ . For the data, we sample  $X = \sqrt{\Sigma}Z$  as in the  
 346 HDA model Example 2.4 where  $Z_{ij}$  are i.i.d standard Gaussian random variables and  $\Sigma$  is as in  
 347 Assumption 2.2. The same set up is used for Figure 1. All code for the experiments are included in  
 348 Appendix K.

349 **E Expanded related works**

350 *Trade-offs in interpolation-based learning.* In addition to Mallinar et al. [2022], Ghosh and Belkin  
 351 [2022], Belkin et al. [2018], previous works have also studied the fundamental trade-off in learning  
 352 algorithms between overparametrization and (Lipschitz) smoothness [Bubeck and Sellke, 2021]  
 353 robustness and smoothness [Zhang et al., 2022].

354 *Power law spectra.* Many works reviewed in this section study the eigenvalues of kernel/gram  
 355 matrices, while we are primarily interested in the covariance matrix spectra. However, we note that  
 356 the covariance matrix have the same eigenvalues. Thus, results regarding the spectra applies to both  
 357 kernel/gram and covariance matrices. Below, we will review works in this area using the term used  
 358 by the original authors.

359 *Power-law spectra datasets.* Synthetic data with artificial power law EVD covariance have been used  
 360 frequently as toy examples [Berthier et al., 2020, Mallinar et al., 2022]. On real datasets, power  
 361 law EVD is often observed to describe neural tangent kernels (NTK) well in practice, including on  
 362 MNIST ([Bahri et al., 2021, Fig. 4] and [Velikanov and Yarotsky, 2022, Fig. 2]), FASHION-MNIST  
 363 [Cui et al., 2021, Fig. 7] CALTECH 101 [Murray et al., 2022, Fig. 1], CIFAR-100 [Wei et al., 2022,  
 364 Fig. 3].

365 *Theoretical machine learning works using power-law spectra.* Bordelon et al. [2020] shows that  
 366 power law EVD implies power law learning curve. Velikanov and Yarotsky [2021, §6.2] computes  
 367 the power law EVD exponent for certain NTKs with ReLU to be  $\alpha = 1 + \frac{1}{d}$ . Murray et al. [2022]  
 368 computes the EVD for NTKs with several different activations. The EVD condition is also known  
 369 as the *capacity condition* in the kernel ridge regression literature. See Bietti et al. [2021] and the  
 370 references there-in.

371 Bartlett et al. [2020, Theorem 6] shows that benign overfitting occurs when the covariance matrix  
372 eigenvalues  $\lambda_i = i^{-1} \log^{-b}(i+1)$  for  $b > 1$ . Mallinar et al. [2022] studies power law decay for  
373  $\alpha \geq 1$  and proposes a taxonomy of overfitting into three categories: catastrophic, tempered and benign.  
374 *Random matrix theory (RMT)*. The signal processing research community have long been using RMT  
375 for theoretical analysis [Couillet and Debbah, 2012]. Increasingly RMT has been applied to machine  
376 learning as well as a key tool for analysis. In particular, Dobriban and Wager [2018], Jacot et al.  
377 [2020] have applied RMT for (kernel) ridge regression. Paquette et al. [2022, 2021] uses the so-called  
378 local Marchenko-Pastur law [Knowles and Yin, 2017] to analyze gradient-based algorithms. Wei  
379 et al. [2022] also applies such local law to analyze the so-called *generalized cross-validation (GCV)*  
380 *estimator*.

## 381 F Proof for Proposition 2.7

382 *Proof of Proposition 2.7.* By definition,  $\hat{\beta}_\varrho$  is feasible for the optimization in Equation (2) and thus  
383  $\|\hat{\beta}_\varrho\|_2 \geq \|\underline{\beta}_\tau\|_2$ . Now, suppose for the sake of contradiction that  $\|\underline{\beta}_\tau\|_2 < \|\hat{\beta}_\rho\|_2$ . Then we have

$$\begin{aligned} & \varrho \|\underline{\beta}_\tau\|_2^2 + \frac{1}{n} \|X \underline{\beta}_\tau - y\|_2^2 \\ & \leq \varrho \|\underline{\beta}_\tau\|_2^2 + \tau \quad \because \underline{\beta}_\tau \text{ is feasible for Equation (2)} \\ & < \varrho \|\hat{\beta}_\varrho\|_2^2 + \tau \quad \because \text{assumption } \|\underline{\beta}_\tau\|_2 < \|\hat{\beta}_\rho\|_2 \\ & = \varrho \|\hat{\beta}_\varrho\|_2^2 + \frac{1}{n} \|X \hat{\beta}_\varrho - y\|_2^2 \quad \because \text{Definition of } \tau \end{aligned}$$

384 This contradicts the minimality of  $\hat{\beta}_\rho$  for Equation (1). Thus, we've shown that  $\|\underline{\beta}_\tau\|_2 = \|\hat{\beta}_\rho\|_2$ .  $\square$

## 385 G Closed-form expression for Proposition 2.9

386 Let  ${}_2F_1(a, b; c; z)$  be the *Gauss hypergeometric function*, implemented in SCIPY as  
387 `scipy.special.hyp2f1`.

388 **Lemma G.1.** Let  $\alpha > 1$  and  $\gamma_* \in \mathbb{R}_{\geq 0}$  be fixed. The functions  $\mathcal{I}, \mathcal{J}$  from Definition 2.8 can be  
389 written in closed form as:

$$\begin{aligned} \mathcal{I}(k) &= \gamma_*^{-1} \times {}_2F_1(1, 1/\alpha; 1 + 1/\alpha; -k\gamma_*^{-\alpha}) \\ \mathcal{J}(k) &= \gamma_*^{-1} \times {}_2F_1(2, 1/\alpha; 1 + 1/\alpha; -k\gamma_*^{-\alpha}). \end{aligned}$$

391 When  $\gamma_* = 0$ , we have

$$\begin{aligned} \mathcal{I}(k) &= \frac{\pi}{\alpha} k^{-1/\alpha} \csc(\pi/\alpha) \\ \mathcal{J}(k) &= \frac{\pi(\alpha-1)}{\alpha^2} k^{-1/\alpha} \csc(\pi/\alpha). \end{aligned}$$

393 The above expressions can be obtained using computer algebra softwares such as MATHEMATICA.  
394 Note that the expression in the  $\gamma_* = 0$  case has appeared previously in [Mallinar et al., 2022,  
395 Eqn. (22)] in a similar context. To the best of our knowledge, the expressions in the  $\gamma_* \neq 0$  case are  
396 new, at least in the machine learning literature.

## 397 H Proofs for supporting lemmas of Theorem 2.3

398 *Proof of Lemma C.2.* Start with the objective function  $\mathcal{F}(\beta) := \frac{1}{n} \|X^\top \beta - y\|_2^2 + \varrho \|\beta\|_2^2$ . Take  
399 derivative with respect to  $\beta$ , we have

$$\frac{1}{2} \nabla_\beta \left( \frac{1}{n} \|X^\top \beta - y\|_2^2 + \varrho \|\beta\|_2^2 \right) = \frac{1}{2} \nabla_\beta \left( \beta^\top (\hat{\Sigma} + \varrho \mathbb{I}_p) \beta - \frac{2}{n} \beta^\top X y \right) = (\hat{\Sigma} + \varrho \mathbb{I}_p) \beta - \frac{1}{n} X y.$$

400 Since  $\nabla_\beta \mathcal{F}(\hat{\beta}_\varrho) = 0$ , we are done.  $\square$

401 **Lemma H.1** (Special case of Woodbury formula). Let  $M \in \mathbb{R}^{p \times n}$  be an arbitrary matrix and  
 402  $\varrho \in (0, \infty)$ . Then

$$(MM^\top + \varrho \mathbb{I}_p)^{-1}M = M(M^\top M + \varrho \mathbb{I}_n)^{-1} \in \mathbb{R}^{n \times p}$$

403 *Proof of Lemma H.1.* It suffices to prove Lemma H.1 for the special case when  $\varrho = 1$ , which we  
 404 assume below. By the Woodbury matrix identity, we have

$$(MM^\top + \mathbb{I}_p)^{-1} = \mathbb{I} - M(M^\top M + \mathbb{I}_n)^{-1}M^\top \quad (4)$$

405 For brevity, let  $P := MM^\top + \mathbb{I}_p$  and let  $N := M^\top M + \mathbb{I}_n$ . To proceed, we have

$$\begin{aligned} P^{-1}M &= M - MN^{-1}M^\top M \quad \because \text{Multiplying (4) by } M \text{ on the right} \\ &= M(\mathbb{I}_n - N^{-1}M^\top M) \quad \because \text{Factoring out } M \text{ on the left} \\ &= M(\mathbb{I}_n - (\mathbb{I}_n - N^{-1})) \quad \because \mathbb{I}_n = N^{-1}N = N^{-1} + N^{-1}M^\top M \\ &= MN^{-1} \end{aligned}$$

406 as desired.  $\square$

407 *Proof of Lemma C.4.* Without the loss of generality, suppose that  $M = \text{diag}(\lambda_1, \dots, \lambda_p)$ . Then we  
 408 have  $f(z) := \text{tr}(z(M + zI_p)^{-1}) = \sum_{i=1}^p \frac{z}{\lambda_i + z}$ . Now, from elementary calculus, we have

$$\frac{d}{dx} \frac{x}{y+x} = (y+x)^{-1} - x(y+x)^{-2} = (y+x)^{-2}((y+x) - x) = \frac{y}{(y+x)^2}.$$

409 From this, we recover the fact that  $\frac{d}{dz} f(z) = \sum_{i=1}^n \frac{\lambda_i}{(\lambda_i + z)^2} = \text{tr}(M(M + z\mathbb{I}_p)^{-2})$ , as desired.  $\square$

410 *Proof of Lemma C.5.* Without the loss of generality, we may assume that  $c = 1$ . Let  $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p$   
 411 be the eigenvalues of  $\hat{\Sigma}$ . Since  $p > n$ , we necessarily have that  $\hat{\lambda}_{n+1} = \dots = \hat{\lambda}_p = 0$ . Moreover,  
 412  $\hat{\lambda}_1, \dots, \hat{\lambda}_n$  are the eigenvalues of  $\check{G}$ . Now, unwinding the definition, we have

$$\mathcal{S}_{\text{esd}(\hat{\Sigma})}(z) = \frac{1}{p} \sum_{i=1}^p \frac{1}{\hat{\lambda}_i - z}$$

413 and

$$\mathcal{S}_{\text{esd}(\check{G})}(z) = \frac{1}{n} \sum_{i=1}^n \frac{1}{\hat{\lambda}_i - z}.$$

414 Thus,

$$\begin{aligned} \mathcal{S}_{\text{esd}(\hat{\Sigma})}(z) &= \frac{1}{p} \left( \sum_{i=1}^n \frac{1}{\hat{\lambda}_i - z} + \sum_{i=n+1}^p \frac{1}{-z} \right) \\ &= \left( \frac{n}{p} \frac{1}{n} \sum_{i=1}^n \frac{1}{\hat{\lambda}_i - z} \right) - \frac{p-n}{p} \frac{1}{z} \\ &= \gamma \cdot \mathcal{S}_{\text{esd}(\check{G})}(z) - \frac{(1-\gamma)}{z} \end{aligned}$$

415 as desired.  $\square$

416 *Proof of Proposition C.1.* Below, for brevity we let  $a := f(X)$  and  $M := (\hat{\Sigma} + \varrho \mathbb{I}_p)^{-1} \frac{1}{n} X$ . We  
 417 recall from the previous lemma that

$$\hat{\beta}_\varrho = (\hat{\Sigma} + \varrho \mathbb{I}_p)^{-1} \frac{1}{n} X y = (\hat{\Sigma} + \varrho \mathbb{I}_p)^{-1} \frac{1}{n} X (f(X) + \varepsilon) = M(a + \varepsilon).$$

418 Thus,

$$\|\hat{\beta}_\varrho\|_2^2 = (a + \varepsilon)^\top M^\top M(a + \varepsilon) \geq \varepsilon^\top M^\top M\varepsilon + 2\varepsilon^\top M^\top Ma$$

419 Note that  $\varepsilon \perp M^\top Ma$  since  $\varepsilon \perp X$ . Thus, since  $\mathbb{E}[\varepsilon] = 0$ , we have

$$\mathbb{E}[\|\hat{\beta}_\varrho\|_2^2] = \mathbb{E}[(a + \varepsilon)^\top M^\top M(a + \varepsilon)] \geq \mathbb{E}[\varepsilon^\top M^\top M\varepsilon] = \mathbb{E}[\text{tr}(M^\top M\varepsilon\varepsilon^\top)]$$

420 Since  $M^\top M \perp \varepsilon\varepsilon^\top$ , we have

$$\mathbb{E}[\text{tr}(M^\top M\varepsilon\varepsilon^\top)] = \text{tr}(\mathbb{E}[M^\top M]\mathbb{E}[\varepsilon\varepsilon^\top]) = \text{tr}(\mathbb{E}[M^\top M]\sigma^2\mathbb{I}_n) = \sigma^2\mathbb{E}[\text{tr}(M^\top M)].$$

421 On the other hand,  $M^\top M = \frac{1}{n}(\hat{\Sigma} + \varrho\mathbb{I}_p)^{-1}\hat{\Sigma}(\hat{\Sigma} + \varrho\mathbb{I}_p)^{-1}$ . Using the cyclic property of trace, we  
422 get the desired inequality.  $\square$

423 *Proof of Proposition C.3.* Recall from Proposition C.1 that  $\mathbb{E}\|\hat{\beta}\|_2^2 \geq n^{-1}\sigma^2\mathbb{E}[\text{tr}((\hat{\Sigma} + \varrho\mathbb{I}_p)^{-2}\hat{\Sigma})]$ .  
424 Below, we analyze the term inside the expectation. By the definition of the Stieltjes transform, we  
425 have

$$\text{tr}(\varrho(\hat{\Sigma} + \varrho\mathbb{I}_p)^{-1}) = \text{tr}(rn^{-\alpha}(\hat{\Sigma} + rn^{-\alpha}\mathbb{I}_p)^{-1}) = \text{tr}(r(n^\alpha\hat{\Sigma} + r\mathbb{I}_p)^{-1}) = pr\mathcal{S}_{\text{esd}(n^\alpha\hat{\Sigma})}(-r).$$

426 Therefore, by Lemma C.4, we have

$$\frac{d}{dr} \left( pr\mathcal{S}_{\text{esd}(n^\alpha\hat{\Sigma})}(-r) \right) = \frac{d}{dr} \text{tr}(\varrho(\hat{\Sigma} + \varrho\mathbb{I}_p)^{-1}) = \frac{d\varrho}{dr} \cdot \frac{d}{d\varrho} \text{tr}(\varrho(\hat{\Sigma} + \varrho\mathbb{I}_p)^{-1}) = n^{-\alpha} \text{tr}((\hat{\Sigma} + \varrho\mathbb{I}_p)^{-2}\hat{\Sigma}).$$

427 By Lemma C.5, we have

$$pr\mathcal{S}_{\text{esd}(n^\alpha\hat{\Sigma})}(-r) = pr \left( \gamma \cdot \mathcal{S}_{\text{esd}(n^\alpha\check{G})}(-r) + \frac{(1-\gamma)}{r} \right) = nr\mathcal{S}_{\text{esd}(n^\alpha\check{G})}(-r) + p(1-\gamma)$$

428 Thus, we have

$$\frac{d}{dr} \left( pr\mathcal{S}_{\text{esd}(n^\alpha\hat{\Sigma})}(-r) \right) = n \frac{d}{dr} \left( r\mathcal{S}_{\text{esd}(n^\alpha\check{G})}(-r) \right)$$

429 from which we conclude that

$$\text{tr}((\hat{\Sigma} + \varrho\mathbb{I}_p)^{-2}\hat{\Sigma}) = n^{\alpha+1} \frac{d}{dr} \left( r\mathcal{S}_{\text{esd}(n^\alpha\check{G})}(-r) \right).$$

430 In view of  $\mathbb{E}\|\hat{\beta}\|_2^2 \geq n^{-1}\sigma^2\mathbb{E}[\text{tr}((\hat{\Sigma} + \varrho\mathbb{I}_p)^{-2}\hat{\Sigma})]$  from Proposition C.1, we get the desired inequality.  
431  $\square$

## 432 I Continued from Appendix C.1

433 Before proceeding, we recall several definitions and notations adapted from Dobriban and Wager  
434 [2018]:

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[ \mathcal{S}_{\text{esd}(n^\alpha\check{G})}(z) \right] = v(z) \quad (5)$$

435 is analogous to the  $v(z)$  defined in the paragraph immediately following [Dobriban and Wager, 2018,  
436 Eqn. (2)]. The difference is our Equation (5) is for the limit of the  $n^\alpha$ -scaled matrices  $n^\alpha\check{G}$ , rather  
437 than for  $\check{G}$  as in Dobriban and Wager [2018].

438 Let  $H = \lim_{n \rightarrow \infty} \text{cdf}[\text{esd}(n^\alpha\Sigma)]$  be the limiting distribution as in Assumption B.2. Plugging in  
439  $z = -r$  into Dobriban and Wager [2018, Eqn. (A.1)], we have

$$-\frac{1}{v(-r)} = -r - \frac{1}{\gamma} \int \frac{tdH(t)}{1 + tv(-r)}.$$

440 Letting  $k \equiv k(r) := \frac{1}{v(-r)}$ , we can rewrite the above as

$$1 = \frac{r}{k} + \frac{1}{\gamma} \int \frac{tdH(t)}{k+t}. \quad (6)$$

441 By construction, we have

$$\frac{1}{\gamma} \int \frac{tdH(t)}{k+t} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^p \frac{1}{1 + kn^{-\alpha} \lambda_i^{-1}}$$

442 where the RHS is as in Assumption B.6. Consequently, the tuple  $r, k$  from Assumption B.6 coincide  
 443 with the earlier definition of  $k := \frac{1}{v(-r)}$  right before Equation (6). Having established the above, we  
 444 now proceed to:

445 *Proof of Lemma B.5.* By the product rule, we have

$$\frac{d}{dr} \left( r \mathcal{S}_{\text{esd}(n^\alpha \check{G})}(-r) \right) = \mathcal{S}_{\text{esd}(n^\alpha \check{G})}(r) - r \mathcal{S}'_{\text{esd}(n^\alpha \check{G})}(-r)$$

446 Now, taking the limit of the above equation on both side, we have

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{E} \left[ \frac{d}{dr} \left( r \mathcal{S}_{\text{esd}(n^\alpha \check{G})}(-r) \right) \right] &= \lim_{n \rightarrow \infty} \mathbb{E} \left[ \mathcal{S}_{\text{esd}(n^\alpha \check{G})}(r) - r \mathcal{S}'_{\text{esd}(n^\alpha \check{G})}(-r) \right] \\ &= v(-r) - r v'(-r) \quad \because \text{Definition of } v \text{ and } v' \\ &= \frac{d}{dr} (rv(-r)) \quad \because \text{Product rule} \\ &= \frac{d}{dr} (k \mathcal{S}_H(-k)) \quad \because \text{Marchenko-Pastur law (Assumption B.8)} \\ &= \frac{dk}{dr} \cdot \frac{d}{dk} (k \mathcal{S}_H(-k)) \quad \because \text{Chain rule} \\ &= \left( \frac{dr}{dk} \right)^{-1} \cdot \frac{d}{dk} (k \mathcal{S}_H(-k)) \quad \because \text{Inverse function theorem} \end{aligned}$$

447 To complete the proof, it suffices to show that both  $\frac{dr}{dk}$  and  $\frac{d}{dk} (k \mathcal{S}_H(-k))$  are positive which will be  
 448 checked in the next two lemmas.  $\square$

449 **Lemma I.1.** The function  $\frac{dr}{dk}$  evaluated at  $k$  is positive.

450 *Proof of Lemma I.1.* Recall that  $k = \frac{1}{v(-r)}$ . Thus, we have

$$\frac{dk}{dr}(r) = (-1) \frac{1}{v(-r)^2} (-1) \cdot v'(-r) = \frac{v'(-r)}{v(-r)^2}.$$

451 From the proof of Silverstein and Choi [1995, Theorem 4.1], we see that  $v'(\cdot) > 0$  for all negative  
 452 inputs. In particular,  $v'(-r) > 0$  which implies that  $\frac{dk}{dr}$  is positive. By the inverse function theorem,  
 453 we have  $\frac{dr}{dk} = \left( \frac{dk}{dr} \right)^{-1}$  is also positive.  $\square$

454 **Lemma I.2.** The quantity  $\frac{d}{dk} (k \mathcal{S}_H(-k))$  is positive.

455 *Proof of Lemma I.2.* Plugging in  $z = -r$  into Dobriban and Wager [2018, Eqn. (3)], we have

$$v(-r) - \frac{1}{r} = \frac{1}{\gamma} \left( m(-r) - \frac{1}{r} \right). \quad (7)$$

456 Now,

$$rm(-r) = \gamma r v(-r) + (1 - \gamma) \quad \because \text{Equation (7)} \quad (8)$$

$$= \gamma \frac{r}{k} + (1 - \gamma) \quad \because \text{Definition of } k \quad (9)$$

$$= \left( \gamma - \int \frac{tdH(t)}{k+t} \right) + (1 - \gamma) \quad \because \text{Equation (6)} \quad (10)$$

$$= 1 - \int \frac{tdH(t)}{k+t} \quad (11)$$

$$= \int \frac{k dH(t)}{k+t} \quad \because 1 = \int dH(t) = \int \frac{k+t}{k+t} dH(t) \quad (12)$$

$$= k \mathcal{S}_H(-k). \quad (13)$$

457 Thus, differentiating under the integral, we have

$$\frac{d}{dk} (k \mathcal{S}_H(-k)) = \int \frac{d}{dk} \left( \frac{k}{k+t} \right) dH(t) = \int \frac{tdH(t)}{(k+t)^2} > 0$$

458 as desired.  $\square$

459 **J Proof of Proposition 2.9**

460 We begin by analyzing the functions defined in Definition 2.8 and prove the items 1 and 2 of the  
461 “Moreover” part of Proposition 2.9:

462 **Proposition J.1.** Let  $\mathcal{I}$  and  $\mathcal{J}$  be functions as defined in Definition 2.8. Under Assumption 2.2 and  
463 Assumption B.6, we have that  $r = \mathcal{R}(k) := k \cdot (1 - \mathcal{I}(k))$  and  $\frac{dr}{dk} = 1 - \mathcal{J}(k)$ .

464 Furthermore, the following holds:

465 1.  $\mathcal{R}(k) \asymp k$  for  $k \gg 0$ ,

466 2. There exists  $k_{\text{crit}} > 0$  such that  $\mathcal{R}(k_{\text{crit}}) = 0$ ,  $\mathcal{R}$  is increasing and positive on  $(k_{\text{crit}}, +\infty)$ .

467 3.  $\mathcal{J}(k) < 1$  for  $k \in (k_{\text{crit}}, +\infty)$  and  $\mathcal{J}(+\infty) = 0$ .

468 *Proof of Proposition J.1.* We begin by proving the first part: that  $r = \mathcal{R}(k) := k \cdot (1 - \mathcal{I}(k))$  and  
469  $\frac{dr}{dk} = 1 - \mathcal{J}(k)$ . Rewrite the limit in Equation (3) as

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^p \frac{1}{1 + kn^{-\alpha} \sigma_i^{-1}} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^{n/\gamma} \frac{1}{1 + k(i/n)^\alpha} = \int_0^{1/\gamma_*} \frac{dx}{1 + kx^\alpha}$$

470 The right-most equality follows from the definition of the (Riemann) integral. If  $\gamma_* = 0$ , then  
471  $1/\gamma_* = +\infty$  and the above is interpreted as an improper Riemann integral. Now, rearranging  
472 Equation (3), we get the desired formula of  $r = \mathcal{R}(k) := k \cdot (1 - \mathcal{I}(k))$ . The formula for  $\frac{dr}{dk}$  follows  
473 from differentiating under the integral theorem. Note that this also proves the assertion made in  
474 Remark B.7.

475 For the first item of the “Furthermore” part, it suffices to show that  $\lim_{k \rightarrow +\infty} \mathcal{I}(k) = 0$ . This follows  
476 from the fact that  $\lim_{k \rightarrow +\infty} \frac{1}{1 + kx^\alpha} = 0$  for all  $x > 0$ , integrability of the function  $(1 + x^\alpha)^{-1}$  over  
477  $\mathbb{R}_{\geq 0}$ , and the dominated convergence theorem. Likewise,  $\lim_{k \rightarrow \infty} \mathcal{J}(k) = 0$  as well.

478 For the second item of the “Furthermore” part, we note that for all  $x$  sufficiently large, we have  
479  $\frac{dr}{dk} > 0$  since  $\lim_{k \rightarrow \infty} \mathcal{J}(k) = 0$ . Now, let  $k_{\text{crit}}$  be the largest real number such that  $\mathcal{R}(k_{\text{crit}}) = 0$ .  
480 Since  $\mathcal{R}(0) = 0$ , we must have  $k_{\text{crit}} \geq 0$ .

481 For all  $k > k_{\text{crit}}$ , we claim that  $\mathcal{I}(k) < 1$ . To see this, assume the contrary. Then by the fact that  
482  $\lim_{k \rightarrow +\infty} \mathcal{I}(k) = 0$  and the intermediate value theorem, there must exist  $k' > k$  such  
483 that  $\mathcal{I}(k') = 1$  which implies that  $\mathcal{R}(k') = 0$ . This contradicts the maximality of  $k_{\text{crit}}$ .

484 Finally, since  $1 + kx^\alpha \leq (1 + kx^\alpha)^2$  for all  $k \geq 0$  and  $x \geq 0$ , we have that  $\mathcal{I}(k) \geq \mathcal{J}(k)$  for all  
485 such  $k$ 's. Thus, by the previous claim, for all  $k > k_{\text{crit}}$ , we have  $1 > \mathcal{I}(k) \geq \mathcal{J}(k)$ . This proves  
486 that  $\frac{dr}{dk} > 0$  for all  $k > k_{\text{crit}}$ , as desired.  $\square$

487 **J.1 Review of the eigenlearning framework**

488 Before proceeding with finishing the proof of Proposition 2.9, we briefly review the eigenlearning  
489 framework. Simon et al. [2022] calculates the test error for the estimator

$$\tilde{\beta}_\delta := X(X^\top X + \delta \mathbb{I}_n)^{-1} y = X(n\check{G} + \delta \mathbb{I}_n)^{-1} y \quad (14)$$

490 for kernel ridge regression using the so-called *eigenlearning equations* [Simon et al., 2022, Section  
491 4.1]. Below, we recall some relevant parts of the framework:

492 **Definition J.2** (Eigenlearning eqn. specialized to setting in Section 2). Suppose that the ground truth  
493 regression function is linear, i.e.,  $f(x) = x^\top \beta^*$  for some  $\beta^* \in \mathbb{R}^p$ . Let  $\delta$  and  $\kappa$  satisfy the equation

$$n = \frac{\delta}{\kappa} + \sum_{i=1}^p \frac{\lambda_i}{\lambda_i + \kappa}. \quad (15)$$

494 Define the following  $n$ -dependent quantities:

495 1. *Overfitting coefficient:*  $\mathcal{E}_{\text{coef}} := n \frac{d\kappa}{d\delta}$

496 2. *Testing error:*  $\mathcal{E}_{\text{test}} := \mathcal{E}_{\text{coef}}(\sigma^2 + C)$  where

$$C = \sum_{i=1}^p (1 - \mathcal{L}_i)(\beta_i^*)^2 \quad \text{and} \quad \mathcal{L}_i := \frac{\lambda_i}{\lambda_i + \kappa}.$$

497 3. *Training error:*  $\mathcal{E}_{\text{train}} := \frac{\delta^2}{n^2 \kappa^2} \mathcal{E}_{\text{test}}$ .

498 **J.2 Completing the proof of Proposition 2.9**

499 Throughout this section, we assume that we are in the situation of Proposition 2.9. Now, Simon  
500 et al. [2022] uses a different scaling for ridge regression than the one we use. We first resolve this  
501 discrepancy. Comparing Equation (14) with the expression in Lemma C.2, if we let  $\delta := n\rho$ , then the  
502 expressions are equivalent, i.e.,  $\check{\beta}_\delta = \hat{\beta}_\rho$ . To see this, note that

$$\begin{aligned}\check{\beta}_\delta &= \check{\beta}_{n\rho} = X(X^\top X + n\rho\mathbb{I}_n)^{-1}y \\ &= (XX^\top + n\rho\mathbb{I}_p)^{-1}Xy \quad \because \text{Lemma H.1} \\ &= (n(n^{-1}XX^\top + \rho\mathbb{I}_p))^{-1}Xy \\ &= (\hat{\Sigma} + \rho\mathbb{I}_p)^{-1}\frac{1}{n}Xy = \hat{\beta}_\rho \quad \because \text{Definition of } \hat{\beta}_\rho\end{aligned}$$

503 Furthermore, we claim that as  $n \rightarrow \infty$ , we have  $r, k$  satisfies Equation (3) if and only if  $(\delta =$   
504  $nrn^{-\alpha}, \kappa = kn^{-\alpha})$  satisfies Equation (15):

$$n = \frac{\delta}{\kappa} + \sum_{i=1}^p \frac{\lambda_i}{\lambda_i + \kappa} \iff n = \frac{nrn^{-\alpha}}{kn^{-\alpha}} + \sum_{i=1}^p \frac{\lambda_i}{\lambda_i + kn^{-\alpha}} \iff 1 = \frac{r}{k} + \frac{1}{n} \sum_{i=1}^p \frac{1}{1 + kn^{-\alpha}\lambda_i^{-1}}.$$

505 Taking limit as  $n \rightarrow \infty$ , we have proved the claim.

506 Next, we show that  $\lim_{n \rightarrow \infty} C = 0$  where  $C$  is as in Definition J.2. We have  $\mathcal{L}_i := \frac{\lambda_i}{\lambda_i + \kappa} =$   
507  $\frac{1}{1 + k(i/n)^\alpha}$ . Note that  $\lim_{n \rightarrow \infty} \mathcal{L}_i = 1$  for all fixed  $i$ . On the other hand, since  $\sup_{n=1,2,\dots} \|\beta^*\|_2 <$   
508  $+\infty$ , dominated convergence theorem implies that  $\lim_{n \rightarrow \infty} C = 0$

509 We claim that the following asymptotic expression for the testing and training error hold:

$$\lim_{n \rightarrow \infty} \mathcal{E}_{\text{test}} = \sigma^2 \cdot \frac{dk}{dr} \quad \text{and} \quad \lim_{n \rightarrow \infty} \mathcal{E}_{\text{train}} = \sigma^2 \cdot \frac{r^2}{k^2} \cdot \frac{dk}{dr} \quad (16)$$

510 where  $r$  and  $k$  are defined as in Assumption B.6.

511 To see this, first note that the overfitting coefficient satisfies

$$\mathcal{E}_{\text{coef}} := n \frac{d\kappa}{d\delta} = n \frac{d\kappa}{d\rho} \frac{d\rho}{d\delta} = n \frac{d\kappa}{d\rho} \frac{1}{n} = \frac{d\kappa}{d\rho} = \frac{dk}{dr}.$$

512 Thus, we obtain the following asymptotic expression

$$\lim_{n \rightarrow \infty} \mathcal{E}_{\text{test}} = \mathcal{E}_{\text{coef}} \cdot \sigma^2 = \sigma^2 \cdot \frac{dk}{dr}.$$

513 On the other hand, the training error is given by

$$\mathcal{E}_{\text{train}} = \frac{\delta^2}{n^2 \kappa^2} \mathcal{E}_{\text{test}} = \frac{\rho^2}{\kappa^2} \mathcal{E}_{\text{test}} = \mathcal{E}_{\text{test}} \cdot \frac{r^2}{k^2}.$$

514 Therefore,  $\lim_{n \rightarrow \infty} \mathcal{E}_{\text{train}} = \sigma^2 \cdot \frac{r^2}{k^2} \cdot \frac{dk}{dr}$ . This proves (16), as desired.

515 Finally, we show that  $\frac{d}{d\alpha} \mathcal{E}_{\text{test}}^* > 0$  for any  $k > 0$ . To this end, we use the expression derived in the  
516 previous step that  $\mathcal{E}_{\text{test}}^* = \sigma^2 \cdot \frac{1}{1 - \mathcal{J}(k)}$ . Taking derivative of both side w.r.t  $\alpha$ , we have

$$\frac{d}{d\alpha} \mathcal{E}_{\text{test}}^* = \sigma^2 \frac{-1}{(1 - \mathcal{J}(k))^2} \frac{d}{d\alpha} \mathcal{J}(k)$$

517 Now, we recall from Definition 2.8 that  $\mathcal{J}(k) := \int_0^{1/\gamma^*} \frac{dx}{(1 + kx^\alpha)^2}$ . Thus, by differentiating under the  
518 integral sign, we have

$$\frac{d}{d\alpha} \mathcal{J}(k) = \int_0^{1/\gamma^*} \frac{-2kx^\alpha \log(x) dx}{(1 + kx^\alpha)^3}.$$

519 Putting it all together, we have

$$\frac{d}{d\alpha} \mathcal{E}_{\text{test}}^* = 2k\sigma^2 \frac{1}{(1 - \mathcal{J}(k))^2} \int_0^{1/\gamma^*} \frac{x^\alpha \log(x) dx}{(1 + kx^\alpha)^3}.$$

520 Since the integrand is positive, the integral is positive as well. Moreover, since  $k > 0$ , we have  
521  $\frac{d}{d\alpha} \mathcal{E}_{\text{test}}^* > 0$  as desired.  $\square$



## 522 K Code

523 Implementation of the  $\mathcal{I}$  and  $\mathcal{J}$  functions from Definition 2.8:

```

524 1 import scipy.special as sc
525 2 gamma = 0.5
526 3 alpha = 1.75
527 4
528 5 # I generator
529 6 I_gen = lambda x,k, alpha: x*sc.hyp2f1(1,(1/alpha), 1 + (1/alpha), -k*
530      x**alpha)
531 7 # J generator
532 8 J_gen = lambda x,k, alpha: x*sc.hyp2f1(2,(1/alpha), 1 + (1/alpha), -k*
533      x**alpha)
534 9
535 10 I = lambda k : I_gen(1/gamma, k, alpha) #\mathcal{I}
536 11 J = lambda k : J_gen(1/gamma, k, alpha) #\mathcal{J}
537 12
538 13 N = lambda k : 1 - I(k) # helper
539 14 D = lambda k : 1 - J(k) # helper
540 15
541 16 Etst = lambda k : 1/D(k) #\mathcal{E}_{\text{test}}/\sigma^2
542 17 Etrn = lambda k : N(k)**2/D(k) #\mathcal{E}_{\text{train}}/\sigma^2
543 18 R = lambda k : k*(1-I(k)) # \mathcal{R}

```

544 For the experiments in Figure 1-Right:

```

545 1 import numpy as np
546 2 gamma = 0.5
547 3 alpha = 1.75
548 4
549 5 k_grid = [ 1.34, 1.99, 2.45, 2.92, 3.44, 4.03, 4.71, 5.5 ,
550      6.44,
551 6      7.55, 8.9 , 10.54, 12.58, 15.15, 18.46, 22.8 , 28.67, 36.87,
552 7      48.82, 67.2 ]
553 8
554 9 n_tst = 1000
555 10 def get_norms(n,r):
556 11     p = int(n/gamma)
557 12
558 13     idx = np.arange(1,p+1) # feature indices
559 14
560 15     pop_evs = idx**(-alpha) # population level eigenvalues
561 16
562 17     X = np.multiply(np.sqrt(pop_evs[:,None]), np.random.normal(size= (
563      p, n)) )
564 18     X_tst = np.multiply(np.sqrt(pop_evs[:,None]), np.random.normal(
565      size= (p, n_tst)) )
566 19
567 20     beta_true = np.sqrt(10)*np.random.normal(size= (p,1))/np.sqrt(p)
568 21
569 22     y = X.T@beta_true + np.random.normal(size= (n,1))
570 23     y_tst = X_tst.T@beta_true + np.random.normal(size= (n_tst,1))
571 24
572 25
573 26     hatSig = (1/n)*X@X.T # sample covariance matrix
574 27
575 28     beta = (1/n)*np.linalg.solve(hatSig + r*n**(-alpha)*np.eye(p), X@y
576      )
577 29
578 30     norm = np.linalg.norm(beta)**2
579 31     Etrn = np.mean(np.square(y-X.T@beta))
580 32     Etst = np.mean(np.square(y_tst-X_tst.T@beta))
581 33     return {"norm": norm, "Etrn":Etrn, "Etst":Etst}
582 34

```

```

583:5 rs = R(np.array(k_grid))
584:6
585:7 n = 1000
586:8
587:9 Etrns = []
588:0 Etsts = []
589:1 for r in rs:
590:2     result = get_norms(n,r)
591:3     Etrns.append(result["Etrn"])
592:4     Etsts.append(result["Etst"])

```

593 For the experiments in Figure 1:

```

594:1 # run the previous block first!
595:2 r = 3.5433549686341
596:3
597:4 ns = np.logspace(1,3.6,num=20)
598:5 categories = ["norm","Etrn","Etst"]
599:6 n_trials = 10
600:7
601:8 results = {cat : [[] for _ in range(n_trials)] for cat in categories}
602:9
603:0 for t in range(n_trials):
604:1     for n in ns:
605:2         out = get_norms(int(n),r)
606:3         for cat in categories:
607:4             results[cat][t].append(out[cat])

```