
Under-Parameterized Double Descent for Ridge Regularized Least Squares Denoising of Data on a Line

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this paper, we present a simple example that provably exhibits double descent in
2 the under-parameterized regime. For simplicity, we look at the ridge regularized
3 least squares denoising problem with data on a line embedded in high-dimension
4 space. By deriving an asymptotically accurate formula for the generalization error,
5 we observe sample-wise and parameter-wise double descent with the peak in the
6 under-parameterized regime rather than at the interpolation point or in the over-
7 parameterized regime. Further, the peak of the sample-wise double descent curve
8 corresponds to a peak in the curve for the norm of the estimator, and adjusting μ ,
9 the strength of the ridge regularization, shifts the location of the peak. We observe
10 that parameter-wise double descent occurs for this model for small μ . For larger
11 values of μ , we observe that the curve for the norm of the estimator has a peak but
12 that this no longer translates to a peak in the generalization error.

13 1 Introduction

14 This paper aims to demonstrate interesting new phenomena that suggest that our understanding of the
15 relationship between the number of data points, the number of parameters, and the generalization
16 error is incomplete, even for simple linear models with data on a line. The classical bias-variance
17 theory postulates that the generalization risk versus the number of parameters for a fixed number
18 of training data points is U-shaped. However, modern machine learning showed that if we keep
19 increasing the number of parameters, the generalization error eventually starts decreasing again [1,
20 2]. This second descent has been termed as *double descent* and occurs in the *over-parameterized*
21 *regime*, that is when the number of parameters exceeds the number of data points. Understanding the
22 location and the cause of such peaks in the generalization error is of significant importance. Hence
23 many recent works have theoretically studied the generalization error for linear regression [3–12]
24 and kernelized regression [13–21] and show that there exists a peak at the boundary between the
25 under and over-parameterized regimes. Further works such as [10, 22–25] show that there can be
26 multiple descents in the over-parameterized regime and [26] shows that any shaped generalization
27 error curve can occur in the over-parameterized regime. However, all prior works assume that the
28 classical bias-variance trade-off is true in the under-parameterized regime.

29 The implicit bias of the learning algorithm is a possible reason that the error decreases in the over-
30 parameterized regime [27–32]. In the under-parameterized regime, there is exactly one solution
31 that minimizes the loss. However, once in the over-parameterized regime, there are many different
32 solutions, and the training algorithm implicitly picks one that generalizes well. For linear models, the
33 generalization error and the variance are very closely related to the norm of the estimator [11, 33].
34 Then, using the well-known fact that the pseudo-inverse solution to the least squares problem is the
35 minimum norm solution, we see that the training algorithm picks solutions with the minimum norm.
36 Hence this learning algorithm minimizes the variance and lowers the generalization error.

Table 1: Table showing various assumptions on the data and the location of the double descent peak for linear regression and denoising. We only present a subset of references for each problem setting.

Noise	Ridge Reg.	Dimension	Peak Location	Reference
Input	Yes	1	Under-parameterized	This paper.
Input	No	Low	Over-parameterized/interpolation point	[33, 37]
Output	No	Full	Over-parameterized/interpolation point	[5, 8, 11]
Output	Yes	Full	Over-parameterized/interpolation point	[11, 24]
Output	No	Low	Over-parameterized/interpolation point	[34, 35]
Output	Yes	Low	Over-parameterized/interpolation point	[36]

37 **Main Contributions.** In contrast with prior work, this paper shows that double descent can occur
 38 in the under-parameterized regime. Specifically, when denoising data on a line embedded in high-
 39 dimensional space using a denoiser obtained as the pseudo-inverse solution for the ridge regularized
 40 least squares problem, we show that a peak in the generalization error curve occurs in the under-
 41 parameterized regime. We also show that changing the ridge regularization strength changes the
 42 location of the peak. The major contributions of this paper are as follows.¹

- 43 • **(Generalization error)** We derive a theoretical formula for the generalization error (Theorem 1).
- 44 • **(Under-parameterized double descent)** We prove (Theorem 2) and empirically demonstrate
 45 that the generalization error versus the number of data points curve has double descent in the
 46 under-parameterized regime.
- 47 • **(Location of the peak)** The peak location depends on the regularization strength. We provide
 48 evidence (Theorem 6) that the peak is near $c = \frac{1}{\mu^2+1}$ for the sample-wise double descent curves.
- 49 • **(Norm of the estimator)** We show that the peak in the curve for the generalization error versus
 50 the number of training data points corresponds to a peak in the norm of the estimator. However,
 51 versus the number of parameters, we show that there is still a peak in the curve for the norm of the
 52 estimator (Theorem 5), but this no longer corresponds to a peak in the generalization error.

53 **Low-Dimensional Data.** It is important to highlight that using low-rank data does not immediately
 54 imply that a peak occurs in the under-parameterized regime. Specifically, [33–37] look at a variety of
 55 different problems with low rank data and see that the peak occurs at the interpolation point or in the
 56 over-parameterized regime. Table 1 compares common assumptions and the location of the peak.

57 2 Background and Model Assumptions

58 Throughout the paper, we assume that noiseless training data x_i live in \mathbb{R}^d and that we have access to
 59 a $d \times N_{trn}$ matrix X_{trn} of training data. Then given new data $X_{tst} \in \mathbb{R}^{d \times N_{tst}}$, we are interested in
 60 the least squares generalization (or test) error. Two scenarios for the generalization error curve are
 61 considered; data scaling and parameter scaling.

62 **Definition 1.** • Data scaling refers to the regime in which we fix the dimension d of the input data
 63 and vary the number of training data points N_{trn} . This is also known as the sample-wise regime.

64 • Parameter scaling refers to the regime in which we fix the number of training data points N_{trn} and
 65 vary the dimension d of the input data. This is also known as the parameter-wise regime.

66 • A linear model is under-parameterized, if $d < N_{trn}$. A linear model is over-parameterized, if
 67 $d > N_{trn}$. The boundary of the under and over-parameterized regimes is when $d = N_{trn}$.

68 • Given N_{trn} , the interpolation point is the smallest d for the which the model has zero training error.

69 • A curve has double descent if the curve has a local maximum or peak.

70 • The aspect ratio of an $m \times n$ matrix is $c := m/n$.

71 **Prior Double Descent** We present a baseline model from prior work on double descent. This is to
 72 highlight prior important phenomena related to double descent in the literature. Concretely, consider
 73 the following simple linear model that is a special case of the general models studied in [5, 8, 11, 24]
 74 amongst many other works. Let $x_i \sim \mathcal{N}(0, I_d)$ and let $\beta \in \mathbb{R}^d$ be a linear model with $\|\beta\| = 1$. Let
 75 $y_i = \beta^T x_i + \xi_i$ where $\xi \sim \mathcal{N}(0, 1)$. Then, let $\beta_{opt} := \arg \min_{\tilde{\beta}} \|\beta^T X_{trn} - \tilde{\beta} X_{trn} + \xi_{trn}\|$, where
 76 $\xi_{trn} \in \mathbb{R}^{N_{trn} \times 1}$. Then the excess risk, when taking the expectation over the new test data point, can
 77 be expressed as $\mathcal{R} = \|\beta - \beta_{opt}\|^2 = \|\beta\|^2 + \|\beta_{opt}\|^2 - 2\beta^T \beta_{opt}$. Let c be the aspect ratio of the

¹All code is available anonymized at [Github Repo]

78 data matrix. That is, $c = d/N_{trn}$. Then it can be shown that²

$$\mathbb{E}_{X_{trn}, \xi_{trn}}[\|\beta_{opt}\|^2] = \begin{cases} 1 + \frac{c}{1-c} & c < 1 \\ \frac{1}{c} + \frac{1}{c-1} & c > 1 \end{cases} \quad \text{and} \quad \mathbb{E}_{X_{trn}, \xi_{trn}}[\beta^T \beta_{opt}] = \begin{cases} 1 & c < 1 \\ \frac{1}{c} & c > 1 \end{cases}$$

79 Then, the excess risk can be expressed as $\mathcal{R} = \begin{cases} \frac{c}{1-c} & c < 1 \\ \frac{c-1}{c} + \frac{1}{c-1} & c > 1 \end{cases}$. There are a few important
80 features that are considered staple in many prior double descent curves that are present in this model.

- 81 1. The peak happens at $c = 1$, on the border between the under and over-parameterized regimes.
- 82 2. Further, at $c = 1$ the training error equals zero. Hence this is the interpolation point.
- 83 3. The peak occurs due to the norm of the estimator β_{opt} blowing up near the interpolation point.

84 **Further, [26] proved risk is monotonic in the under-parameterized regime for the above model.**

85 For the ridge regularized version of the regression problem, as shown in [11, 24], *the peak is always*
86 *at $c = 1$* (see Figure 1 in [24]). Further, as seen in Figure 1 in [24], changing the strength of the
87 regularization changes the magnitude of the peak. *Not the location of the peak.* Building on this, [23]
88 looks at the model where $y_i = f(x_i) + \xi_i$ and shows that triple descent occurs for the random features
89 model [38] in the *over-parameterized* regime. Further [26] shows that by considering a variety of
90 product data distributions, any shaped risk curve can be observed in the *over-parameterized* regime.

91 **Assumptions for Denoising Model** With the context from the previous section in mind, we are
92 now ready to present the assumptions for the input noise model with double descent in the under-
93 parameterized regime. For the denoising problem, let $A_{trn} \in \mathbb{R}^{d \times N_{trn}}$ be the noise matrix, then the
94 ridge regularized least square denoiser W_{opt} is the minimum norm solution to

$$W_{opt} := \arg \min_W \|X_{trn} - W(X_{trn} + A_{trn})\|_F^2 + \mu^2 \|W\|_F^2. \quad (1)$$

95 Given test data X_{tst} , the mean squared generalization error is given by

$$\mathcal{R}(W_{opt}) = \mathbb{E}_{A_{trn}, A_{tst}} \left[\frac{1}{N_{tst}} \|X_{tst} - W_{opt}(X_{tst} + A_{tst})\|_F^2 \right]. \quad (2)$$

96 The reason we consider linear models with the pseudo-inverse solution is that this eliminates other
97 factors, such as the initialization of the network that could be a cause of the double descent [23]. We
98 assume that the data lies on a line embedded in high-dimensional space.

99 **Assumption 1.** Let $\mathcal{U} \subset \mathbb{R}^d$ be a one dimensional space with a unit basis vector u . Then let $X_{trn} =$
100 $\sigma_{trn} w_{trn}^T \in \mathbb{R}^{d \times N_{trn}}$ and $X_{tst} = \sigma_{tst} w_{tst}^T \in \mathbb{R}^{d \times N_{tst}}$ be the respective SVDs for the training data
101 and test data matrices. We further assume that $\sigma_{trn} = O(\sqrt{N_{trn}})$ and $\sigma_{tst} = O(\sqrt{N_{tst}})$.

102 In [26], it was shown that by considering specific data distributions, any shaped generalization error
103 curve could be observed in the over-parameterized regime. Hence to limit the effect of the data, we
104 consider data on a line with norm restrictions.

105 **Assumption 2.** The entries of the noise matrices $A \in \mathbb{R}^{d \times N}$ are I.I.D. from $\mathcal{N}(0, 1/d)$.

106 **Notational note.** One final piece of technical notation is the following function definition.

$$p(\mu) := (4\mu^{15} + 48\mu^{13} + 204\mu^{11} + 352\mu^9 + 192\mu^7) \sqrt{\mu^2 + 4} - (4\mu^{16} + 56\mu^{14} + 292\mu^{12} + 680\mu^{10} + 640\mu^8 + 128\mu^6). \quad (3)$$

107 3 Under-Parameterized Regime Peak

108 We begin by providing a formula for the generalization error given by Equation 2 for the least squares
109 solution given by Equation 1. The over-parameterized case can be found in Appendix F.2. See
110 Appendix A for more discussion. All proofs are in Appendix F.

111 **Theorem 1** (Generalization Error Formula). *Suppose the training data X_{trn} and test data X_{tst}*
112 *satisfy Assumption 1 and the noise A_{trn}, A_{tst} satisfy Assumption 2. Let μ be the regularization*

²The proofs are in Appendix F.1.

113 parameter. Then for the under-parameterized regime (i.e., $c < 1$) for the solution W_{opt} to Problem 1,
 114 the generalization error or risk given by Equation 2 is given by

$$\mathcal{R}(c, \mu) = \tau^{-2} \left(\frac{\sigma_{tst}^2}{N_{tst}} + \frac{c\sigma_{trn}^2(\sigma_{trn}^2 + 1)}{2d} \left(\frac{1 + c + \mu^2 c}{\sqrt{(1 - c + \mu^2 c)^2 + 4\mu^2 c^2}} - 1 \right) \right) + o\left(\frac{1}{d}\right),$$

115 where $\tau^{-1} = \frac{2}{2 + \sigma_{trn}^2(1 + c + \mu^2 c - \sqrt{(1 - c + \mu^2 c)^2 + 4\mu^2 c^2})}$.

116 **Data Scaling.** We prove that the risk curve in Theorem 1 has a peak for $c \in (0, 1)$. Theorem 2
 117 tells us that under certain conditions, we are guaranteed to have a peak in the under-parameterized
 118 regime. This contrasts with prior work such as [3, 5, 8–11, 14, 25]. Further, we conjecture that the
 119 peak occurs near $c = (\mu^2 + 1)^{-1}$ (Appendix B). Figure 1 shows that the theoretically predicted risk
 120 matches the numerical risk. Moreover, the curve has a single peak for $c < 1$. Thus, *verifying that*
 121 *double descent occurs in the under-parameterized regime.* Finally, Figure 1 shows that the location
 122 of the peak is near the conjectured location of $\frac{1}{\mu^2 + 1}$. See Appendix D for the training error curves.

123 **Theorem 2 (Under-Parameterized Peak).** *If $\mu \in \mathbb{R}_{>0}$ is such that $p(\mu) < 0$, $\sigma_{trn}^2 = N_{trn} = d/c$
 124 and $\sigma_{tst}^2 = N_{tst}$, and d is sufficiently large, then the risk $\mathcal{R}(c)$ from Theorem 1, as a function of c ,
 125 has a local maximum in the under-parameterized regime ($c \in (0, 1)$).*

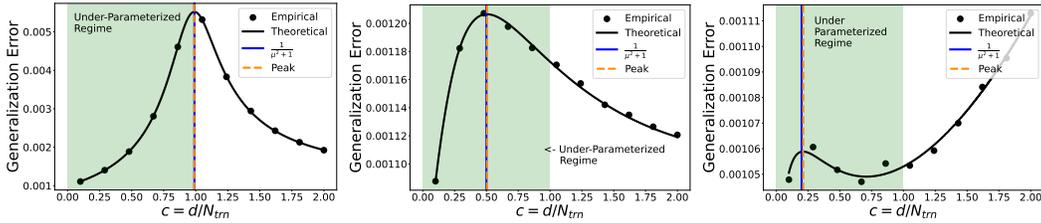


Figure 1: Figure showing the risk curve in the data scaling regime for different values of μ [(L) $\mu = 0.1$, (C) $\mu = 1$, (R) $\mu = 2$]. Here $\sigma_{trn} = \sqrt{N_{trn}}$, $\sigma_{tst} = \sqrt{N_{tst}}$, $d = 1000$, $N_{tst} = 1000$. For each empirical point, we ran at least 100 trials. More details can be found in Appendix G.

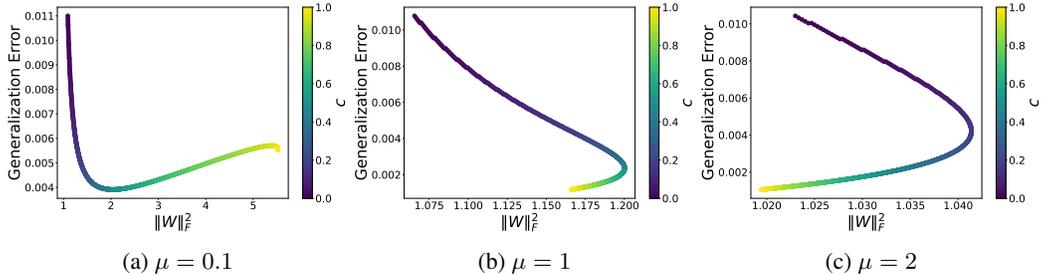


Figure 2: Figure showing generalization error versus $\|W_{opt}\|_F^2$ for the parameter scaling regime for three different values of μ . More details can be found in Appendix B.

126 **Parameter Scaling.** For many prior models, the data and parameter scaling regimes are analogous
 127 in that the shape of the risk is primarily governed by the aspect ratio c of the data matrix. However,
 128 we see significant differences between the parameter scaling and data scaling regimes for our setup.
 129 Figure 2 shows that for small values of μ , double descent occurs in the under-parameterized regime,
 130 for larger values of μ , the risk is monotonically decreasing.³ Further, Figure 2 shows that for larger
 131 values of μ , there is still a peak in the curve for the norm of the estimator $\|W_{opt}\|_F^2$. However, this
 132 does not translate to a peak in the risk curve.

133 **Theorem 3 ($\|W_{opt}\|_F$ Peak).** *If $\sigma_{tst} = \sqrt{N_{tst}}$, $\sigma_{trn} = \sqrt{N_{trn}}$ and μ is such that $p(\mu) < 0$,
 134 then for N_{trn} large enough and $d = cN_{trn}$, we have that $\|W_{opt}\|_F$ has a local maximum in the
 135 under-parameterized regime. Specifically for $c \in ((\mu^2 + 1)^{-1}, 1)$.*

³This is verified for more values of μ in Appendix B.

136 **References**

- 137 [1] Manfred Opper and Wolfgang Kinzel. Statistical Mechanics of Generalization. *Models of*
138 *Neural Networks III: Association, Generalization, and Representation*, 1996 (cited on page 1).
- 139 [2] Mikhail Belkin, Daniel J. Hsu, Siyuan Ma, and Soumik Mandal. Reconciling Modern Machine-
140 Learning Practice and the Classical Bias–Variance Trade-off. *Proceedings of the National*
141 *Academy of Sciences*, 2019 (cited on page 1).
- 142 [3] Madhu S. Advani, Andrew M. Saxe, and Haim Sompolinsky. High-dimensional Dynamics of
143 Generalization Error in Neural Networks. *Neural Networks*, 2020 (cited on pages 1, 4, 9).
- 144 [4] Chen Cheng and Andrea Montanari. Dimension Free Ridge Regression. *arXiv preprint*
145 *arXiv:2210.08571*, 2022 (cited on page 1).
- 146 [5] Edgar Dobriban and Stefan Wager. High-dimensional asymptotics of prediction: ridge regres-
147 sion and classification. *The Annals of Statistics*, 2018 (cited on pages 1, 2, 4, 9).
- 148 [6] Gabriel Mel and Surya Ganguli. A Theory of High Dimensional Regression with Arbitrary
149 Correlations Between Input Features and Target Functions: Sample Complexity, Multiple
150 Descent Curves and a Hierarchy of Phase Transitions. In *Proceedings of the 38th International*
151 *Conference on Machine Learning*, 2021 (cited on page 1).
- 152 [7] Vidya Muthukumar, Kailas Vodrahalli, and Anant Sahai. Harmless Interpolation of Noisy Data
153 in Regression. *2019 IEEE International Symposium on Information Theory (ISIT)*, 2019 (cited
154 on page 1).
- 155 [8] Peter Bartlett, Philip M. Long, Gábor Lugosi, and Alexander Tsigler. Benign Overfitting in
156 Linear Regression. *Proceedings of the National Academy of Sciences*, 2020 (cited on pages 1,
157 2, 4, 9).
- 158 [9] Mikhail Belkin, Daniel J. Hsu, and Ji Xu. Two Models of Double Descent for Weak Features.
159 *SIAM Journal on Mathematics of Data Science*, 2020 (cited on pages 1, 4, 9).
- 160 [10] Michal Dereziński, Feynman T Liang, and Michael W Mahoney. Exact Expressions for Double
161 Descent and Implicit Regularization Via Surrogate Random Design. In *Advances in Neural*
162 *Information Processing Systems*, 2020 (cited on pages 1, 4, 9).
- 163 [11] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J. Tibshirani. Surprises in High-
164 Dimensional Ridgeless Least Squares Interpolation. *The Annals of Statistics*, 2022 (cited on
165 pages 1–4, 9).
- 166 [12] Bruno Loureiro, Gabriele Sicuro, Cedric Gerbelot, Alessandro Pocco, Florent Krzakala, and
167 Lenka Zdeborova. Learning Gaussian Mixtures with Generalized Linear Models: Precise
168 Asymptotics in High-dimensions. In *Advances in Neural Information Processing Systems*,
169 2021 (cited on page 1).
- 170 [13] Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Generalization Error of Random
171 Feature and Kernel Methods: Hypercontractivity and Kernel Matrix Concentration. *Applied*
172 *and Computational Harmonic Analysis*, 2022 (cited on page 1).
- 173 [14] Song Mei and Andrea Montanari. The Generalization Error of Random Features Regression:
174 Precise Asymptotics and the Double Descent Curve. *Communications on Pure and Applied*
175 *Mathematics*, 75, 2021 (cited on pages 1, 4, 9).
- 176 [15] Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A Mean Field View of the Landscape
177 of Two-layer Neural Networks. *Proceedings of the National Academy of Sciences of the United*
178 *States of America*, 2018 (cited on page 1).
- 179 [16] Nilesh Tripuraneni, Ben Adlam, and Jeffrey Pennington. Covariate Shift in High-Dimensional
180 Random Feature Regression. *ArXiv*, 2021 (cited on page 1).
- 181 [17] Federica Gerace, Bruno Loureiro, Florent Krzakala, Marc Mezard, and Lenka Zdeborova.
182 Generalisation Error in Learning with Random Features and the Hidden Manifold Model.
183 In *Proceedings of the 37th International Conference on Machine Learning*, 2020 (cited on
184 page 1).
- 185 [18] Blake Woodworth, Suriya Gunasekar, Jason D. Lee, Edward Moroshko, Pedro Savarese, Itay
186 Golan, Daniel Soudry, and Nathan Srebro. Kernel and Rich Regimes in Overparametrized
187 Models. In *Proceedings of Thirty Third Conference on Learning Theory*, 2020 (cited on
188 page 1).
- 189 [19] Bruno Loureiro, Cedric Gerbelot, Hugo Cui, Sebastian Goldt, Florent Krzakala, Marc Mezard,
190 and Lenka Zdeborova. Learning Curves of Generic Features Maps for Realistic Datasets with
191 a Teacher-Student Model. In *NeurIPS*, 2021 (cited on page 1).

- 192 [20] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Limitations of
193 Lazy Training of Two-layers Neural Network. In *Advances in Neural Information Processing*
194 *Systems*, 2019 (cited on page 1).
- 195 [21] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. When Do
196 Neural Networks Outperform Kernel Methods? In *Advances in Neural Information Processing*
197 *Systems*, 2020 (cited on page 1).
- 198 [22] Ben Adlam and Jeffrey Pennington. The Neural Tangent Kernel in High Dimensions: Triple
199 Descent and a Multi-Scale Theory of Generalization. In *International Conference on Machine*
200 *Learning*, 2020 (cited on page 1).
- 201 [23] Stéphane d’Ascoli, Levent Sagun, and Giulio Biroli. Triple Descent and the Two Kinds of
202 Overfitting: Where and Why Do They Appear? In *Advances in Neural Information Processing*
203 *Systems*, 2020 (cited on pages 1, 3).
- 204 [24] Preetum Nakkiran, Prayaag Venkat, Sham M. Kakade, and Tengyu Ma. Optimal Regularization
205 can Mitigate Double Descent. In *International Conference on Learning Representations*, 2020
206 (cited on pages 1–3, 13).
- 207 [25] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Linearized
208 Two-layers Neural Networks in High Dimension. *The Annals of Statistics*, 2021 (cited on
209 pages 1, 4, 9).
- 210 [26] Lin Chen, Yifei Min, Mikhail Belkin, and Amin Karbasi. Multiple Descent: Design Your Own
211 Generalization Curve. *Advances in Neural Information Processing Systems*, 2021 (cited on
212 pages 1, 3).
- 213 [27] Alnur Ali, Edgar Dobriban, and Ryan J. Tibshirani. The Implicit Regularization of Stochastic
214 Gradient Flow for Least Squares. In *International Conference on Machine Learning*, 2020
215 (cited on page 1).
- 216 [28] Arthur Jacot, Berfin Simsek, Francesco Spadaro, Clement Hongler, and Franck Gabriel.
217 Implicit Regularization of Random Feature Models. In *Proceedings of the 37th International*
218 *Conference on Machine Learning*, 2020 (cited on page 1).
- 219 [29] Ohad Shamir. The Implicit Bias of Benign Overfitting. In *Proceedings of Thirty Fifth Confer-*
220 *ence on Learning Theory*, 2022 (cited on page 1).
- 221 [30] Behnam Neyshabur. Implicit Regularization in Deep Learning. *ArXiv*, abs/1709.01953, 2017
222 (cited on page 1).
- 223 [31] Behnam Neyshabur, Srinadh Bhojanapalli, David Mcallester, and Nati Srebro. Exploring
224 Generalization in Deep Learning. In *Advances in Neural Information Processing Systems*,
225 2017 (cited on page 1).
- 226 [32] Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In Search of the Real Inductive Bias:
227 On the Role of Implicit Regularization in Deep Learning. *CoRR*, abs/1412.6614, 2015 (cited
228 on page 1).
- 229 [33] Rishi Sonthalia and Raj Rao Nadakuditi. Training data size induced double descent for
230 denoising feedforward neural networks and the role of training noise. *Transactions on Machine*
231 *Learning Research*, 2023 (cited on pages 1, 2, 13, 16, 19, 24).
- 232 [34] Ninyuan Huang, David W. Hogg, and Soledad Villar. Dimensionality reduction, regulariza-
233 tion, and generalization in overparameterized regressions. *SIAM Journal on Mathematics of*
234 *Data Science*, 2022 (cited on page 2).
- 235 [35] Ji Xu and Daniel J Hsu. On the number of variables to use in principal component regression.
236 *Advances in neural information processing systems*, 2019 (cited on page 2).
- 237 [36] Denny Wu and Ji Xu. On the Optimal Weighted ℓ_2 Regularization in Overparameterized
238 Linear Regression. *Advances in Neural Information Processing Systems*, 2020 (cited on
239 page 2).
- 240 [37] Risi Sonthalia Chinmaya Kausik Kashvi Srivastva. Generalization Error without Independence:
241 Denoising, Linear Regression, and Transfer Learning, 2023 (cited on page 2).
- 242 [38] Ali Rahimi and Benjamin Recht. Random Features for Large-Scale Kernel Machines. In
243 *Advances in Neural Information Processing Systems*, 2007 (cited on page 3).
- 244 [39] Andrew Ng. CS229 Lecture notes. *CS229 Lecture notes*, 2000 (cited on page 10).
- 245 [40] Carl D. Meyer Jr. Generalized Inversion of Modified Matrices. *SIAM Journal on Applied*
246 *Mathematics*, 1973 (cited on page 17).

- 247 [41] Friedrich Götze and Alexander Tikhomirov. The Rate of Convergence for Spectra of GUE and
248 LUE Matrix Ensembles. *Central European Journal of Mathematics*, 2005 (cited on page 21).
- 249 [42] Friedrich Götze and Alexander Tikhomirov. Rate of Convergence to the Semi-Circular Law.
250 *Probability Theory and Related Fields*, 2003 (cited on page 21).
- 251 [43] Friedrich Götze and Alexander Tikhomirov. Rate of Convergence in Probability to the
252 Marchenko-Pastur Law. *Bernoulli*, 2004 (cited on page 21).
- 253 [44] Vladimir Marcenko and Leonid Pastur. Distribution of Eigenvalues for Some Sets of Random
254 Matrices. *Mathematics of The Ussr-sbornik*, 1967 (cited on page 21).
- 255 [45] Z. Bai, Baiqi. Miao, and Jian-Feng. Yao. Convergence Rates of Spectral Distributions of Large
256 Sample Covariance Matrices. *SIAM Journal on Matrix Analysis and Applications*, 2003 (cited
257 on page 21).

258	Contents	
259	1 Introduction	1
260	2 Background and Model Assumptions	2
261	3 Under-Parameterized Regime Peak	3
262	A Under-Parameterized Regime Peak	9
263	B Peak Location and $\ W_{opt}\ _F$	11
264	B.1 Peak Location for the Data Scaling Regime	11
265	C Generalization error - bias and variance	11
266	D Training Error	11
267	E Regularization Trade-off	13
268	E.1 Optimal Value of μ	14
269	E.2 Trade-off in Parameter Scaling Regime	15
270	F Proofs	16
271	F.1 Linear Regression	16
272	F.2 Proofs for Theorem 1	16
273	F.2.1 Step 1: Decompose the error into bias and variance terms.	16
274	F.2.2 Step 2: Formula for W_{opt}	17
275	F.2.3 Step 3: Decompose the terms into a sum of various trace terms.	18
276	F.2.4 Step 4: Estimate With Random Matrix Theory	19
277	F.2.5 Step 5: Putting it together	27
278	F.3 Proof of Theorem 2	28
279	F.4 Proof of Theorem 6	29
280	F.5 Proof of Theorem 5	31
281	F.6 Proof of Theorem 7	31
282	F.7 Proof of Proposition 1	34
283	G Experiments	35
284	H Technical Assumption on μ	35

285 A Under-Parameterized Regime Peak

286 We begin by providing a formula for the generalization error given by Equation 2 for the least squares
287 solution given by Equation 1. All proofs are in Appendix F.

288 **Theorem 1** (Generalization Error Formula). *Suppose the training data X_{trn} and test data X_{tst}
289 satisfy Assumption 1 and the noise A_{trn}, A_{tst} satisfy Assumption 2. Let μ be the regularization
290 parameter. Then for the under-parameterized regime (i.e., $c < 1$) for the solution W_{opt} to Problem 1,
291 the generalization error or risk given by Equation 2 is given by*

$$\mathcal{R}(c, \mu) = \tau^{-2} \left(\frac{\sigma_{tst}^2}{N_{tst}} + \frac{c\sigma_{trn}^2(\sigma_{trn}^2 + 1)}{2d} \left(\frac{1 + c + \mu^2 c}{\sqrt{(1 - c + \mu^2 c)^2 + 4\mu^2 c^2}} - 1 \right) \right) + o\left(\frac{1}{d}\right),$$

292 where $\tau^{-1} = \frac{2}{2 + \sigma_{trn}^2(1 + c + \mu^2 c - \sqrt{(1 - c + \mu^2 c)^2 + 4\mu^2 c^2})}$.

293 Since the focus is on the under-parameterized regime, Theorem 1 only presents the under-
294 parameterized case. The over-parameterized case can be found in Appendix F.2.

295 **Data Scaling.** Looking at the formula in Theorem 1, the risk curve's shape is unclear. In this
296 section, we prove that the risk curve in Theorem 1 has a peak for $c \in (0, 1)$. Theorem 2 tells us that
297 under certain conditions, we are theoretically guaranteed to have a peak in the under-parameterized
298 regime. This contrasts with prior work such as [3, 5, 8–11, 14, 25] where double descent occurs in
299 the over-parameterized regime or on the boundary between the two regimes.

300 **Theorem 2** (Under-Parameterized Peak). *If $\mu \in \mathbb{R}_{>0}$ is such that $p(\mu) < 0$, $\sigma_{trn}^2 = N_{trn} = d/c$
301 and $\sigma_{tst}^2 = N_{tst}$, and d is sufficiently large, then the risk $\mathcal{R}(c)$ from Theorem 1, as a function of c ,
302 has a local maximum in the under-parameterized regime ($c \in (0, 1)$).*

303 Since the peak no longer occurs at $c = 1$, one important question is to determine the location of the
304 peak. Theorem 6 provides a method for estimating the location of the peak.

305 **Theorem 4** (Peak Location). *If $\mu \in \mathbb{R}_{>0}$ is such that $p(\mu) < 0$, $\sigma_{trn}^2 = N_{trn} = d/c$ and $\sigma_{tst}^2 = N_{tst}$,
306 then the partial derivative with respect to c of the risk $\mathcal{R}(c)$ from Theorem 1 can be written as*

$$\frac{\partial}{\partial c} \mathcal{R}(c, \mu) = \frac{(\mu^2 c + c - 1)P(c, \mu, T(c, \mu), d) + 4d\mu^2 c^2(2\mu^2 c - T(c, \mu))}{Q(c, \mu, T(c, \mu), d)},$$

307 where $T(c, \mu) = \sqrt{(1 - c + \mu^2 c)^2 + 4\mu^2 c^2}$ and P, Q are polynomials in four variables.

308 Here, at $c = (\mu^2 + 1)^{-1}$, the first term in the numerator is zero. Hence we conjecture that the peak of
309 the generalization error curve occurs near $c = (\mu^2 + 1)^{-1}$.

310 **Remark 1.** *Note that as $\mu \rightarrow 0$, we have that $4d\mu^2 c^2(2\mu^2 c - T(c, \mu)) \rightarrow 0$. We also note that, when
311 $\mu = 1$, we have that $2c - T(c, 1) = 0$. Thus, we see that for μ near 0 or 1, we should expect our
312 estimate of the location of the peak to be accurate.*

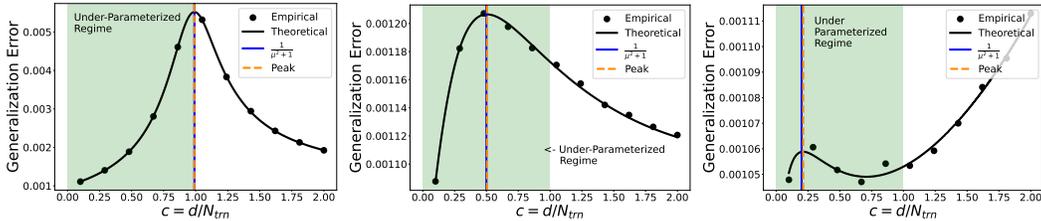


Figure 3: Figure showing the risk curve in the data scaling regime for different values of μ [(L) $\mu = 0.1$, (C) $\mu = 1$, (R) $\mu = 2$]. Here $\sigma_{trn} = \sqrt{N_{trn}}$, $\sigma_{tst} = \sqrt{N_{tst}}$, $d = 1000$, $N_{tst} = 1000$. For each empirical point, we ran at least 100 trials. More details can be found in Appendix G.

313 We numerically verify the predictions from Theorems 1, 2, 6. Figure 1 shows that the theoretically
314 predicted risk matches the numerical risk. Moreover, the curve has a single peak for $c < 1$. Thus,
315 verifying that double descent occurs in the under-parameterized regime. Finally, Figure 3 shows that
316 the location of the peak is near the conjectured location of $\frac{1}{\mu^2 + 1}$. This conjecture is further tested

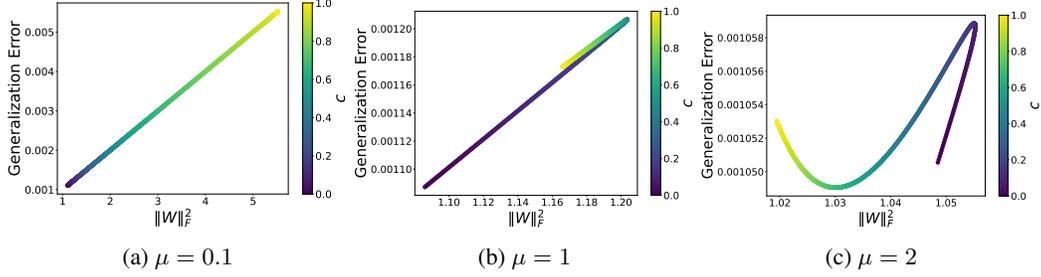


Figure 4: Figure showing generalization error versus $\|W_{opt}\|_F^2$ for the data scaling regime for three different values of μ . More details can be found in Appendix B and G.

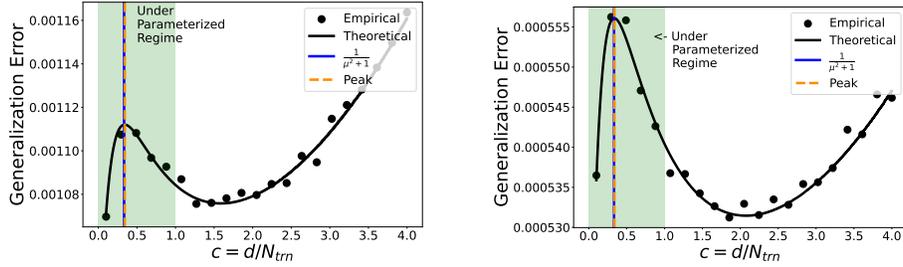


Figure 5: Figure showing that the shape of the risk curve in the data scaling regime depends on d [(L) $d = 1000$, (R) $d = 2000$]. Here $\mu = \sqrt{2}$, $\sigma_{trn} = \sqrt{N_{trn}}$, $\sigma_{tst} = \sqrt{N_{tst}}$, $N_{tst} = 1000$. Each empirical point is an average of at least 200 trials. More details can be found in Appendix G.

317 for a larger range of μ values in Appendix B. One similarity with prior work is that the peak in the
 318 generalization error or risk is corresponds to a peak in the norm of the estimator W_{opt} as seen in
 319 Figure 4 (i.e., the curve passes through the top right corner). The figure further shows, as conjectured
 320 in [39], that the double descent for the generalization error disappears when plotted as a function of
 321 $\|W_{opt}\|_F^2$ and, in some cases, recovers an approximation of the standard U shaped error curve.

322 **Risk curve shape depends on d .** Another interesting aspect of Theorem 2 is that it requires that d
 323 is large enough. Hence the shape of the risk curve depends on d . Most results for the risk are in the
 324 asymptotic regime. While Theorems 1, 2, and 6 are also in the asymptotic regime, we see that the
 325 results are accurate even for (relatively) small values of d , N_{trn} . Figure 5 shows that the shape of the
 326 risk curve depends on the value of d . Both curves still have a peak at the same location.

327 **Parameter Scaling.** For many prior models, the data scaling and parameter scaling regimes are
 328 analogous in that the shape of the risk curve does not depend on which one is scaled. The shape is
 329 primarily governed by the aspect ratio c of the data matrix. However, we see significant differences
 330 between the parameter scaling and data scaling regimes for our setup. Figure 6 shows risk curves
 331 that differ from those in Figure 3. Further, while for small values of μ , double descent occurs in the
 332 under-parameterized regime, for larger values of μ , the risk is monotonically decreasing.⁴

333 Even more astonishing, as shown in Figure 7, is the fact that for larger values of μ , *there is still a*
 334 *peak* in the curve for the norm of the estimator $\|W_{opt}\|_F^2$. However, this *does not* translate to a peak
 335 in the risk curve. Thus, showing that the norm of the estimator increasing cannot solely result in the
 336 generalization error increasing. The following theorem provides a local maximum in the $\|W_{opt}\|_F^2$
 337 versus c curve for $c < 1$.

338 **Theorem 5 ($\|W_{opt}\|_F$ Peak).** *If $\sigma_{tst} = \sqrt{N_{tst}}$, $\sigma_{trn} = \sqrt{N_{trn}}$ and μ is such that $p(\mu) < 0$,
 339 then for N_{trn} large enough and $d = cN_{trn}$, we have that $\|W_{opt}\|_F$ has a local maximum in the
 340 under-parameterized regime. Specifically for $c \in ((\mu^2 + 1)^{-1}, 1)$.*

⁴This is verified for more values of μ in Appendix B.

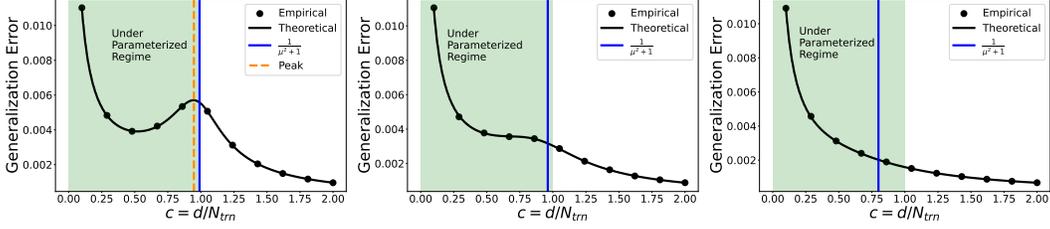


Figure 6: Figure showing the risk curves in the parameter scaling regime for different values of μ [(L) $\mu = 0.1$, (C) $\mu = 0.2$, (R) $\mu = 0.2$]. Here only the $\mu = 0.1$ has a local peak. Here $N_{trn} = N_{tst} = 1000$ and $\sigma_{trn} = \sigma_{tst} = \sqrt{1000}$. Each empirical point is an average of 100 trials.

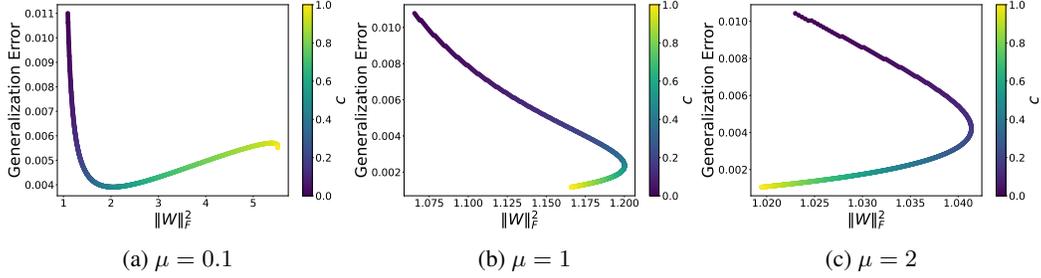


Figure 7: Figure showing generalization error versus $\|W_{opt}\|_F^2$ for the parameter scaling regime for three different values of μ . More details can be found in Appendix B.

341 B Peak Location and $\|W_{opt}\|_F$

342 **Theorem 6** (Peak Location). *If $\mu \in \mathbb{R}_{>0}$ is such that $p(\mu) < 0$, $\sigma_{trn}^2 = N_{trn} = d/c$ and $\sigma_{tst}^2 = N_{tst}$,*
 343 *then the partial derivative with respect to c of the risk $\mathcal{R}(c)$ from Theorem 1 can be written as*

$$\frac{\partial}{\partial c} \mathcal{R}(c, \mu) = \frac{(\mu^2 c + c - 1)P(c, \mu, T(c, \mu), d) + 4d\mu^2 c^2(2\mu^2 c - T(c, \mu))}{Q(c, \mu, T(c, \mu), d)},$$

344 *where $T(c, \mu) = \sqrt{(1 - c + \mu^2 c)^2 + 4\mu^2 c^2}$ and P, Q are polynomials in four variables.*

345 B.1 Peak Location for the Data Scaling Regime

346 We first look at the peak location conjecture for the data scaling regime. For this experiment, for 101
 347 different values of $\mu \in [0.1, 10]$ we compute the generalization error at 101 equally spaced points for

$$c \in \left(\frac{1}{2(\mu^2 + 1)}, \frac{2}{\mu^2 + 1} \right).$$

348 We then pick the c value that has the maximum from amongst these 101 values of c . We notice that
 349 this did not happen at the boundary. Hence it corresponded to a true local maximum. We plot this
 350 value of c on Figure 8 and compare this against $\frac{1}{\mu^2 + 1}$. As we can see from Figure 8, our conjectured
 351 location of the peak is an accurate estimate.

352 C Generalization error - bias and variance

353 For both the data scaling and parameter scaling regimes, Figures 9 and 10 show the bias, $\|W_{opt}\|$ and
 354 the generalization error. Here we see that our estimate is accurate.

355 D Training Error

356 As seen in the prior section, the peak happens in the interior of the under-parameterized regime and
 357 not on the border between the under-parameterized and over-parameterized regimes. In many prior
 358 works, the peak aligns with the interpolation point (i.e., zero training error). Theorem 7 derives a
 359 formula for the training error in the under-parameterized regime. Figure 11 plots the location of

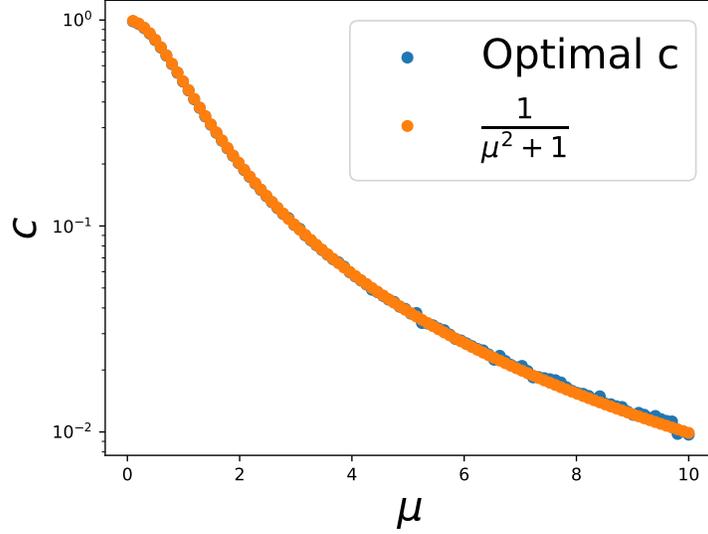


Figure 8: Figure showing the value of c where the peak occurs and the curve $1/(\mu^2 + 1)$

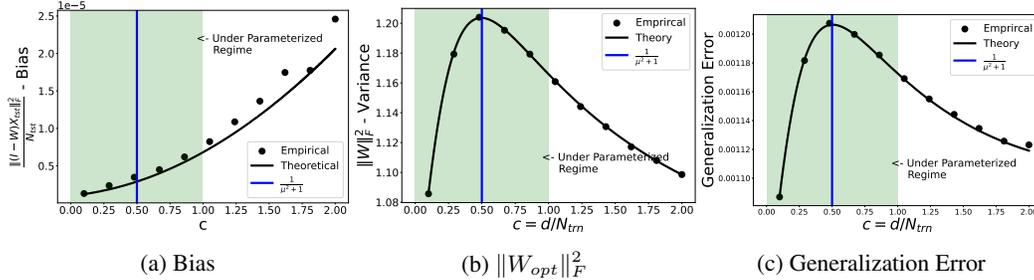


Figure 9: Figure showing the bias, $\|W_{opt}\|_F^2$, and the generalization error in the data scaling regime for $\mu = 1$, $\sigma_{trn} = \sqrt{N_{trn}}$, and $\sigma_{tst} = \sqrt{N_{tst}}$. Here $d = 1000$ and $N_{tst} = 1000$. For each empirical data point, we ran at least 100 trials. More details can be found in Appendix G.

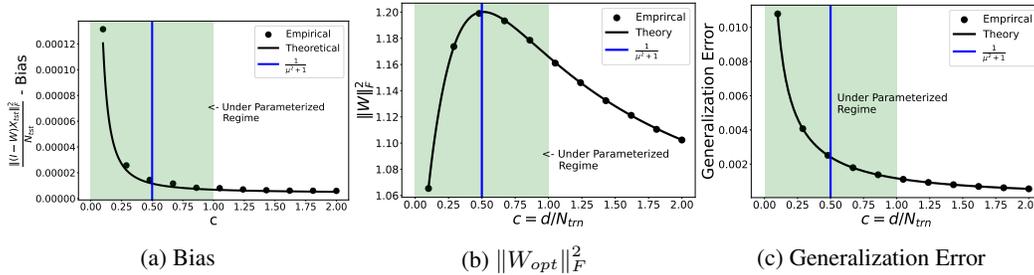


Figure 10: Figure showing the $\|W_{opt}\|_F^2$, and the generalization error in the parameter scaling regime for $\mu = 1$, $\sigma_{trn} = \sqrt{N_{trn}}$, and $\sigma_{tst} = \sqrt{N_{tst}}$. Here $N_{trn} = 1000$ and $N_{tst} = 1000$. For each empirical data point, we ran at least 100 trials. More details can be found in Appendix G.

360 the peak, the training error, and the third derivative of the training error. Here the figure shows that
 361 the training error curve does not signal the location of the peak in the generalization error curve.
 362 However, it shows that for the data scaling regime, the peak roughly corresponds to a local minimum
 363 of the third derivative of the training error.

364 **Theorem 7** (Training Error). *Let τ be as in Theorem 1. The training error for $c < 1$ is given by*

$$\mathbb{E}_{A_{trn}}[\|X_{trn} - W_{opt}(X_{trn} + A_{trn})\|_F^2] = \tau^{-2} (\sigma_{trn}^2 (1 - c \cdot T_1) + \sigma_{trn}^4 T_2) + o(1),$$

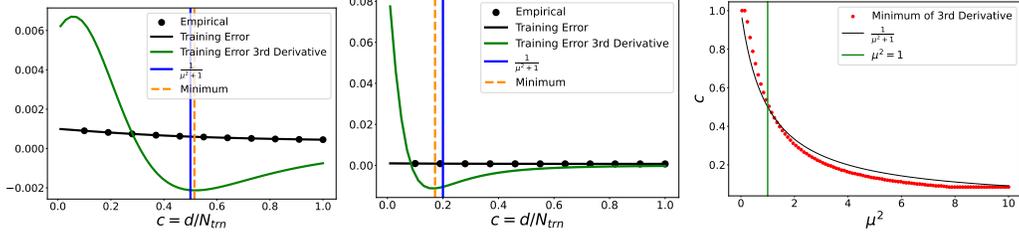


Figure 11: Figure showing the training error, the third derivative of the training error, and the location of the peak of the generalization error for different values of μ [(L) $\mu = 1$, (C) $\mu = 2$] for the data scaling regime. (R) shows the location of the local minimum of the third derivative and $\frac{1}{\mu^2+1}$.

365 where $T_1 = \frac{\mu^2}{2} \left(\frac{1+c+\mu^2c}{\sqrt{(1-c+\mu^2c)^2+4\mu^2c^2}} - 1 \right) + \frac{1}{2} + \frac{1+\mu^2c - \sqrt{(1-c+\mu^2c)^2+4c^2\mu^2}}{2c}$,

366 and $T_2 = \frac{(\mu^2c+c-1 - \sqrt{(1-c+\mu^2c)^2+4c^2\mu^2})^2(\mu^2c+c+1 - \sqrt{(1-c+\mu^2c)^2+4c^2\mu^2})}{2\sqrt{(1-c+\mu^2c)^2+4c^2\mu^2}}$.

367 E Regularization Trade-off

368 We analyze the trade-off between the two regularizers and the generalization error.

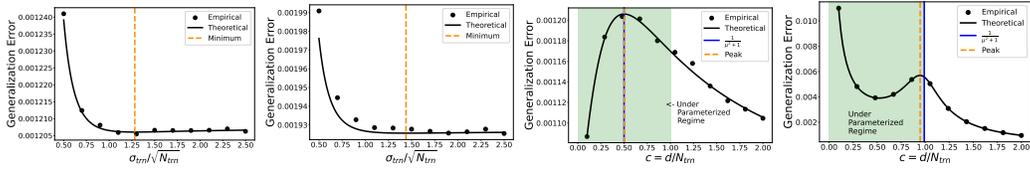


Figure 12: The first two figures show the σ_{trn} versus risk curve for $c = 0.5, \mu = 1$ and $c = 2, \mu = 0.1$ with $d = 1000$. The second two figures show the risk when training using the optimal σ_{trn} for the data scaling and parameter scaling regimes.

369 **Optimal σ_{trn} .** First, we fix μ and determine the optimal σ_{trn} . Figure 12 displays the generalization
 370 error versus σ_{trn}^2 curve. The figure shows that the error is initially large but then decreases until the
 371 optimal generalization error. The generalization error when using the optimal σ_{trn} is also shown in
 372 Figure 12. Here, unlike [24], picking the optimal value of σ_{trn} does not mitigate double descent.

373 **Proposition 1 (Optimal σ_{trn}).** *The optimal value of σ_{trn}^2 for $c < 1$ is given by*

$$\sigma_{trn}^2 = \frac{\sigma_{tst}^2 d [2c(\mu^2 + 1)^2 - 2T(c\mu^2 + c + 1) + 2(c\mu^2 - 2c + 1)] + N_{tst}(\mu^2 c^2 + c^2 + 1 - T)}{N_{tst}(c^3(\mu^2 + 1)^2 - T(\mu^2 c^2 + c^2 - 1) - 2c^2 - 1)}.$$

374 Additionally, it is interesting to determine how the optimal value of σ_{trn} depends on both μ and
 375 c . Figure 13 shows that for small values of μ (0.1, 0.5), as c changes, there exists an (inverted)
 376 double descent curve for the optimal value of σ_{trn} . However, unlike [33], for the data scaling

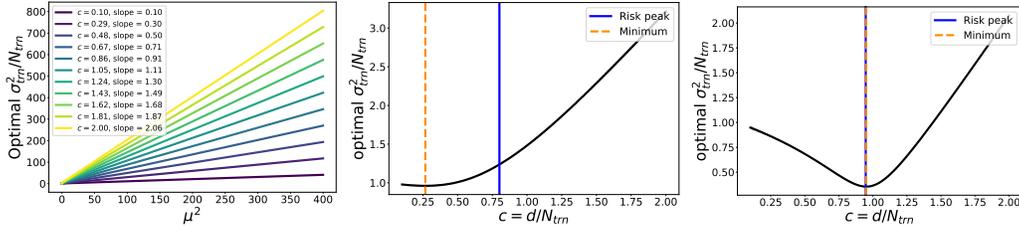


Figure 13: The first figure plots the optimal σ_{trn}^2/N_{trn} versus μ curve. The middle figure plots the optimal σ_{trn}^2/N_{trn} versus c in the data scaling regime for $\mu = 0.5$, and the last figure plots the optimal σ_{trn}^2/N_{trn} versus c in the parameter scaling regime for $\mu = 0.1$.

377 regime, the minimum of this double descent curve *does not match the location for the peak of the*
 378 *generalization error.* Further, as the amount of ridge regularization increases, the optimal amount of
 379 noise regularization decreases proportionally; optimal $\sigma_{trn}^2 \approx d\mu^2$. Thus, for higher values of ridge
 380 regularization, it is preferable to have higher-quality data.

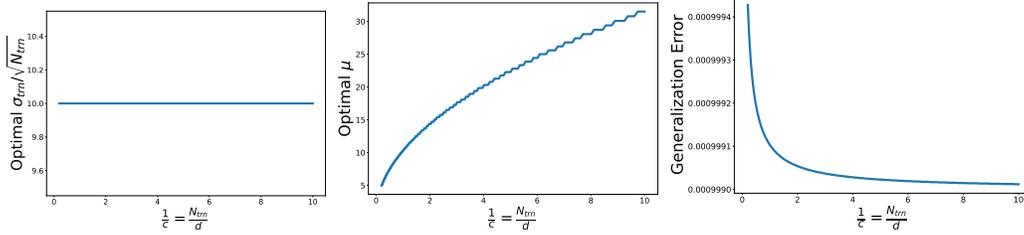


Figure 14: Trade-off between the regularizers. The left column is the optimal σ_{trn} , the central column is the optimal μ , and the right column is the generalization error for these parameter restrictions.

381 **Interaction Between the Regularizers.** The optimal values of μ and σ_{trn} are jointly computed
 382 using grid search for $\mu \in (0, 100]$ and $\sigma_{trn}/\sqrt{N_{trn}} \in (0, 10]$. Figure 14 shows the results. Specifi-
 383 cally, σ_{trn} is at the highest possible value (so best quality data), and then the model regularizes purely
 384 using the ridge regularizer. This results in a monotonically decreasing generalization error curve.
 385 Thus, in the data scaling model, *there is an implicit bias that favors one regularizer over the other.*
 386 Specifically, the model’s implicit bias *is to use higher quality data while using ridge regularization to*
 387 *regularize the model appropriately.* It is surprising that the two regularizers are not balanced.

388 E.1 Optimal Value of μ

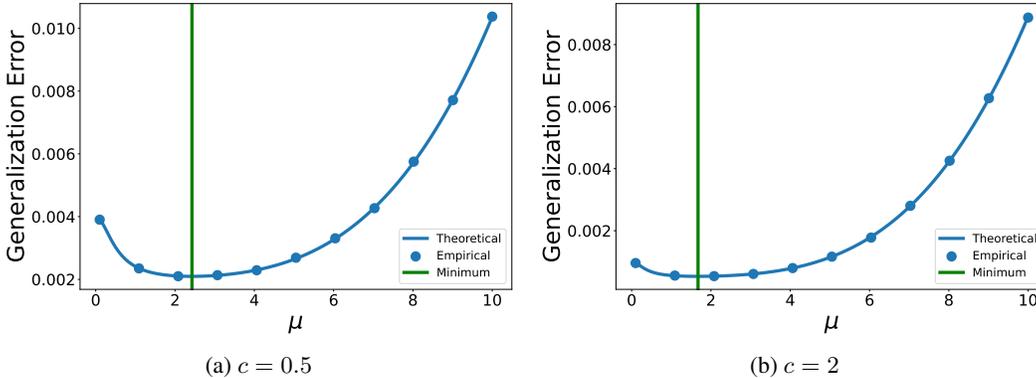


Figure 15: Figure showing the generalization error versus μ for $\sigma_{trn}^2 = N_{trn}$ and $\sigma_{tst}^2 = N_{tst}$.

389 We now explore the effect of fixing σ_{trn} and then changing μ . Figure 15, shows a U shaped curve
 390 for the generalization error versus μ , suggesting that there is an optimal value of μ , which should be
 391 used to minimize the generalization error.

392 Next, we compute the optimal value of μ using grid search and plot it against c . Figure 16 shows
 393 double descent for the optimal value of μ for small values of σ_{trn} . Thus for low SNR data we see
 394 double descent, but we do not for high SNR data.

395 Finally, for a given value of μ and c , we compute the optimal σ_{trn} . We then compute the generalization
 396 error (when using the optimal σ_{trn}) and plot the generalization error versus μ curve. Figure 17
 397 displays a very different trend from Figure 15. Instead of having a U -shaped curve, we have a
 398 monotonically decreasing generalization error curve. *This suggests that we can improve generalization*
 399 *by using higher-quality training while compensating for this by increasing the amount of ridge*
 400 *regularization.*

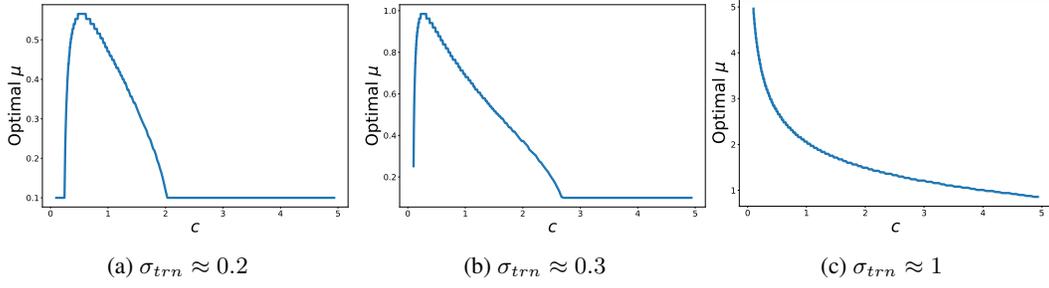


Figure 16: Figure for the optimal value of μ versus for different values of σ_{trn}

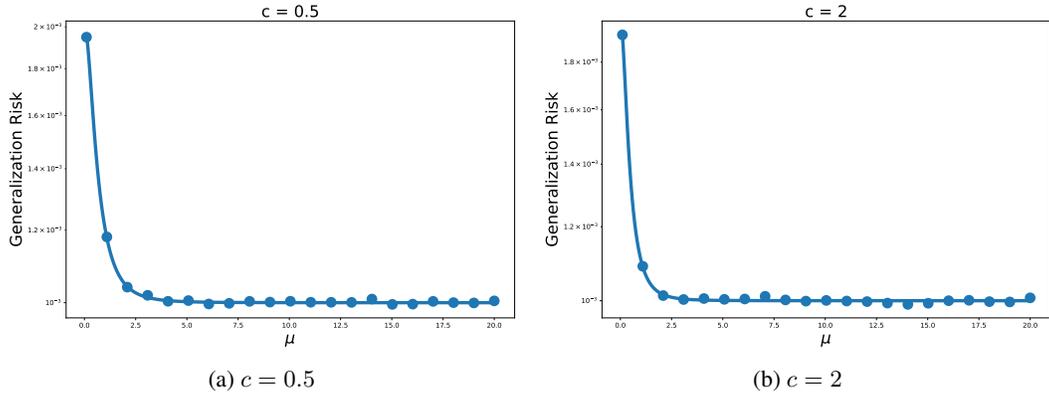


Figure 17: Figure showing the generalization error versus μ for the optimal σ_{trn}^2 and $\sigma_{tst}^2 = N_{tst}$.

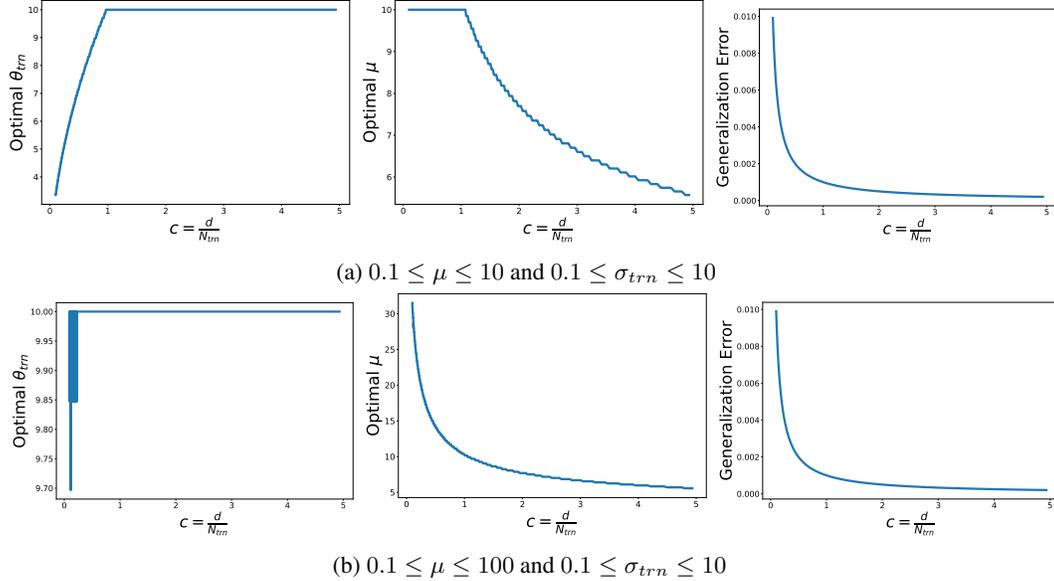


Figure 18: Trade-off between the regularizers. The left column is the optimal σ_{trn} , the central column is the optimal μ , and the right column is the generalization error for these parameter restrictions

401 **E.2 Trade-off in Parameter Scaling Regime**

402 Here we look at the trade-off between σ_{trn} and μ for the parameter scaling regime. We again see
 403 that the model implicitly prefers regularizing via ridge regularization and not via input data noise
 404 regularizer.

405 **F Proofs**

406 **F.1 Linear Regression**

407 We begin by noting,

$$\beta^T = (\beta_{opt}^T X + \xi_{trn}) X_{trn}^\dagger.$$

408 Thus, we have,

$$\begin{aligned} \|\beta\|^2 &= \text{Tr}(\beta^T \beta) \\ &= \text{Tr}(\beta_{opt}^T X_{trn} X_{trn}^\dagger (X_{trn}^\dagger)^T X_{trn} \beta_{opt}) + \text{Tr}(\xi_{trn} X_{trn}^\dagger (X_{trn}^\dagger)^T \xi_{trn}^T) + 2 \text{Tr}(\beta_{opt}^T X_{trn} X_{trn}^\dagger X_{trn}^\dagger)^T \xi_{trn}^T. \end{aligned}$$

409 Taking the expectation, with respect to ξ_{trn} , we see that the last term vanishes.

410 Letting $X_{trn} = U_X \Sigma_X V_X^T$. We see that using the rotational invariance of X , U_X , V_X are independent and uniformly random. Thus, $s := \beta_{opt}^T U_X$ is a uniformly random unit vector.

412 Thus, we see,

$$\mathbb{E}_{X_{trn}, \xi_{trn}} \left[\text{Tr}(\beta_{opt}^T X_{trn} X_{trn}^\dagger (X_{trn}^\dagger)^T X_{trn} \beta_{opt}) \right] = \sum_{i=1}^{\min(d, N_{trn})} \mathbb{E}[s_i^2] = \min\left(1, \frac{1}{c}\right)$$

413 Similarly, we see,

$$\mathbb{E}_{X_{trn}, \xi_{trn}} \left[\xi_{trn} X_{trn}^\dagger (X_{trn}^\dagger)^T \xi_{trn}^T \right] = \sum_{i=1}^{\min(d, N_{trn})} \mathbb{E} \left[\frac{1}{\sigma_i(X_{trn})^2} \right]$$

414 Multiplying and dividing by d , normalizes the singular values squared of X_{trn} so that the limiting distribution is the Marchenko Pastur distribution with shape c . Thus, we can estimate using Lemma 5 from Sonthalia and Nadakuditi [33] to get,

$$\begin{cases} \frac{c}{1-c} + o(1) & c < 1 \\ \frac{1}{c-1} + o(1) & c > 1 \end{cases}.$$

417 Finally, the cross-term has an expectation equal to zero. Thus,

$$\mathbb{E}_{X_{trn}, \xi_{trn}} [\|\beta_{opt}\|^2] = \begin{cases} 1 + \frac{c}{1-c} & c < 1 \\ \frac{1}{c} + \frac{1}{c-1} & c > 1 \end{cases}$$

418 Then we have,

$$\beta^T \beta_{opt} = \beta_{opt}^T X_{trn} X_{trn}^\dagger \beta_{opt} + \xi_{trn} X_{trn}^\dagger \beta_{opt}$$

419 The second term has an expectation equal to zero, and the first term is similar to before and has an expectation equal to $\min\left(1, \frac{1}{c}\right)$.

421 **F.2 Proofs for Theorem 1**

422 The proof structure closely follows that of [33].

423 **F.2.1 Step 1: Decompose the error into bias and variance terms.**

424 First, we decompose the error. Since we are not in the supervised learning setup, we do not have standard definitions of bias/variance. However, we will call the following terms the bias/variance of the model. First, we recall the following from [33].

427 **Lemma 1** (Sonthalia and Nadakuditi [33]). *If A_{tst} has mean 0 entries and A_{tst} is independent of X_{tst} and W , then*

$$\mathbb{E}_{A_{tst}} [\|X_{tst} - WY_{tst}\|_F^2] = \underbrace{\mathbb{E}_{A_{tst}} [\|X_{tst} - WX_{tst}\|_F^2]}_{\text{Bias}} + \underbrace{\mathbb{E}_{A_{tst}} [\|WA_{tst}\|_F^2]}_{\text{Variance}}. \quad (4)$$

429 **F.2.2 Step 2: Formula for W_{opt}**

430 Here, we compute the explicit formula for W_{opt} in Problem 1. Let $\hat{A}_{trn} = [A_{trn} \quad \mu I]$, $\hat{X}_{trn} =$
431 $[X_{trn} \quad 0]$, and $\hat{Y}_{trn} = \hat{X}_{trn} + \hat{A}_{trn}$. Then solving $\arg \min_W \|X_{trn} - WY_{trn}\|_F^2 + \mu^2 \|W\|_F^2$ is
432 equivalent to solving $\arg \min_W \|\hat{X}_{trn} - W\hat{Y}_{trn}\|_F^2$. Thus, $W_{opt} = \arg \min_W \|\hat{X}_{trn} - W\hat{Y}_{trn}\|_F^2 =$
433 $\hat{X}_{trn}\hat{Y}_{trn}^\dagger$. Expanding this out, we get the following formula for \hat{W} . Let \hat{u} be the left singular
434 vector and \hat{v}_{trn} be the right singular vectors of \hat{X}_{trn} . Note that the left singular does not change
435 after ridge regularization, so $\hat{u} = u$. Let $\hat{h} = \hat{v}_{trn}^T \hat{A}_{trn}^\dagger$, $\hat{k} = \hat{A}_{trn}^\dagger u$, $\hat{s} = (I - \hat{A}_{trn} \hat{A}_{trn}^\dagger)u$,
436 $\hat{t} = \hat{v}_{trn}(I - \hat{A}_{trn} \hat{A}_{trn}^\dagger)$, $\hat{\gamma} = 1 + \sigma_{trn} \hat{v}_{trn}^T \hat{A}_{trn}^\dagger u$, $\hat{\tau} = \sigma_{trn}^2 \|\hat{t}\|^2 \|\hat{k}\|^2 + \hat{\gamma}^2$.

437 **Proposition 2.** *If $\hat{\gamma} \neq 0$ and A_{trn} has full rank then*

$$W_{opt} = \frac{\sigma_{trn} \hat{\gamma}}{\hat{\tau}} u \hat{h} + \frac{\sigma_{trn}^2 \|\hat{t}\|^2}{\hat{\tau}} u \hat{k}^T \hat{A}_{trn}^\dagger.$$

438 *Proof.* Here we know that u is arbitrary. We have that \hat{A}_{trn} has full rank. Thus, the rank of \hat{A}_{trn} is
439 d , and the range of \hat{A}_{trn} is the whole space. Thus, u lives in the range of \hat{A}_{trn} . In this case, we want
440 Theorem 3 from [40]. We define

$$\hat{p} = -\frac{\sigma_{trn}^2 \|\hat{k}\|^2}{\hat{\gamma}} \hat{t}^T - \sigma_{trn} \hat{k} \text{ and } \hat{q}^T = -\frac{\sigma_{trn} \|\hat{t}\|^2}{\hat{\gamma}} \hat{k}^T \hat{A}_{trn}^\dagger - \hat{h}.$$

441 Then we have,

$$(\hat{A}_{trn} + \sigma_{trn} u \hat{v}_{trn}^T)^\dagger = \hat{A}_{trn}^\dagger + \frac{\sigma_{trn}}{\hat{\gamma}} \hat{t}^T \hat{k}^T \hat{A}_{trn}^\dagger - \frac{\hat{\gamma}}{\hat{\tau}} \hat{p} \hat{q}^T.$$

442 Note that, by our assumptions, we have $\hat{t} = \hat{v}_{trn}(I - \hat{A}_{trn} \hat{A}_{trn}^\dagger)$, and $(I - \hat{A}_{trn} \hat{A}_{trn}^\dagger)$ is a projection
443 matrix, thus

$$\begin{aligned} \hat{v}_{trn}^T \hat{t}^T &= \hat{v}_{trn}^T (I - \hat{A}_{trn} \hat{A}_{trn}^\dagger)^T \hat{v}_{trn} \\ &= \hat{v}_{trn}^T (I - \hat{A}_{trn} \hat{A}_{trn}^\dagger)^T (I - \hat{A}_{trn} \hat{A}_{trn}^\dagger)^T \hat{v}_{trn}^T. \end{aligned}$$

444 To compute $W_{opt} = \hat{X}_{trn}(\hat{X}_{trn} + \hat{A}_{trn})^\dagger = \sigma_{trn} u \hat{v}_{trn}^T (\hat{A}_{trn} + \sigma_{trn} u \hat{v}_{trn}^T)^\dagger$, using $\hat{\gamma} - 1 =$
445 $\sigma_{trn} \hat{v}_{trn}^T \hat{A}_{trn}^\dagger u = \sigma_{trn} \hat{h} u$, we multiply this through.

$$\begin{aligned} \sigma_{trn} u \hat{v}_{trn}^T (\hat{A}_{trn} + \sigma_{trn} u \hat{v}_{trn}^T)^\dagger &= \sigma_{trn} u \hat{v}_{trn}^T (\hat{A}_{trn}^\dagger + \frac{\sigma_{trn}}{\hat{\gamma}} \hat{t}^T \hat{k}^T \hat{A}_{trn}^\dagger - \frac{\hat{\gamma}}{\hat{\tau}} \hat{p} \hat{q}^T) \\ &= \sigma_{trn} u \hat{h} + \frac{\sigma_{trn}^2 \|\hat{t}\|^2}{\hat{\gamma}} u \hat{k}^T \hat{A}_{trn}^\dagger \\ &\quad + \frac{\sigma_{trn} \hat{\gamma}}{\hat{\tau}} u \hat{v}_{trn}^T \left(\frac{\sigma_{trn}^2 \|\hat{k}\|^2}{\hat{\gamma}} \hat{t}^T + \sigma_{trn} \hat{k} \right) \hat{q}^T \\ &= \sigma_{trn} u \hat{h} + \frac{\sigma_{trn}^2 \|\hat{t}\|^2}{\hat{\gamma}} u \hat{k}^T \hat{A}_{trn}^\dagger + \frac{\sigma_{trn}^3 \|\hat{k}\|^2 \|\hat{t}\|^2}{\hat{\tau}} u \hat{q}^T \\ &\quad + \frac{\sigma_{trn} \hat{\gamma} (\hat{\gamma} - 1)}{\hat{\tau}} u \hat{q}^T. \end{aligned}$$

446 Then we have,

$$\begin{aligned} \frac{\sigma_{trn}^3 \|\hat{k}\|^2 \|\hat{t}\|^2}{\hat{\tau}} u \hat{q}^T &= \frac{\sigma_{trn}^3 \|\hat{k}\|^2 \|\hat{t}\|^2}{\hat{\tau}} u \left(-\frac{\sigma_{trn} \|\hat{t}\|^2}{\hat{\gamma}} \hat{k}^T \hat{A}_{trn}^\dagger - \hat{h} \right) \\ &= -\frac{\sigma_{trn}^4 \|\hat{k}\|^2 \|\hat{t}\|^4}{\hat{\tau} \hat{\gamma}} u \hat{k}^T \hat{A}_{trn}^\dagger - \frac{\sigma_{trn}^3 \|\hat{k}\|^2 \|\hat{t}\|^2}{\hat{\tau}} u \hat{h} \end{aligned}$$

447 and

$$\begin{aligned} \frac{\sigma_{trn}\hat{\gamma}(\hat{\gamma}-1)}{\hat{\tau}}u\hat{q}^T &= \frac{\sigma_{trn}\hat{\gamma}(\hat{\gamma}-1)}{\hat{\tau}}u\left(-\frac{\sigma_{trn}\|\hat{t}\|^2}{\hat{\gamma}}\hat{k}^T\hat{A}_{trn}^\dagger - \hat{h}\right) \\ &= -\frac{\sigma_{trn}^2\|\hat{t}\|^2(\hat{\gamma}-1)}{\hat{\tau}}u\hat{k}^T\hat{A}_{trn}^\dagger - \frac{\sigma_{trn}\hat{\gamma}(\hat{\gamma}-1)}{\hat{\tau}}u\hat{h}. \end{aligned}$$

448 Substituting back in and collecting like terms, we get,

$$\begin{aligned} \sigma_{trn}u\hat{v}_{trn}^T(\hat{A}_{trn} + \sigma_{trn}u\hat{v}_{trn}^T)^\dagger &= \sigma_{trn}\left(1 - \frac{\sigma_{trn}^2\|\hat{k}\|^2\|\hat{t}\|^2}{\hat{\tau}} - \frac{\hat{\gamma}(\hat{\gamma}-1)}{\hat{\tau}}\right)u\hat{h} + \\ &\quad \sigma_{trn}^2\left(\frac{\|\hat{t}\|^2}{\hat{\gamma}} - \frac{\sigma_{trn}^2\|\hat{k}\|^2\|\hat{t}\|^4}{\hat{\tau}\hat{\gamma}} - \frac{\|\hat{t}\|^2(\hat{\gamma}-1)}{\hat{\tau}}\right)u\hat{k}^T\hat{A}_{trn}^\dagger. \end{aligned}$$

449 We can then simplify the constants as follows.

$$1 - \frac{\sigma_{trn}^2\|\hat{k}\|^2\|\hat{t}\|^2}{\hat{\tau}} - \frac{\hat{\gamma}(\hat{\gamma}-1)}{\hat{\tau}} = \frac{\hat{\tau} - \sigma_{trn}^2\|\hat{k}\|^2\|\hat{t}\|^2 - \hat{\gamma}^2 + \hat{\gamma}}{\hat{\tau}} = \frac{\hat{\gamma}}{\hat{\tau}}$$

450 and

$$\frac{\|\hat{t}\|^2}{\hat{\gamma}} - \frac{\sigma_{trn}^2\|\hat{k}\|^2\|\hat{t}\|^4}{\hat{\tau}\hat{\gamma}} - \frac{\|\hat{t}\|^2(\hat{\gamma}-1)}{\hat{\tau}} = \frac{\|\hat{t}\|^2\left(\hat{\tau} - \sigma_{trn}^2\|\hat{k}\|^2\|\hat{t}\|^2 - \hat{\gamma}^2 + \hat{\gamma}\right)}{\hat{\tau}\hat{\gamma}} = \frac{\|\hat{t}\|^2}{\hat{\tau}}.$$

451 This gives us the result. □

452 **F.2.3 Step 3: Decompose the terms into a sum of various trace terms.**

453 For the bias and variance terms, we have the following two Lemmas.

454 **Lemma 2.** *If W_{opt} is the solution to Equation 1, then*

$$X_{tst} - W_{opt}X_{tst} = \frac{\hat{\gamma}}{\hat{\tau}}X_{tst}.$$

455 *Proof.* To see this, note that we have $N_{trn} + M > M$.

$$\begin{aligned} X_{tst} - W_{opt}X_{tst} &= X_{tst} - \frac{\sigma_{trn}\hat{\gamma}}{\hat{\tau}}u\hat{h}u_{tst}^T - \frac{\sigma_{trn}^2\|\hat{t}\|^2}{\hat{\tau}}u\hat{k}^T\hat{A}_{trn}^\dagger u_{tst}^T \\ &= X_{tst} - \frac{\hat{\sigma}_{trn}\hat{\gamma}}{\hat{\tau}}u\hat{v}_{trn}^T\hat{A}_{trn}^\dagger u_{tst}^T - \frac{\sigma_{trn}^2\|\hat{t}\|^2}{\hat{\tau}}u\hat{k}^T\hat{A}_{trn}^\dagger u_{tst}^T. \end{aligned}$$

456 Note that $\hat{\gamma} = 1 + \sigma_{trn}\hat{v}_{trn}^T\hat{A}_{trn}^\dagger u$. Thus, we have that $\sigma_{trn}\hat{v}_{trn}^T\hat{A}_{trn}^\dagger u = \hat{\gamma} - 1$. Substituting this
457 into the second term, we get,

$$X_{tst} - W_{opt}X_{tst} = X_{tst} - \frac{\hat{\gamma}(\hat{\gamma}-1)}{\hat{\tau}}u_{tst}^T - \frac{\sigma_{trn}^2\|\hat{t}\|^2}{\hat{\tau}}u\hat{k}^T\hat{A}_{trn}^\dagger u_{tst}^T.$$

458 For the third term, since $\hat{k} = \hat{A}_{trn}^\dagger u$, $\hat{k}^T\hat{A}_{trn}^\dagger u = \hat{k}^T\hat{k} = \|\hat{k}\|^2$. Substituting this into the expression,
459 we get that

$$X_{tst} - W_{opt}X_{tst} = X_{tst} - \frac{\hat{\gamma}(\hat{\gamma}-1)}{\hat{\tau}}u_{tst}^T - \frac{\sigma_{trn}^2\|\hat{t}\|^2\|\hat{k}\|^2}{\hat{\tau}}u_{tst}^T.$$

460 Since $X_{tst} = u_{tst}^T$, we get,

$$X_{tst} - W_{opt}X_{tst} = X_{tst} \left(1 - \frac{\hat{\gamma}(\hat{\gamma}-1)}{\hat{\tau}} - \frac{\sigma_{trn}^2\|\hat{t}\|^2\|\hat{k}\|^2}{\hat{\tau}}\right).$$

461 Simplify the constants using $\hat{\tau} = \sigma_{trn}^2 \|\hat{t}\|^2 \|\hat{k}\|^2 + \hat{\gamma}^2$, we get,

$$\frac{\hat{\tau} + \hat{\gamma} - \hat{\gamma}^2 - \sigma_{trn}^2 \|\hat{t}\|^2 \|\hat{k}\|^2}{\hat{\tau}} = \frac{\hat{\gamma}}{\hat{\tau}}.$$

462

□

463 **Lemma 3** (Sonthalia and Nadakuditi [33]). *If the entries of A_{tst} are independent with mean 0, and*
 464 *variance $1/d$, then we have that $\mathbb{E}_{A_{tst}}[\|W_{opt} A_{tst}\|^2] = \frac{N_{tst}}{d} \|W_{opt}\|^2$.*

465 **Lemma 4.** *If $\hat{\gamma} \neq 0$ and A_{trn} has full rank, then we have that*

$$\|W_{opt}\|_F^2 = \frac{\sigma_{trn}^2 \hat{\gamma}^2}{\hat{\tau}^2} \text{Tr}(\hat{h}^T \hat{h}) + 2 \frac{\sigma_{trn}^3 \|\hat{t}\|^2 \hat{\gamma}}{\hat{\tau}^2} \text{Tr}(\hat{h}^T \hat{k}^T \hat{A}_{trn}^\dagger) + \frac{\sigma_{trn}^4 \|\hat{t}\|^4}{\hat{\tau}^2} \underbrace{\text{Tr}((\hat{A}_{trn}^\dagger)^T \hat{k} \hat{k}^T \hat{A}_{trn}^\dagger)}_{\rho}.$$

466 *Proof.* We have

$$\begin{aligned} \|W_{opt}\|_F^2 &= \text{Tr}(W_{opt}^T W_{opt}) \\ &= \text{Tr} \left(\left(\frac{\sigma_{trn} \hat{\gamma}}{\hat{\tau}} u \hat{h} + \frac{\sigma_{trn}^2 \|\hat{t}\|^2}{\hat{\tau}} u \hat{k}^T \hat{A}_{trn}^\dagger \right)^T \left(\frac{\sigma_{trn} \hat{\gamma}}{\hat{\tau}} u \hat{h} + \frac{\sigma_{trn}^2 \|\hat{t}\|^2}{\hat{\tau}} u \hat{k}^T \hat{A}_{trn}^\dagger \right) \right) \\ &= \frac{\sigma_{trn}^2 \hat{\gamma}^2}{\hat{\tau}^2} \text{Tr}(\hat{h}^T u^T u \hat{h}) + 2 \frac{\sigma_{trn}^3 \|\hat{t}\|^2 \hat{\gamma}}{\hat{\tau}^2} \text{Tr}(\hat{h}^T u^T u \hat{k}^T \hat{A}_{trn}^\dagger) \\ &\quad + \frac{\sigma_{trn}^4 \|\hat{t}\|^4}{\hat{\tau}^2} \text{Tr}((\hat{A}_{trn}^\dagger)^T \hat{k} u^T u \hat{k}^T \hat{A}_{trn}^\dagger) \\ &= \frac{\sigma_{trn}^2 \hat{\gamma}^2}{\hat{\tau}^2} \text{Tr}(\hat{h}^T \hat{h}) + 2 \frac{\sigma_{trn}^3 \|\hat{t}\|^2 \hat{\gamma}}{\hat{\tau}^2} \text{Tr}(\hat{h}^T \hat{k}^T \hat{A}_{trn}^\dagger) + \frac{\sigma_{trn}^4 \|\hat{t}\|^4}{\hat{\tau}^2} \text{Tr}((\hat{A}_{trn}^\dagger)^T \hat{k} \hat{k}^T \hat{A}_{trn}^\dagger). \end{aligned}$$

467 Where the last inequality is true due to the fact that $\|u\|^2 = 1$.

□

468 F.2.4 Step 4: Estimate With Random Matrix Theory

469 **Lemma 5.** *Let A be a $p \times q$ matrix and let $\hat{A} = [A \quad \mu I] \in \mathbb{R}^{p \times q+p}$. Suppose $A = U \Sigma V^T$ be the*
 470 *singular value decomposition of A . If $\hat{A} = \hat{U} \hat{\Sigma} \hat{V}^T$ is the singular value decomposition of \hat{A} , then*
 471 *$\hat{U} = U$ and if $p < q$*

$$\hat{\Sigma} = \begin{bmatrix} \sqrt{\sigma_1(A)^2 + \mu^2} & 0 & \cdots & 0 \\ 0 & \sqrt{\sigma_2(A)^2 + \mu^2} & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & \sqrt{\sigma_p(A)^2 + \mu^2} \end{bmatrix} \in \mathbb{R}^{p \times p},$$

472 and

$$\hat{V} = \begin{bmatrix} V_{1:p} \Sigma \hat{\Sigma}^{-1} \\ \mu U \hat{\Sigma}^{-1} \end{bmatrix} \in \mathbb{R}^{q+p \times p}.$$

473 Here $V_{1:p}$ are the first p columns of V .

474 *Proof.* Since $p < q$, we have that $U \in \mathbb{R}^{p \times p}$, $\Sigma \in \mathbb{R}^{p \times p}$ are invertible. Here also consider the form
 475 of the SVD in which $V^T \in \mathbb{R}^{p \times q}$.

476 We start by noting that $\hat{U} \hat{\Sigma}^2 \hat{U}^T = \hat{A} \hat{A}^T = A A^T + \mu^2 I = U(\Sigma^2 + \mu^2 I_p) U^T$. Thus, we immediately
 477 see that $\sigma_i(\hat{A})^2 = \sigma_i(A)^2 + \mu^2$ and that $\hat{U} = U$.

478 Finally, we see,

$$\hat{V}^T = \hat{\Sigma}^{-1} U^T \hat{A} = [\hat{\Sigma}^{-1} \Sigma V_{1:p}^T \quad \mu \hat{\Sigma}^{-1} U^T]$$

479

□

480 **Lemma 6.** Let A be a $p \times q$ matrix and let $\hat{A} = [A \quad \mu I] \in \mathbb{R}^{p \times q+p}$. Suppose $A = U\Sigma V^T$ be the
481 singular value decomposition of A . If $\hat{A} = \hat{U}\hat{\Sigma}\hat{V}^T$ is the singular value decomposition of \hat{A} , then
482 $\hat{U} = U$ and if $p > q$

$$\hat{\Sigma} = \begin{bmatrix} \sqrt{\sigma_1(A)^2 + \mu^2} & 0 & \cdots & 0 & \cdots & 0 \\ 0 & \sqrt{\sigma_2(A)^2 + \mu^2} & & 0 & & \\ \vdots & & \ddots & \vdots & & \vdots \\ 0 & 0 & \cdots & \sqrt{\sigma_q(A)^2 + \mu^2} & & 0 \\ \vdots & & & & \mu & \\ 0 & 0 & \cdots & 0 & \cdots & 0 \\ \vdots & & & & & \ddots \\ 0 & 0 & \cdots & 0 & \cdots & 0 \\ & & & & & \mu \end{bmatrix} \in \mathbb{R}^{p \times p}.$$

483 Here we will denote the upper left $q \times q$ block by C . Further,

$$\hat{V} = \begin{bmatrix} V\Sigma_{1:q,1:q}^T C^{-1} & 0 \\ \mu U_{1:q} C^{-1} & U_{q+1:p} \end{bmatrix} \in \mathbb{R}^{q+p \times p}.$$

484 *Proof.* Since $p > q$, we have that $U \in \mathbb{R}^{p \times p}$ and we have that $\Sigma \in \mathbb{R}^{p \times q}$. Here $V^T \in \mathbb{R}^{q \times q}$ is
485 invertible.

486 We start with nothing,

$$\hat{U}\hat{\Sigma}^2\hat{U}^T = \hat{A}\hat{A}^T = AA^T + \mu^2 I = U \left(\begin{bmatrix} \Sigma_{1:q,1:q}^2 & 0 \\ 0 & 0_{q-p} \end{bmatrix} + \mu^2 I_q \right) U^T.$$

487 Thus, we immediately see that for $i = 1, \dots, p$ $\sigma_i(\hat{A})^2 = \sigma_i(A)^2 + \mu^2$ and for $i = p+1, \dots, q$, we
488 have that $\sigma_i(\hat{A})^2 = \mu^2$ and that $\hat{U} = U$.

489 Then, we see,

$$\hat{V}^T = \hat{\Sigma}^{-1} U^T \hat{A} = \begin{bmatrix} \hat{\Sigma}^{-1} \Sigma V^T & \mu \hat{\Sigma}^{-1} U^T \end{bmatrix}.$$

490 Note that Σ has 0 for the last $p - q$ entries. Thus,

$$\hat{\Sigma}^{-1} \Sigma V = \begin{bmatrix} C^{-1} \Sigma_{1:q,1:q} V \\ 0_{q-p,q} \end{bmatrix}.$$

491 Similarly, due to the structure of $\hat{\Sigma}$, we see,

$$\mu \hat{\Sigma}^{-1} U^T = [\mu C^{-1} U_{1:q}^T \quad \mu \frac{1}{\mu} U_{q+1:p}^T].$$

492 □

493 **Lemma 7.** Suppose A is an p by q matrix such that $p < q$, the entries of A are independent and have
494 mean 0, variance $1/p$, and bounded fourth moment. Let $c = p/q$. Let $\hat{A} = [A \quad \mu I] \in \mathbb{R}^{p \times q+p}$. Let
495 $W_p = \hat{A}\hat{A}^T$ and let $W_q = \hat{A}^T \hat{A}$. Suppose λ_p is a random non-zero eigenvalue from the largest p
496 eigenvalues of W_p , and λ_q is a random non-zero eigenvalue of W_q . Then

497 1. $\mathbb{E} \left[\frac{1}{\lambda_p} \right] = \mathbb{E} \left[\frac{1}{\lambda_q} \right] = \frac{\sqrt{(1+\mu^2 c - c)^2 + 4\mu^2 c^2} - 1 - \mu^2 c + c}{2\mu^2 c} + o(1).$

498 2. $\mathbb{E} \left[\frac{1}{\lambda_p^2} \right] = \mathbb{E} \left[\frac{1}{\lambda_q^2} \right] = \frac{\mu^2 c^2 + c^2 + \mu^2 c - 2c + 1}{2\mu^4 c \sqrt{4\mu^2 c^2 + (1 - c + \mu^2 c)^2}} + \frac{1}{2\mu^4} \left(1 - \frac{1}{c} \right) + o(1).$

499 *Proof.* First, we note that the non-zero eigenvalues of W_p and W_q are the same. Hence we focus on
500 W_p . W_p is nearly a Wishart matrix but is not normalized by the correct value. However, cW_p does
501 have the correct normalization.

502 Due to the assumptions on A , we have that the eigenvalues of cAA^T converge to the Marchenko-
503 Pastur. Hence since the eigenvalues of cW_p are

$$(c\lambda_p)_i = c\sigma_i(A)^2 + c\mu^2,$$

504 we can estimate them by estimating $c\sigma_i(A)^2$ with the Marchenko-Pastur [41–45]. In particular, we
 505 want the expectation of the inverse. We need to use the Stieljes transform. We know that if $m_c(z)$ is
 506 the Stieljes transform for the Marchenko-Pastur with shape parameter c , then if λ is sampled from the
 507 Marchenko-Pastur distribution, then

$$m_c(z) = \mathbb{E}_\lambda \left[\frac{1}{\lambda - z} \right].$$

508 Thus, we have that the expected inverse of the eigenvalue can be approximated $m(-c\mu^2)$. We know
 509 that the Stieljes transform:

$$m_c(z) = -\frac{1 - z - c - \sqrt{(1 - z - c)^2 - 4cz}}{-2zc}.$$

510 Thus, we have,

$$\mathbb{E} \left[\frac{1}{c\lambda_p} \right] = m(-c\mu^2) = \frac{\sqrt{(1 + \mu^2c - c)^2 + 4\mu^2c^2} - 1 - \mu^2c + c}{2\mu^2c^2}.$$

511 Canceling $1/c$ from both sides, we get,

$$\mathbb{E} \left[\frac{1}{\lambda_p} \right] = \frac{\sqrt{(1 + \mu^2c - c)^2 + 4\mu^2c^2} - 1 - \mu^2c + c}{2\mu^2c}.$$

512 Then for the estimate of $\mathbb{E} [1/\lambda_p^2]$, we need to compute the derivative of the $m_c(z)$ and evaluate it at
 513 $-c\mu^2$. Hence, we see,

$$m'_c(z) = \frac{(c - z + \sqrt{-4cz + (1 - c - z)^2} - 1)(c + z + \sqrt{-4cz + (1 - c - z)^2} - 1)}{4cz^2 \sqrt{-4cz + (1 - c - z)^2}}.$$

514 Thus,

$$\begin{aligned} \mathbb{E} \left[\frac{1}{c^2\lambda_p^2} \right] &= m'_c(-c\mu^2) \\ &= \frac{(c + \mu^2c + \sqrt{4\mu^2c^2 + (1 - c + \mu^2c)^2} - 1)(c - \mu^2c + \sqrt{4\mu^2c^2 + (1 - c + \mu^2c)^2} - 1)}{4\mu^4c^3 \sqrt{4\mu^2c^2 + (1 - c + \mu^2c)^2}}. \end{aligned}$$

515 Canceling the $1/c^2$ from both sides, we get,

$$\mathbb{E} \left[\frac{1}{\lambda_p^2} \right] = \frac{(c + \mu^2c + \sqrt{4\mu^2c^2 + (1 - c + \mu^2c)^2} - 1)(c - \mu^2c + \sqrt{4\mu^2c^2 + (1 - c + \mu^2c)^2} - 1)}{4\mu^4c \sqrt{4\mu^2c^2 + (1 - c + \mu^2c)^2}}.$$

516 Multiplying out and simplifying

$$\mathbb{E} \left[\frac{1}{\lambda_p^2} \right] = \frac{\mu^2c^2 + c^2 + \mu^2c - 2c + 1}{2\mu^4c \sqrt{4\mu^2c^2 + (1 - c + \mu^2c)^2}} + \frac{1}{2\mu^4} \left(1 - \frac{1}{c} \right).$$

517 □

518 **Lemma 8.** Suppose A is an p by q matrix such that $p > q$, the entries of A are independent and have
 519 mean 0, variance $1/p$, and bounded fourth moment. Let $c = p/q$. Let $\hat{A} = [A \quad \mu I] \in \mathbb{R}^{p \times q+p}$. Let
 520 $W_p = \hat{A}\hat{A}^T$ and let $W_q = \hat{A}^T\hat{A}$. Suppose λ_p is a random non-zero eigenvalue of W_p , and λ_q is a
 521 random eigenvalue from the largest q eigenvalues of W_q . Then

522 1. $\mathbb{E} \left[\frac{1}{\lambda_q} \right] = \mathbb{E} \left[\frac{1}{\lambda_p} \right] = \frac{\sqrt{4\mu^2c + (1 - c + \mu^2c)^2} - c - \mu^2c + 1}{2\mu^2} + o(1).$

523 2. $\mathbb{E} \left[\frac{1}{\lambda_q^2} \right] = \mathbb{E} \left[\frac{1}{\lambda_p^2} \right] = \frac{1 - 2c + c^2 + \mu^2c + \mu^2c^2}{2\mu^4 \sqrt{4\mu^2c + (1 - c + \mu^2c)^2}} + (1 - c) \frac{1}{2\mu^4} + o(1).$

524 *Proof.* First, we note that the non-zero eigenvalues of W_p and W_q are the same. Hence we focus on
 525 W_p . Due to the assumptions on A , we have that the eigenvalues of $A^T A$ converge to the Marchenko-
 526 Pastur with shape c^{-1} . Hence if λ_p is one of the first q eigenvalues of W_p , we see,

$$\mathbb{E} \left[\frac{1}{\lambda_p} \right] = m_{c^{-1}}(\mu^2) = \frac{\sqrt{(1 + \mu^2 - 1/c)^2 + 4\mu^2/c} - 1 - \mu^2 + 1/c}{2\mu^2/c}.$$

527 Then for the estimate of $\mathbb{E} [1/\lambda_p^2]$, we need to compute the derivative of the $m_{c^{-1}}(z)$ and evaluate it
 528 at $-\mu^2$. Hence, we see,

$$\begin{aligned} \mathbb{E} \left[\frac{1}{\lambda_p^2} \right] &= \frac{(1/c + \mu^2 + \sqrt{4\mu^2/c + (1 - 1/c + \mu^2)^2} - 1)(1/c - \mu^2 + \sqrt{4\mu^2/c + (1 - 1/c + \mu^2)^2} - 1)}{4\mu^4/c\sqrt{4\mu^2/c + (1 - 1/c + \mu^2)^2}} \\ &= \frac{(1 + \mu^2c + c\sqrt{4\mu^2/c + (1 - 1/c + \mu^2)^2} - c)(1 - \mu^2c + c\sqrt{4\mu^2/c + (1 - 1/c + \mu^2)^2} - c)}{4\mu^4c\sqrt{4\mu^2/c + (1 - 1/c + \mu^2)^2}} \\ &= \frac{(1 + \mu^2c + \sqrt{4\mu^2c + (-1 + c + \mu^2c)^2} - c)(1 - \mu^2c + \sqrt{4\mu^2c + (-1 + c + \mu^2c)^2} - c)}{4\mu^4\sqrt{4\mu^2c + (-1 + c + \mu^2c)^2}} \end{aligned}$$

529 This can be further simplified to

$$\frac{1 - 2c + c^2 + \mu^2c + \mu^2c^2}{2\mu^4\sqrt{4\mu^2c + (-1 + c + \mu^2c)^2}} + (1 - c)\frac{1}{2\mu^4} + o(1)$$

530

□

531 We will also need to estimate some other terms.

532 **Lemma 9.** *Suppose A is an p by q matrix such that the entries of A are independent and have mean*
 533 *0, variance $1/p$, and bounded fourth moment. Let $\hat{A} = [A \quad \mu I] \in \mathbb{R}^{p \times q+p}$. Let $W_p = \hat{A}\hat{A}^T$ and let*
 534 *$W_q = \hat{A}^T\hat{A}$. Suppose λ_p, λ_q are random non-zero eigenvalues of W_p, W_q from the largest $\min(p, q)$*
 535 *eigenvalues of W_p, W_q . Then*

536 1. If $p > q$, $\mathbb{E} \left[\frac{\lambda_p}{\lambda_p + \mu^2} \right] = c \left(\frac{1}{2} + \frac{1 + \mu^2c - \sqrt{(-1 + c + \mu^2c)^2 + 4\mu^2c}}{2c} \right) + o(1).$

537 2. If $p < q$, $\mathbb{E} \left[\frac{\lambda_q}{\lambda_q + \mu^2} \right] = \frac{1}{2} + \frac{1 + \mu^2c - \sqrt{(1 - c + \mu^2c)^2 + 4c^2\mu^2}}{2c} + o(1).$

538 3. If $p > q$, $\mathbb{E} \left[\frac{\lambda_p}{(\lambda_p + \mu^2)^2} \right] = c \left(\frac{1 + c + \mu^2c}{2\sqrt{(-1 + c + \mu^2c)^2 + 4\mu^2c}} - \frac{1}{2} \right) + o(1).$

539 4. If $p < q$, $\mathbb{E} \left[\frac{\lambda_q}{(\lambda_q + \mu^2)^2} \right] = \frac{1 + c + \mu^2c}{2\sqrt{(1 - c + \mu^2c)^2 + 4c^2\mu^2}} - \frac{1}{2} + o(1).$

540 *Proof.* Notice that

$$\frac{\lambda}{\lambda + \mu^2} = 1 - \frac{\mu^2}{\lambda + \mu^2} \quad \text{and} \quad \frac{\lambda}{(\lambda + \mu^2)^2} = \frac{1}{\lambda + \mu^2} - \frac{\mu^2}{(\lambda + \mu^2)^2}$$

541 Then use Lemmas 7, and 8 to finish the proof. □

542 **Bounding the Variance.**

543 **Lemma 10.** *Let η_n be a uniform measure on n numbers a_1, \dots, a_n such that $\eta^n \rightarrow \eta$ weakly in*
 544 *probability. Then for any bounded continuous function f*

$$\frac{1}{n} \sum_{i=1}^{n-1} f(a_i) \rightarrow \mathbb{E}_{x \sim \eta}[f(x)].$$

545 *Proof.* Using weak convergence

$$\frac{1}{n} \sum_{i=1}^n f(a_i) \rightarrow \mathbb{E}_{x \sim \eta}[f(x)].$$

546 Then using the boundedness of f , we get,

$$\frac{1}{n} \sum_{i=1}^{n-1} f(a_i) - \frac{1}{n} \sum_{i=1}^n f(a_i) = -\frac{1}{n} f(a_n) \rightarrow 0.$$

547

□

548 **Lemma 11.** Let η_n be a uniform measure on n numbers a_1, \dots, a_n such that $\eta_n \rightarrow \eta$ weakly in
 549 probability. Let s be a uniformly random unit vector in \mathbb{R}^m independent of η_n . Suppose $n/m \rightarrow \zeta \in$
 550 $(0, 1]$. Then for any bounded function f ,

$$\mathbb{E}_s \left[\sum_{i=1}^n s_i^2 f(a_i) \right] \rightarrow \zeta \mathbb{E}_{x \sim \eta}[f(x)]$$

551 and

$$\mathbb{E}_s \left[\left(\sum_{i=1}^n s_i^2 f(a_i) \right)^2 \right] - \mathbb{E}_s \left[\sum_{i=1}^n s_i^2 f(a_i) \right]^2 \rightarrow 0.$$

552 *Proof.* The first limit comes directly from weak convergence.

553 For the second, notice,

$$\left(\sum_{i=1}^n s_i^2 f(a_i) \right)^2 = \sum_{i=1}^n s_i^4 f(a_i)^2 + \sum_{i \neq j} s_i^2 s_j^2 f(a_i) f(a_j) = \sum_{i=1}^n s_i^4 f(a_i)^2 + \sum_{i=1}^n s_i^2 f(a_i) \sum_{j \neq i} s_j^2 f(a_j).$$

554 Taking the expectation with respect to s we get,

$$\mathbb{E}_s \left[\left(\sum_{i=1}^n s_i^2 f(a_i) \right)^2 \right] = \frac{1}{m^2 + O(m)} \sum_{i=1}^n f(a_i)^2 + \frac{1}{m^2 + O(m)} \sum_{i=1}^n f(a_i) \sum_{j \neq i} f(a_j)$$

555 Then using Lemma 10 for any fixed i , we have,

$$\frac{1}{m} \sum_{j \neq i} f(a_j) \rightarrow \zeta \mathbb{E}_{x \sim \eta}[f(x)].$$

556 Thus, as $n \rightarrow \infty$, we have,

$$\mathbb{E}_s \left[\left(\sum_{i=1}^n s_i^2 f(a_i) \right)^2 \right] \rightarrow \zeta^2 \mathbb{E}_{x \sim \eta}[f(x)]^2.$$

557 Then since

$$\mathbb{E}_s \left[\sum_{i=1}^n s_i^2 f(a_i) \right]^2 \rightarrow \zeta^2 \mathbb{E}_{x \sim \eta}[f(x)]^2.$$

558 Thus, the variance goes to zero. □

559 The interpretation of the above Lemma is that the variance of the sum decays to zero as $m \rightarrow \infty$.

560 **Lemma 12.** Suppose A is an p by q matrix such that the entries of A are independent and have
 561 mean 0, variance $1/p$, and bounded fourth moment. Let $\hat{A} = [A \quad \mu I] \in \mathbb{R}^{p \times q+p}$. Let $x \in \mathbb{R}^p$ and
 562 $\hat{y} \in \mathbb{R}^{p+q}$ be unit norm vectors such that $\hat{y}^T = [y^T \quad 0_p]$. Then

563 1. If $p < q$, then $\mathbb{E}[\text{Tr}(x^T (\hat{A} \hat{A}^T)^\dagger x)] = \frac{\sqrt{(1-c+\mu^2 c)^2 + 4\mu^2 c^2} - 1 - \mu^2 c + c}{2\mu^2 c} + o(1)$.

564 2. If $p > q$, then $\mathbb{E}[\text{Tr}(x^T(\hat{A}\hat{A}^T)^\dagger x)] = \frac{\sqrt{(-1+c+\mu^2c)^2+4\mu^2c}-1-\mu^2c+c}{2\mu^2c} + o(1)$.

565 3. If $p < q$, then $\mathbb{E}[\text{Tr}(\hat{y}^T(\hat{A}^T\hat{A})^\dagger \hat{y})] = c \left(\frac{1+c+\mu^2c}{2\sqrt{(1-c+\mu^2c)^2+4c^2\mu^2}} - \frac{1}{2} \right) + o(1)$.

566 4. If $p > q$, then $\mathbb{E}[\text{Tr}(\hat{y}^T(\hat{A}^T\hat{A})^\dagger \hat{y})] = c \left(\frac{1+c+\mu^2c}{2\sqrt{(-1+c+\mu^2c)^2+4\mu^2c}} - \frac{1}{2} \right) + o(1)$.

567 The variance of each above is $o(1)$.

568 *Proof.* Let us start with $p < q$.

569 Let $\hat{A} = \hat{U}\hat{\Sigma}\hat{V}^T$, where $\hat{\Sigma}$ is $p \times p$. Then we see,

$$(\hat{A}\hat{A}^T)^\dagger = \hat{U}\hat{\Sigma}^{-2}\hat{U}^T.$$

570 Where \hat{U} is uniformly random. Thus similar to [33], we can use Lemma 7 to get,

$$\mathbb{E}[\text{Tr}(x^T(\hat{A}\hat{A}^T)^\dagger x)] = \frac{\sqrt{(1+\mu^2c-c)^2+4\mu^2c^2}-1-\mu^2c+c}{2\mu^2c} + o(1).$$

571 On the other hand, for $p > q$, we have that only the first q eigenvalues have the expectation in Lemma
572 8 The other $p - q$ are equal to $\frac{1}{\mu^2}$. Thus, we see,

$$\begin{aligned} \mathbb{E}[\text{Tr}(x^T(\hat{A}\hat{A}^T)^\dagger x)] &= \frac{1}{c} \left(\frac{\sqrt{4\mu^2c+(-1+c+\mu^2c)^2}-c-\mu^2c+1}{2\mu^2} + o(1) \right) + \left(1 - \frac{1}{c}\right) \frac{1}{\mu^2} \\ &= \frac{\sqrt{4\mu^2c+(-1+c+\mu^2c)^2}+c-\mu^2c-1}{2c\mu^2}. \end{aligned}$$

573 Again let us first consider the case when $p < q$. Then we have,

$$(\hat{A}^T\hat{A})^\dagger = \hat{V}\hat{\Sigma}^{-2}\hat{V}^T = \begin{bmatrix} V_{1:p}\Sigma\hat{\Sigma}^{-1} \\ \mu U\hat{\Sigma}^{-1} \end{bmatrix} \hat{\Sigma}^{-2} \begin{bmatrix} \hat{\Sigma}^{-1}\Sigma V_{1:p}^T & \mu\hat{\Sigma}^{-1}U^T \end{bmatrix}.$$

574 Since \hat{y} has zeros in the last p coordinates, we see,

$$\hat{y}^T(\hat{A}^T\hat{A})^\dagger \hat{y} = y^T V_{1:p}\Sigma\hat{\Sigma}^{-4}\Sigma V_{1:p}^T y.$$

575 Thus, we can use Lemma 9 to estimate this as,

$$c \left(\frac{1+c+\mu^2c}{2\sqrt{(1-c+c\mu^2)^2+4c^2\mu^2}} - \frac{1}{2} \right) + o(1).$$

576 The extra factor of c comes from the sum of p coordinates of a uniformly unit vector in q dimensional
577 space. And for $p > q$, we have that the estimate is

$$\frac{1+c+\mu^2c}{2\sqrt{(1+\mu^2-1/c)^2+4\mu^2/c}} - \frac{c}{2} + o(1).$$

578 For the variance term, use Lemma 11. For three of the cases, the limiting distribution is the Marchenko-
579 Pastur distribution. For the other case, the limiting measure is a mixture of the Marchenko-Pastur and
580 a dirac delta at $1/\mu^2$. \square

581 The rest of the lemmas in this section are used to compute the mean and variance of the various terms
582 that appear in the formula of W_{opt} .

583 **Lemma 13.** *We have that*

$$\mathbb{E}_{A_{trn}} \left[\|\hat{h}\|^2 \right] = \begin{cases} c \left(\frac{1+c+\mu^2c}{2\sqrt{(1-c+\mu^2c)^2+4\mu^2c^2}} - \frac{1}{2} \right) + o(1) & c < 1 \\ c \left(\frac{1+c+\mu^2c}{2\sqrt{(-1+c+\mu^2c)^2+4\mu^2c}} - \frac{1}{2} \right) + o(1) & c > 1 \end{cases}$$

584 and that $\mathbb{V}(\|\hat{h}\|^2) = o(1)$.

585 *Proof.* Here we see that

$$\|\hat{h}\|^2 = \text{Tr}(\hat{v}_{trn}^T (\hat{A}_{trn}^T \hat{A}_{trn})^\dagger \hat{v}_{trn}^T).$$

586 Thus, using the Lemma 12 we get that if $c < 1$

$$\mathbb{E}[\|\hat{h}\|^2] = c \left(\frac{1 + c + \mu^2 c}{2\sqrt{(1 - c + \mu^2 c)^2 + 4\mu^2 c^2}} - \frac{1}{2} \right) + o(1)$$

587 and if $c > 1$

$$\mathbb{E}[\|\hat{h}\|^2] = c \left(\frac{1 + c + \mu^2 c}{2\sqrt{(-1 + c + \mu^2 c)^2 + 4\mu^2 c}} - \frac{1}{2} \right) + o(1).$$

588

□

589 **Lemma 14.** *We have*

$$\mathbb{E}_{A_{trn}} [\|\hat{k}\|^2] = \begin{cases} \frac{\sqrt{(1-c+\mu^2c)^2+4\mu^2c^2}-1-\mu^2c+c}{2\mu^2c} + o(1) & c < 1 \\ \frac{\sqrt{(-1+c+\mu^2c)^2+4\mu^2c}-1-\mu^2c+c}{2\mu^2c} + o(1) & c > 1 \end{cases}$$

590 and that $\mathbb{V}(\|\hat{k}\|^2) = o(1)$.

591 *Proof.* Since $\hat{k} = \hat{A}_{trn}^\dagger u$, we have that

$$\|\hat{k}\|^2 = \text{Tr}(u^T (\hat{A}_{trn} \hat{A}_{trn}^\dagger)^\dagger u).$$

592 According to the Lemma 12, if $c < 1$

$$\mathbb{E}[\|\hat{k}\|^2] = \frac{\sqrt{(1-c+\mu^2c)^2+4\mu^2c^2}-1-\mu^2c+c}{2\mu^2c} + o(1)$$

593 and if $c > 1$

$$\mathbb{E}[\|\hat{k}\|^2] = \frac{\sqrt{(-1+c+\mu^2c)^2+4\mu^2c}-1-\mu^2c+c}{2\mu^2c} + o(1).$$

594

□

595 **Lemma 15.** *We have that*

$$\mathbb{E}_{A_{trn}} [\|\hat{t}\|^2] = \begin{cases} \frac{1}{2} \left(1 - c - \mu^2 c + \sqrt{(1 - c + \mu^2 c)^2 + 4c^2 \mu^2} \right) + o(1) & c < 1 \\ \frac{1}{2} \left(1 - c - \mu^2 c + \sqrt{(-1 + c + \mu^2 c)^2 + 4\mu^2 c} \right) + o(1) & c > 1 \end{cases}$$

596 and we have that $\mathbb{V}(\|\hat{t}\|^2) = o(1)$

597 *Proof.* Here we see that $\hat{t} = \hat{v}_{trn} (I - \hat{A}_{trn}^\dagger \hat{A}_{trn})$. Thus, we see that

$$\|\hat{t}\|^2 = \|\hat{v}_{trn}\|^2 - \hat{v}_{trn}^T \hat{A}_{trn}^\dagger \hat{A}_{trn} \hat{v}_{trn} = 1 - \hat{v}_{trn}^T \hat{A}_{trn}^\dagger \hat{A}_{trn} \hat{v}_{trn}.$$

598 If $\hat{V} \in \mathbb{R}^{p+q \times p+q}$, we have that

$$\hat{A}_{trn}^\dagger \hat{A}_{trn} = \hat{V} \begin{bmatrix} I_p & 0 \\ 0 & 0_q \end{bmatrix} \hat{V}^T.$$

599 Then if $p < q$ using Lemma 6 and the fact that the last p coordinates of \hat{v}_{trn} are 0, we see that

$$\hat{v}_{trn}^T \hat{A}_{trn}^\dagger \hat{A}_{trn} \hat{v}_{trn} = \hat{v}_{trn}^T V_{1:p} \Sigma \hat{\Sigma}^{-2} \Sigma V_{1:p}^T \hat{v}_{trn}.$$

600 Then using Lemma 9 to estimate the middle diagonal matrix, we get that

$$\begin{aligned} \mathbb{E}[\|\hat{t}\|^2] &= 1 - c \left(\frac{1}{2} + \frac{1 + \mu^2 c - \sqrt{(1 + \mu^2 c - c)^2 + 4c^2 \mu^2}}{2c} \right) \\ &= \frac{1}{2} \left(1 - c - \mu^2 c + \sqrt{(1 - c + \mu^2 c)^2 + 4c^2 \mu^2} \right) + o(1). \end{aligned}$$

601 Similarly for $c > 1$, we have that

$$\begin{aligned}\mathbb{E}[\|\hat{t}\|^2] &= 1 - \left(\frac{1}{2} + \frac{c + \mu^2 c - c\sqrt{(1 + \mu^2 - 1/c)^2 + 4\mu^2/c}}{2} \right) + o(1) \\ &= \frac{1}{2} \left(1 - c - \mu^2 c + \sqrt{(-1 + c + \mu^2 c)^2 + 4\mu^2 c} \right) + o(1).\end{aligned}$$

602 The variance of $\hat{A}_{trn}^\dagger \hat{A}_{trn}$ is also $o(1)$ using Lemma 11. \square

603 **Lemma 16.** We have that $\mathbb{E}_{A_{trn}}[\hat{\gamma}] = 1$ and $\mathbb{V}(\gamma) = O(\sigma_{trn}^2/d)$.

604 *Proof.* Noting that $\hat{A} = U\hat{\Sigma}\hat{V}^T$, we have that

$$\hat{\gamma} = 1 + \sigma_{trn} \hat{v}_{trn}^T \hat{A}_{trn}^\dagger u = 1 + \sigma_{trn} \sum_{i=1}^{\min(N_{trn}, d)} \sigma_i(\hat{A})^{-1} \hat{a}_i b_i.$$

605 Here $\hat{a}^T = \hat{v}_{trn}^T \hat{V}$ and $b = U^T u$. U is a uniformly random rotation matrix that is independent of $\hat{\Sigma}$
606 and \hat{V} . Thus, taking the expectation with respect to A_{trn} , we get that the expectation is equal to zero.

607 For the variance, let us first consider the case when $c < 1$. For this case, we have that

$$\hat{V} = \begin{bmatrix} V_{1:d} \hat{\Sigma}^{-1} \\ \mu U \hat{\Sigma}^{-1} \end{bmatrix}.$$

608 Thus, letting $a^T = v_{trn}^T V_{1:d}$, we get that

$$\hat{\gamma} = 1 + \sum_{i=1}^d \frac{\sigma_i(A)}{\sigma_i^2(A) + \mu^2} a_i b_i.$$

609 Squaring and taking the expectation, we see that

$$\mathbb{E}[\gamma^2] = 1 + \frac{\sigma_{trn}^2}{N_{trn}} \mathbb{E}_{\lambda \sim \mu_c} \left[\frac{\lambda}{(\lambda + \mu^2)^2} \right] + o\left(\frac{\sigma_{trn}^2}{N_{trn}}\right).$$

610 Similarly for $c > 1$, we have that

$$\mathbb{E}[\gamma^2] = 1 + \frac{\sigma_{trn}^2}{d} \mathbb{E}_{\lambda \sim \mu_c} \left[\frac{\lambda}{(\lambda + \mu^2)^2} \right] + o\left(\frac{\sigma_{trn}^2}{d}\right).$$

611 \square

612 **Lemma 17.** We have that

$$\mathbb{E} \left[\text{Tr}((\hat{A}_{trn}^\dagger)^T \hat{k} \hat{k}^T \hat{A}_{trn}^\dagger) \right] = \mathbb{E}[\rho] = \begin{cases} \frac{\mu^2 c^2 + c^2 + \mu^2 c - 2c + 1}{2\mu^4 c \sqrt{4\mu^2 c^2 + (1 - c + \mu^2 c)^2}} + \frac{1}{2\mu^4} \left(1 - \frac{1}{c}\right) + o(1) & c < 1 \\ \frac{1 - 2c + c^2 + \mu^2 c + \mu^2 c^2}{2\mu^4 c \sqrt{4\mu^2 c + (-1 + c + \mu^2 c)^2}} + \left(1 - \frac{1}{c}\right) \frac{1}{2\mu^4} + o(1) \end{cases}$$

613 and that $\mathbb{V}(\rho) = o(1)$.

614 *Proof.* Here we have that

$$\rho = \text{Tr}(\hat{k}^T (\hat{A}_{trn}^\dagger \hat{A}_{trn})^\dagger \hat{k}) = \text{Tr}(u^T (\hat{A}_{trn} \hat{A}_{trn}^T)^\dagger (\hat{A}_{trn} \hat{A}_{trn}^T)^\dagger u).$$

615 We first notice that

$$(\hat{A}_{trn} \hat{A}_{trn}^T)^\dagger (\hat{A}_{trn} \hat{A}_{trn}^T)^\dagger = \hat{U}^T \hat{\Sigma}^2 \hat{U}.$$

616 Thus using Lemmas 7 and 8, we see that if $c < 1$

$$\mathbb{E}[\rho] = \frac{\mu^2 c^2 + c^2 + \mu^2 c - 2c + 1}{2\mu^4 c \sqrt{4\mu^2 c^2 + (1 - c + \mu^2 c)^2}} + \frac{1}{2\mu^4} \left(1 - \frac{1}{c}\right)$$

617 and if $c > 1$

$$\begin{aligned}\mathbb{E}[\rho] &= \frac{1}{c} \left(\frac{1 - 2c + c^2 + \mu^2 c + \mu^2 c^2}{2\mu^4 \sqrt{4\mu^2 c + (-1 + c + \mu^2 c)^2}} + (1 - c) \frac{1}{2\mu^4} \right) + \left(1 - \frac{1}{c}\right) \frac{1}{\mu^4} \\ &= \frac{1 - 2c + c^2 + \mu^2 c + \mu^2 c^2}{2\mu^4 c \sqrt{4\mu^2 c + (-1 + c + \mu^2 c)^2}} + \left(1 - \frac{1}{c}\right) \frac{1}{2\mu^4}.\end{aligned}$$

618 The variance being $o(1)$ comes from Lemma 11 again. \square

619 **Lemma 18.** *We have that*

$$\mathbb{E}_{A_{trn}} \left[\text{Tr}(\hat{h}^T \hat{k}^T \hat{A}_{trn}^\dagger) \right] = 0$$

620 *and the variance is $o(1)$.*

621 *Proof.* Letting $\hat{A} = U \hat{\Sigma} \hat{V}^T$, we get that

$$\text{Tr}(\hat{h}^T \hat{k}^T \hat{A}^T) = u^T U \hat{\Sigma}^{-3} \hat{V}^T \hat{v}_{trn}^T.$$

622 Then again since U is uniformly random and independent of $\hat{\Sigma}$ and \hat{V} , the expectation is equal to
623 zero. The variance is computed similarly to Lemma 16. \square

624 **F.2.5 Step 5: Putting it together**

625 **Lemma 19.** *We have that*

$$\mathbb{E} \left[\frac{\tau}{\sigma_{trn}^2} \right] = \begin{cases} \frac{1}{\sigma_{trn}^2} + \frac{1}{2} \left(1 + \mu^2 c + c - \sqrt{(1 - c + \mu^2 c)^2 + 4\mu^2 c^2} \right) + o(1) & c < 1 \\ \frac{1}{\sigma_{trn}^2} + \frac{1}{2} \left(1 + \mu^2 c + c - \sqrt{(-1 + c + \mu^2 c)^2 + 4\mu^2 c} \right) + o(1) & c > 1 \end{cases}$$

626 *and that $\mathbb{V}(\tau/\sigma_{trn}^2) = o(1)$.*

627 *Proof.* Using the fact that all of the quantities concentrate, we can use the previous estimates.
628 Specifically, we use that

$$|\mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]| \leq \sqrt{\mathbb{V}[X]\mathbb{V}[Y]}.$$

629 Thus, since our variances decay, we can use the product of the expectations. Further,

$$\begin{aligned} |\mathbb{V}[XY]| &= |\mathbb{V}[X]\mathbb{V}[Y] + \mathbb{E}[X]^2\mathbb{V}[Y] + \mathbb{E}[Y]^2\mathbb{V}[X] - 2\mathbb{E}[X]\mathbb{E}[Y]\text{Cov}(X, Y) + \text{Cov}(X^2, Y^2) - \text{Cov}(X, Y)^2| \\ &\leq |\mathbb{V}[X]\mathbb{V}[Y] + \mathbb{E}[X]^2\mathbb{V}[Y] + \mathbb{E}[Y]^2\mathbb{V}[X]| + 2|\mathbb{E}[X]\mathbb{E}[Y]|\sqrt{\mathbb{V}[X]\mathbb{V}[Y]} + |\mathbb{V}[X]\mathbb{V}[Y]| + |\sqrt{\mathbb{V}[X^2]\mathbb{V}[Y^2]}|. \end{aligned}$$

630 Thus, since the variances individually go to 0, we see that the variance of the product also goes to 0.
631 Then using Lemma 15 and 14, we have that if $c < 1$

$$\mathbb{E} \left[\|\hat{t}\|^2 \|\hat{k}\|^2 \right] = \frac{1}{2} \left(1 + \mu^2 c + c - \sqrt{(1 - c + \mu^2 c)^2 + 4\mu^2 c^2} \right) + o(1)$$

632 *and $\mathbb{V}(\|\hat{t}\|^2 \|\hat{k}\|^2) = o(1)$. Then since*

$$|\mathbb{V}[X + Y]| \leq |\mathbb{V}[X] + \mathbb{V}[Y]| + 2\sqrt{\mathbb{V}[X]\mathbb{V}[Y]}$$

633 *we have that using Lemma 16, that if $c < 1$*

$$\mathbb{E} \left[\frac{\tau}{\sigma_{trn}^2} \right] = \frac{1}{\sigma_{trn}^2} + \frac{1}{2} \left(1 + \mu^2 c + c - \sqrt{(1 - c + \mu^2 c)^2 + 4\mu^2 c^2} \right) + o(1)$$

634 *and that that variance is $o(1)$. If $c > 1$*

$$\mathbb{E} \left[\frac{\tau}{\sigma_{trn}^2} \right] = \frac{1}{\sigma_{trn}^2} + \frac{1}{2} \left(1 + \mu^2 c + c - \sqrt{(-1 + c + \mu^2 c)^2 + 4\mu^2 c} \right) + o(1).$$

635 \square

636 **Lemma 20.** *We have that*

$$\mathbb{E}_{A_{trn}} \left[\frac{1}{\sigma_{trn}^2} \|\hat{h}\|^2 + \|\hat{t}\|^4 \rho \right] = \begin{cases} \frac{c(1+\sigma_{trn}^{-2})}{2} \left(\frac{\mu^2 c + c + 1}{\sqrt{(1 - c + \mu^2 c)^2 + 4\mu^2 c^2}} - 1 \right) + o(1) & c < 1 \\ \frac{c(1+\sigma_{trn}^{-2})}{2} \left(\frac{\mu^2 c + c + 1}{\sqrt{(-1 + c + \mu^2 c)^2 + 4\mu^2 c}} - 1 \right) + o(1) & c > 1 \end{cases}$$

637 *and that the variance is $o(1)$.*

638 *Proof.* Similar to Lemma 19, we can multiply the expectations since the variances are small. For
 639 $c < 1$, simplifying, we get that

$$\mathbb{E}_{A_{trn}} \left[\frac{1}{\sigma_{trn}^2} \|\hat{h}\|^2 + \|\hat{t}\|^4 \rho \right] = \frac{c(1 + \sigma_{trn}^{-2})}{2} \left(\frac{\mu^2 c + c + 1}{\sqrt{(1 - c + \mu^2 c)^2 + 4\mu^2 c^2}} - 1 \right) + o(1)$$

640 and if $c > 1$, we get that

$$\mathbb{E}_{A_{trn}} \left[\frac{1}{\sigma_{trn}^2} \|\hat{h}\|^2 + \|\hat{t}\|^4 \rho \right] = \frac{c(1 + \sigma_{trn}^{-2})}{2} \left(\frac{\mu^2 c + c + 1}{\sqrt{(-1 + c + \mu^2 c)^2 + 4\mu^2 c}} - 1 \right) + o(1)$$

641 and the variance decays since the variances decay individually. \square

642 **Lemma 21.** *We have that*

$$\mathbb{E}_{A_{trn}} [\|W_{opt}\|_F^2] = \frac{\sigma_{trn}^4}{\tau^2} \begin{cases} \frac{c(1 + \sigma_{trn}^{-2})}{2} \left(\frac{\mu^2 c + c + 1}{\sqrt{(1 - c + \mu^2 c)^2 + 4\mu^2 c^2}} - 1 \right) + o(1) & c < 1 \\ \frac{c(1 + \sigma_{trn}^{-2})}{2} \left(\frac{\mu^2 c + c + 1}{\sqrt{(-1 + c + \mu^2 c)^2 + 4\mu^2 c}} - 1 \right) + o(1) & c > 1 \end{cases}$$

643 and that $\mathbb{V}(\|W_{opt}\|_F^2) = o(1)$.

644 *Proof.* Follows immediately from Lemmas 4, 17, 18, and 20. \square

645 **Theorem 1** (Generalization Error Formula). *Suppose the training data X_{trn} and test data X_{tst}*
 646 *satisfy Assumption 1 and the noise A_{trn}, A_{tst} satisfy Assumption 2. Let μ be the regularization*
 647 *parameter. Then for the under-parameterized regime (i.e., $c < 1$) for the solution W_{opt} to Problem 1,*
 648 *the generalization error or risk given by Equation 2 is given by*

$$\mathcal{R}(c, \mu) = \tau^{-2} \left(\frac{\sigma_{tst}^2}{N_{tst}} + \frac{c\sigma_{trn}^2(\sigma_{trn}^2 + 1)}{2d} \left(\frac{1 + c + \mu^2 c}{\sqrt{(1 - c + \mu^2 c)^2 + 4\mu^2 c^2}} - 1 \right) \right) + o\left(\frac{1}{d}\right),$$

649 where $\tau^{-1} = \frac{2}{2 + \sigma_{trn}^2(1 + c + \mu^2 c - \sqrt{(1 - c + \mu^2 c)^2 + 4\mu^2 c^2})}$.

650 *Proof.* Rewriting $\frac{\hat{\gamma}^2}{\tau^2}$ as $\frac{\hat{\gamma}^2/\sigma_{trn}^4}{\tau^2/\sigma_{trn}^4}$, we can the concentration from Lemmas 16 and 19. Then using
 651 Lemma 21 we get the needed result. \square

652 **Theorem 8.** *For the over-parameterized case, we have that the generalization error is given by*

$$\mathcal{R}(c, \mu) = \tau^{-2} \left(\frac{\sigma_{tst}^2}{N_{tst}} + \frac{c\sigma_{trn}^2(\sigma_{trn}^2 + 1)}{2d} \left(\frac{1 + c + \mu^2 c}{\sqrt{(-1 + c + \mu^2 c)^2 + 4\mu^2 c}} - 1 \right) \right) + o\left(\frac{1}{d}\right),$$

653 where $\tau^{-1} = \frac{2}{2 + \sigma_{trn}^2(1 + c + \mu^2 c - \sqrt{(-1 + c + \mu^2 c)^2 + 4\mu^2 c})}$.

654 *Proof.* Rewriting $\frac{\hat{\gamma}^2}{\tau^2}$ as $\frac{\hat{\gamma}^2/\sigma_{trn}^4}{\tau^2/\sigma_{trn}^4}$, we can the concentration from Lemmas 16 and 19. Then using
 655 Lemma 21 we get the needed result. \square

656 F.3 Proof of Theorem 2

657 **Theorem 2** (Under-Parameterized Peak). *If $\mu \in \mathbb{R}_{>0}$ is such that $p(\mu) < 0$, $\sigma_{trn}^2 = N_{trn} = d/c$*
 658 *and $\sigma_{tst}^2 = N_{tst}$, and d is sufficiently large, then the risk $\mathcal{R}(c)$ from Theorem 1, as a function of c ,*
 659 *has a local maximum in the under-parameterized regime ($c \in (0, 1)$).*

679 *Proof.* To begin, we note that the derivative is,

$$\partial_c \mathcal{R}(c, \mu) = \frac{P(c, \mu, d, T)}{Q(c, \mu, d, T)}.$$

680 Where

$$\begin{aligned} P(c, \mu, d, T) = & -4T^2(-Tc^3d^2\mu^6 - 3Tc^3d^2\mu^4 - 3Tc^3d^2\mu^2 - Tc^3d^2 - Tc^3d\mu^4 \\ & - 5Tc^3d\mu^2 - 4Tc^3d - Tc^2d^2\mu^4 - Tc^2d^2\mu^2 - 2Tc^2d\mu^2 + 5Tc^2d + Tcd^2\mu^2 - Tcd \\ & + Td^2 + c^4d^2\mu^8 + 4c^4d^2\mu^6 + 6c^4d^2\mu^4 + 4c^4d^2\mu^2 + c^4d^2 + c^4d\mu^6 + 2c^4d\mu^4 \\ & + c^4d\mu^2 + 2c^4\mu^2 + 2c^4 + 2c^3d^2\mu^6 + 3c^3d^2\mu^4 - c^3d^2 + 3c^3d\mu^4 + 5c^3d - 2c^3 \\ & + 3c^2d\mu^2 - 6c^2d - 2cd^2\mu^2 + cd^2 + cd - d^2), \end{aligned}$$

681

$$Q(c, \mu, d, T) = T^7 (-Td + cd\mu^2 + cd + 2c + d)^3,$$

682 and

$$T = \sqrt{c^2\mu^4 + 2c^2\mu^2 + c^2 + 2c\mu^2 - 2c + 1}.$$

683 Then if a critical point exists, it must be the case that $P(c, \mu, d, T) = 0$. This happens either if

684 $T^2 = 0$ or $\hat{P} = P/(-4T^2) = 0$. Note we can simplify T^2 as

$$c^2(\mu^2 + 1)^2 + 2(\mu^2 - 1)c + 1$$

685 Then since this is a quadratic, we get that,

$$c = \frac{-2(\mu^2 - 1) \pm \sqrt{4(\mu^2 - 1)^2 - 4(\mu^2 + 1)^2}}{2(\mu^2 + 1)^2} = \frac{-2(\mu^2 - 1) \pm \sqrt{-16\mu^4}}{2(\mu^2 + 1)^2}.$$

686 Thus, the solutions live in \mathbb{C} and not in \mathbb{R} . Since we want to find a root in $(0, 1)$, we can discard this
687 factor and focus on \hat{P} .

688 Looking at \hat{P} , we see that

$$\hat{P} = \hat{P}_1 + \hat{P}_2 + \hat{P}_3 + \hat{P}_4 + \hat{P}_5,$$

689 where

$$\hat{P}_1 = -d^2T(\mu^2c + c - 1)(\mu^4c^2 + 2\mu^2c^2 + 2\mu^2c + c^2 + c + 1).$$

$$\hat{P}_2 = -dTc(\mu^4c^2 + 5\mu^2c^2 + 2\mu^2c + 4c^2 - 5c + 1).$$

$$\hat{P}_3 = d^2(\mu^2c + c - 1)(\mu^2c + c + 1)(\mu^4c^2 + 2\mu^2c^2 + 2\mu^2c + c^2 - c + 1).$$

$$\hat{P}_4 = dc(\mu^6c^3 + 2\mu^4c^3 + 3\mu^4c^2 + \mu^2c^3 + 3\mu^2c + 5c^2 - 6c + 1).$$

$$\hat{P}_5 = 2c^3(\mu^2c + c - 1).$$

690 Here we see that $\mu^2c + c - 1$ is a factor for three of the five polynomials. Hence, the hope is that a
691 multiple of $\mu^2c + c - 1$ can approximate the sum of the other two. Dividing \hat{P}_2, \hat{P}_4 by $\mu^2c + c - 1$,
692 we get that

$$\hat{P}_2 = -dTc(\mu^2c + c - 1)(\mu^2c + 4c - 1) - 4dT\mu^2c^2.$$

$$\hat{P}_4 = dc(\mu^2c + c - 1)(\mu^4c^2 + \mu^2c^2 + 4\mu^2c - 1) - 3d\mu^2c^3 + 8d\mu^2c^2 + 5dc^3 - 5dC^2.$$

693 Now we see that for some \tilde{P}

$$\hat{P} = (\mu^2c + c - 1)\tilde{P} - 4dT\mu^2c^2 - 3d\mu^2c^3 + 8d\mu^2c^2 + 5dc^3 - 5dC^2.$$

694 We further simplify this by dividing the remainder again by $\mu^2c + c - 1$ to get that

$$-4dT\mu^2c^2 - 3d\mu^2c^3 + 8d\mu^2c^2 + 5dc^3 - 5dC^2 = dc^2(\mu^2c + c - 1)(5\mu^2c - 8\mu^2) + 4d\mu^2c^2(2\mu^2c - T).$$

695 Thus, redefining \tilde{P} , we get that

$$\hat{P} = (\mu^2c + c - 1)\tilde{P} + 4d\mu^2c^2(2\mu^2c - T),$$

696 with

$$\begin{aligned} \tilde{P} = & -Tc^2d^2\mu^4 - 2Tc^2d^2\mu^2 - Tc^2d^2 - Tc^2d\mu^2 - 4Tc^2d - 2Tcd^2\mu^2 - Tcd^2 + Tcd \\ & - Td^2 + c^3d^2\mu^6 + 3c^3d^2\mu^4 + 3c^3d^2\mu^2 + c^3d^2 + c^3d\mu^4 + c^3d\mu^2 + 2c^3 \\ & + 3c^2d^2\mu^4 + 3c^2d^2\mu^2 - 4c^2d\mu^2 + 5c^2d + 3cd^2\mu^2 - cd + d^2. \end{aligned}$$

697 Thus, we have the needed result.

698

□

699 **E.5 Proof of Theorem 5**

700 **Theorem 5** ($\|W_{opt}\|_F$ Peak). If $\sigma_{tst} = \sqrt{N_{tst}}$, $\sigma_{trn} = \sqrt{N_{trn}}$ and μ is such that $p(\mu) < 0$,
 701 then for N_{trn} large enough and $d = cN_{trn}$, we have that $\|W_{opt}\|_F$ has a local maximum in the
 702 under-parameterized regime. Specifically for $c \in ((\mu^2 + 1)^{-1}, 1)$.

703 *Proof.* Here we note that the expression for the norm of W_{opt} is given by Lemma 21. Differentiating
 704 with respect to c , we get that the derivative is given by

$$\begin{aligned} & \frac{c\sigma_{trn}^4 \left(-1 + \frac{c\mu^2+c+1}{\sqrt{4c^2\mu^2+(c\mu^2-c+1)^2}} \right) (\sigma_{trn}^2 + 1) \left(\mu^2 - \frac{4c\mu^2 + \frac{(2\mu^2-2)(c\mu^2-c+1)}{2}}{\sqrt{4c^2\mu^2+(c\mu^2-c+1)^2}} + 1 \right)}{2 \left(\frac{\sigma_{trn}^2 (c\mu^2+c-\sqrt{4c^2\mu^2+(c\mu^2-c+1)^2}+1)}{2} + 1 \right)^3} \\ & + \frac{c\sigma_{trn}^2 (\sigma_{trn}^2 + 1) \left(\frac{\mu^2+1}{\sqrt{4c^2\mu^2+(c\mu^2-c+1)^2}} + \frac{\left(-4c\mu^2 - \frac{(2\mu^2-2)(c\mu^2-c+1)}{2} \right) (c\mu^2+c+1)}{(4c^2\mu^2+(c\mu^2-c+1)^2)^{\frac{3}{2}}} \right)}{2 \left(\frac{\sigma_{trn}^2 (c\mu^2+c-\sqrt{4c^2\mu^2+(c\mu^2-c+1)^2}+1)}{2} + 1 \right)^2} \\ & + \frac{\sigma_{trn}^2 \left(-1 + \frac{c\mu^2+c+1}{\sqrt{4c^2\mu^2+(c\mu^2-c+1)^2}} \right) (\sigma_{trn}^2 + 1)}{2 \left(\frac{\sigma_{trn}^2 (c\mu^2+c-\sqrt{4c^2\mu^2+(c\mu^2-c+1)^2}+1)}{2} + 1 \right)^2}. \end{aligned}$$

705 At $c = \frac{1}{\mu^2+1}$, this has value

$$\frac{2\sigma_{trn}^2 (\mu^2 + 1)^{\frac{3}{2}} (-256\mu^7 + 256\mu^6\sqrt{\mu^2+1}) (\sigma_{trn}^2 + 1)}{\left(\frac{4\mu^4}{(\mu^2+1)^2} + \frac{4\mu^2}{(\mu^2+1)^2} \right)^{\frac{7}{2}} (-2\mu\sigma_{trn}^2 + 2\sigma_{trn}^2\sqrt{\mu^2+1} + 2\sqrt{\mu^2+1})^3 (\mu^6\sqrt{\mu^2+1} + 3\mu^4\sqrt{\mu^2+1} + 3\mu^2\sqrt{\mu^2+1} + \sqrt{\mu^2+1})}$$

706 Then since $\sqrt{\mu^2+1} > \mu$, we have that the derivative is positive at this point. Next, we compute the
 707 limit of the derivative as $c \rightarrow 1^-$ and see that this is given by

$$\frac{\sigma_{trn}^2 (\sigma_{trn}^2 + 1) (\sigma_{trn}^2 p(\mu) + 4\mu^{14} + 56\mu^{12} + 280\mu^{10} + 576\mu^8 + 384\mu^6 - (4\mu^{13} + 48\mu^{11} + 192\mu^9 + 256\mu^7)\sqrt{\mu^2+4})}{(\mu^4 + 4\mu^2)^{\frac{7}{2}} (\sigma_{trn}^2 (\mu^2 - \mu\sqrt{\mu^2+4} + 2) + 2)^3}$$

708 Then we see that the denominator is positive. Hence the sign is determined by the numerator. Again,
 709 we assumed $p(\mu) < 0$. Hence the leading coefficient in term of σ_{trn}^2 is negative. Since $\sigma_{trn}^2 = N_{trn}$.
 710 If N_{trn} is sufficiently large the derivative is negative near $c = 1$. Thus, we have a peak. \square

711 **E.6 Proof of Theorem 7**

712 **Theorem 7** (Training Error). Let τ be as in Theorem 1. The training error for $c < 1$ is given by

$$\mathbb{E}_{A_{trn}} [\|X_{trn} - W_{opt}(X_{trn} + A_{trn})\|_F^2] = \tau^{-2} (\sigma_{trn}^2 (1 - c \cdot T_1) + \sigma_{trn}^4 T_2) + o(1),$$

713 where $T_1 = \frac{\mu^2}{2} \left(\frac{1 + c + \mu^2 c}{\sqrt{(1 - c + \mu^2 c)^2 + 4\mu^2 c^2}} - 1 \right) + \frac{1}{2} + \frac{1 + \mu^2 c - \sqrt{(1 - c + \mu^2 c)^2 + 4c^2 \mu^2}}{2c}$,

714 and $T_2 = \frac{(\mu^2 c + c - 1 - \sqrt{(1 - c + \mu^2 c)^2 + 4c^2 \mu^2})^2 (\mu^2 c + c + 1 - \sqrt{(1 - c + \mu^2 c)^2 + 4c^2 \mu^2})}{2\sqrt{(1 - c + \mu^2 c)^2 + 4c^2 \mu^2}}$.

715 *Proof.* Note that we have:

$$\begin{aligned}
\mathbb{E}_{A_{trn}} \left[\frac{\|X_{trn} - W_{opt} Y_{trn}\|_F^2}{N_{trn}} \right] &= \frac{1}{N_{trn}} \mathbb{E}_{A_{trn}} [\|X_{trn} - W_{opt}(X_{trn} + A_{trn})\|_F^2] \\
&= \frac{1}{N_{trn}} \mathbb{E}[\|X_{trn} - W_{opt} X_{trn}\|^2] + \frac{1}{N_{trn}} \mathbb{E}[\|W_{opt} A_{trn}\|^2] \\
&\quad + \frac{2}{N_{trn}} \mathbb{E}[\text{Tr}((X_{trn} - W_{opt} X_{trn})^T W_{opt} A_{trn})].
\end{aligned}$$

716 First, by Lemma 2, we have $X_{trn} - W_{opt} X_{trn} = \frac{\hat{\gamma}}{\hat{\tau}} X_{trn}$. Then, $\mathbb{E}[\|X_{trn} - W_{opt} X_{trn}\|^2] =$
717 $\frac{\hat{\gamma}^2}{\hat{\tau}^2} \mathbb{E}[\|X_{trn}\|^2] = \frac{\hat{\gamma}^2 \sigma_{trn}^2}{\hat{\tau}^2}$. Then, let us look at the $\mathbb{E}_{A_{trn}}[\|W_{opt} A_{trn}\|_F^2]$ term.

$$\begin{aligned}
\mathbb{E}_{A_{trn}}[\|W_{opt} A_{trn}\|_F^2] &= \mathbb{E}[\text{Tr}(A_{trn}^T W_{opt}^T W_{opt} A_{trn})] \\
&= \frac{\sigma_{trn}^2 \hat{\gamma}^2}{\hat{\tau}^2} \mathbb{E}[\text{Tr}(A_{trn}^T \hat{h}^T u^T u \hat{h} A_{trn})] \\
&\quad + \frac{\sigma_{trn}^3 \hat{\gamma} \|\hat{t}\|^2}{\hat{\tau}^2} \mathbb{E}[\text{Tr}(A_{trn}^T \hat{h}^T u^T u \hat{k}^T \hat{A}_{trn}^\dagger A_{trn})] \\
&\quad + \frac{\sigma_{trn}^3 \hat{\beta} \|\hat{t}\|^2}{\hat{\tau}^2} \mathbb{E}[\text{Tr}(A_{trn}^T (\hat{A}_{trn}^\dagger)^T \hat{k} u^T u \hat{h} A_{trn})] \\
&\quad + \frac{\sigma_{trn}^4 \|\hat{t}\|^4}{\hat{\tau}^2} \mathbb{E}[\text{Tr}(A_{trn}^T (\hat{A}_{trn}^\dagger)^T \hat{k} u^T u \hat{k}^T \hat{A}_{trn}^\dagger A_{trn})] \\
&= \frac{\sigma_{trn}^2 \hat{\gamma}^2}{\hat{\tau}^2} \mathbb{E}[\text{Tr}(\hat{h} A_{trn} A_{trn}^T \hat{h}^T)] \\
&\quad + \frac{\sigma_{trn}^3 \hat{\gamma} \|\hat{t}\|^2}{\hat{\tau}^2} \mathbb{E}[\text{Tr}(\hat{k}^T \hat{A}_{trn}^\dagger A_{trn} A_{trn}^T \hat{h}^T)] \\
&\quad + \frac{\sigma_{trn}^3 \hat{\gamma} \|\hat{t}\|^2}{\hat{\tau}^2} \mathbb{E}[\text{Tr}(\hat{h} A_{trn} A_{trn}^T (\hat{A}_{trn}^\dagger)^T \hat{k})] \\
&\quad + \frac{\sigma_{trn}^4 \|\hat{t}\|^4}{\hat{\tau}^2} \mathbb{E}[\text{Tr}(\hat{k}^T \hat{A}_{trn}^\dagger A_{trn} A_{trn}^T (\hat{A}_{trn}^\dagger)^T \hat{k})] \\
&= \frac{\sigma_{trn}^2 \hat{\gamma}^2}{\hat{\tau}^2} \mathbb{E}[\text{Tr}(\hat{v}_{trn}^T \hat{A}_{trn}^\dagger A_{trn} A_{trn}^T (\hat{A}_{trn}^\dagger)^T \hat{v}_{trn})] \\
&\quad + \frac{\sigma_{trn}^3 \hat{\gamma} \|\hat{t}\|^2}{\hat{\tau}^2} \mathbb{E}[\text{Tr}(u^T (\hat{A}_{trn}^\dagger)^T \hat{A}_{trn}^\dagger A_{trn} A_{trn}^T (\hat{A}_{trn}^\dagger)^T \hat{v}_{trn}^T)] \\
&\quad + \frac{\sigma_{trn}^3 \hat{\gamma} \|\hat{t}\|^2}{\hat{\tau}^2} \mathbb{E}[\text{Tr}(\hat{v}_{trn}^T \hat{A}_{trn}^\dagger A_{trn} A_{trn}^T (\hat{A}_{trn}^\dagger)^T \hat{A}_{trn}^\dagger u)] \\
&\quad + \frac{\sigma_{trn}^4 \|\hat{t}\|^4}{\hat{\tau}^2} \mathbb{E}[\text{Tr}(u^T (\hat{A}_{trn}^\dagger)^T \hat{A}_{trn}^\dagger A_{trn} A_{trn}^T (\hat{A}_{trn}^\dagger)^T \hat{A}_{trn}^\dagger u)] \\
&= \frac{\sigma_{trn}^2 \hat{\gamma}^2}{\hat{\tau}^2} \mathbb{E}[\text{Tr}(\hat{v}_{trn}^T \hat{A}_{trn}^\dagger A_{trn} A_{trn}^T (\hat{A}_{trn}^\dagger)^T \hat{v}_{trn})] \\
&\quad + \frac{\sigma_{trn}^4 \|\hat{t}\|^4}{\hat{\tau}^2} \mathbb{E}[\text{Tr}(u^T (\hat{A}_{trn}^\dagger)^T \hat{A}_{trn}^\dagger A_{trn} A_{trn}^T (\hat{A}_{trn}^\dagger)^T \hat{A}_{trn}^\dagger u)].
\end{aligned}$$

718 Then, we look at the $\text{Tr}((X_{trn} - W_{opt}X_{trn})^T W_{opt}A_{trn})$ term. By Lemma 2, we have $X_{trn} -$
719 $W_{opt}X_{trn} = \frac{\hat{\gamma}}{\hat{\tau}}X_{trn}$. Then,

$$\begin{aligned}
\frac{\hat{\gamma}}{\hat{\tau}} \text{Tr}(X_{trn}^T W_{opt}A_{trn}) &= \frac{\hat{\gamma}}{\hat{\tau}} \text{Tr} \left(X_{trn}^T \left(\frac{\sigma_{trn}\hat{\gamma}}{\hat{\tau}}u\hat{h} + \frac{\sigma_{trn}^2\|\hat{t}\|^2}{\hat{\tau}}u\hat{k}^T \hat{A}_{trn}^\dagger \right) A_{trn} \right) \\
&= \frac{\sigma_{trn}\hat{\gamma}^2}{\hat{\tau}^2} \text{Tr} \left(X_{trn}^T u\hat{h}A_{trn} \right) \\
&\quad + \frac{\sigma_{trn}^2\hat{\gamma}\|\hat{t}\|^2}{\hat{\tau}^2} \text{Tr} \left(X_{trn}^T u\hat{k}^T \hat{A}_{trn}^\dagger A_{trn} \right) \\
&= \frac{\sigma_{trn}\hat{\gamma}^2}{\hat{\tau}^2} \text{Tr} \left(\sigma_{trn}v_{trn}\hat{v}_{trn}^T \hat{A}_{trn}^\dagger A_{trn} \right) \\
&\quad + \frac{\sigma_{trn}^2\hat{\gamma}\|\hat{t}\|^2}{\hat{\tau}^2} \text{Tr} \left(\sigma_{trn}v_{trn}u^T (\hat{A}_{trn}^\dagger)^T \hat{A}_{trn}^\dagger A_{trn} \right) \\
&= \frac{\sigma_{trn}^2\hat{\gamma}^2}{\hat{\tau}^2} \text{Tr} \left(\hat{v}_{trn}^T \hat{A}_{trn}^\dagger A_{trn} v_{trn} \right) \\
&\quad + \frac{\sigma_{trn}^3\hat{\gamma}\|\hat{t}\|^2}{\hat{\tau}^2} \text{Tr} \left(u^T (\hat{A}_{trn}^\dagger)^T \hat{A}_{trn}^\dagger A_{trn} v_{trn} \right) \\
&= \frac{\sigma_{trn}^2\hat{\gamma}^2}{\hat{\tau}^2} \text{Tr} \left(\hat{v}_{trn}^T \hat{A}_{trn}^\dagger A_{trn} v_{trn} \right).
\end{aligned}$$

720 In conclusion, we have the training error:

$$\begin{aligned}
\mathbb{E}_{A_{trn}} \left[\frac{\|X_{trn} - W_{opt}Y_{trn}\|_F^2}{N_{trn}} \right] &= \frac{\hat{\gamma}^2\sigma_{trn}^2}{N_{trn}\hat{\tau}^2} + \frac{\sigma_{trn}^2\hat{\gamma}^2}{N_{trn}\hat{\tau}^2} \mathbb{E}[\text{Tr}(\hat{v}_{trn}^T \hat{A}_{trn}^\dagger A_{trn} A_{trn}^T (\hat{A}_{trn}^\dagger)^T \hat{v}_{trn})] \\
&\quad + \frac{\sigma_{trn}^4\|\hat{t}\|^4}{N_{trn}\hat{\tau}^2} \mathbb{E}[\text{Tr}(u^T (\hat{A}_{trn}^\dagger)^T \hat{A}_{trn}^\dagger A_{trn} A_{trn}^T (\hat{A}_{trn}^\dagger)^T \hat{A}_{trn}^\dagger u)] \\
&\quad + 2 \frac{\sigma_{trn}^2\hat{\gamma}^2}{N_{trn}\hat{\tau}^2} \mathbb{E} \left[\text{Tr} \left(\hat{v}_{trn}^T \hat{A}_{trn}^\dagger A_{trn} v_{trn} \right) \right].
\end{aligned}$$

721 Now we estimate the above terms using random matrix theory. Here we focus on the $c < 1$ case. For
722 $c < 1$, we note that

$$\hat{A}_{trn}^\dagger A_{trn} A_{trn}^T (\hat{A}_{trn}^\dagger)^T = \hat{V} \hat{\Sigma}^{-1} \Sigma \Sigma^T \hat{\Sigma}^{-1} \hat{V}^T.$$

723 Thus, for $c < 1$

$$\hat{v}_{trn}^T \hat{A}_{trn}^\dagger A_{trn} A_{trn}^T (\hat{A}_{trn}^\dagger)^T \hat{v}_{trn} = \sum_{i=1}^d a_i^2 \frac{\sigma_i(A)^4}{(\sigma_i(A)^2 + \mu^2)^2}$$

724 where $a^T = v_{trn}^T V_{1:d}$. Taking the expectation, and using Lemma 9 we get that

$$\begin{aligned}
\mathbb{E}_{A_{trn}} \left[\hat{v}_{trn}^T \hat{A}_{trn}^\dagger A_{trn} A_{trn}^T (\hat{A}_{trn}^\dagger)^T \hat{v}_{trn} \right] &= \\
c \left(\frac{1}{2} + \frac{1 + \mu^2 c - \sqrt{(1 - c + \mu^2 c)^2 + 4c^2 \mu^2}}{2c} + \mu^2 \left(\frac{1 + c + \mu^2 c}{2\sqrt{(1 - c + c\mu^2)^2 + 4c^2 \mu^2}} - \frac{1}{2} \right) \right) &+ o(1).
\end{aligned}$$

725 Using Lemma 11, we see that the variance is $o(1)$. Similarly, we have that

$$(\hat{A}_{trn}^\dagger)^T \hat{A}_{trn}^\dagger A_{trn} A_{trn}^T (\hat{A}_{trn}^\dagger)^T \hat{A}_{trn}^\dagger = U \hat{\Sigma}^{-2} \Sigma \Sigma^T \hat{\Sigma}^{-2} U^T.$$

726 Thus, again, using a similar argument, we see that

$$\mathbb{E}_{A_{trn}} \left[\text{Tr}(u^T (\hat{A}_{trn}^\dagger)^T \hat{A}_{trn}^\dagger A_{trn} A_{trn}^T (\hat{A}_{trn}^\dagger)^T \hat{A}_{trn}^\dagger u) \right] = \frac{1 + c + \mu^2 c}{2\sqrt{(1 - c + c\mu^2)^2 + 4c^2 \mu^2}} - \frac{1}{2} + o(1)$$

727 and again using Lemma 11, the variance is $o(1)$. Finally,

$$\hat{A}_{trn}^\dagger A_{trn} = \hat{V} \hat{\Sigma}^{-1} \Sigma V.$$

728 Thus,

$$\mathrm{Tr}(\hat{v}_{trn}^T \hat{A}_{trn}^\dagger A_{trn} v_{trn}) = \sum_{i=1}^d a_i^2 \frac{\sigma_i(A)^2}{\sigma_i(A)^2 + \mu^2}.$$

729 Thus, using Lemma 9, we get that

$$\mathbb{E}_{A_{trn}} \left[\mathrm{Tr}(\hat{v}_{trn}^T \hat{A}_{trn}^\dagger A_{trn} v_{trn}) \right] = \frac{1}{2} + \frac{1 + \mu^2 c - \sqrt{(1 - c + \mu^2 c)^2 + 4c^2 \mu^2}}{2c} + o(1)$$

730 and using Lemma 11, the variance is $o(1)$. Then similar to the proof of Theorem 1, we can simplify
731 the above expression to get the final result. \square

732 E.7 Proof of Proposition 1

733 **Proposition 1** (Optimal σ_{trn}). *The optimal value of σ_{trn}^2 for $c < 1$ is given by*

$$\sigma_{trn}^2 = \frac{\sigma_{tst}^2 d [2c(\mu^2 + 1)^2 - 2T(c\mu^2 + c + 1) + 2(c\mu^2 - 2c + 1)] + N_{tst}(\mu^2 c^2 + c^2 + 1 - T)}{N_{tst}(c^3(\mu^2 + 1)^2 - T(\mu^2 c^2 + c^2 - 1) - 2c^2 - 1)}.$$

734 *Proof.* Let $\sigma := \sigma_{trn}^2$ and

$$F = \tau^{-2} \left(\frac{\sigma_{tst}^2}{N_{tst}} + \frac{1}{d} (\sigma \|\hat{h}\|_2^2 + \sigma^2 \|\hat{t}\|_{2\rho}^4) \right).$$

735 Notice that only τ is a function of σ , $\|\hat{h}\|_2^2$, $\|\hat{t}\|_2^2$, and $\|\hat{k}\|_2^2$ are all functions of μ . Then

$$\begin{aligned} \frac{\partial F}{\partial \sigma} &= \tau^{-2} \frac{1}{d} (\|\hat{h}\|_2^2 + 2\sigma \|\hat{t}\|_{2\rho}^4) - 2\tau^{-3} \frac{\partial \tau}{\partial \sigma} \left(\frac{\sigma_{tst}^2}{N_{tst}} + \frac{1}{d} (\sigma \|\hat{h}\|_2^2 + \sigma^2 \|\hat{t}\|_{2\rho}^4) \right) \\ &= \tau^{-2} \frac{1}{d} (\|\hat{h}\|_2^2 + 2\sigma \|\hat{t}\|_{2\rho}^4) - 2\tau^{-3} \|\hat{t}\|_2^2 \|\hat{k}\|_2^2 \left(\frac{\sigma_{tst}^2}{N_{tst}} + \frac{1}{d} (\sigma \|\hat{h}\|_2^2 + \sigma^2 \|\hat{t}\|_{2\rho}^4) \right) \\ &= \tau^{-2} \left(\frac{1}{d} (\|\hat{h}\|_2^2 + 2\sigma \|\hat{t}\|_{2\rho}^4) - 2\tau^{-1} \|\hat{t}\|_2^2 \|\hat{k}\|_2^2 \left(\frac{\sigma_{tst}^2}{N_{tst}} + \frac{1}{d} (\sigma \|\hat{h}\|_2^2 + \sigma^2 \|\hat{t}\|_{2\rho}^4) \right) \right). \end{aligned}$$

736 The optimal σ^* satisfies $\frac{\partial F}{\partial \sigma} \Big|_{\sigma=\sigma^*} = 0$. Thus, we can solve the equation

$$\tau^{-2} = 0 \quad \text{or} \quad \frac{1}{d} (\|\hat{h}\|_2^2 + 2\sigma \|\hat{t}\|_{2\rho}^4) - 2\tau^{-1} \|\hat{t}\|_2^2 \|\hat{k}\|_2^2 \left(\frac{\sigma_{tst}^2}{N_{tst}} + \frac{1}{d} (\sigma \|\hat{h}\|_2^2 + \sigma^2 \|\hat{t}\|_{2\rho}^4) \right).$$

737 Let $\alpha := \|\hat{t}\|_2^2 \|\hat{k}\|_2^2$, $\delta := d \frac{\sigma_{tst}^2}{N_{tst}}$. Then

$$\tau^{-2} = 0 \implies \sigma = -\frac{1}{\|\hat{t}\|_2^2 \|\hat{k}\|_2^2}.$$

738 Notice that $\sigma < 0$ implies σ_{trn} is an imaginary number, something we don't want. Thus, we look at
739 the other expression.

$$\begin{aligned} 0 &= \frac{1}{d} (\|\hat{h}\|_2^2 + 2\sigma \|\hat{t}\|_{2\rho}^4) - 2\tau^{-1} \|\hat{t}\|_2^2 \|\hat{k}\|_2^2 \left(\frac{\sigma_{tst}^2}{N_{tst}} + \frac{1}{d} (\sigma \|\hat{h}\|_2^2 + \sigma^2 \|\hat{t}\|_{2\rho}^4) \right) \\ &= \frac{1}{d} (\|\hat{h}\|_2^2 + 2\sigma \|\hat{t}\|_{2\rho}^4) - 2\tau^{-1} \alpha \left(\frac{\delta}{d} + \frac{1}{d} (\sigma \|\hat{h}\|_2^2 + \sigma^2 \|\hat{t}\|_{2\rho}^4) \right). \quad [\alpha = \|\hat{t}\|_2^2 \|\hat{k}\|_2^2] \end{aligned}$$

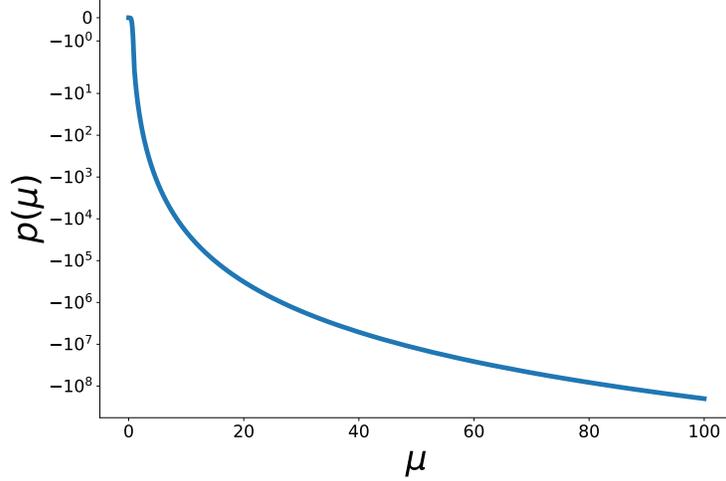


Figure 19: Figure showing the value of $p(\mu)$

740 Then multiplying through by d and τ

$$\begin{aligned}
0 &= (1 + \alpha\sigma)(\|\hat{h}\|_2^2 + 2\sigma\|\hat{t}\|_2^4\rho) - 2\alpha(\delta + \sigma\|\hat{h}\|_2^2 + \sigma^2\|\hat{t}\|_2^4\rho) & [\tau = 1 + \alpha\sigma] \\
&= \|\hat{h}\|_2^2 + 2\|\hat{t}\|_2^4\rho\sigma + \alpha\|\hat{h}\|_2^2\sigma + 2\alpha\|\hat{t}\|_2^4\rho\sigma^2 - 2\alpha\delta - 2\alpha\|\hat{h}\|_2^2\sigma - 2\alpha\|\hat{t}\|_2^4\rho\sigma^2 \\
&= \|\hat{h}\|_2^2 + 2\|\hat{t}\|_2^4\rho\sigma + \alpha\|\hat{h}\|_2^2\sigma - 2\alpha\delta - 2\alpha\|\hat{h}\|_2^2\sigma.
\end{aligned}$$

741 Then solving for σ , we get that

$$\sigma = \frac{2\alpha\delta - \|\hat{h}\|_2^2}{2\|\hat{t}\|_2^4\rho - \alpha\|\hat{h}\|_2^2} = \frac{2d\|\hat{t}\|_2^2\|\hat{k}\|_2^2\sigma_{tst}^2 - \|\hat{h}\|_2^2N_{tst}}{N_{tst}(2\|\hat{t}\|_2^4\rho - \|\hat{t}\|_2^2\|\hat{k}\|_2^2\|\hat{h}\|_2^2)}.$$

742 Then we use the random matrix theory lemmas to estimate this quantity. \square

743 G Experiments

744 All experiments were conducted using Pytorch and run on Google Colab using an A100 GPU. For
745 each empirical data point, we did at least 100 trials. The maximum number of trials for any experiment
746 was 20000 trials.

747 For each configuration of the parameters, N_{trn} , N_{tst} , d , σ_{trn} , σ_{tst} , and μ . For each trial, we sampled
748 u , v_{trn} , v_{tst} uniformly at random from the appropriate dimensional sphere. We also sampled new
749 training and test noise for each trial.

750 For the data scaling regime, we kept $d = 1000$ and for the parameter scaling regime, we kept
751 $N_{trn} = 1000$. For all experiments, $N_{tst} = 1000$.

752 H Technical Assumption on μ

753 Notice that we had this assumption that $p(\mu) < 0$. We compute $p(\mu)$ for a million equally spaced
754 points in $(0, 100]$ and see that $p(\mu) < 0$. Here we use Mpmath with a precision of 1000. The result is
755 shown in Figure 19. Hence we see that the assumption is satisfied for $\mu \in (0, 100]$.