
A random matrix analysis of random Fourier features: beyond the Gaussian kernel, a precise phase transition, and the corresponding double descent

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This article characterizes the exact asymptotics of random Fourier feature (RFF)
2 regression, in the realistic setting where the number of data samples n , their
3 dimension p , and the dimension of feature space N are all large and comparable.
4 In this regime, the random RFF Gram matrix no longer converges to the well-
5 known limiting Gaussian kernel matrix (as it does when $N \rightarrow \infty$ alone), but it
6 still has a tractable behavior that is captured by our analysis. This analysis also
7 provides accurate estimates of training and test regression errors for large n, p, N .
8 Based on these estimates, a precise characterization of two qualitatively different
9 phases of learning, including the phase transition between them, is provided;
10 and the corresponding double descent test error curve is derived from this phase
11 transition behavior. These results do not depend on strong assumptions on the data
12 distribution, and they perfectly match empirical results on real-world data sets.

13 1 Introduction

14 For a machine learning system having N parameters, trained on a data set of size n , asymptotic
15 analysis as used in classical statistical learning theory typically either focuses on the (statistical)
16 population $n \rightarrow \infty$ limit, for N fixed, or the over-parameterized $N \rightarrow \infty$ limit, for a given n . These
17 two settings are technically more convenient to work with, yet less practical, as they essentially
18 assume that one of the two dimensions is negligibly small compared to the other, and this is rarely the
19 case in practice. Indeed, with a factor of 2 or 10 more data, one typically works with a more complex
20 model. This has been highlighted perhaps most prominently in recent work on neural network models,
21 in which model complexity and data size increase together. For this reason, the “double asymptotic”
22 regime where $n, N \rightarrow \infty$, with $N/n \rightarrow \alpha$, a constant, is a particularly interesting (and likely more
23 realistic) limit, despite being technically more challenging [46, 49, 22, 16, 36, 30, 31, 5]. In particular,
24 working in this regime allows for a finer quantitative assessment of machine learning systems, as
25 a function of their *relative* complexity N/n , as well as for a precise description of the under- to
26 over-parameterized “phase transition” (that does not appear in the $N \rightarrow \infty$ alone analysis). This
27 transition is largely hidden in the usual style of statistical learning theory analysis [47], but it is
28 well-known in the statistical mechanics approach to learning theory [46, 49, 22, 16, 30], and empirical
29 signatures of it have received attention recently under the name “double descent” phenomena [1, 7].

30 This article addresses the example of the systematically exploited yet, as we shall see, quite inappropriate
31 large- N alone asymptotics of random Fourier features [42] and more generally random feature
32 maps, which may be viewed also as a single-hidden-layer neural network model. More precisely,
33 let $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$ denote the data matrix of size n with data vectors $\mathbf{x}_i \in \mathbb{R}^p$ as column
34 vectors. The random feature matrix $\Sigma_{\mathbf{X}}$ of \mathbf{X} is generated by multiplying some random matrix

35 $\mathbf{W} \in \mathbb{R}^{N \times p}$ having i.i.d. entries and then passing through some *entry-wise* nonlinear function $\sigma(\cdot)$,
 36 i.e., $\Sigma_{\mathbf{X}} \equiv \sigma(\mathbf{W}\mathbf{X}) \in \mathbb{R}^{N \times n}$. Popularly used random feature techniques such as random Fourier
 37 features (RFFs) [42] and homogeneous kernel maps [48], however, rarely involve a single nonlinearity.
 38 The popular RFF maps are built with cosine and sine nonlinearities, so that $\Sigma_{\mathbf{X}} \in \mathbb{R}^{2N \times n}$ is obtained
 39 by cascading the random features of both, i.e., $\Sigma_{\mathbf{X}}^T \equiv [\cos(\mathbf{W}\mathbf{X})^T, \sin(\mathbf{W}\mathbf{X})^T]$. Note that, by
 40 combining both nonlinearities, RFFs generated from $\mathbf{W} \in \mathbb{R}^{N \times p}$ are of dimension $2N$.

41 The large N asymptotics of random feature maps is closely related to their limiting kernel $\mathbf{K}_{\mathbf{X}}$.
 42 In the case of RFF, it was shown in [42] that *entry-wise* the Gram matrix $\Sigma_{\mathbf{X}}^T \Sigma_{\mathbf{X}}/N$ converges
 43 to the Gaussian kernel matrix $\mathbf{K}_{\mathbf{X}} \equiv \{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/2)\}_{i,j=1}^n$, as $N \rightarrow \infty$. This fol-
 44 lows from $\frac{1}{N}[\Sigma_{\mathbf{X}}^T \Sigma_{\mathbf{X}}]_{ij} = \frac{1}{N} \sum_{t=1}^N \cos(\mathbf{x}_i^T \mathbf{w}_t) \cos(\mathbf{w}_t^T \mathbf{x}_j) + \sin(\mathbf{x}_i^T \mathbf{w}_t) \sin(\mathbf{w}_t^T \mathbf{x}_j)$, for \mathbf{w}_t in-
 45 dependent Gaussian random vectors, so that by the strong law of large numbers, for fixed n, p ,
 46 $[\Sigma_{\mathbf{X}}^T \Sigma_{\mathbf{X}}/N]_{ij}$ goes to its expectation (with respect to $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$) almost surely as $N \rightarrow \infty$, i.e.,
 47 $[\Sigma_{\mathbf{X}}^T \Sigma_{\mathbf{X}}/N]_{ij} \xrightarrow{\text{a.s.}} \mathbb{E}_{\mathbf{w}} [\cos(\mathbf{x}_i^T \mathbf{w}) \cos(\mathbf{w}^T \mathbf{x}_j) + \sin(\mathbf{x}_i^T \mathbf{w}) \sin(\mathbf{w}^T \mathbf{x}_j)] \equiv \mathbf{K}_{\cos} + \mathbf{K}_{\sin}$, with

$$\mathbf{K}_{\cos} + \mathbf{K}_{\sin} \equiv e^{-\frac{1}{2}(\|\mathbf{x}_i\|^2 + \|\mathbf{x}_j\|^2)} (\cosh(\mathbf{x}_i^T \mathbf{x}_j) + \sinh(\mathbf{x}_i^T \mathbf{x}_j)) = e^{-\frac{1}{2}(\|\mathbf{x}_i - \mathbf{x}_j\|^2)} \equiv [\mathbf{K}_{\mathbf{X}}]_{ij}. \quad (1)$$

48 While this result holds in the $N \rightarrow \infty$ limit, recent advances in random matrix theory [28] suggest
 49 that, in the more practical setting where N is not much larger than n, p and $n, p, N \rightarrow \infty$ at the same
 50 pace, the situation is more subtle. In particular, the above entry-wise convergence remains valid, but
 51 the convergence $\|\Sigma_{\mathbf{X}}^T \Sigma_{\mathbf{X}}/N - \mathbf{K}_{\mathbf{X}}\| \rightarrow 0$ no longer holds in spectral norm, due to the factor n , now
 52 large, in the norm inequality $\|\mathbf{A}\|_{\infty} \leq \|\mathbf{A}\| \leq n\|\mathbf{A}\|_{\infty}$ for $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $\|\mathbf{A}\|_{\infty} \equiv \max_{ij} |\mathbf{A}_{ij}|$.
 53 This implies that, in this large n, p, N regime, the assessment of the behavior of $\Sigma_{\mathbf{X}}^T \Sigma_{\mathbf{X}}/N$ via $\mathbf{K}_{\mathbf{X}}$
 54 may result in a spectral norm error that blows up. As a consequence, for various machine learning
 55 algorithms [10], the performance guarantee offered by the limiting Gaussian kernel is less likely to
 56 agree with empirical observations in real-world large-scale problems, when n, p are large.¹

57 1.1 Our Main Contributions

58 We consider the RFF model in the more realistic large n, p, N limit. While, in this setting, the RFF
 59 empirical Gram matrix does *not* converge to the Gaussian kernel matrix, we can characterize the
 60 Gram matrix behavior as $n, p, N \rightarrow \infty$ and provide *asymptotic performance guarantees* for RFF on
 61 large-scale problems. We also identify a phase transition as a function of the ratio N/n , including the
 62 corresponding double descent phenomenon. In more detail, our contributions are the following.

63 1. We provide a *precise* characterization of the asymptotics of the RFF empirical Gram matrix, in the
 64 large n, p, N limit (Theorem 1). This is accomplished by constructing a deterministic equivalent for
 65 the resolvent of the RFF Gram matrix. Based on this, the behavior of the RFF model is (asymptot-
 66 ically) accessible through a fixed-point equation, that can be interpreted in terms of an angle-like
 67 correction induced by the non-trivial large n, p, N limit (relative to the $N \rightarrow \infty$ alone limit).

68 2. We derive the asymptotic training and test mean squared errors (MSEs) of RFF ridge regression, as
 69 a function of the ratio N/n , regularization penalty λ , training as well as test sets (Theorem 2 and 3,
 70 respectively). We identify precisely the under- to over-parameterization phase transition, as a function
 71 of the relative model complexity N/n ; we prove the existence of a “singular” peak of test error at the
 72 $N/n = 1/2$ boundary; and we characterize the corresponding *double descent* behavior. Importantly,
 73 our result is valid *with almost no specific assumption* on the data distribution. This is a significant
 74 improvement over existing double descent analyses, which fundamentally rely on the knowledge of
 75 the data distribution (often assumed to be multivariate Gaussian for simplicity) [21, 35].

76 3. We provide a detailed empirical evaluation of our theoretical results, demonstrating that the theory
 77 closely matches empirical results on a range of real-world data sets (Section 3 and Section F in the
 78 supplementary material). This includes the correction due to the large n, p, N limit, sharp transitions
 79 (as a function of N/n) in angle-like quantities, and the corresponding double descent. This also
 80 includes an evaluation of the impact of training-test similarity and the effect of different data sets,
 81 thus confirming, as stated in 2., that (unlike in prior work) the phase transition and double descent
 82 hold with almost no specific assumption on the data distribution.

¹For readers not familiar with the impact of spectral norm error in learning, or with the random matrix theory techniques that we will use in our analysis, such as resolvent analysis and the use of deterministic equivalents, see Appendix A for a warm-up discussion.

1.2 Related Work

Here, we provide a brief review of related previous efforts.

Random features and limiting kernels. In most RFF work [43, 4, 3, 44], non-asymptotic bounds are given, on the number of random features N needed for a predefined approximation error, for a given kernel matrix with fixed n, p . A more recent line of work [2, 15, 23, 9] has focused on the over-parameterized $N \rightarrow \infty$ limit of large neural networks by studying the corresponding *neural tangent kernels*. Here, we position ourselves in the more practical regime where n, p, N are all large and comparable, and provide *asymptotic performance guarantees* that better fit large-scale problems.

Random matrix theory. From a random matrix theory perspective, nonlinear Gram matrices of the type $\Sigma_{\mathbf{X}}^T \Sigma_{\mathbf{X}}$ have recently received an unprecedented research interests, due to their close connection to neural networks [40, 38, 8, 37], with a particular focus on the associated eigenvalue distribution. Here we propose a deterministic equivalent [11, 20] analysis for the resolvent matrix that provides access, not only to the eigenvalue distribution, but also to the regression error of central interest in this article. While most existing deterministic equivalent analyses are performed on linear models, here we focus on the *nonlinear* RFF model. From a technical perspective, the most relevant work is [28, 35]. We improve their results by considering *generic* data model on the nonlinear RFF model.

Statistical mechanics of learning. A long history of connections between statistical mechanics and machine learning models (such as neural networks) exists, including a range of techniques to establish generalization bounds [46, 49, 22, 16], and recently there has been renewed interest [30, 31, 33, 32, 34, 5]. Their relevance to our results lies in the use of the thermodynamic limit (akin to the large n, p, N limit), rather than the classical limits more commonly used in statistical learning theory, where uniform convergence bounds and related techniques can be applied.

Double descent in large-scale learning systems. The large n, N asymptotics of statistical models has received considerable research interests in machine learning [39, 21, 13], resulting in a (somehow) counterintuitive phenomenon referred to as the “double descent.” Instead of focusing on different “phases of learning” [46, 49, 22, 16, 30], the “double descent” phenomenon focuses on an empirical manifestation of the phase boundary and refers to the empirical observations of the test error curve as a function of the model complexity, which differs from the usual textbook description of the bias-variance tradeoff [1, 7, 18]. Theoretical investigation into this phenomenon mainly focuses on various regression models [14, 13, 6, 12, 26, 21, 35]. In most cases, quite specific (and rather strong) assumptions are imposed on the input data distribution. In this respect, our work extends the analysis in [35] to handle the RFF model and its phase structure *on real-world data sets*.

2 Main Technical Results

In this section, we present our main theoretical results. To investigate the large n, p, N asymptotics of the RFF model, we shall technically position ourselves under the following assumption.

Assumption 1. As $n \rightarrow \infty$, we have

1. $0 < \liminf_n \min\{\frac{p}{n}, \frac{N}{n}\} \leq \limsup_n \max\{\frac{p}{n}, \frac{N}{n}\} < \infty$; or, *practically speaking*, the ratios p/n and N/n are only moderately large or moderately small.

2. $\limsup_n \|\mathbf{X}\| < \infty$ and $\limsup_n \|\mathbf{y}\|_\infty < \infty$, i.e., they are normalized with respect to n .

Under Assumption 1, we consider the RFF regression model. For training data $\mathbf{X} \in \mathbb{R}^{p \times n}$ of size n , the associated random Fourier features, $\Sigma_{\mathbf{X}} \in \mathbb{R}^{2N \times n}$, are obtained by computing $\mathbf{W}\mathbf{X} \in \mathbb{R}^{N \times n}$, for standard Gaussian random matrix $\mathbf{W} \in \mathbb{R}^{N \times p}$, and then applying entry-wise cosine and sine nonlinearities on $\mathbf{W}\mathbf{X}$, i.e., $\Sigma_{\mathbf{X}}^T = [\cos(\mathbf{W}\mathbf{X})^T, \sin(\mathbf{W}\mathbf{X})^T]$ with $\mathbf{W}_{ij} \sim \mathcal{N}(0, 1)$. Given this setup, the RFF ridge regressor $\beta \in \mathbb{R}^{2N}$ is given by, for $\lambda \geq 0$,

$$\beta \equiv \frac{1}{n} \Sigma_{\mathbf{X}} \left(\frac{1}{n} \Sigma_{\mathbf{X}}^T \Sigma_{\mathbf{X}} + \lambda \mathbf{I}_n \right)^{-1} \mathbf{y} \cdot \mathbf{1}_{2N \geq n} + \left(\frac{1}{n} \Sigma_{\mathbf{X}} \Sigma_{\mathbf{X}}^T + \lambda \mathbf{I}_{2N} \right)^{-1} \frac{1}{n} \Sigma_{\mathbf{X}} \mathbf{y} \cdot \mathbf{1}_{2N < n}. \quad (2)$$

The two forms of β in (2) are equivalent for any $\lambda > 0$ and minimize the (ridge-regularized) squared loss $\frac{1}{n} \|\mathbf{y} - \Sigma_{\mathbf{X}}^T \beta\|^2 + \lambda \|\beta\|^2$ on the training set (\mathbf{X}, \mathbf{y}) . Our objective is to characterize the large

129 n, p, N asymptotics of both the *training MSE*, E_{train} , and the *test MSE*, E_{test} , defined as

$$E_{\text{train}} = \frac{1}{n} \|\mathbf{y} - \Sigma_{\mathbf{X}}^{\top} \beta\|^2, \quad E_{\text{test}} = \frac{1}{\hat{n}} \|\hat{\mathbf{y}} - \Sigma_{\hat{\mathbf{X}}}^{\top} \beta\|^2, \quad (3)$$

130 with $\Sigma_{\hat{\mathbf{X}}}^{\top} \equiv [\cos(\mathbf{W}\hat{\mathbf{X}})^{\top}, \sin(\mathbf{W}\hat{\mathbf{X}})^{\top}] \in \mathbb{R}^{\hat{n} \times 2N}$ on a test set $(\hat{\mathbf{X}}, \hat{\mathbf{y}})$ of size \hat{n} .

131 2.1 Asymptotic Deterministic Equivalent

132 To start, we observe that the training MSE, E_{train} , in (3), can be written as $E_{\text{train}} = \frac{\lambda^2}{n} \|\mathbf{Q}(\lambda) \mathbf{y}\|^2 =$
 133 $-\frac{\lambda^2}{n} \mathbf{y}^{\top} \partial \mathbf{Q}(\lambda) \mathbf{y} / \partial \lambda$, which depends on the quadratic form $\mathbf{y}^{\top} \mathbf{Q}(\lambda) \mathbf{y}$ of

$$\mathbf{Q}(\lambda) \equiv \left(\frac{1}{n} \Sigma_{\mathbf{X}}^{\top} \Sigma_{\mathbf{X}} + \lambda \mathbf{I}_n \right)^{-1} \in \mathbb{R}^{n \times n}, \quad (4)$$

134 the *resolvent* of $\frac{1}{n} \Sigma_{\mathbf{X}}^{\top} \Sigma_{\mathbf{X}}$ (also denoted \mathbf{Q} when there is no ambiguity) with $\lambda > 0$.

135 In order to assess the asymptotic training MSE, it suffices to find a deterministic equivalent for $\mathbf{Q}(\lambda)$
 136 (i.e., a *deterministic* matrix that captures the asymptotic behavior of the latter). One possibility is
 137 the expectation $\mathbb{E}_{\mathbf{W}}[\mathbf{Q}(\lambda)]$. Informally, if the training MSE E_{train} (that is random due to random
 138 \mathbf{W}) is “close to” some deterministic quantity \bar{E}_{train} , in the large n, p, N limit, then \bar{E}_{train} must have
 139 the same limit as $\mathbb{E}_{\mathbf{W}}[E_{\text{train}}] = -\frac{\lambda^2}{n} \partial \mathbf{y}^{\top} \mathbb{E}_{\mathbf{W}}[\mathbf{Q}(\lambda)] \mathbf{y} / \partial \lambda$ for $n, p, N \rightarrow \infty$. However, $\mathbb{E}_{\mathbf{W}}[\mathbf{Q}]$
 140 involves integration (with no closed form due to the matrix inverse), and it is not a convenient quantity
 141 with which to work. Our objective is to find an asymptotic “alternative” for $\mathbb{E}_{\mathbf{W}}[\mathbf{Q}]$ that is (i) close
 142 to $\mathbb{E}_{\mathbf{W}}[\mathbf{Q}]$ in the large $n, p, N \rightarrow \infty$ limit and (ii) numerically more accessible.

143 In the following theorem (proved in Section B of the appendix), we introduce an asymptotic equivalent
 144 for $\mathbb{E}_{\mathbf{W}}[\mathbf{Q}]$. Instead of being directly related to the Gaussian kernel $\mathbf{K}_{\mathbf{X}} = \mathbf{K}_{\text{cos}} + \mathbf{K}_{\text{sin}}$ as suggested
 145 by (1) in the large- N limit, it depends on the two components $\mathbf{K}_{\text{cos}}, \mathbf{K}_{\text{sin}}$ in a more involved manner.

146 **Theorem 1** (Asymptotic equivalent for $\mathbb{E}_{\mathbf{W}}[\mathbf{Q}]$). *Under Assumption 1, for \mathbf{Q} defined in (4) and*
 147 *$\lambda > 0$, we have, as $n \rightarrow \infty$*

$$\|\mathbb{E}_{\mathbf{W}}[\mathbf{Q}] - \bar{\mathbf{Q}}\| \rightarrow 0$$

148 *for $\bar{\mathbf{Q}} \equiv \left(\frac{N}{n} \left(\frac{\mathbf{K}_{\text{cos}}}{1+\delta_{\text{cos}}} + \frac{\mathbf{K}_{\text{sin}}}{1+\delta_{\text{sin}}} \right) + \lambda \mathbf{I}_n \right)^{-1}$, $\mathbf{K}_{\text{cos}} \equiv \mathbf{K}_{\text{cos}}(\mathbf{X}, \mathbf{X})$, $\mathbf{K}_{\text{sin}} \equiv \mathbf{K}_{\text{sin}}(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{n \times n}$ and*

$$\mathbf{K}_{\text{cos}}(\mathbf{X}, \mathbf{X}')_{ij} = e^{-\frac{\|\mathbf{x}_i\|^2 + \|\mathbf{x}'_j\|^2}{2}} \cosh(\mathbf{x}_i^{\top} \mathbf{x}'_j), \quad \mathbf{K}_{\text{sin}}(\mathbf{X}, \mathbf{X}')_{ij} = e^{-\frac{\|\mathbf{x}_i\|^2 + \|\mathbf{x}'_j\|^2}{2}} \sinh(\mathbf{x}_i^{\top} \mathbf{x}'_j), \quad (5)$$

149 *where $(\delta_{\text{cos}}, \delta_{\text{sin}})$ is the unique positive solution to $\delta_{\text{cos}} = \frac{1}{n} \text{tr}(\mathbf{K}_{\text{cos}} \bar{\mathbf{Q}})$, $\delta_{\text{sin}} = \frac{1}{n} \text{tr}(\mathbf{K}_{\text{sin}} \bar{\mathbf{Q}})$.*

150 **Remark 1** (Correction to large- N behavior). Taking $N/n \rightarrow \infty$, one has $\delta_{\text{cos}} \rightarrow 0$, $\delta_{\text{sin}} \rightarrow 0$ so
 151 that $\frac{\mathbf{K}_{\text{cos}}}{1+\delta_{\text{cos}}} + \frac{\mathbf{K}_{\text{sin}}}{1+\delta_{\text{sin}}} \rightarrow \mathbf{K}_{\text{cos}} + \mathbf{K}_{\text{sin}} = \mathbf{K}$ and $\bar{\mathbf{Q}} \sim \frac{n}{N} \mathbf{K}^{-1}$, for $\lambda > 0$, in accordance with the
 152 large- N asymptotic prediction. In this sense, the pair $(\delta_{\text{cos}}, \delta_{\text{sin}})$ introduced in Theorem 1 accounts
 153 for the “correction” due to the non-trivial n/N , as opposed to the $N \rightarrow \infty$ alone analysis. Also, in
 154 the $N/n \rightarrow \infty$ limit, when the number of features is large, the regularization effect of λ flattens out
 155 and $\bar{\mathbf{Q}}$ behaves like (a scaled version of) the inverse Gaussian kernel matrix \mathbf{K}^{-1} .

156 **Remark 2** (Geometric interpretation). Since $\bar{\mathbf{Q}}$ shares the same eigenspace with $\frac{\mathbf{K}_{\text{cos}}}{1+\delta_{\text{cos}}} + \frac{\mathbf{K}_{\text{sin}}}{1+\delta_{\text{sin}}}$, one
 157 can geometrically interpret $(\delta_{\text{cos}}, \delta_{\text{sin}})$ as a sort of “angle” between the eigenspace of $\mathbf{K}_{\text{cos}}, \mathbf{K}_{\text{sin}} \in$
 158 $\mathbb{R}^{n \times n}$ and that of $\frac{\mathbf{K}_{\text{cos}}}{1+\delta_{\text{cos}}} + \frac{\mathbf{K}_{\text{sin}}}{1+\delta_{\text{sin}}}$, weighted by the associated eigenvalues. For fixed n , as $N \rightarrow \infty$,
 159 we have $\frac{1}{N} \sum_{t=1}^N \cos(\mathbf{X}^{\top} \mathbf{w}_t) \cos(\mathbf{w}_t^{\top} \mathbf{X}) \rightarrow \mathbf{K}_{\text{cos}}$, $\frac{1}{N} \sum_{t=1}^N \sin(\mathbf{X}^{\top} \mathbf{w}_t) \sin(\mathbf{w}_t^{\top} \mathbf{X}) \rightarrow \mathbf{K}_{\text{sin}}$, the
 160 eigenspaces of which are “orthogonal” to each other, so that $\delta_{\text{cos}}, \delta_{\text{sin}} \rightarrow 0$. On the other hand, as
 161 $N, n \rightarrow \infty$, the eigenspaces of \mathbf{K}_{cos} and \mathbf{K}_{sin} “intersect” with each other, captured by the non-trivial
 162 correction $(\delta_{\text{cos}}, \delta_{\text{sin}})$.

163 2.2 Asymptotic Training Performance

164 Theorem 1 provides an asymptotically more tractable approximation of $\mathbb{E}_{\mathbf{W}}[\mathbf{Q}]$ under the form
 165 of a fixed-point equation. Together with some additional concentration arguments (e.g., from [28,
 166 Theorem 2]), this permits us to provide a complete description of the bilinear forms $\mathbf{a}^{\top} \mathbf{Q} \mathbf{b}$, for
 167 $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ of bounded Euclidean norms, such that $\mathbf{a}^{\top} \mathbf{Q} \mathbf{b} - \mathbf{a}^{\top} \bar{\mathbf{Q}} \mathbf{b} \xrightarrow{\text{a.s.}} 0$, as $n, p, N \rightarrow \infty$. This,
 168 together with the fact that $E_{\text{train}} = \frac{\lambda^2}{n} \mathbf{y}^{\top} \mathbf{Q}(\lambda)^2 \mathbf{y} = -\frac{\lambda^2}{n} \mathbf{y}^{\top} \partial \mathbf{Q}(\lambda) \mathbf{y} / \partial \lambda$, leads to the following
 169 result on the asymptotic training error, the proof of which is given in Section C of the appendix.

170 **Theorem 2** (Asymptotic training performance). *Under Assumption 1, for training MSE, E_{train}*
 171 *defined in (3), as $n \rightarrow \infty$*

$$E_{\text{train}} - \bar{E}_{\text{train}} \xrightarrow{\text{a.s.}} 0, \quad \bar{E}_{\text{train}} = \frac{\lambda^2}{n} \|\bar{\mathbf{Q}}\mathbf{y}\|^2 + \frac{N}{n} \frac{\lambda^2}{n^2} \begin{bmatrix} \frac{\text{tr}(\bar{\mathbf{Q}}\mathbf{K}_{\cos}\bar{\mathbf{Q}})}{(1+\delta_{\cos})^2} & \frac{\text{tr}(\bar{\mathbf{Q}}\mathbf{K}_{\sin}\bar{\mathbf{Q}})}{(1+\delta_{\sin})^2} \end{bmatrix} \Omega \begin{bmatrix} \mathbf{y}^\top \bar{\mathbf{Q}}\mathbf{K}_{\cos}\bar{\mathbf{Q}}\mathbf{y} \\ \mathbf{y}^\top \bar{\mathbf{Q}}\mathbf{K}_{\sin}\bar{\mathbf{Q}}\mathbf{y} \end{bmatrix}$$

172 *over the randomness of \mathbf{W} , for $\bar{\mathbf{Q}}$ defined in Theorem 1 and*

$$\Omega^{-1} \equiv \mathbf{I}_2 - \frac{N}{n} \begin{bmatrix} \frac{1}{n} \frac{\text{tr}(\bar{\mathbf{Q}}\mathbf{K}_{\cos}\bar{\mathbf{Q}}\mathbf{K}_{\cos})}{(1+\delta_{\cos})^2} & \frac{1}{n} \frac{\text{tr}(\bar{\mathbf{Q}}\mathbf{K}_{\cos}\bar{\mathbf{Q}}\mathbf{K}_{\sin})}{(1+\delta_{\sin})^2} \\ \frac{1}{n} \frac{\text{tr}(\bar{\mathbf{Q}}\mathbf{K}_{\cos}\bar{\mathbf{Q}}\mathbf{K}_{\sin})}{(1+\delta_{\cos})^2} & \frac{1}{n} \frac{\text{tr}(\bar{\mathbf{Q}}\mathbf{K}_{\sin}\bar{\mathbf{Q}}\mathbf{K}_{\sin})}{(1+\delta_{\sin})^2} \end{bmatrix}. \quad (6)$$

173 One can show that for a given n and fixed $\lambda > 0$, \bar{E}_{train} decreases as the model size N increases;
 174 and for a given ratio N/n , \bar{E}_{train} increases as the regularization penalty λ grows large, as expected.

175 2.3 Asymptotic Test Performance

176 Theorem 2 holds without any restriction on the training set, (\mathbf{X}, \mathbf{y}) , except for Assumption 1, since
 177 only the randomness of \mathbf{W} is involved, and thus one can simply treat (\mathbf{X}, \mathbf{y}) as known in this
 178 result. This is no longer the case for the test error. Intuitively, the test data $\hat{\mathbf{X}}$ cannot be chosen
 179 arbitrarily, and one must ensure that the test data “behave” statistically like the training data, in
 180 a “well-controlled” manner, so that the test MSE is asymptotically deterministic and bounded as
 181 $n, \hat{n}, p, N \rightarrow \infty$. Following this intuition, we work under the following assumption.

182 **Assumption 2** (Data as concentrated random vectors [27]). *The training data $\mathbf{x}_i \in \mathbb{R}^p, i \in$*
 183 *$\{1, \dots, n\}$, are independently drawn from one of $K > 0$ distribution classes² μ_1, \dots, μ_K . There*
 184 *exist constants $C, \sigma, q > 0$ such that for any $\mathbf{x}_i \sim \mu_k, k \in \{1, \dots, K\}$ and any 1-Lipschitz function*
 185 *$f : \mathbb{R}^p \rightarrow \mathbb{R}$, we have the concentration*

$$\mathbb{P}(|f(\mathbf{x}_i) - \mathbb{E}[f(\mathbf{x}_i)]| > t) \leq Ce^{-(t/\sigma)^q}, \quad t \geq 0. \quad (7)$$

186 *The test data $\hat{\mathbf{x}}_i \sim \mu_k, i \in \{1, \dots, \hat{n}\}$ are mutually independent, but may depend on training data \mathbf{X} .*

187 To facilitate the discussion of the phase transition and the double descent, we do not assume indepen-
 188 dence between training data and test data (but we do assume independence between different columns
 189 within \mathbf{X} and $\hat{\mathbf{X}}$). In this respect, Assumption 2 is weaker than the classical i.i.d. assumption, and it
 190 permits us to illustrate the impact of training-test mismatch on the model performance in Section 3.3.

191 A first example of concentrated random vectors satisfying (7) is the Gaussian vector $\mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ [25].
 192 Moreover, since the concentration property in (7) is stable over Lipschitz transformations [27], it
 193 holds, for any 1-Lipschitz mapping $g : \mathbb{R}^d \rightarrow \mathbb{R}^p$ and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$, that $g(\mathbf{z})$ also satisfies (7).
 194 In this respect, Assumption 2, although seemingly quite restrictive, represents a large family of
 195 “generative models”, including notably the “fake images” generated by modern generative adversarial
 196 networks (GANs) that are, by construction, Lipschitz transformations of large random Gaussian
 197 vectors [19, 45]. As such, from a practical consideration, Assumption 2 provides a more realistic and
 198 flexible statistical model for real-world data.

199 With Assumption 2, we have the following result on the asymptotic test error, proved in Section D.

200 **Theorem 3** (Asymptotic test performance). *Under Assumptions 1 and 2, we have, for test MSE*
 201 *E_{test} defined in (3) and test data $(\hat{\mathbf{X}}, \hat{\mathbf{y}})$ satisfying $\limsup_{\hat{n}} \|\hat{\mathbf{X}}\| < \infty, \limsup_{\hat{n}} \|\hat{\mathbf{y}}\|_\infty < \infty$ with*
 202 *$\hat{n}/n \in (0, \infty)$ that, as $n \rightarrow \infty$*

$$E_{\text{test}} - \bar{E}_{\text{test}} \xrightarrow{\text{a.s.}} 0, \quad \bar{E}_{\text{test}} = \frac{1}{\hat{n}} \|\hat{\mathbf{y}} - \frac{N}{n} \hat{\Phi} \bar{\mathbf{Q}}\mathbf{y}\|^2 + \frac{N^2}{n^2} \frac{1}{\hat{n}} \begin{bmatrix} \frac{\Theta_{\cos}}{(1+\delta_{\cos})^2} & \frac{\Theta_{\sin}}{(1+\delta_{\sin})^2} \end{bmatrix} \Omega \begin{bmatrix} \mathbf{y}^\top \bar{\mathbf{Q}}\mathbf{K}_{\cos}\bar{\mathbf{Q}}\mathbf{y} \\ \mathbf{y}^\top \bar{\mathbf{Q}}\mathbf{K}_{\sin}\bar{\mathbf{Q}}\mathbf{y} \end{bmatrix},$$

203 *over the randomness of \mathbf{W}, \mathbf{X} and $\hat{\mathbf{X}}$, for Ω defined in (6),*

$$\Theta_\sigma = \frac{1}{N} \text{tr} \mathbf{K}_\sigma(\hat{\mathbf{X}}, \hat{\mathbf{X}}) + \frac{N}{n} \frac{1}{n} \text{tr} \bar{\mathbf{Q}} \hat{\Phi}^\top \hat{\Phi} \bar{\mathbf{Q}} \mathbf{K}_\sigma - \frac{2}{n} \text{tr} \bar{\mathbf{Q}} \hat{\Phi}^\top \mathbf{K}_\sigma(\hat{\mathbf{X}}, \mathbf{X}), \quad \sigma \in \{\cos, \sin\}, \quad (8)$$

204 *and $\hat{\Phi} \equiv \frac{\mathbf{K}_{\cos}}{1+\delta_{\cos}} + \frac{\mathbf{K}_{\sin}}{1+\delta_{\sin}}$, $\hat{\Phi} \equiv \frac{\mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})}{1+\delta_{\cos}} + \frac{\mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})}{1+\delta_{\sin}}$, with $\mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X}), \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X}) \in \mathbb{R}^{\hat{n} \times n}$*

205 *and $\mathbf{K}_{\cos}(\hat{\mathbf{X}}, \hat{\mathbf{X}}), \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \hat{\mathbf{X}}) \in \mathbb{R}^{\hat{n} \times \hat{n}}$ defined in (5).*

² $K \geq 2$ is included to cover (multi-class) classification problems; and K should remain fixed as $n, p \rightarrow \infty$.

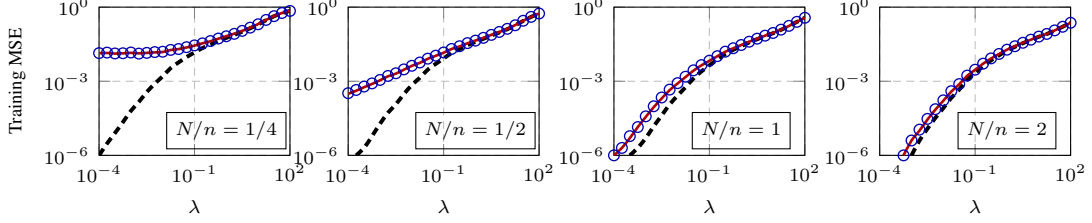


Figure 1: Training MSEs of RFF ridge regression on MNIST data (class 3 versus 7), as a function of regression penalty λ , for $p = 784$, $n = 1000$, $N = 250, 500, 1000, 2000$. Empirical results displayed in blue circles; Gaussian kernel predictions (assuming $N \rightarrow \infty$ alone) in black dashed lines; and Theorems 2 and 3 in red solid lines. Results obtained by averaging over 30 runs.

206 Taking $(\hat{\mathbf{X}}, \hat{\mathbf{y}}) = (\mathbf{X}, \mathbf{y})$, one gets $\bar{E}_{\text{test}} = \bar{E}_{\text{train}}$, as expected. From this perspective, Theorem 3
 207 can be seen as an extension of Theorem 2, with the “interaction” between training and test data (e.g.,
 208 test-versus-test $\mathbf{K}_{\sigma}(\hat{\mathbf{X}}, \hat{\mathbf{X}})$ and test-versus-train $\mathbf{K}_{\sigma}(\hat{\mathbf{X}}, \mathbf{X})$ interaction matrices) summarized in the
 209 scalar parameter Θ_{σ} defined in (8), for $\sigma \in \{\cos, \sin\}$.

210 3 Practical Implications

211 In this section, we provide a detailed empirical evaluation, including a discussion of the behavior
 212 of the fixed point equation in Theorem 1, and its consequences in Theorem 2 and Theorem 3. In
 213 particular, we describe the behavior of $(\delta_{\cos}, \delta_{\sin})$ that characterizes the necessary correction in the
 214 large n, p, N regime, as a function of the regularization λ and the ratio N/n . This explains: (i) the
 215 mismatch between empirical regression errors from the Gaussian kernel prediction (Figure 1); and
 216 (ii) the behavior of $(\delta_{\cos}, \delta_{\sin})$ as a function of N/n , which clearly indicates two phases of learning
 217 (Figure 3) and the corresponding double descent test error curves (Figure 4).

218 3.1 Correction due to the Large n, p, N Regime

219 The nonlinear Gram matrix $\frac{1}{n} \Sigma_{\mathbf{X}}^T \Sigma_{\mathbf{X}}$ is *not* close to the classical Gaussian kernel matrix \mathbf{K} in the
 220 large n, p, N regime; and, as a consequence, its resolvent $\bar{\mathbf{Q}}$, as well the training and test MSE, E_{train}
 221 and E_{test} (that are functions of $\bar{\mathbf{Q}}$), behave quite differently from the Gaussian kernel predictions.
 222 As already discussed in Remark 1, for $\lambda > 0$, the pair $(\delta_{\cos}, \delta_{\sin})$ characterizes the correction when
 223 considering n, p, N all large, compared to the large- N only asymptotic behavior:

$$\delta_{\cos} = \frac{1}{n} \text{tr} \mathbf{K}_{\cos} \bar{\mathbf{Q}}, \quad \delta_{\sin} = \frac{1}{n} \text{tr} \mathbf{K}_{\sin} \bar{\mathbf{Q}}, \quad \bar{\mathbf{Q}} = \left(\frac{N}{n} \left(\frac{\mathbf{K}_{\cos}}{1 + \delta_{\cos}} + \frac{\mathbf{K}_{\sin}}{1 + \delta_{\sin}} \right) + \lambda \mathbf{I}_n \right)^{-1}. \quad (9)$$

224 To start, Figure 1 compares the training MSE of RFF ridge regression to the predictions from Gaussian
 225 kernel regression and to the predictions from our Theorem 2, on the popular MNIST data set [24].
 226 Observe that there is a huge gap for training errors between empirical results and the classical
 227 Gaussian kernel predictions, especially when $N/n < 1$, while our theory *consistently* fits empirical
 228 observations almost perfectly.

229 Next, from (9) we know that both δ_{\cos} and δ_{\sin} are decreasing functions of λ . (See Lemma 7 in
 230 Appendix E for a proof of this fact.) Figure 2 shows that: (i) over a range of different N/n , both δ_{\cos}
 231 and δ_{\sin} decrease monotonically as λ increases; (ii) the behavior for $N/n < 1$, which is decreasing
 232 from an initial value of $\delta \gg 1$, is very different from the behavior for $N/n \gtrsim 1$, with an initially flat
 233 region where $\delta < 1$ for all values of λ ; and (iii) the impact of regularization λ becomes less significant
 234 as the ratio N/n becomes large. This is in accordance with the limiting behavior of $\bar{\mathbf{Q}} \sim \frac{n}{N} \mathbf{K}^{-1}$ in
 235 Remark 1 that is *independent* of λ as $N/n \rightarrow \infty$.

236 Note also that, while δ_{\cos} and δ_{\sin} can be geometrically interpreted as a sort of weighted “angle”
 237 between different kernels, and therefore one might expect to have $\delta \in [0, 1]$, this is not the case for
 238 the leftmost plot with $N/n = 1/4$. For $N/n = 1/4$, for small values of λ (say $\lambda \lesssim 0.1$), both δ_{\cos}
 239 and δ_{\sin} scale like λ^{-1} , while they are observed to saturate to a fixed $O(1)$ value for $N/n = 1, 4, 16$.
 240 This corresponds to two different phases of learning in the ridgeless $\lambda \rightarrow 0$ limit, as discussed below.

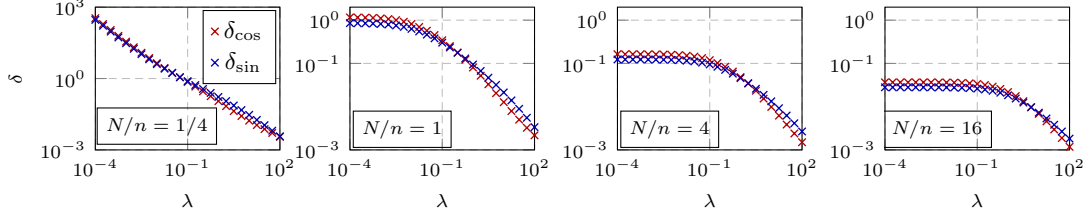


Figure 2: Behavior of $(\delta_{\cos}, \delta_{\sin})$ in (9) on MNIST data set (class 3 versus 7), as a function of the regularization parameter λ , for $p = 784$, $n = 1\,000$, $N = 250, 1\,000, 4\,000, 16\,000$.

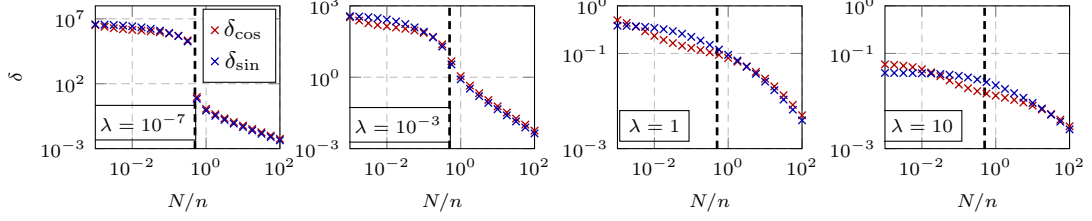


Figure 3: $(\delta_{\cos}, \delta_{\sin})$ in (9) on MNIST data set (class 3 versus 7), as a function of N/n , for $p = 784$, $n = 1\,000$, $\lambda = 10^{-7}, 10^{-3}, 1, 10$. The **black** dashed line is the interpolation threshold $2N = n$.

241 3.2 Phase Transition and Corresponding Double Descent

242 Both δ_{\cos} and δ_{\sin} in (9) are decreasing functions of N . This is depicted empirically in Figure 3.
 243 (See Lemma 6 in Appendix E for a proof.) More importantly, Figure 3 also illustrates that δ_{\cos} and
 244 δ_{\sin} exhibit qualitatively different behavior, depending on the ratio N/n . For λ not too small ($\lambda = 1$
 245 or 10), both δ_{\cos} and δ_{\sin} decrease *smoothly*, as N/n grows large. However, for λ relatively small
 246 ($\lambda = 10^{-3}$ and 10^{-7}), we observe a “phase transition” on two sides of the interpolation threshold
 247 $2N = n$. (Note that the scale of the y-axis is different in different subfigures.) More precisely, in the
 248 leftmost plot with $\lambda = 10^{-7}$, δ_{\cos} and δ_{\sin} “jump” from order $O(1)$ (when $2N > n$) to much higher
 249 values of the order of λ^{-1} (when $2N < n$). A similar behavior is also observed for $\lambda = 10^{-3}$.

250 This phase transition can be theoretically justified by considering the “ridgeless” $\lambda \rightarrow 0$ limit in
 251 Theorem 1. First note that, for $\lambda = 0$ and $2N < n$, the (random) resolvent $\mathbf{Q}(\lambda = 0)$ in (4) is simply
 252 undefined, as it involves inverting a singular matrix $\Sigma_{\mathbf{X}}^T \Sigma_{\mathbf{X}} \in \mathbb{R}^{n \times n}$ that is of rank at most $2N < n$.
 253 As a consequence, we expect to see $\mathbf{Q} \sim \lambda^{-1}$ as $\lambda \rightarrow 0$ for $2N < n$, while for $2N > n$ this is no
 254 longer the case. As a consequence, we have the following two phases:

- 255 1. *Under-parameterized* with $2N < n$. Here, \mathbf{Q} is not well defined (indeed $\mathbf{Q} \sim \lambda^{-1}$) and one must
 256 consider instead the properly scaled $\lambda \delta_{\cos}$, $\lambda \delta_{\sin}$ and $\lambda \bar{\mathbf{Q}}$ as $\lambda \rightarrow 0$.
- 257 2. *Over-parameterized* with $2N > n$, where one can take $\lambda \rightarrow 0$ in (9) to obtain δ_{\cos} , δ_{\sin} and $\bar{\mathbf{Q}}$.

258 **Remark 3** (Double descent test error curve). On account of these two phases of learning, it is not
 259 surprising to observe a “singular” behavior at $2N = n$, when no regularization is applied. More
 260 precisely, we consider the (asymptotic) test MSE in Theorem 3 in the *ridgeless* $\lambda \rightarrow 0$ limit and
 261 focus here on the situation where the test data $\hat{\mathbf{X}}$ is sufficiently different from the training data \mathbf{X}
 262 (see more discussions on this point in Section 3.3 below). Then, the two-by-two matrix Ω defined
 263 in (6) diverges to infinity at $2N = n$ as $\lambda \rightarrow 0$. (Indeed, the determinant $\det(\Omega^{-1})$ scales as λ , per
 264 Lemma 5 in Section E of the supplementary material.) As a consequence, we have $\bar{E}_{\text{test}} \rightarrow \infty$ as
 265 $N/n \rightarrow 1/2$, resulting in a sharp deterioration in the test performance around $2N = n$. It is also
 266 interesting to note that, while Ω also appears in \bar{E}_{train} , we still obtain (asymptotically) zero training
 267 MSE at $2N = n$, despite the divergence of Ω , due to the prefactor λ^2 in \bar{E}_{train} .

268 Figure 4 depicts the empirical and theoretical test MSEs with different λ . In particular, for $\lambda = 10^{-7}$
 269 and $\lambda = 10^{-3}$, a double descent-type behavior is observed, with a singularity at $2N = n$, while for
 270 larger values of λ ($\lambda = 1, 10$), a smoother and monotonically decreasing test error curve is observed,
 271 as a function of N/n , in accordance with the observations in [35] on Gaussian data.

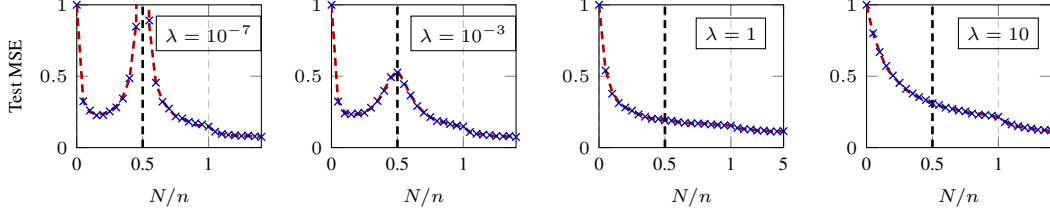


Figure 4: Empirical (blue crosses) and theoretical (red dashed lines) test errors of RFF regression as a function of the ratio N/n , on MNIST data set (class 3 versus 7), for $p = 784$, $n = 500$, $\lambda = 10^{-7}, 10^{-3}, 1, 10$. The black dashed line is the interpolation threshold $2N = n$.

Remark 4 (Double descent as a consequence of phase transition). While the double descent phenomenon has received considerable attention recently, our analysis makes it clear that in this model (and presumably many others) it is a natural consequence of the phase transition between two qualitatively different phases of learning [30].

3.3 Impact of Training-test Similarity

We see that the (asymptotic) test error behaves entirely differently, depending on whether $\hat{\mathbf{X}}$ is “close to” \mathbf{X} or not. For $\hat{\mathbf{X}} = \mathbf{X}$, one has $\bar{E}_{\text{test}} = \bar{E}_{\text{train}}$ that decreases monotonically as N grows large; while for $\hat{\mathbf{X}}$ sufficiently different from \mathbf{X} , \bar{E}_{test} diverges at $2N = n$. To have a more quantitative assessment of the impact of training-test similarity on the model performance, we consider here the special case $\hat{\mathbf{y}} = \mathbf{y}$. Since in the ridgeless $\lambda \rightarrow 0$ limit, Ω scales as λ^{-1} (Remark 3), one must then have $\Theta_\sigma \sim \lambda$ so that \bar{E}_{test} does not diverge at $2N = n$ as $\lambda \rightarrow 0$. A first example is the case where the test data is a small perturbation of the training data. In Figure 5, the test data are generated by adding Gaussian white noise of variance σ^2 to the training data, i.e.,

$$\hat{\mathbf{x}}_i = \mathbf{x}_i + \sigma \varepsilon_i \quad (10)$$

for independent $\varepsilon_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p/p)$. In Figure 5, we observe that (i) below the threshold $\sigma^2 = \lambda$, the test error coincides with the training error and both are close to zero; and (ii) as soon as $\sigma^2 > \lambda$, the test error diverges from the training error and grows large (but linearly in σ^2) as the noise level increases. Note also from the two rightmost plots of Figure 5 that the training-to-test “transition” at $\sigma^2 \sim \lambda$ is sharp only for relatively small values of λ , as predicted by our theory.

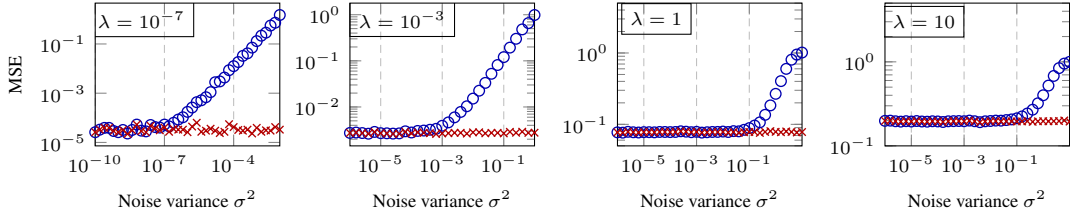


Figure 5: Empirical training (red crosses) and test (blue circles) errors of RFF ridge regression on MNIST data (class 3 versus 7), as a function of the noise level σ^2 , for $N = 512$, $p = 784$, $n = \hat{n} = 1024 = 2N$, $\lambda = 10^{-7}, 10^{-3}, 1, 10$. Results obtained by averaging over 30 runs.

4 Conclusion

We have established a precise description of the resolvent of RFF Gram matrices, and provided asymptotic training and test performance guarantees for RFF ridge regression, in the $n, p, N \rightarrow \infty$ limit. We have also discussed the under- and over-parameterized regimes, where the resolvent behaves dramatically differently. These observations involve only mild regularity assumptions on the data, yielding phase transition behavior and double descent test error curves for RFF regression that closely match experiments on real-world data. Extended to a (technically more involved) multi-layer setting in the more practical large n, p, N setting as in [17], our analysis may shed new light on the theoretical understanding of modern deep neural nets, beyond the large- N alone neural tangent kernel limit.

Broader Impact

In this article, we perform theoretical assessment of the popular random Fourier features (RFFs), in the practical setting where n, p, N are all large and comparable. Asymptotic performance guarantees are provided for RFF ridge regression in this $n, p, N \rightarrow \infty$ limit, as an important positive impact of this work in the development of more reliable large-scale machine learning systems. The developed theoretical framework in this article presents fair and non-offensive societal consequence.

References

- [1] Madhu S Advani and Andrew M Saxe. High-dimensional dynamics of generalization error in neural networks. *arXiv preprint arXiv:1710.03667*, 2017.
- [2] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In *International Conference on Machine Learning*, pages 242–252, 2019.
- [3] Haim Avron, Michael Kapralov, Cameron Musco, Christopher Musco, Ameya Velingker, and Amir Zandieh. Random fourier features for kernel ridge regression: Approximation bounds and statistical guarantees. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 253–262. JMLR. org, 2017.
- [4] Francis Bach. On the equivalence between kernel quadrature rules and random feature expansions. *The Journal of Machine Learning Research*, 18(1):714–751, 2017.
- [5] Y. Bahri, J. Kadmon, J. Pennington, S. Schoenholz, J. Sohl-Dickstein, and S. Ganguli. Statistical mechanics of deep learning. *Annual Review of Condensed Matter Physics*, 11:501–528, 2020.
- [6] Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *arXiv preprint arXiv:1906.11300*, 2019.
- [7] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine learning and the bias-variance trade-off. *arXiv preprint arXiv:1812.11118*, 2018.
- [8] Lucas Benigni and Sandrine Péché. Eigenvalue distribution of nonlinear models of random matrices. *arXiv preprint arXiv:1904.03090*, 2019.
- [9] Lenaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. In *Advances in Neural Information Processing Systems*, pages 2933–2943, 2019.
- [10] Corinna Cortes, Mehryar Mohri, and Ameet Talwalkar. On the impact of kernel approximation on learning accuracy. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 113–120, 2010.
- [11] Romain Couillet and Merouane Debbah. *Random matrix methods for wireless communications*. Cambridge University Press, 2011.
- [12] Zeyu Deng, Abba Kammoun, and Christos Thrampoulidis. A model of double descent for high-dimensional binary linear classification. *arXiv preprint arXiv:1911.05822*, 2019.
- [13] Michał Dereziński, Feynman Liang, and Michael W Mahoney. Exact expressions for double descent and implicit regularization via surrogate random design. *arXiv preprint arXiv:1912.04533*, 2019.
- [14] Edgar Dobriban, Stefan Wager, et al. High-dimensional asymptotics of prediction: Ridge regression and classification. *The Annals of Statistics*, 46(1):247–279, 2018.
- [15] Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, pages 1675–1685, 2019.
- [16] A. Engel and C. P. L. Van den Broeck. *Statistical mechanics of learning*. Cambridge University Press, New York, NY, USA, 2001.
- [17] Zhou Fan and Zhichao Wang. Spectra of the conjugate kernel and neural tangent kernel for linear-width neural networks. *arXiv preprint arXiv:2005.11879*, 2020.
- [18] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*, volume 1. Springer series in statistics New York, 2001.

- [19] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [20] Walid Hachem, Philippe Loubaton, Jamal Najim, et al. Deterministic equivalents for certain functionals of large random matrices. *The Annals of Applied Probability*, 17(3):875–930, 2007.
- [21] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. *arXiv preprint arXiv:1903.08560*, 2019.
- [22] D. Haussler, M. Kearns, H. S. Seung, and N. Tishby. Rigorous learning curve bounds from statistical mechanics. *Machine Learning*, 25(2):195–236, 1996.
- [23] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018.
- [24] Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [25] Michel Ledoux. *The concentration of measure phenomenon*. Number 89. American Mathematical Soc., 2005.
- [26] Tengyuan Liang and Alexander Rakhlin. Just interpolate: Kernel “ridgeless” regression can generalize. *arXiv preprint arXiv:1808.00387*, 2018.
- [27] Cosme Louart and Romain Couillet. Concentration of measure and large random matrices with an application to sample covariance matrices. *arXiv preprint arXiv:1805.08295*, 2018.
- [28] Cosme Louart, Zhenyu Liao, and Romain Couillet. A random matrix approach to neural networks. *The Annals of Applied Probability*, 28(2):1190–1248, 2018.
- [29] Vladimir A Marčenko and Leonid Andreevich Pastur. Distribution of eigenvalues for some sets of random matrices. *Mathematics of the USSR-Sbornik*, 1(4):457, 1967.
- [30] C. H. Martin and M. W. Mahoney. Rethinking generalization requires revisiting old ideas: statistical mechanics approaches and complex learning behavior. Technical Report Preprint: arXiv:1710.09553, 2017.
- [31] C. H. Martin and M. W. Mahoney. Implicit self-regularization in deep neural networks: Evidence from random matrix theory and implications for learning. Technical Report Preprint: arXiv:1810.01075, 2018.
- [32] C. H. Martin and M. W. Mahoney. Statistical mechanics methods for discovering knowledge from modern production quality neural networks. In *Proceedings of the 25th Annual ACM SIGKDD Conference*, pages 3239–3240, 2019.
- [33] C. H. Martin and M. W. Mahoney. Traditional and heavy-tailed self regularization in neural network models. In *Proceedings of the 36th International Conference on Machine Learning*, pages 4284–4293, 2019.
- [34] C. H. Martin and M. W. Mahoney. Heavy-tailed Universality predicts trends in test accuracies for very large pre-trained deep neural networks. In *Proceedings of the 20th SIAM International Conference on Data Mining*, 2020.
- [35] Song Mei and Andrea Montanari. The generalization error of random features regression: Precise asymptotics and double descent curve. *arXiv preprint arXiv:1908.05355*, 2019.
- [36] M. Mézard and A. Montanari. *Information, Physics, and Computation*. Oxford University Press, New York, 2009.
- [37] Leonid Pastur. On random matrices arising in deep neural networks. gaussian case. *arXiv preprint arXiv:2001.06188*, 2020.
- [38] Jeffrey Pennington, Samuel Schoenholz, and Surya Ganguli. Resurrecting the sigmoid in deep learning through dynamical isometry: theory and practice. In *Advances in neural information processing systems*, pages 4785–4795, 2017.
- [39] Jeffrey Pennington, Samuel Schoenholz, and Surya Ganguli. The emergence of spectral universality in deep networks. In *International Conference on Artificial Intelligence and Statistics*, pages 1924–1932, 2018.

- 399 [40] Jeffrey Pennington and Pratik Worah. Nonlinear random matrix theory for deep learning. In
400 *Advances in Neural Information Processing Systems*, pages 2634–2643, 2017.
- 401 [41] Vinay Uday Prabhu. Kannada-mnist: A new handwritten digits dataset for the kannada language.
402 *arXiv preprint arXiv:1908.01242*, 2019.
- 403 [42] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances*
404 *in neural information processing systems*, pages 1177–1184, 2008.
- 405 [43] Ali Rahimi and Benjamin Recht. Weighted sums of random kitchen sinks: Replacing minimiza-
406 tion with randomization in learning. In *Advances in neural information processing systems*,
407 pages 1313–1320, 2009.
- 408 [44] Alessandro Rudi and Lorenzo Rosasco. Generalization properties of learning with random
409 features. In *Advances in Neural Information Processing Systems*, pages 3218–3228, 2017.
- 410 [45] Mohamed El Amine Seddik, Cosme Louart, Mohamed Tamaazousti, and Romain Couillet.
411 Random matrix theory proves that deep learning representations of GAN-data behave as gaussian
412 mixtures. *arXiv preprint arXiv:2001.08370*, 2020.
- 413 [46] H. S. Seung, H. Sompolinsky, and N. Tishby. Statistical mechanics of learning from examples.
414 *Physical Review A*, 45(8):6056–6091, 1992.
- 415 [47] V.N. Vapnik. *Statistical Learning Theory*. John Wiley & Sons, New York, 1998.
- 416 [48] Andrea Vedaldi and Andrew Zisserman. Efficient additive kernels via explicit feature maps.
417 *IEEE transactions on pattern analysis and machine intelligence*, 34(3):480–492, 2012.
- 418 [49] T. L. H. Watkin, A. Rau, and M. Biehl. The statistical mechanics of learning a rule. *Rev. Mod.*
419 *Phys.*, 65(2):499–556, 1993.
- 420 [50] Christopher KI Williams. Computing with infinite networks. *Advances in neural information*
421 *processing systems*, pages 295–301, 1997.
- 422 [51] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset for
423 benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 424 [52] Roy D. Yates. A framework for uplink power control in cellular radio systems. *IEEE Journal*
425 *on selected areas in communications*, 13(7):1341–1347, 1995.

A Warm-up Example: Sample Covariance and the Marčenko-Pastur Equation

Consider the sample covariance matrix $\hat{\mathbf{C}} = \frac{1}{n} \mathbf{X} \mathbf{X}^\top$ from some data $\mathbf{X} \in \mathbb{R}^{p \times n}$ composed of n i.i.d. $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$ with nonnegative definite $\mathbf{C} \in \mathbb{R}^{p \times p}$. In this zero-mean Gaussian setting, the sample covariance $\hat{\mathbf{C}}$, despite being the maximum likelihood estimator of the *population covariance* \mathbf{C} and providing *entry-wise* consistent estimate for it, is an extremely poor estimator of \mathbf{C} in a *spectral norm* sense, for n, p large. More precisely, $\|\hat{\mathbf{C}} - \mathbf{C}\| \not\rightarrow 0$ as $n, p \rightarrow \infty$ with $p/n \in (0, \infty)$. Indeed, one has $\|\hat{\mathbf{C}} - \mathbf{C}\|/\|\mathbf{C}\| \approx 20\%$, even with $n = 100p$, in the simple $\mathbf{C} = \mathbf{I}_p$ setting.

In the regression analysis (such as ridge regression) based on \mathbf{X} , of more immediate interest is the *resolvent* $\mathbf{Q}_{\hat{\mathbf{C}}}(\lambda) \equiv (\hat{\mathbf{C}} + \lambda \mathbf{I}_p)^{-1}$, $\lambda > 0$ of the sample covariance $\hat{\mathbf{C}}$, and more concretely, the bilinear forms of the type $\mathbf{a}^\top \mathbf{Q}_{\hat{\mathbf{C}}}(\lambda) \mathbf{b}$ for $\mathbf{a}, \mathbf{b} \in \mathbb{R}^p$. As a result of the spectral norm inconsistency $\|\hat{\mathbf{C}} - \mathbf{C}\| \not\rightarrow 0$ in the large n, p regime, it is unlikely that for most \mathbf{a}, \mathbf{b} , the convergence $\mathbf{a}^\top \mathbf{Q}_{\hat{\mathbf{C}}}(\lambda) \mathbf{b} - \mathbf{a}^\top (\mathbf{C} + \lambda \mathbf{I}_p)^{-1} \mathbf{b} \rightarrow 0$ would still hold.

While the *random* variable $\mathbf{a}^\top \mathbf{Q}_{\hat{\mathbf{C}}}(\lambda) \mathbf{b}$ is not getting close to $\mathbf{a}^\top (\mathbf{C} + \lambda \mathbf{I}_p)^{-1} \mathbf{b}$ as $n, p \rightarrow \infty$, it does exhibit a tractable asymptotically *deterministic* behavior, described by the Marčenko-Pastur equation [29] for $\mathbf{C} = \mathbf{I}_p$. Notably, for $\mathbf{a}, \mathbf{b} \in \mathbb{R}^p$ deterministic vectors of bounded Euclidean norms, we have, as $n, p \rightarrow \infty$ and $p/n \rightarrow c \in (0, \infty)$,

$$\mathbf{a}^\top \mathbf{Q}_{\hat{\mathbf{C}}}(\lambda) \mathbf{b} - m(\lambda) \mathbf{a}^\top \mathbf{b} \xrightarrow{\text{a.s.}} 0,$$

with $m(\lambda)$ the unique positive solution to the following Marčenko-Pastur equation [29]

$$c\lambda m^2(\lambda) + (1 + \lambda - c)m(\lambda) - 1 = 0. \quad (11)$$

In a sense, $\bar{\mathbf{Q}}(\lambda) \equiv m(\lambda) \mathbf{I}_p$ can be seen as a *deterministic equivalent* [20, 11] for the *random* $\mathbf{Q}_{\hat{\mathbf{C}}}(\lambda)$ that asymptotically characterizes the behavior of the latter, when bilinear forms are considered.

BLUE add more discussion for deterministic equivalents!

B Proof of Theorem 1

Our objective is to prove, under Assumption 1, the asymptotic equivalence between the expectation (over \mathbf{W} , omitted from now on) $\mathbb{E}[\mathbf{Q}]$ and

$$\bar{\mathbf{Q}} \equiv \left(\frac{N}{n} \left(\frac{\mathbf{K}_{\cos}}{1 + \delta_{\cos}} + \frac{\mathbf{K}_{\sin}}{1 + \delta_{\sin}} \right) + \lambda \mathbf{I}_n \right)^{-1}$$

for $\mathbf{K}_{\cos} \equiv \mathbf{K}_{\cos}(\mathbf{X}, \mathbf{X})$, $\mathbf{K}_{\sin} \equiv \mathbf{K}_{\sin}(\mathbf{X}, \mathbf{X}) \in \mathbb{R}^{n \times n}$ defined in (5), with $(\delta_{\cos}, \delta_{\sin})$ the unique positive solution to

$$\delta_{\cos} = \frac{1}{n} \text{tr}(\mathbf{K}_{\cos} \bar{\mathbf{Q}}), \quad \delta_{\sin} = \frac{1}{n} \text{tr}(\mathbf{K}_{\sin} \bar{\mathbf{Q}}).$$

The existence and uniqueness of the above fixed-point equation is standard in random matrix literature and can be reached for instance with the standard interference function framework [52].

The asymptotic equivalence should be announced in the sense that $\|\mathbb{E}[\mathbf{Q}] - \bar{\mathbf{Q}}\| \rightarrow 0$ as $n, p, N \rightarrow \infty$ at the same pace. We shall proceed by introducing an intermediary resolvent $\hat{\mathbf{Q}}$ (see definition in (13)) and show subsequently that

$$\|\mathbb{E}[\mathbf{Q}] - \hat{\mathbf{Q}}\| \rightarrow 0, \quad \|\hat{\mathbf{Q}} - \bar{\mathbf{Q}}\| \rightarrow 0.$$

We start by introducing the following lemma.

Lemma 1 (Expectation of $\sigma_1(\mathbf{x}_i^\top \mathbf{w}) \sigma_2(\mathbf{w}^\top \mathbf{x}_j)$). *For $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ and $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^p$ we have (per Definition in (5))*

$$\begin{aligned} \mathbb{E}_{\mathbf{w}}[\cos(\mathbf{x}_i^\top \mathbf{w}) \cos(\mathbf{w}^\top \mathbf{x}_j)] &= e^{-\frac{1}{2}(\|\mathbf{x}_i\|^2 + \|\mathbf{x}_j\|^2)} \cosh(\mathbf{x}_i^\top \mathbf{x}_j) \equiv [\mathbf{K}_{\cos}(\mathbf{X}, \mathbf{X})]_{ij} \equiv [\mathbf{K}_{\cos}]_{ij} \\ \mathbb{E}_{\mathbf{w}}[\sin(\mathbf{x}_i^\top \mathbf{w}) \sin(\mathbf{w}^\top \mathbf{x}_j)] &= e^{-\frac{1}{2}(\|\mathbf{x}_i\|^2 + \|\mathbf{x}_j\|^2)} \sinh(\mathbf{x}_i^\top \mathbf{x}_j) \equiv [\mathbf{K}_{\sin}(\mathbf{X}, \mathbf{X})]_{ij} \equiv [\mathbf{K}_{\sin}]_{ij} \\ \mathbb{E}_{\mathbf{w}}[\cos(\mathbf{x}_i^\top \mathbf{w}) \sin(\mathbf{w}^\top \mathbf{x}_j)] &= 0. \end{aligned}$$

460 *Proof of Lemma 1.* The proof follows the integration tricks in [50, 28]. Note in particular that
 461 the third equality holds in the case of (\cos, \sin) nonlinearity but in general not true for arbitrary
 462 (σ_1, σ_2) . \square

463 Let us focus on the resolvent $\mathbf{Q} \equiv \left(\frac{1}{n} \Sigma_{\mathbf{X}}^T \Sigma_{\mathbf{X}} + \lambda \mathbf{I}_n \right)^{-1}$ of $\frac{1}{n} \Sigma_{\mathbf{X}}^T \Sigma_{\mathbf{X}} \in \mathbb{R}^{n \times n}$, for random Fourier
 464 feature matrix $\Sigma_{\mathbf{X}} \equiv \begin{bmatrix} \cos(\mathbf{W}\mathbf{X}) \\ \sin(\mathbf{W}\mathbf{X}) \end{bmatrix}$ that can be rewritten as

$$\Sigma_{\mathbf{X}}^T = [\cos(\mathbf{X}^T \mathbf{w}_1), \dots, \cos(\mathbf{X}^T \mathbf{w}_N), \sin(\mathbf{X}^T \mathbf{w}_1), \dots, \sin(\mathbf{X}^T \mathbf{w}_N)] \quad (12)$$

465 for $\mathbf{w}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$, $i = 1, \dots, N$, that is at the core of our analysis. Note from (12) that we have

$$\Sigma_{\mathbf{X}}^T \Sigma_{\mathbf{X}} = \sum_{i=1}^N (\cos(\mathbf{X}^T \mathbf{w}_i) \cos(\mathbf{w}_i^T \mathbf{X}) + \sin(\mathbf{X}^T \mathbf{w}_i) \sin(\mathbf{w}_i^T \mathbf{X})) = \sum_{i=1}^N \mathbf{U}_i \mathbf{U}_i^T$$

466 with $\mathbf{U}_i = \begin{bmatrix} \cos(\mathbf{X}^T \mathbf{w}_i) & \sin(\mathbf{X}^T \mathbf{w}_i) \end{bmatrix} \in \mathbb{R}^{n \times 2}$.

467 Letting

$$\hat{\mathbf{Q}} \equiv \left(\frac{N}{n} \frac{\mathbf{K}_{\cos}}{1 + \alpha_{\cos}} + \frac{N}{n} \frac{\mathbf{K}_{\sin}}{1 + \alpha_{\sin}} + \lambda \mathbf{I}_n \right)^{-1} \quad (13)$$

468 with

$$\alpha_{\cos} = \frac{1}{n} \text{tr}(\mathbf{K}_{\cos} \mathbb{E}[\mathbf{Q}]), \quad \alpha_{\sin} = \frac{1}{n} \text{tr}(\mathbf{K}_{\sin} \mathbb{E}[\mathbf{Q}]) \quad (14)$$

469 we have, with the resolvent identity $(\mathbf{A}^{-1} - \mathbf{B}^{-1} = \mathbf{A}^{-1}(\mathbf{B} - \mathbf{A})\mathbf{B}^{-1})$ for invertible \mathbf{A}, \mathbf{B} that

$$\begin{aligned} \mathbb{E}[\mathbf{Q}] - \hat{\mathbf{Q}} &= \mathbb{E} \left[\mathbf{Q} \left(\frac{N}{n} \frac{\mathbf{K}_{\cos}}{1 + \alpha_{\cos}} + \frac{N}{n} \frac{\mathbf{K}_{\sin}}{1 + \alpha_{\sin}} - \frac{1}{n} \Sigma_{\mathbf{X}}^T \Sigma_{\mathbf{X}} \right) \right] \hat{\mathbf{Q}} \\ &= \mathbb{E}[\mathbf{Q}] \frac{N}{n} \left(\frac{\mathbf{K}_{\cos}}{1 + \alpha_{\cos}} + \frac{\mathbf{K}_{\sin}}{1 + \alpha_{\sin}} \right) \hat{\mathbf{Q}} - \frac{N}{n} \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{Q} \mathbf{U}_i \mathbf{U}_i^T] \hat{\mathbf{Q}} \\ &= \mathbb{E}[\mathbf{Q}] \frac{N}{n} \left(\frac{\mathbf{K}_{\cos}}{1 + \alpha_{\cos}} + \frac{\mathbf{K}_{\sin}}{1 + \alpha_{\sin}} \right) \hat{\mathbf{Q}} - \frac{N}{n} \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{Q}_{-i} \mathbf{U}_i (\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_i^T \mathbf{Q}_{-i} \mathbf{U}_i)^{-1} \mathbf{U}_i^T] \hat{\mathbf{Q}}, \end{aligned}$$

470 for $\mathbf{Q}_{-i} \equiv \left(\frac{1}{n} \Sigma_{\mathbf{X}}^T \Sigma_{\mathbf{X}} - \frac{1}{n} \mathbf{U}_i \mathbf{U}_i^T + \lambda \mathbf{I}_n \right)^{-1}$ that is **independent** of \mathbf{U}_i (and thus \mathbf{w}_i), where we
 471 applied the following Woodbury identity with $\mathbf{U} = \frac{1}{\sqrt{n}} \mathbf{U}_i = \begin{bmatrix} \cos(\mathbf{X}^T \mathbf{w}_i) & \sin(\mathbf{X}^T \mathbf{w}_i) \end{bmatrix}$ and
 472 $\mathbf{A} = \mathbf{Q}_{-i}^{-1}$.

473 **Lemma 2** (Woodbury). *For $\mathbf{A}, \mathbf{A} + \mathbf{U}\mathbf{U}^T \in \mathbb{R}^{p \times p}$ both invertible and $\mathbf{U} \in \mathbb{R}^{p \times n}$, we have*

$$(\mathbf{A} + \mathbf{U}\mathbf{U}^T)^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{U} (\mathbf{I}_n + \mathbf{U}^T \mathbf{A}^{-1} \mathbf{U})^{-1} \mathbf{U}^T \mathbf{A}^{-1}$$

474 so that in particular $(\mathbf{A} + \mathbf{U}\mathbf{U}^T)^{-1} \mathbf{U} = \mathbf{A}^{-1} \mathbf{U} (\mathbf{I}_n + \mathbf{U}^T \mathbf{A}^{-1} \mathbf{U})^{-1}$.

475 Consider now the two-by-two matrix

$$\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_i^T \mathbf{Q}_{-i} \mathbf{U}_i = \begin{bmatrix} 1 + \frac{1}{n} \cos(\mathbf{w}_i^T \mathbf{X}) \mathbf{Q}_{-i} \cos(\mathbf{X} \mathbf{w}_i) & \frac{1}{n} \cos(\mathbf{w}_i^T \mathbf{X}) \mathbf{Q}_{-i} \sin(\mathbf{X} \mathbf{w}_i) \\ \frac{1}{n} \sin(\mathbf{w}_i^T \mathbf{X}) \mathbf{Q}_{-i} \cos(\mathbf{X} \mathbf{w}_i) & 1 + \frac{1}{n} \sin(\mathbf{w}_i^T \mathbf{X}) \mathbf{Q}_{-i} \sin(\mathbf{X} \mathbf{w}_i) \end{bmatrix}$$

476 which, according to the following lemma, is expected to be close to $\begin{bmatrix} 1 + \alpha_{\cos} & 0 \\ 0 & 1 + \alpha_{\sin} \end{bmatrix}$ as defined
 477 in (14).

478 **Lemma 3** (Concentration of quadratic forms). *Under Assumption 1, for $\sigma_1(\cdot), \sigma_2(\cdot)$ two real Lipschitz
 479 functions, $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ and $\mathbf{A} \in \mathbb{R}^{n \times n}$ independent of \mathbf{w} with $\|\mathbf{A}\| \leq 1$, then*

$$\mathbb{P} \left(\left| \frac{1}{n} \sigma_1(\mathbf{w}^T \mathbf{X}) \mathbf{A} \sigma_2(\mathbf{X}^T \mathbf{w}) - \frac{1}{n} \text{tr}(\mathbf{A} \mathbb{E}_{\mathbf{w}} [\sigma_2(\mathbf{X}^T \mathbf{w}) \sigma_1(\mathbf{w}^T \mathbf{X})]) \right| > t \right) \leq C e^{-c n \min(t, t^2)}$$

480 for some universal constants $C, c > 0$.

481 *Proof of Lemma 3.* Lemma 3 is a trivial extension of Lemma 1 in [28], where one observes the proof
 482 actually holds when different types of nonlinear functions $\sigma_1(\cdot), \sigma_2(\cdot)$ (and in particular \cos and \sin)
 483 are considered. \square

484 As a consequence, we continue to write, with again the resolvent identity, that

$$\begin{aligned} & (\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_i^\top \mathbf{Q}_{-i} \mathbf{U}_i)^{-1} - \begin{bmatrix} 1 + \alpha_{\cos} & 0 \\ 0 & 1 + \alpha_{\sin} \end{bmatrix}^{-1} \\ &= \begin{bmatrix} 1 + \frac{1}{n} \cos(\mathbf{w}_i^\top \mathbf{X}) \mathbf{Q}_{-i} \cos(\mathbf{X} \mathbf{w}_i) & \frac{1}{n} \cos(\mathbf{w}_i^\top \mathbf{X}) \mathbf{Q}_{-i} \sin(\mathbf{X} \mathbf{w}_i) \\ \frac{1}{n} \sin(\mathbf{w}_i^\top \mathbf{X}) \mathbf{Q}_{-i} \cos(\mathbf{X} \mathbf{w}_i) & 1 + \frac{1}{n} \sin(\mathbf{w}_i^\top \mathbf{X}) \mathbf{Q}_{-i} \sin(\mathbf{X} \mathbf{w}_i) \end{bmatrix}^{-1} - \begin{bmatrix} 1 + \alpha_{\cos} & 0 \\ 0 & 1 + \alpha_{\sin} \end{bmatrix}^{-1} \\ &= (\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_i^\top \mathbf{Q}_{-i} \mathbf{U}_i)^{-1} \begin{bmatrix} \alpha_{\cos} - \frac{1}{n} \cos(\mathbf{w}_i^\top \mathbf{X}) \mathbf{Q}_{-i} \cos(\mathbf{X} \mathbf{w}_i) & -\frac{1}{n} \cos(\mathbf{w}_i^\top \mathbf{X}) \mathbf{Q}_{-i} \sin(\mathbf{X} \mathbf{w}_i) \\ -\frac{1}{n} \sin(\mathbf{w}_i^\top \mathbf{X}) \mathbf{Q}_{-i} \cos(\mathbf{X} \mathbf{w}_i) & \alpha_{\sin} - \frac{1}{n} \sin(\mathbf{w}_i^\top \mathbf{X}) \mathbf{Q}_{-i} \sin(\mathbf{X} \mathbf{w}_i) \end{bmatrix} \\ &\times \begin{bmatrix} \frac{1}{1 + \alpha_{\cos}} & 0 \\ 0 & \frac{1}{1 + \alpha_{\sin}} \end{bmatrix} \equiv (\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_i^\top \mathbf{Q}_{-i} \mathbf{U}_i)^{-1} D_i \begin{bmatrix} \frac{1}{1 + \alpha_{\cos}} & 0 \\ 0 & \frac{1}{1 + \alpha_{\sin}} \end{bmatrix}, \end{aligned}$$

485 where we note from Lemma 3 (with $\|\mathbf{Q}_{-i}\| \leq \lambda^{-1}$) that the matrix D_i should be of spectral norm
 486 $O(n^{-\frac{1}{2}})$ with high probability. So that

$$\begin{aligned} \mathbb{E}[\mathbf{Q}] - \hat{\mathbf{Q}} &= \mathbb{E}[\mathbf{Q}] \frac{N}{n} \left(\frac{\mathbf{K}_{\cos}}{1 + \alpha_{\cos}} + \frac{\mathbf{K}_{\sin}}{1 + \alpha_{\sin}} \right) \hat{\mathbf{Q}} - \frac{N}{n} \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{Q}_{-i} \mathbf{U}_i (\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_i^\top \mathbf{Q}_{-i} \mathbf{U}_i)^{-1} \mathbf{U}_i^\top] \hat{\mathbf{Q}} \\ &= \mathbb{E}[\mathbf{Q}] \frac{N}{n} \left(\frac{\mathbf{K}_{\cos}}{1 + \alpha_{\cos}} + \frac{\mathbf{K}_{\sin}}{1 + \alpha_{\sin}} \right) \hat{\mathbf{Q}} - \frac{N}{n} \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{Q}_{-i} \mathbf{U}_i \begin{bmatrix} \frac{1}{1 + \alpha_{\cos}} & 0 \\ 0 & \frac{1}{1 + \alpha_{\sin}} \end{bmatrix} \mathbf{U}_i^\top] \hat{\mathbf{Q}} \\ &\quad - \frac{N}{n} \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{Q}_{-i} \mathbf{U}_i (\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_i^\top \mathbf{Q}_{-i} \mathbf{U}_i)^{-1} D_i \begin{bmatrix} \frac{1}{1 + \alpha_{\cos}} & 0 \\ 0 & \frac{1}{1 + \alpha_{\sin}} \end{bmatrix} \mathbf{U}_i^\top] \hat{\mathbf{Q}} \\ &= (\mathbb{E}[\mathbf{Q}] - \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{Q}_{-i}]) \frac{N}{n} \left(\frac{\mathbf{K}_{\cos}}{1 + \alpha_{\cos}} + \frac{\mathbf{K}_{\sin}}{1 + \alpha_{\sin}} \right) \hat{\mathbf{Q}} - \frac{N}{n} \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{Q} \mathbf{U}_i D_i \begin{bmatrix} \frac{1}{1 + \alpha_{\cos}} & 0 \\ 0 & \frac{1}{1 + \alpha_{\sin}} \end{bmatrix} \mathbf{U}_i^\top] \hat{\mathbf{Q}}, \end{aligned}$$

487 where we used $\mathbb{E}_{\mathbf{w}_i}[\mathbf{U}_i \mathbf{U}_i^\top] = \mathbf{K}_{\cos} + \mathbf{K}_{\sin}$ by Lemma 1 and then Lemma 2 in reverse for the last
 488 equality. Moreover, since

$$\mathbb{E}[\mathbf{Q}] - \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{Q}_{-i}] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{Q} - \mathbf{Q}_{-i}] = -\frac{1}{n} \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{Q} \mathbf{U}_i (\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_i^\top \mathbf{Q}_{-i} \mathbf{U}_i)^{-1} \mathbf{U}_i^\top \mathbf{Q}]$$

so that with the fact $\frac{1}{\sqrt{n}} \|\mathbf{Q} \Sigma_X^\top\| \leq \sqrt{\mathbf{Q}_n^\top \Sigma_X^\top \Sigma_X \mathbf{Q}_n} \leq \lambda^{-\frac{1}{2}}$ we have for the first term

$$\|\mathbb{E}[\mathbf{Q}] - \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\mathbf{Q}_{-i}]\| = O(n^{-1}).$$

489 It thus remains to treat the second term, which, with the relation $\mathbf{A} \mathbf{B}^\top + \mathbf{B} \mathbf{A}^\top \preceq \mathbf{A} \mathbf{A}^\top + \mathbf{B} \mathbf{B}^\top$ (in
 490 the sense of symmetric matrices), and the same line of arguments as above, can be shown to have
 491 vanishing spectral norm (of order $O(n^{-\frac{1}{2}})$) as $n, p, N \rightarrow \infty$.

492 We thus have $\|\mathbb{E}[\mathbf{Q}] - \hat{\mathbf{Q}}\| = O(n^{-\frac{1}{2}})$, which concludes the first part of the proof of Theorem 1.

493 We shall show next that $\|\hat{\mathbf{Q}} - \bar{\mathbf{Q}}\| \rightarrow 0$ as $n, p, N \rightarrow \infty$. First note from previous derivation that
 494 $\alpha_\sigma - \frac{1}{n} \text{tr} \mathbf{K}_\sigma \hat{\mathbf{Q}} = O(n^{-\frac{1}{2}})$ for $\sigma = \cos, \sin$. To compare $\hat{\mathbf{Q}}$ and $\bar{\mathbf{Q}}$, it follows again from the
 495 resolvent identity that

$$\hat{\mathbf{Q}} - \bar{\mathbf{Q}} = \hat{\mathbf{Q}} \left(\frac{N}{n} \frac{\mathbf{K}_{\cos}(\alpha_{\cos} - \delta_{\cos})}{(1 + \delta_{\cos})(1 + \alpha_{\cos})} + \frac{N}{n} \frac{\mathbf{K}_{\sin}(\alpha_{\sin} - \delta_{\sin})}{(1 + \delta_{\sin})(1 + \alpha_{\sin})} \right) \bar{\mathbf{Q}}$$

496 so that the control of $\|\hat{\mathbf{Q}} - \bar{\mathbf{Q}}\|$ boils down to the control of $\max\{|\alpha_{\cos} - \delta_{\cos}|, |\alpha_{\sin} - \delta_{\sin}|\}$. To
 497 this end, it suffices to write

$$\alpha_{\cos} - \delta_{\cos} = \frac{1}{n} \text{tr} \mathbf{K}_{\cos} (\mathbb{E}[\mathbf{Q}] - \bar{\mathbf{Q}}) = \frac{1}{n} \text{tr} \mathbf{K}_{\cos} (\hat{\mathbf{Q}} - \bar{\mathbf{Q}}) + O(n^{-\frac{1}{2}})$$

where we used $|\text{tr}(\mathbf{AB})| \leq \|\mathbf{A}\| \text{tr}(\mathbf{B})$ for nonnegative definite \mathbf{B} , together with the fact that $\frac{1}{n} \text{tr} \mathbf{K}_\sigma$ is (uniformly) bounded under Assumption 1, for $\sigma = \cos, \sin$.
As a consequence, we have

$$|\alpha_{\cos} - \delta_{\cos}| \leq |\alpha_{\cos} - \delta_{\cos}| \frac{N \frac{1}{n} \text{tr}(\mathbf{K}_{\cos} \hat{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}})}{n (1 + \delta_{\cos})(1 + \alpha_{\cos})} + O(n^{-\frac{1}{2}}).$$

It thus remains to show

$$\frac{N \frac{1}{n} \text{tr}(\mathbf{K}_{\cos} \hat{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}})}{n (1 + \delta_{\cos})(1 + \alpha_{\cos})} < 1$$

or alternatively, by the Cauchy–Schwarz inequality, to show

$$\frac{N \frac{1}{n} \text{tr}(\mathbf{K}_{\cos} \hat{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}})}{n (1 + \delta_{\cos})(1 + \alpha_{\cos})} \leq \sqrt{\frac{N \frac{1}{n} \text{tr}(\mathbf{K}_{\cos} \bar{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}})}{n (1 + \delta_{\cos})^2} \cdot \frac{N \frac{1}{n} \text{tr}(\mathbf{K}_{\cos} \hat{\mathbf{Q}} \mathbf{K}_{\cos} \hat{\mathbf{Q}})}{n (1 + \alpha_{\cos})^2}} < 1.$$

To treat the first right-hand side term (the second can be done similarly), it unfolds from $|\text{tr}(\mathbf{AB})| \leq \|\mathbf{A}\| \text{tr}(\mathbf{B})$ for nonnegative definite \mathbf{B} that

$$\frac{N \frac{1}{n} \text{tr}(\mathbf{K}_{\cos} \bar{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}})}{n (1 + \delta_{\cos})^2} \leq \left\| \frac{N \mathbf{K}_{\cos} \bar{\mathbf{Q}}}{n (1 + \delta_{\cos})} \right\| \frac{\frac{1}{n} \text{tr}(\mathbf{K}_{\cos} \bar{\mathbf{Q}})}{1 + \delta_{\cos}} = \left\| \frac{N \mathbf{K}_{\cos} \bar{\mathbf{Q}}}{n (1 + \delta_{\cos})} \right\| \frac{\gamma_{\cos}}{1 + \delta_{\cos}} \leq \frac{\gamma_{\cos}}{1 + \delta_{\cos}} < 1$$

where we used the fact that $\frac{N \mathbf{K}_{\cos} \bar{\mathbf{Q}}}{n (1 + \delta_{\cos})} = \mathbf{I}_n - \frac{N \mathbf{K}_{\sin} \bar{\mathbf{Q}}}{n (1 + \delta_{\sin})} - \lambda \bar{\mathbf{Q}}$. This concludes the proof of Theorem 1.
■

C Proof of Theorem 2

To prove Theorem 2, it indeed suffices to prove the following lemma.

Lemma 4 (Asymptotic behavior of $\mathbb{E}[\mathbf{QAQ}]$). *Under Assumption 1, for \mathbf{Q} defined in (4) and symmetric nonnegative definite $\mathbf{A} \in \mathbb{R}^{n \times n}$ of bounded spectral norm, we have*

$$\left\| \mathbb{E}[\mathbf{QAQ}] - \left(\bar{\mathbf{Q}} \mathbf{A} \bar{\mathbf{Q}} + \frac{N}{n} \begin{bmatrix} \frac{1}{n} \text{tr}(\mathbf{QAQK}_{\cos}) & \frac{1}{n} \text{tr}(\mathbf{QAQK}_{\sin}) \\ \frac{1}{n} \text{tr}(\mathbf{QAQK}_{\cos}) & \frac{1}{n} \text{tr}(\mathbf{QAQK}_{\sin}) \end{bmatrix} \Omega \begin{bmatrix} \bar{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}} \\ \bar{\mathbf{Q}} \mathbf{K}_{\sin} \bar{\mathbf{Q}} \end{bmatrix} \right) \right\| \rightarrow 0$$

almost surely as $n \rightarrow \infty$, with $\Omega^{-1} \equiv \mathbf{I}_2 - \frac{N}{n} \begin{bmatrix} \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}} \mathbf{K}_{\cos}) & \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}} \mathbf{K}_{\sin}) \\ \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \mathbf{K}_{\sin} \bar{\mathbf{Q}} \mathbf{K}_{\cos}) & \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \mathbf{K}_{\sin} \bar{\mathbf{Q}} \mathbf{K}_{\sin}) \end{bmatrix}$. In particular,

we have

$$\left\| \mathbb{E} \begin{bmatrix} \mathbf{QK}_{\cos} \mathbf{Q} \\ \mathbf{QK}_{\sin} \mathbf{Q} \end{bmatrix} - \Omega \begin{bmatrix} \bar{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}} \\ \bar{\mathbf{Q}} \mathbf{K}_{\sin} \bar{\mathbf{Q}} \end{bmatrix} \right\| \rightarrow 0.$$

Proof of Lemma 4. The proof of Lemma 4 essentially follows the same line of arguments as that of Theorem 1. Writing

$$\begin{aligned} \mathbb{E}[\mathbf{QAQ}] &= \mathbb{E}[\bar{\mathbf{Q}} \mathbf{A} \bar{\mathbf{Q}}] + \mathbb{E}[(\mathbf{Q} - \bar{\mathbf{Q}}) \mathbf{A} \mathbf{Q}] \\ &\simeq \bar{\mathbf{Q}} \mathbf{A} \bar{\mathbf{Q}} + \mathbb{E} \left[\mathbf{Q} \left(\frac{N}{n} \frac{\mathbf{K}_{\cos}}{1 + \delta_{\cos}} + \frac{N}{n} \frac{\mathbf{K}_{\sin}}{1 + \delta_{\sin}} - \frac{1}{n} \Sigma_{\mathbf{X}}^{\top} \Sigma_{\mathbf{X}} \right) \bar{\mathbf{Q}} \mathbf{A} \mathbf{Q} \right] \\ &= \bar{\mathbf{Q}} \mathbf{A} \bar{\mathbf{Q}} + \frac{N}{n} \mathbb{E}[\mathbf{Q} \Phi \bar{\mathbf{Q}} \mathbf{A} \mathbf{Q}] - \frac{1}{n} \sum_{i=1}^N \mathbb{E}[\mathbf{Q} \mathbf{U}_i \mathbf{U}_i^{\top} \bar{\mathbf{Q}} \mathbf{A} \mathbf{Q}] \end{aligned}$$

515 where we note \simeq by ignoring matrices with spectral norm of order $O(n^{-\frac{1}{2}})$ and recall the shortcut
 516 $\Phi \equiv \frac{\mathbf{K}_{\cos}}{1+\delta_{\cos}} + \frac{\mathbf{K}_{\sin}}{1+\delta_{\sin}}$. Developing rightmost term with Lemma 2 as

$$\begin{aligned} \mathbb{E}[\mathbf{Q}\mathbf{U}_i\mathbf{U}_i^\top\bar{\mathbf{Q}}\mathbf{A}\mathbf{Q}] &= \mathbb{E}\left[\mathbf{Q}_{-i}\mathbf{U}_i(\mathbf{I}_2 + \frac{1}{n}\mathbf{U}_i^\top\mathbf{Q}_{-i}\mathbf{U}_i)^{-1}\mathbf{U}_i^\top\bar{\mathbf{Q}}\mathbf{A}\mathbf{Q}\right] \\ &= \mathbb{E}\left[\mathbf{Q}_{-i}\mathbf{U}_i(\mathbf{I}_2 + \frac{1}{n}\mathbf{U}_i^\top\mathbf{Q}_{-i}\mathbf{U}_i)^{-1}\mathbf{U}_i^\top\bar{\mathbf{Q}}\mathbf{A}\mathbf{Q}_{-i}\right] \\ &\quad - \frac{1}{n}\mathbb{E}\left[\mathbf{Q}_{-i}\mathbf{U}_i(\mathbf{I}_2 + \frac{1}{n}\mathbf{U}_i^\top\mathbf{Q}_{-i}\mathbf{U}_i)^{-1}\mathbf{U}_i^\top\bar{\mathbf{Q}}\mathbf{A}\mathbf{Q}_{-i}\mathbf{U}_i(\mathbf{I}_2 + \frac{1}{n}\mathbf{U}_i^\top\mathbf{Q}_{-i}\mathbf{U}_i)^{-1}\mathbf{U}_i^\top\mathbf{Q}_{-i}\right] \\ &\simeq \mathbb{E}[\mathbf{Q}_{-i}\Phi\bar{\mathbf{Q}}\mathbf{A}\mathbf{Q}_{-i}] \\ &= \mathbb{E}\left[\mathbf{Q}_{-i}\mathbf{U}_i\begin{bmatrix} \frac{1}{1+\delta_{\cos}} & 0 \\ 0 & \frac{1}{1+\delta_{\sin}} \end{bmatrix}\begin{bmatrix} \frac{1}{n}\text{tr}(\bar{\mathbf{Q}}\mathbf{A}\bar{\mathbf{Q}}\mathbf{K}_{\cos}) & 0 \\ 0 & \frac{1}{n}\text{tr}(\bar{\mathbf{Q}}\mathbf{A}\bar{\mathbf{Q}}\mathbf{K}_{\sin}) \end{bmatrix}\begin{bmatrix} \frac{1}{1+\delta_{\cos}} & 0 \\ 0 & \frac{1}{1+\delta_{\sin}} \end{bmatrix}\mathbf{U}_i^\top\mathbf{Q}_{-i}\right] \end{aligned}$$

517 so that

$$\begin{aligned} \mathbb{E}[\mathbf{Q}\mathbf{A}\mathbf{Q}] &\simeq \bar{\mathbf{Q}}\mathbf{A}\bar{\mathbf{Q}} + \frac{N}{n}\mathbb{E}\left[\mathbf{Q}\left(\frac{\frac{1}{n}\text{tr}(\bar{\mathbf{Q}}\mathbf{A}\bar{\mathbf{Q}}\mathbf{K}_{\cos})}{(1+\delta_{\cos})^2}\mathbf{K}_{\cos} + \frac{\frac{1}{n}\text{tr}(\bar{\mathbf{Q}}\mathbf{A}\bar{\mathbf{Q}}\mathbf{K}_{\sin})}{(1+\delta_{\sin})^2}\mathbf{K}_{\sin}\right)\mathbf{Q}\right] \\ &= \bar{\mathbf{Q}}\mathbf{A}\bar{\mathbf{Q}} + \frac{N}{n}\begin{bmatrix} \frac{\frac{1}{n}\text{tr}(\bar{\mathbf{Q}}\mathbf{A}\bar{\mathbf{Q}}\mathbf{K}_{\cos})}{(1+\delta_{\cos})^2} & \frac{\frac{1}{n}\text{tr}(\bar{\mathbf{Q}}\mathbf{A}\bar{\mathbf{Q}}\mathbf{K}_{\sin})}{(1+\delta_{\sin})^2} \end{bmatrix}\mathbb{E}\begin{bmatrix} \mathbf{Q}\mathbf{K}_{\cos}\mathbf{Q} \\ \mathbf{Q}\mathbf{K}_{\sin}\mathbf{Q} \end{bmatrix} \end{aligned} \quad (15)$$

518 by taking $\mathbf{A} = \mathbf{K}_{\cos}$ or \mathbf{K}_{\sin} , we result in

$$\begin{aligned} \mathbb{E}[\mathbf{Q}\mathbf{K}_{\cos}\mathbf{Q}] &\simeq \frac{c}{ac-bd}\bar{\mathbf{Q}}\mathbf{K}_{\cos}\bar{\mathbf{Q}} + \frac{b}{ac-bd}\bar{\mathbf{Q}}\mathbf{K}_{\sin}\bar{\mathbf{Q}} \\ \mathbb{E}[\mathbf{Q}\mathbf{K}_{\sin}\mathbf{Q}] &\simeq \frac{a}{ac-bd}\bar{\mathbf{Q}}\mathbf{K}_{\sin}\bar{\mathbf{Q}} + \frac{d}{ac-bd}\bar{\mathbf{Q}}\mathbf{K}_{\cos}\bar{\mathbf{Q}} \end{aligned}$$

519 with $a = 1 - \frac{N}{n}\frac{\frac{1}{n}\text{tr}(\bar{\mathbf{Q}}\mathbf{K}_{\cos}\bar{\mathbf{Q}}\mathbf{K}_{\cos})}{(1+\delta_{\cos})^2}$, $b = \frac{N}{n}\frac{\frac{1}{n}\text{tr}(\bar{\mathbf{Q}}\mathbf{K}_{\cos}\bar{\mathbf{Q}}\mathbf{K}_{\sin})}{(1+\delta_{\sin})^2}$, $c = 1 - \frac{N}{n}\frac{\frac{1}{n}\text{tr}(\bar{\mathbf{Q}}\mathbf{K}_{\sin}\bar{\mathbf{Q}}\mathbf{K}_{\sin})}{(1+\delta_{\sin})^2}$ and
 520 $d = \frac{N}{n}\frac{\frac{1}{n}\text{tr}(\bar{\mathbf{Q}}\mathbf{K}_{\sin}\bar{\mathbf{Q}}\mathbf{K}_{\cos})}{(1+\delta_{\cos})^2}$ such that $(1+\delta_{\sin})^2b = (1+\delta_{\cos})^2d$.

$$\mathbb{E}\begin{bmatrix} \mathbf{Q}\mathbf{K}_{\cos}\mathbf{Q} \\ \mathbf{Q}\mathbf{K}_{\sin}\mathbf{Q} \end{bmatrix} \simeq \begin{bmatrix} a & -b \\ -d & c \end{bmatrix}^{-1} \begin{bmatrix} \bar{\mathbf{Q}}\mathbf{K}_{\cos}\bar{\mathbf{Q}} \\ \bar{\mathbf{Q}}\mathbf{K}_{\sin}\bar{\mathbf{Q}} \end{bmatrix} \equiv \Omega \begin{bmatrix} \bar{\mathbf{Q}}\mathbf{K}_{\cos}\bar{\mathbf{Q}} \\ \bar{\mathbf{Q}}\mathbf{K}_{\sin}\bar{\mathbf{Q}} \end{bmatrix}$$

521 for $\Omega \equiv \begin{bmatrix} a & -b \\ -d & c \end{bmatrix}^{-1}$. Plugging back into (15) we conclude the proof of Lemma 4. \square

522 Theorem 2 can be achieved by considering the concentration of (the bilinear form) $\frac{1}{n}\mathbf{y}^\top\mathbf{Q}^2\mathbf{y}$ around
 523 its expectation $\frac{1}{n}\mathbf{y}^\top\mathbb{E}[\mathbf{Q}^2]\mathbf{y}$ (with for instance Lemma 3 in [28]), together with Lemma 4. This
 524 concludes the proof of Theorem 2. \blacksquare

525 D Proof of Theorem 3

526 Recall the definition of $E_{\text{test}} = \frac{1}{\hat{n}}\|\hat{\mathbf{y}} - \Sigma_{\hat{\mathbf{X}}}^\top\beta\|^2$ from (3) with $\Sigma_{\hat{\mathbf{X}}} = \begin{bmatrix} \cos(\mathbf{W}\hat{\mathbf{X}}) \\ \sin(\mathbf{W}\hat{\mathbf{X}}) \end{bmatrix} \in \mathbb{R}^{2N \times \hat{n}}$ on a
 527 test set $(\hat{\mathbf{X}}, \hat{\mathbf{y}})$ of size \hat{n} , and first focus on the case $2N > n$ where $\beta = \frac{1}{n}\Sigma_{\mathbf{X}}\mathbf{Q}\mathbf{y}$ as per (2). By
 528 (12), we have

$$E_{\text{test}} = \frac{1}{\hat{n}}\left\|\hat{\mathbf{y}} - \frac{1}{n}\Sigma_{\hat{\mathbf{X}}}^\top\Sigma_{\mathbf{X}}\mathbf{Q}\mathbf{y}\right\|^2 = \frac{1}{\hat{n}}\left\|\hat{\mathbf{y}} - \frac{1}{n}\sum_{i=1}^N\hat{\mathbf{U}}_i\mathbf{U}_i^\top\mathbf{Q}\mathbf{y}\right\|^2$$

529 where, similar to the notation $\mathbf{U}_i = [\cos(\mathbf{X}^\top\mathbf{w}_i) \quad \sin(\mathbf{X}^\top\mathbf{w}_i)] \in \mathbb{R}^{n \times 2}$ as in the proof of Theo-
 530 rem 1, we denote

$$\hat{\mathbf{U}}_i \equiv [\cos(\hat{\mathbf{X}}^\top\mathbf{w}_i) \quad \sin(\hat{\mathbf{X}}^\top\mathbf{w}_i)] \in \mathbb{R}^{\hat{n} \times 2}.$$

531 As a consequence, we further get

$$\begin{aligned}
\mathbb{E}[E_{\text{test}}] &= \frac{1}{\hat{n}} \|\hat{\mathbf{y}}\|^2 - \frac{2}{n\hat{n}} \sum_{i=1}^N \hat{\mathbf{y}}^\top \mathbb{E}[\hat{\mathbf{U}}_i \mathbf{U}_i^\top \mathbf{Q}] \mathbf{y} + \frac{1}{n^2 \hat{n}} \sum_{i,j=1}^N \mathbf{y}^\top \mathbb{E}[\mathbf{Q} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \hat{\mathbf{U}}_j \mathbf{U}_j^\top \mathbf{Q}] \mathbf{y} \\
&= \frac{1}{\hat{n}} \|\hat{\mathbf{y}}\|^2 - \frac{2}{n\hat{n}} \sum_{i=1}^N \hat{\mathbf{y}}^\top \mathbb{E} \left[\hat{\mathbf{U}}_i (\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_i^\top \mathbf{Q}_{-i} \mathbf{U}_i)^{-1} \mathbf{U}_i^\top \mathbf{Q}_{-i} \right] \mathbf{y} + \frac{1}{n^2 \hat{n}} \sum_{i,j=1}^N \mathbf{y}^\top \mathbb{E}[\mathbf{Q} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \hat{\mathbf{U}}_j \mathbf{U}_j^\top \mathbf{Q}] \mathbf{y} \\
&\simeq \frac{1}{\hat{n}} \|\hat{\mathbf{y}}\|^2 - \frac{2}{n\hat{n}} \sum_{i=1}^N \hat{\mathbf{y}}^\top \mathbb{E} \left[\hat{\mathbf{U}}_i \begin{bmatrix} \frac{1}{1+\delta_{\cos}} & 0 \\ 0 & \frac{1}{1+\delta_{\sin}} \end{bmatrix} \mathbf{U}_i^\top \mathbf{Q}_{-i} \right] \mathbf{y} + \frac{1}{n^2 \hat{n}} \sum_{i,j=1}^N \mathbf{y}^\top \mathbb{E}[\mathbf{Q} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \hat{\mathbf{U}}_j \mathbf{U}_j^\top \mathbf{Q}] \mathbf{y} \\
&\simeq \frac{1}{\hat{n}} \|\hat{\mathbf{y}}\|^2 - \frac{2}{\hat{n}} \hat{\mathbf{y}}^\top \left(\frac{N}{n} \frac{\mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})}{1 + \delta_{\cos}} + \frac{N}{n} \frac{\mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})}{1 + \delta_{\sin}} \right) \bar{\mathbf{Q}} \mathbf{y} + \frac{1}{n^2 \hat{n}} \sum_{i,j=1}^N \mathbf{y}^\top \mathbb{E}[\mathbf{Q} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \hat{\mathbf{U}}_j \mathbf{U}_j^\top \mathbf{Q}] \mathbf{y}
\end{aligned}$$

532 where we similarly denote

$$\begin{aligned}
\mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X}) &\equiv \left\{ e^{-\frac{1}{2}(\|\hat{\mathbf{x}}_i\|^2 + \|\mathbf{x}_j\|^2)} \cosh(\hat{\mathbf{x}}_i^\top \mathbf{x}_j) \right\}_{i,j=1}^{\hat{n},n} \\
\mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X}) &\equiv \left\{ e^{-\frac{1}{2}(\|\hat{\mathbf{x}}_i\|^2 + \|\mathbf{x}_j\|^2)} \sinh(\hat{\mathbf{x}}_i^\top \mathbf{x}_j) \right\}_{i,j=1}^{\hat{n},n} \in \mathbb{R}^{\hat{n} \times n}.
\end{aligned}$$

533 Note that, different from the proof of Theorem 1 and 2 where we constantly use the fact that
534 $\|\mathbf{Q}\| \leq \lambda^{-1}$ and

$$\frac{1}{n} \Sigma_{\mathbf{X}}^\top \Sigma_{\mathbf{X}} \mathbf{Q} = \mathbf{I}_n - \lambda \mathbf{Q}$$

535 so that $\|\frac{1}{n} \Sigma_{\mathbf{X}}^\top \Sigma_{\mathbf{X}} \mathbf{Q}\| \leq 1$, we do not have in general a simple control for $\|\frac{1}{n} \Sigma_{\hat{\mathbf{X}}}^\top \Sigma_{\mathbf{X}} \mathbf{Q}\|$, when
536 arbitrary $\hat{\mathbf{X}}$ is considered. Intuitively speaking, this is due to the loss-of-control for $\|\frac{1}{n} (\Sigma_{\hat{\mathbf{X}}} -$
537 $\Sigma_{\mathbf{X}})^\top \Sigma_{\mathbf{X}} \mathbf{Q}\|$ when $\hat{\mathbf{X}}$ can be chosen arbitrarily with respect to \mathbf{X} . It was remarked in [28] that in
538 general only a $O(\sqrt{n})$ upper bound can be derived for $\|\frac{1}{\sqrt{n}} \Sigma_{\mathbf{X}}\|$ or $\|\frac{1}{\sqrt{n}} \Sigma_{\hat{\mathbf{X}}}\|$. Nonetheless, this
539 problem can be resolved with the additional Assumption 2 by mimicking the same construction in
540 Section 1.2.2 of [27].

541 It thus remains to handle the last term (noted \mathbf{Z}) as follows

$$\begin{aligned}
\mathbf{Z} &\equiv \frac{1}{n^2 \hat{n}} \sum_{i,j=1}^N \mathbf{y}^\top \mathbb{E}[\mathbf{Q} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \hat{\mathbf{U}}_j \mathbf{U}_j^\top \mathbf{Q}] \mathbf{y} \\
&= \frac{1}{n^2 \hat{n}} \sum_{i=1}^N \mathbf{y}^\top \mathbb{E}[\mathbf{Q} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \hat{\mathbf{U}}_i \mathbf{U}_i^\top \mathbf{Q}] \mathbf{y} + \frac{1}{n^2 \hat{n}} \sum_{i=1}^N \sum_{j \neq i} \mathbf{y}^\top \mathbb{E}[\mathbf{Q} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \hat{\mathbf{U}}_j \mathbf{U}_j^\top \mathbf{Q}] \mathbf{y} = \mathbf{Z}_1 + \mathbf{Z}_2
\end{aligned}$$

542 where \mathbf{Z}_1 term can be treated as

$$\begin{aligned}
\mathbf{Z}_1 &\equiv \frac{1}{n^2 \hat{n}} \sum_{i=1}^N \mathbf{y}^\top \mathbb{E}[\mathbf{Q} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \hat{\mathbf{U}}_i \mathbf{U}_i^\top \mathbf{Q}] \mathbf{y} \\
&= \frac{1}{n\hat{n}} \sum_{i=1}^N \mathbf{y}^\top \mathbb{E}[\mathbf{Q}_{-i} \mathbf{U}_i (\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_i^\top \mathbf{Q}_{-i} \mathbf{U}_i)^{-1} \frac{1}{n} \hat{\mathbf{U}}_i^\top \hat{\mathbf{U}}_i (\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_i^\top \mathbf{Q}_{-i} \mathbf{U}_i)^{-1} \mathbf{U}_i^\top \mathbf{Q}_{-i}] \mathbf{y} \\
&\simeq \frac{1}{n\hat{n}} \sum_{i=1}^N \mathbf{y}^\top \mathbb{E}[\mathbf{Q}_{-i} \mathbf{U}_i \begin{bmatrix} \frac{1}{1+\delta_{\cos}} & 0 \\ 0 & \frac{1}{1+\delta_{\sin}} \end{bmatrix} \begin{bmatrix} \frac{1}{n} \text{tr} \hat{\mathbf{K}}_{\cos} & 0 \\ 0 & \frac{1}{n} \text{tr} \hat{\mathbf{K}}_{\sin} \end{bmatrix} \begin{bmatrix} \frac{1}{1+\delta_{\cos}} & 0 \\ 0 & \frac{1}{1+\delta_{\sin}} \end{bmatrix} \mathbf{U}_i^\top \mathbf{Q}_{-i}] \mathbf{y} \\
&\simeq \frac{N}{n} \frac{1}{\hat{n}} \mathbf{y}^\top \mathbb{E} \left[\mathbf{Q} \left(\frac{\frac{1}{n} \text{tr} \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \hat{\mathbf{X}})}{(1 + \delta_{\cos})^2} \mathbf{K}_{\cos} + \frac{\frac{1}{n} \text{tr} \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \hat{\mathbf{X}})}{(1 + \delta_{\sin})^2} \mathbf{K}_{\sin} \right) \mathbf{Q} \right] \mathbf{y} \\
&\simeq \frac{N}{n} \frac{1}{\hat{n}} \begin{bmatrix} \frac{1}{n} \text{tr} \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \hat{\mathbf{X}}) & \frac{1}{n} \text{tr} \frac{1}{n} \text{tr} \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \hat{\mathbf{X}}) \end{bmatrix} \Omega \begin{bmatrix} \mathbf{y}^\top \bar{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}} \mathbf{y} \\ \mathbf{y}^\top \bar{\mathbf{Q}} \mathbf{K}_{\sin} \bar{\mathbf{Q}} \mathbf{y} \end{bmatrix}
\end{aligned}$$

543 where we apply Lemma 4 and recall

$$\mathbf{K}_{\cos}(\hat{\mathbf{X}}, \hat{\mathbf{X}}) \equiv \left\{ e^{-\frac{1}{2}(\|\hat{\mathbf{x}}_i\|^2 + \|\hat{\mathbf{x}}_j\|^2)} \cosh(\hat{\mathbf{x}}_i^\top \hat{\mathbf{x}}_j) \right\}_{i,j=1}^{\hat{n}}, \quad \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \hat{\mathbf{X}}) \equiv \left\{ e^{-\frac{1}{2}(\|\hat{\mathbf{x}}_i\|^2 + \|\hat{\mathbf{x}}_j\|^2)} \sinh(\hat{\mathbf{x}}_i^\top \hat{\mathbf{x}}_j) \right\}_{i,j=1}^{\hat{n}}$$

544 Moving on to \mathbf{Z}_2 and we write

$$\begin{aligned} \mathbf{Z}_2 &\equiv \frac{1}{n^2 \hat{n}} \mathbb{E} \sum_{i=1}^N \sum_{j \neq i} \mathbf{y}^\top \mathbf{Q} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \hat{\mathbf{U}}_j \mathbf{U}_j^\top \mathbf{Q} \mathbf{y} \\ &= \frac{1}{n^2 \hat{n}} \mathbb{E} \sum_{i=1}^N \sum_{j \neq i} \mathbf{y}^\top \mathbf{Q}_{-j} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \hat{\mathbf{U}}_j (\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_j^\top \mathbf{Q}_{-j} \mathbf{U}_j)^{-1} \mathbf{U}_j^\top \mathbf{Q}_{-j} \mathbf{y} \\ &\quad - \frac{1}{n^2 \hat{n}} \mathbb{E} \sum_{i=1}^N \sum_{j \neq i} \mathbf{y}^\top \mathbf{Q}_{-j} \mathbf{U}_j (\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_j^\top \mathbf{Q}_{-j} \mathbf{U}_j)^{-1} \mathbf{U}_j^\top \mathbf{Q}_{-j} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \hat{\mathbf{U}}_j (\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_j^\top \mathbf{Q}_{-j} \mathbf{U}_j)^{-1} \mathbf{U}_j^\top \mathbf{Q}_{-j} \mathbf{y} \\ &\simeq \frac{1}{n \hat{n}} \mathbb{E} \sum_{i=1}^N \sum_{j \neq i} \mathbf{y}^\top \mathbf{Q}_{-j} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \left(\frac{\mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})}{1 + \delta_{\cos}} + \frac{\mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})}{1 + \delta_{\sin}} \right) \mathbf{Q}_{-j} \mathbf{y} \\ &\quad - \frac{1}{n^2 \hat{n}} \mathbb{E} \sum_{i=1}^N \sum_{j \neq i} \mathbf{y}^\top \mathbf{Q}_{-j} \mathbf{U}_j \begin{bmatrix} \frac{1}{1 + \delta_{\cos}} & 0 \\ 0 & \frac{1}{1 + \delta_{\sin}} \end{bmatrix} \begin{bmatrix} \frac{1}{n} \text{tr}(\mathbf{Q}_{-j} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})) & 0 \\ 0 & \frac{1}{n} \text{tr}(\mathbf{Q}_{-j} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})) \end{bmatrix} \\ &\quad \begin{bmatrix} \frac{1}{1 + \delta_{\cos}} & 0 \\ 0 & \frac{1}{1 + \delta_{\sin}} \end{bmatrix} \mathbf{U}_j^\top \mathbf{Q}_{-j} \mathbf{y} \equiv \mathbf{Z}_{21} - \mathbf{Z}_{22}. \end{aligned}$$

545 For the term \mathbf{Z}_{21} , note that $\mathbf{Q}_{-j} \simeq \mathbf{Q}$ and **depends** on \mathbf{U}_i (and $\hat{\mathbf{U}}_i$), such that

$$\begin{aligned} \mathbf{Z}_{21} &\equiv \frac{1}{n^2 \hat{n}} \mathbb{E} \sum_{i=1}^N \sum_{j \neq i} \mathbf{y}^\top \mathbf{Q}_{-j} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \left(\frac{\mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})}{1 + \delta_{\cos}} + \frac{\mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})}{1 + \delta_{\sin}} \right) \mathbf{Q}_{-j} \mathbf{y} \\ &\simeq \frac{N}{n} \frac{1}{n \hat{n}} \mathbb{E} \sum_{i=1}^N \mathbf{y}^\top \mathbf{Q} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \left(\frac{\mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})}{1 + \delta_{\cos}} + \frac{\mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})}{1 + \delta_{\sin}} \right) \mathbf{Q} \mathbf{y} \\ &= \frac{N}{n} \frac{1}{n \hat{n}} \mathbb{E} \sum_{i=1}^N \mathbf{y}^\top \mathbf{Q}_{-i} \mathbf{U}_i (\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_i^\top \mathbf{Q}_{-i} \mathbf{U}_i)^{-1} \hat{\mathbf{U}}_i^\top \hat{\mathbf{\Phi}} \mathbf{Q}_{-i} \mathbf{y} \\ &\quad - \frac{N}{n} \frac{1}{n \hat{n}} \mathbb{E} \sum_{i=1}^N \mathbf{y}^\top \mathbf{Q}_{-i} \mathbf{U}_i (\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_i^\top \mathbf{Q}_{-i} \mathbf{U}_i)^{-1} \hat{\mathbf{U}}_i^\top \hat{\mathbf{\Phi}} \mathbf{Q}_{-i} \mathbf{U}_i (\mathbf{I}_2 + \frac{1}{n} \mathbf{U}_i^\top \mathbf{Q}_{-i} \mathbf{U}_i)^{-1} \mathbf{U}_i^\top \mathbf{Q}_{-i} \mathbf{y} \\ &\simeq \frac{N}{n} \frac{1}{n \hat{n}} \mathbb{E} \sum_{i=1}^N \mathbf{y}^\top \mathbf{Q}_{-i} \left(\frac{\mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})}{1 + \delta_{\cos}} + \frac{\mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})}{1 + \delta_{\sin}} \right)^\top \hat{\mathbf{\Phi}} \mathbf{Q}_{-i} \mathbf{y} \\ &\quad - \frac{N}{n} \frac{1}{\hat{n}} \mathbb{E} \sum_{i=1}^N \mathbf{y}^\top \mathbf{Q}_{-i} \mathbf{U}_i \begin{bmatrix} \frac{1}{1 + \delta_{\cos}} & 0 \\ 0 & \frac{1}{1 + \delta_{\sin}} \end{bmatrix} \frac{1}{n} \hat{\mathbf{U}}_i^\top \hat{\mathbf{\Phi}} \mathbf{Q}_{-i} \mathbf{U}_i \begin{bmatrix} \frac{1}{1 + \delta_{\cos}} & 0 \\ 0 & \frac{1}{1 + \delta_{\sin}} \end{bmatrix} \mathbf{U}_i^\top \mathbf{Q}_{-i} \mathbf{y} \end{aligned}$$

546 where we recall the shortcut $\Phi \equiv \frac{\mathbf{K}_{\cos}}{1+\delta_{\cos}} + \frac{\mathbf{K}_{\sin}}{1+\delta_{\sin}}$ and similarly $\hat{\Phi} \equiv \frac{\mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})}{1+\delta_{\cos}} + \frac{\mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})}{1+\delta_{\sin}} \in \mathbb{R}^{\hat{n} \times n}$.
 547 As a consequence, we further have, with Lemma 4 that

$$\begin{aligned}
 \mathbf{Z}_{21} &\simeq \left(\frac{N}{n}\right)^2 \frac{1}{\hat{n}} \mathbf{y}^\top \mathbb{E} \left[\mathbf{Q} \hat{\Phi}^\top \hat{\Phi} \mathbf{Q} \right] \mathbf{y} \\
 &\quad - \frac{N}{n} \frac{1}{\hat{n}} \mathbb{E} \sum_{i=1}^N \mathbf{y}^\top \mathbf{Q}_{-i} \mathbf{U}_i \begin{bmatrix} \frac{1}{1+\delta_{\cos}} & 0 \\ 0 & \frac{1}{1+\delta_{\sin}} \end{bmatrix} \begin{bmatrix} \frac{1}{n} \text{tr}(\hat{\Phi} \bar{\mathbf{Q}} \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})^\top) & 0 \\ 0 & \frac{1}{n} \text{tr}(\hat{\Phi} \bar{\mathbf{Q}} \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})^\top) \end{bmatrix} \\
 &\quad \times \begin{bmatrix} \frac{1}{1+\delta_{\cos}} & 0 \\ 0 & \frac{1}{1+\delta_{\sin}} \end{bmatrix} \mathbf{U}_i^\top \mathbf{Q}_{-i} \mathbf{y} \\
 &\simeq \left(\frac{N}{n}\right)^2 \frac{1}{\hat{n}} \mathbf{y}^\top \mathbb{E} \left[\mathbf{Q} \hat{\Phi}^\top \hat{\Phi} \mathbf{Q} \right] \mathbf{y} \\
 &\quad - \left(\frac{N}{n}\right)^2 \frac{1}{\hat{n}} \mathbf{y}^\top \mathbf{Q} \left(\frac{1}{n} \text{tr}(\hat{\Phi} \bar{\mathbf{Q}} \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})^\top) \frac{\mathbf{K}_{\cos}}{(1+\delta_{\cos})^2} + \frac{1}{n} \text{tr}(\hat{\Phi} \bar{\mathbf{Q}} \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})^\top) \frac{\mathbf{K}_{\sin}}{(1+\delta_{\sin})^2} \right) \mathbf{Q} \mathbf{y} \\
 &\simeq \left(\frac{N}{n}\right)^2 \frac{1}{\hat{n}} \mathbf{y}^\top \mathbb{E} \left[\mathbf{Q} \hat{\Phi}^\top \hat{\Phi} \mathbf{Q} \right] \mathbf{y} - \left(\frac{N}{n}\right)^2 \frac{1}{\hat{n}} \mathbf{y}^\top \left(\begin{bmatrix} \frac{1}{n} \text{tr}(\hat{\Phi} \bar{\mathbf{Q}} \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})^\top) & \frac{1}{n} \text{tr}(\hat{\Phi} \bar{\mathbf{Q}} \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})^\top) \\ \frac{1}{n} \text{tr}(\hat{\Phi} \bar{\mathbf{Q}} \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})^\top) & \frac{1}{n} \text{tr}(\hat{\Phi} \bar{\mathbf{Q}} \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})^\top) \end{bmatrix} \mathbb{E} \begin{bmatrix} \mathbf{Q} \mathbf{K}_{\cos} \mathbf{Q} \\ \mathbf{Q} \mathbf{K}_{\sin} \mathbf{Q} \end{bmatrix} \right) \mathbf{y} \\
 &\simeq \left(\frac{N}{n}\right)^2 \frac{1}{\hat{n}} \mathbf{y}^\top \bar{\mathbf{Q}} \Phi^\top \hat{\Phi} \bar{\mathbf{Q}} \mathbf{y} \\
 &\quad + \left(\frac{N}{n}\right)^2 \frac{1}{\hat{n}} \begin{bmatrix} \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \frac{N}{n} \hat{\Phi}^\top \hat{\Phi} \bar{\mathbf{Q}} \mathbf{K}_{\cos} - \frac{1}{n} \text{tr} \bar{\mathbf{Q}} \hat{\Phi} \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})) & \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \frac{N}{n} \hat{\Phi}^\top \hat{\Phi} \bar{\mathbf{Q}} \mathbf{K}_{\sin} - \frac{1}{n} \text{tr} \bar{\mathbf{Q}} \hat{\Phi} \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})) \\ \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \frac{N}{n} \hat{\Phi}^\top \hat{\Phi} \bar{\mathbf{Q}} \mathbf{K}_{\cos} - \frac{1}{n} \text{tr} \bar{\mathbf{Q}} \hat{\Phi} \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})) & \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \frac{N}{n} \hat{\Phi}^\top \hat{\Phi} \bar{\mathbf{Q}} \mathbf{K}_{\sin} - \frac{1}{n} \text{tr} \bar{\mathbf{Q}} \hat{\Phi} \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})) \end{bmatrix} \Omega \begin{bmatrix} \mathbf{y}^\top \bar{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}} \mathbf{y} \\ \mathbf{y}^\top \bar{\mathbf{Q}} \mathbf{K}_{\sin} \bar{\mathbf{Q}} \mathbf{y} \end{bmatrix}
 \end{aligned}$$

548 The last term \mathbf{Z}_{22} can be similarly treated as

$$\mathbf{Z}_{22} \simeq \frac{1}{n^2 \hat{n}} \mathbb{E} \sum_{i=1}^N \sum_{j \neq i} \mathbf{y}^\top \mathbf{Q}_{-j} \mathbf{U}_j \begin{bmatrix} \frac{\frac{1}{n} \text{tr}(\mathbf{Q} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X}))}{(1+\delta_{\cos})^2} & 0 \\ 0 & \frac{\frac{1}{n} \text{tr}(\mathbf{Q} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X}))}{(1+\delta_{\sin})^2} \end{bmatrix} \mathbf{U}_j^\top \mathbf{Q}_{-j} \mathbf{y}$$

549 where by Lemma 2 we deduce

$$\begin{aligned}
 \frac{1}{n} \text{tr}(\mathbf{Q} \mathbf{U}_i \hat{\mathbf{U}}_i^\top \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})) &\simeq \frac{1}{n} \text{tr} \left(\mathbf{Q}_{-i} \mathbf{U}_i (\mathbf{I}_2 + \mathbf{U}_i^\top \mathbf{Q}_{-i} \mathbf{U}_i)^{-1} \hat{\mathbf{U}}_i^\top \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X}) \right) \\
 &\simeq \frac{1}{n} \text{tr} \left(\mathbf{Q}_{-i} \mathbf{U}_i \begin{bmatrix} \frac{1}{1+\delta_{\cos}} & 0 \\ 0 & \frac{1}{1+\delta_{\sin}} \end{bmatrix} \hat{\mathbf{U}}_i^\top \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X}) \right) \simeq \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \hat{\Phi}^\top \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X}))
 \end{aligned}$$

550 so that by again Lemma 4

$$\begin{aligned}
 \mathbf{Z}_{22} &\simeq \frac{N}{n} \frac{1}{n \hat{n}} \mathbb{E} \sum_{j=1}^N \mathbf{y}^\top \mathbf{Q}_{-j} \mathbf{U}_j \begin{bmatrix} \frac{\frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \hat{\Phi}^\top \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X}))}{(1+\delta_{\cos})^2} & 0 \\ 0 & \frac{\frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \hat{\Phi}^\top \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X}))}{(1+\delta_{\sin})^2} \end{bmatrix} \mathbf{U}_j^\top \mathbf{Q}_{-j} \mathbf{y} \\
 &\simeq \left(\frac{N}{n}\right)^2 \frac{1}{\hat{n}} \mathbf{y}^\top \mathbb{E} \left[\mathbf{Q} \left(\frac{\frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \hat{\Phi}^\top \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X}))}{(1+\delta_{\cos})^2} \mathbf{K}_{\cos} + \frac{\frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \hat{\Phi}^\top \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X}))}{(1+\delta_{\sin})^2} \mathbf{K}_{\sin} \right) \mathbf{Q} \right] \mathbf{y} \\
 &\simeq \left(\frac{N}{n}\right)^2 \frac{1}{\hat{n}} \mathbf{y}^\top \left(\bar{\mathbf{Q}} \Xi \bar{\mathbf{Q}} + \frac{N}{n} \begin{bmatrix} \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \Xi \bar{\mathbf{Q}} \mathbf{K}_{\cos}) & \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \Xi \bar{\mathbf{Q}} \mathbf{K}_{\sin}) \\ \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \Xi \bar{\mathbf{Q}} \mathbf{K}_{\cos}) & \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \Xi \bar{\mathbf{Q}} \mathbf{K}_{\sin}) \end{bmatrix} \Omega \begin{bmatrix} \bar{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}} \\ \bar{\mathbf{Q}} \mathbf{K}_{\sin} \bar{\mathbf{Q}} \end{bmatrix} \right) \mathbf{y} \\
 &\simeq \left(\frac{N}{n}\right)^2 \frac{1}{\hat{n}} \begin{bmatrix} \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \hat{\Phi}^\top \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})) & \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \hat{\Phi}^\top \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})) \\ \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \hat{\Phi}^\top \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})) & \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \hat{\Phi}^\top \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})) \end{bmatrix} \Omega \begin{bmatrix} \mathbf{y}^\top \bar{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}} \mathbf{y} \\ \mathbf{y}^\top \bar{\mathbf{Q}} \mathbf{K}_{\sin} \bar{\mathbf{Q}} \mathbf{y} \end{bmatrix}.
 \end{aligned}$$

551 Assembling the estimates for \mathbf{Z}_1 , \mathbf{Z}_{21} and \mathbf{Z}_{22} , we get

$$\begin{aligned}
 \mathbb{E}[E_{\text{test}}] &\simeq \frac{1}{\hat{n}} \|\hat{\mathbf{y}}\|^2 - \frac{2}{\hat{n}} \hat{\mathbf{y}}^\top \frac{N}{n} \hat{\Phi} \bar{\mathbf{Q}} \mathbf{y} + \frac{1}{\hat{n}} \mathbf{y}^\top \left(\frac{N^2}{n^2} \bar{\mathbf{Q}} \hat{\Phi}^\top \hat{\Phi} \bar{\mathbf{Q}} \right) \mathbf{y} + \left(\frac{N}{n}\right)^2 \frac{1}{n \hat{n}} \times \\
 &\quad \left[\frac{\frac{n}{N} \text{tr} \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \hat{\mathbf{X}}) + \frac{N}{n} \text{tr} \bar{\mathbf{Q}} \hat{\Phi}^\top \hat{\Phi} \bar{\mathbf{Q}} \mathbf{K}_{\cos} - 2 \text{tr} \bar{\mathbf{Q}} \hat{\Phi}^\top \mathbf{K}_{\cos}(\hat{\mathbf{X}}, \mathbf{X})}{(1+\delta_{\cos})^2} \quad \frac{\frac{n}{N} \text{tr} \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \hat{\mathbf{X}}) + \frac{N}{n} \text{tr} \bar{\mathbf{Q}} \hat{\Phi}^\top \hat{\Phi} \bar{\mathbf{Q}} \mathbf{K}_{\sin} - 2 \text{tr} \bar{\mathbf{Q}} \hat{\Phi}^\top \mathbf{K}_{\sin}(\hat{\mathbf{X}}, \mathbf{X})}{(1+\delta_{\sin})^2} \right] \\
 &\quad \times \Omega \begin{bmatrix} \mathbf{y}^\top \bar{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}} \mathbf{y} \\ \mathbf{y}^\top \bar{\mathbf{Q}} \mathbf{K}_{\sin} \bar{\mathbf{Q}} \mathbf{y} \end{bmatrix}
 \end{aligned}$$

552 which, up to further simplifications, concludes the proof of Theorem 3.

553 E Several Useful Lemmas

554 **Lemma 5** (Some useful properties of Ω). *For any $\lambda > 0$ and Ω defined in (6), we have*

555 *1. all entries of Ω are positive;*

556 *2. for $2N = n$, $\det(\Omega^{-1})$, as well as the entries of Ω , scales like λ as $\lambda \rightarrow 0$;*

557 *Proof.* Developing the inverse we obtain

$$\Omega = \begin{bmatrix} 1 - \frac{N}{n} \frac{\frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}} \mathbf{K}_{\cos})}{(1+\delta_{\cos})^2} & -\frac{N}{n} \frac{\frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}} \mathbf{K}_{\sin})}{(1+\delta_{\sin})^2} \\ -\frac{N}{n} \frac{\frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \mathbf{K}_{\sin} \bar{\mathbf{Q}} \mathbf{K}_{\cos})}{(1+\delta_{\cos})^2} & 1 - \frac{N}{n} \frac{\frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \mathbf{K}_{\sin} \bar{\mathbf{Q}} \mathbf{K}_{\sin})}{(1+\delta_{\sin})^2} \end{bmatrix}^{-1}$$

558 we have $[\Omega^{-1}]_{11} = \frac{1}{1+\delta_{\cos}} + \frac{\lambda}{n} \text{tr} \bar{\mathbf{Q}} \frac{\mathbf{K}_{\cos}}{1+\delta_{\cos}} \bar{\mathbf{Q}} + \frac{N}{n} \frac{1}{n} \text{tr} \bar{\mathbf{Q}} \frac{\mathbf{K}_{\cos}}{1+\delta_{\cos}} \bar{\mathbf{Q}} \frac{\mathbf{K}_{\sin}}{1+\delta_{\sin}} > 0$, $[\Omega^{-1}]_{12} < 0$, and

559 similarly $[\Omega^{-1}]_{21} < 0$, $[\Omega^{-1}]_{22} > 0$. Furthermore, the determinant writes

$$\begin{aligned} \det(\Omega^{-1}) &= \left(1 - \frac{1}{n} \text{tr} \bar{\mathbf{Q}} \frac{\mathbf{K}_{\cos}}{1+\delta_{\cos}} + \frac{\lambda}{n} \text{tr} \bar{\mathbf{Q}} \frac{\mathbf{K}_{\cos}}{1+\delta_{\cos}} \bar{\mathbf{Q}}\right) \left(1 - \frac{1}{n} \text{tr} \bar{\mathbf{Q}} \frac{\mathbf{K}_{\sin}}{1+\delta_{\sin}} + \frac{\lambda}{n} \text{tr} \bar{\mathbf{Q}} \frac{\mathbf{K}_{\sin}}{1+\delta_{\sin}} \bar{\mathbf{Q}}\right) \\ &+ \left(1 - \frac{1}{n} \text{tr} \bar{\mathbf{Q}} \frac{\mathbf{K}_{\cos}}{1+\delta_{\cos}} + 1 - \frac{1}{n} \text{tr} \bar{\mathbf{Q}} \frac{\mathbf{K}_{\sin}}{1+\delta_{\sin}} + \frac{\lambda}{n} \text{tr} \bar{\mathbf{Q}} \left(\frac{\mathbf{K}_{\cos}}{1+\delta_{\cos}} + \frac{\mathbf{K}_{\sin}}{1+\delta_{\sin}}\right) \bar{\mathbf{Q}}\right) \\ &\times \frac{N}{n} \frac{1}{n} \text{tr} \bar{\mathbf{Q}} \frac{\mathbf{K}_{\cos}}{1+\delta_{\cos}} \bar{\mathbf{Q}} \frac{\mathbf{K}_{\sin}}{1+\delta_{\sin}} \end{aligned}$$

560 where we constantly use the fact that $\bar{\mathbf{Q}} \frac{N}{n} \left(\frac{\mathbf{K}_{\cos}}{1+\delta_{\cos}} + \frac{\mathbf{K}_{\sin}}{1+\delta_{\sin}}\right) = \mathbf{I}_n - \lambda \bar{\mathbf{Q}}$. Note that

$$\begin{aligned} 1 - \frac{1}{n} \text{tr} \bar{\mathbf{Q}} \frac{\mathbf{K}_{\cos}}{1+\delta_{\cos}} &= \frac{1}{1+\delta_{\cos}} > 0, \quad 1 - \frac{1}{n} \text{tr} \bar{\mathbf{Q}} \frac{\mathbf{K}_{\sin}}{1+\delta_{\sin}} = \frac{1}{1+\delta_{\sin}} > 0 \\ \frac{1}{1+\delta_{\cos}} + \frac{1}{1+\delta_{\sin}} &= 2 - \frac{n}{N} + \frac{\lambda}{N} \text{tr} \bar{\mathbf{Q}} > 0 \end{aligned}$$

561 so that 1) $\det(\Omega^{-1}) > 0$ and 2) for $2N = n$, $\det(\Omega^{-1})$ scales like λ as $\lambda \rightarrow 0$. □

562 **Lemma 6** (Derivatives with respect to N). *Let Assumption 1 holds, for any $\lambda > 0$ and*

$$\begin{cases} \delta_{\cos} = \frac{1}{n} \text{tr}(\mathbf{K}_{\cos} \bar{\mathbf{Q}}) = \frac{1}{n} \text{tr} \mathbf{K}_{\cos} \left(\frac{N}{n} \left(\frac{\mathbf{K}_{\cos}}{1+\delta_{\cos}} + \frac{\mathbf{K}_{\sin}}{1+\delta_{\sin}} \right) + \lambda \mathbf{I}_n \right)^{-1} \\ \delta_{\sin} = \frac{1}{n} \text{tr}(\mathbf{K}_{\sin} \bar{\mathbf{Q}}) = \frac{1}{n} \text{tr} \mathbf{K}_{\sin} \left(\frac{N}{n} \left(\frac{\mathbf{K}_{\cos}}{1+\delta_{\cos}} + \frac{\mathbf{K}_{\sin}}{1+\delta_{\sin}} \right) + \lambda \mathbf{I}_n \right)^{-1} \end{cases}$$

563 *defined in Theorem 1, we have that $(\delta_{\cos}, \delta_{\sin})$ and $\|\bar{\mathbf{Q}}\|$ are all decreasing functions of N . Note in*
564 *particular that the same conclusion holds for $2N > n$ as $\lambda \rightarrow 0$.*

565 *Proof.* We write

$$\begin{bmatrix} \frac{\partial \delta_{\cos}}{\partial N} \\ \frac{\partial \delta_{\sin}}{\partial N} \end{bmatrix} = -\frac{1}{n} \Omega \begin{bmatrix} \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \Phi \bar{\mathbf{Q}} \mathbf{K}_{\cos}) \\ \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \Phi \bar{\mathbf{Q}} \mathbf{K}_{\sin}) \end{bmatrix} = -\frac{n}{N} \frac{1}{n} \Omega \begin{bmatrix} \delta_{\cos} - \frac{\lambda}{n} \text{tr}(\bar{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}}) \\ \delta_{\sin} - \frac{\lambda}{n} \text{tr}(\bar{\mathbf{Q}} \mathbf{K}_{\sin} \bar{\mathbf{Q}}) \end{bmatrix} \quad (16)$$

566 for Ω defined in (6) and $\Phi = \frac{\mathbf{K}_{\cos}}{1+\delta_{\cos}} + \frac{\mathbf{K}_{\sin}}{1+\delta_{\sin}}$, which, together with Lemma 5, allows us to conclude
567 that $\frac{\partial \delta_{\cos}}{\partial N}, \frac{\partial \delta_{\sin}}{\partial N} < 0$. Further note that

$$\frac{\partial \bar{\mathbf{Q}}}{\partial N} = -\frac{1}{n} \bar{\mathbf{Q}} \left(\Phi - \frac{\mathbf{K}_{\cos}}{(1+\delta_{\cos})^2} N \frac{\partial \delta_{\cos}}{\partial N} - \frac{\mathbf{K}_{\sin}}{(1+\delta_{\sin})^2} N \frac{\partial \delta_{\sin}}{\partial N} \right) \bar{\mathbf{Q}}$$

568 which concludes the proof. □

569 **Lemma 7** (Derivative with respect to λ). *For any $\lambda > 0$, $(\delta_{\cos}, \delta_{\sin})$ and $\|\bar{\mathbf{Q}}\|$ defined in Theorem 1*
570 *decrease as λ grows large.*

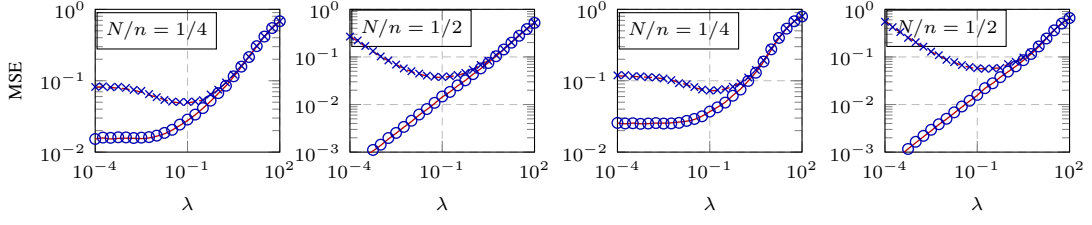


Figure 6: MSEs of RFF regression on Fashion-MNIST (**left two**) and Kannada-MNIST (**right two**) data (class 5 versus 6), as a function of regression parameter λ , for $p = 784$, $n = \hat{n} = 1\,024$, $N = 256$ and 512 . Empirical results displayed in **blue** (circles for training and crosses for test); and the asymptotics from Theorem 2 and 3 displayed in **red** (solid lines for training and dashed for test). Results obtained by averaging over 30 runs.

571 *Proof.* Taking the derivative of $(\delta_{\cos}, \delta_{\sin})$ with respect to $\lambda > 0$, we have explicitly

$$\begin{bmatrix} \frac{\partial \delta_{\cos}}{\partial \lambda} \\ \frac{\partial \delta_{\sin}}{\partial \lambda} \end{bmatrix} = -\Omega \begin{bmatrix} \frac{1}{q} \text{tr}(\bar{\mathbf{Q}} \mathbf{K}_{\cos} \bar{\mathbf{Q}}) \\ \frac{1}{n} \text{tr}(\bar{\mathbf{Q}} \mathbf{K}_{\sin} \bar{\mathbf{Q}}) \end{bmatrix} \quad (17)$$

572 which, together with the fact that all entries of Ω are positive (Lemma 5), allows us to conclude that
 573 $\frac{\partial \delta_{\cos}}{\partial \lambda}, \frac{\partial \delta_{\sin}}{\partial \lambda} < 0$. Further considering

$$\frac{\partial \bar{\mathbf{Q}}}{\partial \lambda} = \bar{\mathbf{Q}} \left(\frac{N}{n} \frac{\mathbf{K}_{\cos}}{(1 + \delta_{\cos})^2} \frac{\partial \delta_{\cos}}{\partial \lambda} + \frac{N}{n} \frac{\mathbf{K}_{\sin}}{(1 + \delta_{\sin})^2} \frac{\partial \delta_{\sin}}{\partial \lambda} - \mathbf{I}_n \right) \bar{\mathbf{Q}}$$

574 and thus the conclusion for $\bar{\mathbf{Q}}$.

575 □

576 F Additional Real-world Data sets

577 We have presented results in detail for one particular real-world data set, but we have extensive
 578 empirical results demonstrating that similar conclusions hold more broadly. As an example of this,
 579 here we present a numerical evaluation of our results on several other real-world image data sets. We
 580 consider the classification task on two MNIST-like data sets composed of 28×28 grayscale images:
 581 the Fashion-MNIST [51] and the Kannada-MNIST [41] data sets. Each image is represented as a
 582 $p = 784$ -dimensional vector and the output targets $\mathbf{y}, \hat{\mathbf{y}}$ are taken to have $-1, +1$ entries depending
 583 on the image class. As a consequence, both the training and test MSEs in (3) are approximately 1 for
 584 $N = 0$ and significantly small λ , as observed in Figure 4 (and Figure 8). For each data set, images
 585 were jointly centered and scaled so to fall close to the setting of Assumption 1 on \mathbf{X} and $\hat{\mathbf{X}}$.

586 In Figure 6, we compare the empirical training and test errors with their limiting behaviors derived
 587 from Theorem 2 and 3, as a function of the penalty parameter λ , on a training set of size $n = 1\,024$
 588 (512 images from class 5 and 512 images from class 6) with feature dimension $N = 256$, on both
 589 data sets. A close fit between theory and practice is observed, for moderately large values of n, p, N ,
 590 demonstrating thus a wide practical applicability of the proposed asymptotic analyses, particularly
 591 compared to the (limiting) Gaussian kernel predictions per Figure 1.

592 In Figure 7, we report the behavior of the pair $(\delta_{\cos}, \delta_{\sin})$ for small values of $\lambda = 10^{-7}$ and 10^{-3} .
 593 Similar to the two leftmost plots in Figure 3 for MNIST, a jump from the under- to over-parameterized
 594 regime occurs at the interpolation threshold $2N = n$, in both Fashion- and Kannada-MNIST data
 595 sets, clearly indicating the two phases of learning and the phase transition between them.

596 In Figure 8, we report the empirical and theoretical test errors as a function of the ratio N/n , on a
 597 training test of size $n = 500$ (250 images from class 8 and 250 images from class 9), by varying
 598 feature dimension N . An exceedingly small regularization $\lambda = 10^{-7}$ is applied to mimic the
 599 “ridgeless” limiting behavior as $\lambda \rightarrow 0$. On both data sets, the corresponding double descent curve
 600 is observed where the test errors goes down and up, with a singular peak around $2N = n$, and then
 601 goes down monotonically as N continues to increase when $2N > n$.

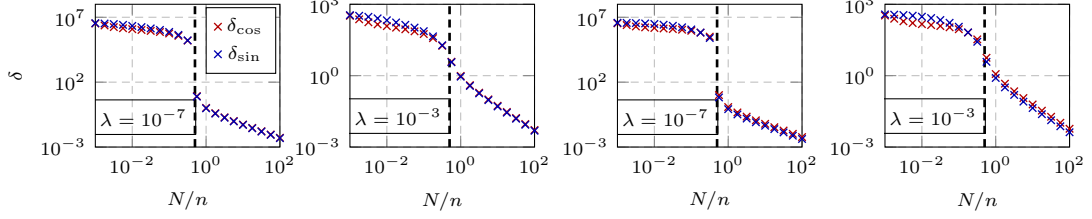


Figure 7: Behavior of $(\delta_{\cos}, \delta_{\sin})$ in (9), on Fashion-MNIST (**left two**) and Kannada-MNIST (**right two**) data (class 8 versus 9), for $p = 784$, $n = 1000$, $\lambda = 10^{-7}$ and 10^{-3} . The **black** dashed line is the interpolation threshold $2N = n$.

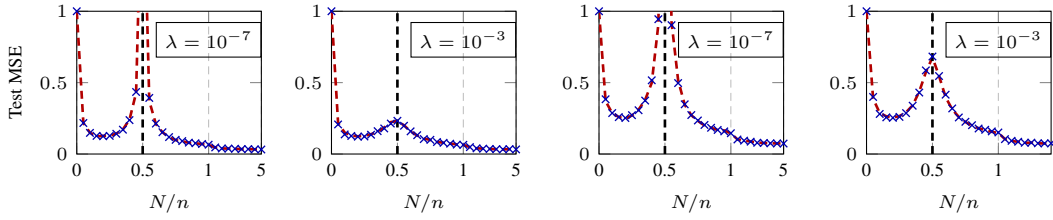


Figure 8: Empirical (crosses) and theoretical (dashed lines) test error of RFF regression, as a function of the ratio N/n , on Fashion-MNIST (**left two**) and Kannada-MNIST (**right two**) data (class 8 versus 9), for $p = 784$, $n = 500$, $\lambda = 10^{-7}$ and 10^{-3} . The **black** dashed line is the interpolation threshold $2N = n$. Results obtained by averaging over 30 runs.