

Concentration of Random Feature Matrices in High-Dimensions

Zhijun Chen¹, Hayden Schaeffer¹ and Rachel Ward²

¹Carnegie Mellon University

²The University of Texas at Austin

Abstract

The spectra of random feature matrices provide essential information on the conditioning of the linear system used in random feature regression problems and are thus connected to the consistency and generalization of random feature models. Random feature matrices are asymmetric rectangular nonlinear matrices depending on two input variables, the data and the weights, which can make their characterization challenging. We consider two settings for the two input variables, either both are random variables or one is a random variable and the other is well-separated, i.e. there is a minimum distance between points. With conditions on the dimension, the complexity ratio, and the sampling variance, we show that the singular values of these matrices concentrate near their full expectation and near one with high-probability. In particular, since the dimension depends only on the logarithm of the number of random weights or the number of data points, our complexity bounds can be achieved even in moderate dimensions for many practical setting. The theoretical results are verified with numerical experiments.

1 Introduction

Kernel methods are some of the most popular approaches in machine learning and have been applied to image processing, classification, and data-based regression problems. Suppose $K(x, x')$ is the kernel, which measures the similarity between two input samples x and x' , then the kernel approach generates an approximation by using a weighted sum of the kernel applied to the data. This is justified since minimizers of kernel training problems within the reproducing kernel Hilbert space associated with the kernel K are guaranteed to be of the form $\sum_{j=1}^m c_j K(x, x_j)$ by the representer theorem [20], where $\{x_j\}_{j=1}^m$ are data samples. Standard kernel methods require the formation of the kernel matrix which uses all of the data points and thus can be intractable for large problems. However, if the kernel is symmetric and positive definite, i.e. $K : \mathbb{R}^{2d} \rightarrow \mathbb{R}$, $K(x, x') = K(x', x)$, and $\langle \mathbf{c}, \mathbf{K}\mathbf{c} \rangle > 0$ where \mathbf{K} is the matrix whose elements are $K(x_i, x_j)$ where x_i and x_j are data points, then $K(x, x') = \langle \psi(x), \psi(x') \rangle$ with feature map ψ . Rather than forming and computing solutions using the kernel matrix directly, the random feature method (RFM) [26–28] uses a randomized approximation to the inner product. That is, the kernel can be approximated by $\Phi(x)^T \Phi(x')$, where $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^N$ and N is the number of randomized features used. When N is not large, the randomized method leads to a significant reduction in the computational cost.

In this work, we examine the singular values (and thus the condition number) of the random feature matrix. In particular, consider the following formulation of the random feature regression problem. Let f be the target function and suppose we are given samples (\mathbf{x}_j, y_j) where $y_j = f(\mathbf{x}_j) + e_j$ with noise e_j for $j \in [m]$. We approximate the target function by the finite sum

$$f^\sharp(\mathbf{x}) = \sum_{k=1}^N c_k^\sharp \phi(\mathbf{x}, \boldsymbol{\omega}_k).$$

where $\phi(\mathbf{x}, \boldsymbol{\omega}) = \phi(\langle \mathbf{x}, \boldsymbol{\omega} \rangle)$ is a feature map parameterized by a random variable $\boldsymbol{\omega}$ drawn from some user defined probability $\rho(\boldsymbol{\omega})$. The function ϕ is an activation function and is often chosen to be the complex exponential function or ReLU. In this way, a random feature model is a two-layer neural network where the hidden weight layer is randomized by some prescribed process and not trained. If we define the output vector $\mathbf{y} = [y_1, \dots, y_m]^T \in \mathbb{C}^m$ and the random feature matrix $\mathbf{A} \in \mathbb{C}^{m \times N}$ whose elements are defined as $\mathbf{A}_{j,k} = \phi(\langle \mathbf{x}_j, \boldsymbol{\omega}_k \rangle)$ for $j \in [m]$ and $k \in [N]$, then the random feature ridge regression problem is

$$\min_{\mathbf{c} \in \mathbb{C}^N} \|\mathbf{y} - \mathbf{A}\mathbf{c}\|_2^2 + \lambda \|\mathbf{c}\|_2^2$$

for ridge parameter $\lambda \geq 0$. In the case where $\lambda = 0$ we refer to the problem as *ridgeless* regression. The solution to the regression problem and its generalization error depends on the spectrum of the Gram matrix $\mathbf{A}^* \mathbf{A}$, which is also the kernel matrix, see for example [1, 6, 15, 17, 31].

The earlier results on spectra of kernel matrices, at least from the random matrix theory perspective, focused on square inner-product kernels that take the form $K_{i,j} = \phi(\langle \mathbf{x}_i, \mathbf{x}_j \rangle)$ and thus only depend on one random variable. In [8], the spectra of kernel matrices that depend on $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$ or $\|\mathbf{x}_i - \mathbf{x}_j\|^2$ and are consistent with their linearization were shown to converge to a Marchenko-Pastur distribution. The limiting distribution of square inner-product kernels in the high-dimensional and large dataset size setting were derived in [7]. The convergence of the spectral norm and the extreme eigenvalues of these matrices were determined in [9], specifically, that the spectral norm converges almost surely to the edge of the limiting spectrum.

In this work, we consider the concentration of the spectrum of asymmetric rectangular (non-linear) random matrices $\mathbf{A}_{j,k} = \phi(\langle \mathbf{x}_j, \boldsymbol{\omega}_k \rangle)$ which depends on two variables, motivated by the random matrices used in random feature models. Specifically, these matrices are the (non-fixed) dictionaries that are obtained by randomizing and not training hidden layers in a neural network. The spectrum of the random feature matrices, in particular the spread of the singular values, can be used to understand the numerical and theoretical properties of various random feature algorithms. Recent results in the literature consider the asymptotic distribution of the the Gram matrix formed from this type of asymmetric rectangular random matrix. The limiting spectral density of $\mathbf{A}^T \mathbf{A}$ when the samples and weights are Gaussian was computed in [25] using the moment method and in [14] using an approach from [7], see also the related work in [23, 24]. The authors of [5] extended the previous results to the case where \mathbf{x} and $\boldsymbol{\omega}$ are both subgaussian. In the setting where \mathbf{x} is deterministic and $\boldsymbol{\omega}$ is random, [18] determined the empirical spectrum in the large data and dimension limit. In [16], the authors provide a precise characterization of the Gram matrix in the setting where the dimension, the size of the feature space, and the size of the dataset are large and comparable. In particular, in this setting they showed that the limit of the Gram matrix generated from random Fourier features is not the Gaussian kernel. In [35], non-asymptotic estimates are

given in the overparameterized setting using a pairwise approximate orthogonality condition [10]. Also, [22] considers the behavior of the random Fourier features model in high dimensions. For the setting with random data and fixed frequencies, Theorem 6.1 in [22] shows high probability bounds similar to our Theorem 4.1. The article uses the bounds to characterize performance for sparse recovery with Fourier features, which we became aware of after publication of this work.

Random feature regression also exhibits the double descent phenomena [6, 16, 19], in which the risk is low in the underparameterized and overparameterized regions but peaks at the interpolation threshold [2–4]. This behavior is intrinsically dependent on the characterization of the spectrum of the random feature matrix. In [19], a detailed analysis of the double descent behavior of the random feature regression problem is shown, with the assumption that one has m data samples from the d -dimensional sphere \mathbb{S}^{d-1} and N random features, with $N, m, d \rightarrow \infty$ but comparable. In particular, they showed that overparameterization is necessary to obtain the optimal test error in certain settings. In [6], the singular values of the random feature matrix with Gaussian data samples and weights are shown to concentration around 1 with high probability when the complexity ratio $\frac{N}{m}$ scales like $\log^{-1}(N)$ (underparameterized) or $\log(m)$ (overparameterized). In addition, they showed that the condition number becomes unbounded for N close to m (i.e. a double descent phenomena for the condition number), which also provided a mechanism for the double descent in the generalization error associated with ridgeless random feature regression.

1.1 Our Contributions

In this work, we derive concentration bounds on the spectrum of asymmetric rectangular (nonlinear) random matrices whose entries are of the form $\mathbf{A}_{j,k} = \phi(\langle \mathbf{x}_j, \boldsymbol{\omega}_k \rangle)$. Throughout this work, we set the activation function to be the complex exponential, i.e. $\phi(z) = \exp(iz)$. We expect similar results to hold for other activation functions.

We consider two settings on the two variables, either they are both random variables (like in [6, 14, 19, 23–25]) or one is a random variable and the other is *well-separated* (see Section 4 for the precise statement). Similar to [6], we focus on the finite m and N setting; however, in this work we provide a new characterization for the spectrum as a function of the parameters. Our results improve and generalize [6] by only requiring one of the variables to be Gaussian and separating the variance and dimensional parameters in the theory. Our results also complement the previous work in the literature highlighted in the introduction, for example, by considering the dimensional scaling in the sampling process [19], subgaussian random variables [5], and incorporating more general data sampling processes.

Rectangular matrices whose entries are sampled i.i.d. from a Gaussian (or subgaussian) distribution are close to isometries when the ratio of the number of rows and columns scale logarithmically. Intuitively, this should imply that a random feature matrix built from a random rectangular matrix (i.e. the random weights) should be well-conditioned if the dimension scales like the number of features and log factors. We show that the conditional expectation of the random feature matrix concentrates quickly to the expectation as the dimension of the input increases. Thus, one can relax conditions in previous works when the dimension is sufficiently large (this is made precise in the theorem statements). One reason this is important is that these results provide more practical parameter regimes in which one should expect (with high-probability) to obtain well-conditioned linear systems and trained models that generalize to new data. Our results hold for a finite number of data samples and weights, which differs from the asymptotic analysis provided in other works.

1.2 Notation

Let $i = \sqrt{-1}$ be the imaginary unit. For an integer N , the set $[N]$ is defined as $[N] = \{1, 2, \dots, N\}$. We use bold letters to denote vectors and matrices. The $d \times d$ identity matrix is denoted \mathbf{I}_d . We denote the ℓ^p -norm of a vector $\mathbf{x} \in \mathbb{C}^d$ by $\|\mathbf{x}\|_p$ and the induced ℓ^p norm of a matrix $\mathbf{A} \in \mathbb{C}^{m \times N}$ by $\|\mathbf{A}\|_p$. The transpose of a matrix $\mathbf{A} \in \mathbb{C}^{m \times N}$ is denoted by \mathbf{A}^T , and the conjugate transpose is denoted by \mathbf{A}^* . We use $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ to denote the Gaussian distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$.

2 Summary of Main Results

Let $\{\mathbf{x}_j\}_{j \in [m]} \subset \mathbb{R}^d$ be data points sampled from a distribution μ and $\{\boldsymbol{\omega}_k\}_{k \in [N]} \subset \mathbb{R}^d$ be feature weights sampled from another distribution ρ . Define the random feature matrix \mathbf{A} component-wise by $\mathbf{A}_{j,k} = \exp(i\langle \mathbf{x}_j, \boldsymbol{\omega}_k \rangle)$. In this section, we give the conditions on the dimension d , the number of data points m , and the number of feature weights N so that the normalized Gram matrix $m^{-1}\mathbf{A}^*\mathbf{A}$ (or $N^{-1}\mathbf{A}\mathbf{A}^*$) is close to both the identity matrix \mathbf{I} and its expectation.

2.1 Concentration with Randomized Inputs

We consider the setting where the covariance matrix of μ is a multiple of $\frac{1}{d}\mathbf{I}_d$ so that the expectation of $\|\mathbf{x}_j\|_2$ is a constant independent of d and thus the data concentration are the sphere with a fixed radius. This is a weaker version of the assumptions used in [19]. The motivation for this is to avoid the data becoming arbitrary large in high-dimensions. The main theorems when both \mathbf{x} and $\boldsymbol{\omega}$ are random variables are stated below.

Theorem 2.1 (Concentration in the Underparameterized Setting). *Suppose that the random feature matrix \mathbf{A} is defined component-wise by $\mathbf{A}_{j,k} = \exp(i\langle \mathbf{x}_j, \boldsymbol{\omega}_k \rangle)$, $\{\mathbf{x}_j\}_{j \in [m]} \subset \mathbb{R}^d$ are data points sampled from $\mathcal{N}(\mathbf{0}, \frac{\gamma^2}{d}\mathbf{I})$, and $\{\boldsymbol{\omega}_k\}_{k \in [N]} \subset \mathbb{R}^d$ are feature weights such that the components of $\boldsymbol{\omega}_k$ are independent mean-zero subgaussian random variables with the same variance σ^2 and the same subgaussian parameters β, κ (see Section 4.2). Then there exist a constant $C_1 > 0$ (depending only on the subgaussian parameters) and a universal constant $C_2 > 0$ such that if the following conditions hold*

$$d \geq C_1 \log \left(\frac{N}{\delta} \right) \tag{2.1}$$

$$\gamma^2 \sigma^2 \geq 4 \log \left(\frac{2N}{\eta} \right) \tag{2.2}$$

$$m \geq C_2 \eta^{-2} N \log \left(\frac{2N}{\delta} \right), \tag{2.3}$$

for $\delta, \eta \in (0, 1)$, then we have

$$\left\| \frac{1}{m} \mathbf{A}^* \mathbf{A} - \mathbf{I}_N \right\|_2 \leq 2\eta, \tag{2.4}$$

with probability at least $1 - 3\delta$. Moreover, if $\eta \geq 2\delta$ (which holds for practical η and δ), then we simultaneously have

$$\left\| \frac{1}{m} \mathbf{A}^* \mathbf{A} - \mathbb{E}_{\mathbf{x}, \omega} \left[\frac{1}{m} \mathbf{A}^* \mathbf{A} \right] \right\|_2 \leq 2\eta. \quad (2.5)$$

The main idea in the proof of Theorem 2.1 is to bound the difference in (2.4) by

$$\left\| \frac{1}{m} \mathbf{A}^* \mathbf{A} - \mathbf{I}_N \right\|_2 \leq \frac{1}{m} \|\mathbf{A}^* \mathbf{A} - \mathbb{E}_{\mathbf{x}} [\mathbf{A}^* \mathbf{A}]\|_2 + \left\| \mathbb{E}_{\mathbf{x}} \left[\frac{1}{m} \mathbf{A}^* \mathbf{A} \right] - \mathbf{I}_N \right\|_2, \quad (2.6)$$

and the difference in (2.5) by

$$\left\| \frac{1}{m} \mathbf{A}^* \mathbf{A} - \mathbb{E}_{\mathbf{x}, \omega} \left[\frac{1}{m} \mathbf{A}^* \mathbf{A} \right] \right\|_2 \leq \frac{1}{m} \|\mathbf{A}^* \mathbf{A} - \mathbb{E}_{\mathbf{x}} [\mathbf{A}^* \mathbf{A}]\|_2 + \frac{1}{m} \|\mathbb{E}_{\mathbf{x}} [\mathbf{A}^* \mathbf{A}] - \mathbb{E}_{\mathbf{x}, \omega} [\mathbf{A}^* \mathbf{A}]\|_2. \quad (2.7)$$

While entries of these matrices are not i.i.d., the first term on the right-hand side of (2.6) (and (2.7)) can be decomposed as the summation of independent matrices which allows us to leverage stronger matrix concentration inequalities. The remaining terms use a weaker result on large deviations which, surprisingly, does not change the overall complexity bounds. In the proofs, we will show that each of the terms on the right-hand side of (2.6) and (2.7) are bounded by η simultaneously when the conditions (2.1), (2.2) and (2.3) are satisfied (see Section 4). Condition (2.1) ensures that the dimension is large enough so that points separate and thus provides a minimal distance condition between random weight vectors in high dimensions. Condition (2.2) resembles an uncertainty principle between the spread of the samples and the weights. Lastly, condition (2.3) is a complexity relation between the number of samples and the number of features. Similar conditions are imposed in the other theorems.

By the symmetry of the input variables, we also have the bounds for the overparameterized case $m < N$.

Theorem 2.2 (Concentration in the Overparameterized Setting). *Suppose that the random feature matrix \mathbf{A} is defined component-wise by $\mathbf{A}_{j,k} = \exp(i\langle \mathbf{x}_j, \boldsymbol{\omega}_k \rangle)$, $\{\boldsymbol{\omega}_k\}_{k \in [N]} \subset \mathbb{R}^d$ are feature weights sampled from $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d)$, and $\{\mathbf{x}_j\}_{j \in [m]} \subset \mathbb{R}^d$ are data points such that the components of \mathbf{x}_j are independent mean-zero subgaussian random variables with the same variance γ^2/d and the same subgaussian parameters β, κ . Then there exist a constant $C_1 > 0$ (depending only on subgaussian parameters) and a universal constant $C_2 > 0$ such that if the following conditions hold*

$$\begin{aligned} d &\geq C_1 \log \left(\frac{m}{\delta} \right) \\ \gamma^2 \sigma^2 &\geq 4 \log \left(\frac{2m}{\eta} \right) \\ N &\geq C_2 \eta^{-2} m \log \left(\frac{2m}{\delta} \right), \end{aligned}$$

for $\delta, \eta \in (0, 1)$, then we have

$$\left\| \frac{1}{N} \mathbf{A} \mathbf{A}^* - \mathbf{I}_m \right\|_2 \leq 2\eta,$$

with probability at least $1 - 3\delta$. Moreover, if $\eta \geq 2\delta$ (which holds for practical η and δ), then we simultaneously have

$$\left\| \frac{1}{N} \mathbf{A} \mathbf{A}^* - \mathbb{E}_{\mathbf{x}, \boldsymbol{\omega}} \left[\frac{1}{N} \mathbf{A} \mathbf{A}^* \right] \right\|_2 \leq 2\eta.$$

A direct consequence of Theorem 2.1 is that all of the eigenvalues of $m^{-1} \mathbf{A}^* \mathbf{A}$ are close to 1. Specifically, if the conditions in Theorem 2.1 are satisfied, then

$$\left| \lambda_k \left(\frac{1}{m} \mathbf{A}^* \mathbf{A} \right) - 1 \right| \leq 2\eta,$$

with probability at least $1 - 3\delta$. Here $\lambda_k(\mathbf{B})$ is the k -th eigenvalue of the matrix \mathbf{B} . Similar results also hold for $N^{-1} \mathbf{A} \mathbf{A}^*$ if the conditions in Theorem 2.2 are satisfied. This provides an upper bound for the largest eigenvalue and a lower bound for the smallest eigenvalue. Therefore, we can conclude that the matrix \mathbf{A} has small condition number when the complexity conditions in the theorems are satisfied.

2.2 Concentration to the Gaussian Kernel

Theorem 2.2 is actually a consequence of the more general results which only requires that the data is well-separated, i.e. there is a minimum distance between data samples.

Theorem 2.3 (Concentration to the Kernel). *Suppose that the random feature matrix \mathbf{A} is defined (component-wise) by $\mathbf{A}_{j,k} = \exp(i\langle \mathbf{x}_j, \boldsymbol{\omega}_k \rangle)$, the feature weights $\{\boldsymbol{\omega}_k\}_{k \in [N]} \subset \mathbb{R}^d$ are sampled from $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d)$, and that for the data points $\{\mathbf{x}_j\}_{j \in [m]} \subset \mathbb{R}^d$ there is a constant $R > 0$ such that $\|\mathbf{x}_j - \mathbf{x}_k\|_2^2 \geq R$ for all $j, k \in [m]$ with $j \neq k$. If the following conditions hold*

$$N \geq C\eta^{-2} m \log \left(\frac{2m}{\delta} \right) \tag{2.8}$$

$$\sigma^2 \geq \frac{2}{R} \log \left(\frac{m}{\eta} \right), \tag{2.9}$$

for some $\delta, \eta \in (0, 1)$, where $C > 0$ is a universal constant. Then with probability at least $1 - \delta$ we have

$$\left\| \frac{1}{N} \mathbf{A} \mathbf{A}^* - \mathbb{E}_{\boldsymbol{\omega}} \left[\frac{1}{N} \mathbf{A} \mathbf{A}^* \right] \right\|_2 \leq \eta.$$

The term $\mathbb{E}_{\boldsymbol{\omega}} \left[\frac{1}{N} \mathbf{A} \mathbf{A}^* \right]$ is in fact the associated kernel matrix for this problem, i.e. the Gaussian kernel matrix. Although conditions (2.8) and (2.9) do not directly use the dimension of the data, the results implicitly improve in high-dimensions where the distance between data points can be larger.

2.3 Comparison with Similar Results

There are several important differences between Theorems 2.1 and 2.2 as compared to Theorem 3.1(a-b) from [6]. For this discussion, we will ignore the universal constants since they may not be

optimal in the theorems discussed (and come from difference sampling distributions) and instead focus on the parameter dependencies. The first difference is the dependence between the failure rate δ , the concentration parameter η , and the number of random features N . Theorems 2.1 and 2.2 uncovers a refined relationship between the parameters when the dimension is sufficiently large, in particular, N is controlled by

$$\min(\delta \exp(d), \eta \exp(\gamma^2 \sigma^2))$$

i.e. larger dimensions lead to a smaller failure rate, while in [6] (when d is sufficiently large) N is controlled by

$$\sqrt{\delta \eta} \exp(\gamma^2 \sigma^2)$$

which implies that the product of the variances must be used to compensate for the failure rate and concentration parameter. We also extend the results to include subgaussian sampling or weights, which only change the universal constants in our bounds. In addition, we show that the concentration relies on the separation between points and does not require both variables to be random (see Theorem 2.3).

While the theorems and proofs in Section 4 use similar concentration techniques to those found in the compressive sensing (CS) literature, they differ in several key places. Both consider rectangular nonlinear random matrices; however, in this work the two inputs for the matrix can both be random variables. The bounded orthonormal systems (BOS) matrices found in compressed sensing assumes that basis parameters (ω for example) are fixed [11]. Some standard examples of BOS include the Fourier basis on $[0, 1]$ where \mathbf{x} is sampled from the uniform distribution on $[0, 1]$, the tensorized Legendre polynomial basis where \mathbf{x} is sampled from the uniform distribution on $[-1, 1]^d$, or the multivariate Hermite polynomials where \mathbf{x} is sampled from the multivariate Gaussian distribution in d -dimensions [29]. In each of these cases, the basis is generated so that it is orthogonal with respect to the sampling density for \mathbf{x} which is not the case for random feature matrices. In addition, our results hold for unbounded data (e.g. Gaussian), which has been a long-standing question for the analysis of the restricted isometry property for BOS matrices [12]. This is an example where the random feature matrices are more robust to the data sampling than orthogonality-based methods.

3 Numerical Experiments

In this section, we verify some of the results numerically and show that the condition on the dimension, i.e. $d \geq C_1 \log(\frac{N}{\delta})$, may have a favorable constant in practice. In Figure 3.1, we display the concentration of the singular values as a function of the dimension by plotting the maximum and minimum singular values. The curves envelop the range of the singular values. The random feature matrix is constructed using the complex exponential function with $m = 100$ random samples drawn from the normal distribution with $\gamma = 1$ and $N = 5000$ random weights. The matrices are normalized so that each column has unit ℓ^2 norm. We chose the scaling between m and N so that Equation (2.3) would hold. The weights are drawn from the normal distribution with the standard deviation σ specified in Figure 3.1. For each dimension, we used 10 trials to calculate the mean of the extreme singular values (the solid curves) and one standard deviation (the shaded regions). The plots indicate that as σ increases or as the dimension d increases, the singular values concentrate quickly. In this case as σ increases, the potential range decreases up

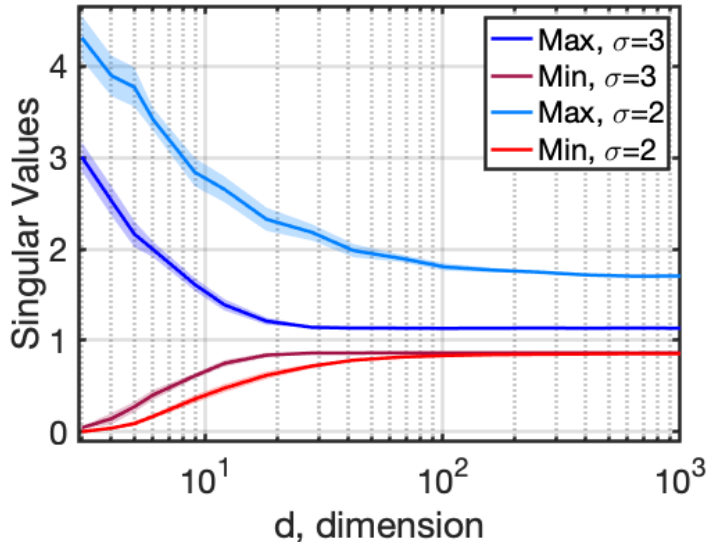


Figure 3.1: **Extreme Singular Values versus dimension:** The plot displays the maximum and minimum singular values for various dimensions d . For each dimension, 10 trials are used to calculate the mean value (the solid curves) and one standard deviation (shaded regions). The random feature matrix is the complex exponential with $m = 100$, $N = 5000$, and $\gamma = 1$. The standard deviation of the weights σ are specified for each curve. Even with a moderate dimension, i.e. $d \geq 10$, the condition number is already less than 15 for both examples.

to the range indicated by $\sigma = 3$, i.e. all pairs of curves for $\sigma \geq 3$ will have the same upper and the same lower plateaus as the curves generated with $\sigma = 3$. Additionally, for $\sigma = 3$, once the dimension exceeds 3 the matrix is well-conditioned.

We study the distribution of the singular value for different parameters in Figure 3.2 and Figure 3.3. The theory shows that the distribution should be close to the line $y = 1$. In Figure 3.2, the dimension d and the number of random features N are varied and the distribution of the singular values are plotted in ascending order (thus the x -axis is the sorted index). The random feature matrix is the complex exponential (with normalized columns) with $m = 100$, $\gamma = 1$, and $\sigma = 3$. The left plot in Figure 3.2 has $N = 500$ random features and the right plot in Figure 3.2 has $N = 5000$. In both plots, when the dimension is 3, the minimum singular values are close to zero, thus leading to ill-conditioning of \mathbf{A} . For $d \geq 6$ the singular values range from about 0.5 to 2 with noticeable concentration for $d = 12$. Note that the concentration is more significant for larger N while retaining a similar dependency on d . In Figure 3.3, we consider the case where the dimension d and the standard deviation of the random features σ are varied. The random feature matrix uses the same setup as Figure 3.2 except that $N = 5000$ while the standard deviation σ is varied. The left plot in the Figure 3.3 uses $\sigma = 2$ and the right plot in Figure 3.3 uses $\sigma = 4$, showing that larger values of σ and d lead to better concentration. Note that there is an asymmetry in the distributions of the singular values around 1 as seen in Figure 3.2 and Figure 3.3.

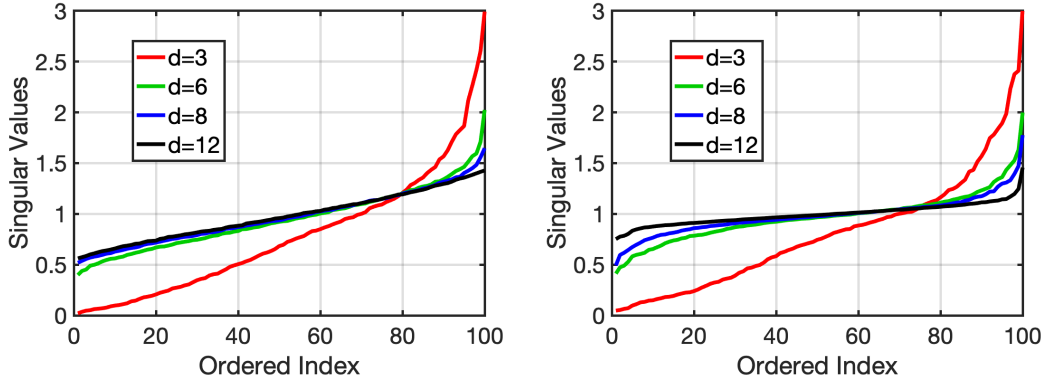


Figure 3.2: **Singular Value Distribution with Different d and N** : The figure shows the distribution of the singular values in ascending as a function of the dimension d and the number of random features N . The random feature matrix is the complex exponential with $m = 100$, $\gamma = 1$, and $\sigma = 3$. The plot on the left uses $N = 500$ and the plot on the right uses $N = 5000$. This experiment shows that the singular values concentration around 1 quickly in dimension and N .

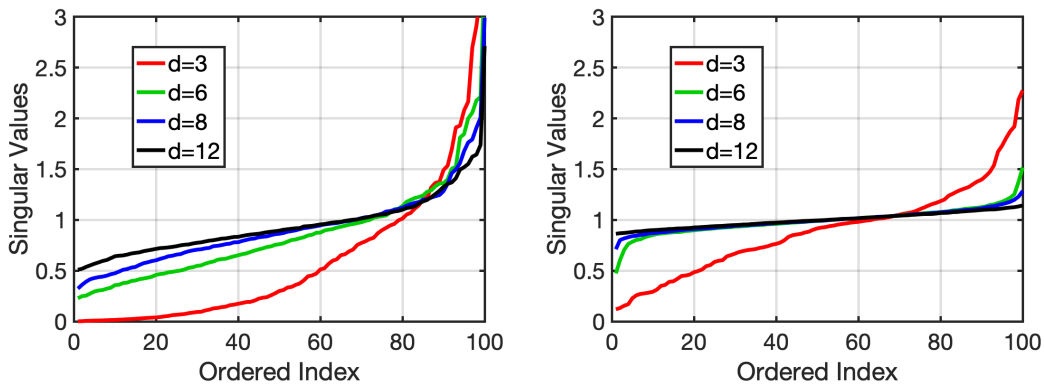


Figure 3.3: **Singular Value Distribution with Different d and σ** : The figure shows the distribution of the singular values in ascending as a function of the dimension d and the standard deviation of the random weights σ . The random feature matrix is the complex exponential with $m = 100$, $N = 5000$, and $\gamma = 1$. The plot on the left uses $\sigma = 2$ and the plot on the right uses $\sigma = 4$. The plots indicate that a small increase in σ has a dramatic effect on the concentration once d is above some value.

4 Theoretical Analysis

In this section, we present the key results for bounding the differences in (2.4) and (2.5). These will then lead to the proof of Theorems 2.1, 2.2, and 2.3. The main arguments are split into multiple theorems with shorter proofs for clarity and modularity. In all of our results we assume that the random feature matrix \mathbf{A} is defined as $\mathbf{A}_{j,k} = \exp(i\langle \mathbf{x}_j, \boldsymbol{\omega}_k \rangle)$.

4.1 Concentration with Separated Weights

Using the matrix Bernstein inequality (Lemma A.1 in the Appendix), we show that the Gram matrix is close to its expectation for feature weights that are well-separated.

Theorem 4.1 (Concentration with Separated Weights). *Let $\{\mathbf{x}_j\}_{j \in [m]} \subset \mathbb{R}^d$ be the data sampled from $\mathcal{N}(\mathbf{0}, \frac{\gamma^2}{d} \mathbf{I}_d)$. Suppose that $\{\boldsymbol{\omega}_k\}_{k \in [N]} \subset \mathbb{R}^d$ is a set of feature weights, and there is a constant $R > 0$ such that $\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2^2 \geq Rd$ for all $j, k \in [N]$ with $j \neq k$. If the following conditions hold*

$$m \geq C\eta^{-2}N \log\left(\frac{2N}{\delta}\right) \quad (4.1)$$

$$\gamma^2 \geq \frac{2}{R} \log\left(\frac{N}{\eta}\right) \quad (4.2)$$

for some $\delta, \eta \in (0, 1)$, where $C > 0$ is a universal constant. Then with probability at least $1 - \delta$ we have

$$\left\| \frac{1}{m} \mathbf{A}^* \mathbf{A} - \mathbb{E}_{\mathbf{x}} \left[\frac{1}{m} \mathbf{A}^* \mathbf{A} \right] \right\|_2 \leq \eta.$$

Proof. Let \mathbf{X}_ℓ be ℓ -th column of \mathbf{A}^* . Defined the random matrices $\{\mathbf{Y}_\ell\}_{\ell \in [m]}$ as

$$\mathbf{Y}_\ell = \mathbf{X}_\ell \mathbf{X}_\ell^* - \mathbb{E}_{\mathbf{x}}[\mathbf{X}_\ell \mathbf{X}_\ell^*].$$

Then $(\mathbf{Y}_\ell)_{j,j} = 0$ and $(\mathbf{Y}_\ell)_{j,k} = \exp(i\langle \mathbf{x}_\ell, \boldsymbol{\omega}_k - \boldsymbol{\omega}_j \rangle) - \exp(-\gamma^2 \|\boldsymbol{\omega}_k - \boldsymbol{\omega}_j\|_2^2 / (2d))$ for $j, k \in [N]$ with $j \neq k$. Note that \mathbf{Y}_ℓ is self-adjoint and its induced ℓ^2 norm is bounded by its largest eigenvalue. By Gershgorin's disk theorem and condition (4.2),

$$\|\mathbf{Y}_\ell\|_2 \leq \max_{j \in [N]} \sum_{k \neq j} \left| e^{i\langle \mathbf{x}_\ell, \boldsymbol{\omega}_k - \boldsymbol{\omega}_j \rangle} - e^{-\frac{\gamma^2}{2d} \|\boldsymbol{\omega}_k - \boldsymbol{\omega}_j\|_2^2} \right| \leq N \left(1 + e^{-\frac{\gamma^2}{2} R} \right) \leq N + \eta.$$

The variance parameter in Lemma A.1 is bounded by

$$\begin{aligned} \left\| \sum_{\ell=1}^m \mathbb{E}_{\mathbf{x}}[\mathbf{Y}_\ell^2] \right\|_2 &\leq \sum_{\ell=1}^m \|\mathbb{E}_{\mathbf{x}}[\mathbf{Y}_\ell^2]\|_2 = \sum_{\ell=1}^m \|N\mathbb{E}_{\mathbf{x}}[\mathbf{X}_\ell \mathbf{X}_\ell^*] - (\mathbb{E}_{\mathbf{x}}[\mathbf{X}_\ell \mathbf{X}_\ell^*])^2\|_2 \\ &\leq m[N(1 + \eta) + (1 + \eta)^2]. \end{aligned}$$

Here we use the fact that \mathbf{X}_ℓ is a vector with $\|\mathbf{X}_\ell\|_2 = \sqrt{N}$, which implies $\mathbf{X}_\ell \mathbf{X}_\ell^* \mathbf{X}_\ell \mathbf{X}_\ell^* = N \mathbf{X}_\ell \mathbf{X}_\ell^*$, and $\mathbb{E}_{\mathbf{x}}[\mathbf{X}_\ell \mathbf{X}_\ell^*]$ is self-adjoint whose ℓ^2 norm is bounded by

$$\|\mathbb{E}_{\mathbf{x}}[\mathbf{X}_\ell \mathbf{X}_\ell^*]\|_2 \leq 1 + \max_{j \in [N]} \sum_{k \neq j} \left| e^{-\frac{\gamma^2}{2d} \|\boldsymbol{\omega}_k - \boldsymbol{\omega}_j\|_2^2} \right| \leq 1 + N \exp\left(-\frac{\gamma^2}{2} R\right) \leq 1 + \eta,$$

by Gershgorin's disk theorem. Since $\{\mathbf{Y}_\ell\}_{\ell \in [m]}$ are independent mean-zero self-adjoint matrices,

applying Lemma A.1 with $K = N + \eta$ and $\sigma^2 = m[N(1 + \eta) + (1 + \eta)^2]$ then gives

$$\begin{aligned} \mathbb{P}\left(\left\|\frac{1}{m}\mathbf{A}^*\mathbf{A} - \mathbb{E}_{\mathbf{x}}\left[\frac{1}{m}\mathbf{A}^*\mathbf{A}\right]\right\|_2 \geq \eta\right) &= \mathbb{P}\left(\left\|\sum_{\ell=1}^m \mathbf{Y}_{\ell}\right\|_2 \geq m\eta\right) \\ &\leq 2N \exp\left(-\frac{m\eta^2/2}{N(1 + \eta) + (1 + \eta)^2 + (N + \eta)\eta/3}\right) \\ &\leq 2N \exp\left(-\frac{m\eta^2}{5N + 9}\right). \end{aligned}$$

The left-hand term is less than δ , provided condition (4.1) is satisfied with $C = 6$ (assuming that $N \geq 9$ and $\eta < 1$). This completes the proof. \square

4.2 Separation of Subgaussian Weights

Theorem 4.1 in the previous section requires that the weights $\{\boldsymbol{\omega}_k\}_{k \in [N]}$ are sufficiently separated, but does not place a restriction on the sampling process. Next, we show that if $\{\boldsymbol{\omega}_k\}_{k \in [N]}$ are sampled independently from a subgaussian distribution then they are separated with high probability. Recall that a random variable X is called subgaussian if there exist $\beta, \kappa > 0$ such that

$$\mathbb{P}(|X| \geq t) \leq \beta e^{-\kappa t^2} \quad \text{for all } t > 0.$$

We call a random vector $\mathbf{X} = (X_1, \dots, X_d) \in \mathbb{R}^d$ subgaussian if X_i are mean-zero independent subgaussian with the same subgaussian parameters. Using a concentration inequality for ℓ^2 norm of subgaussian vectors (Lemma A.2 in the Appendix), we have the following result.

Theorem 4.2 (Separation of Subgaussian Weights). *Suppose $\{\boldsymbol{\omega}_k\}_{k \in [N]} \subset \mathbb{R}^d$ is a set of random vectors such that the components of $\boldsymbol{\omega}_k$ are independent mean-zero subgaussian random variables with variance 1 and the same subgaussian parameters β, κ . If the dimension d satisfies the following condition*

$$d \geq Ct^{-2} \log\left(\frac{N}{\delta}\right),$$

for $\delta, t \in (0, 1)$, where $C > 0$ is a constant depends on the subgaussian parameters, then

$$\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2^2 \geq (2 - 2t)d$$

with probability at least $1 - \delta$.

Proof. We use a result for subgaussian matrices to estimate the squared distance between $\boldsymbol{\omega}_j$ and $\boldsymbol{\omega}_k$. Denote by \mathbf{W} the matrix which has $\boldsymbol{\omega}_j/\sqrt{d}$ as its j -th column. The s -th restricted isometry property (RIP) constant $\delta_s = \delta_s(\mathbf{W})$ of the matrix \mathbf{W} is the smallest $\Delta \geq 0$ such that

$$(1 - \Delta)\|\mathbf{x}\|_2^2 \leq \|\mathbf{W}\mathbf{x}\|_2^2 \leq (1 + \Delta)\|\mathbf{x}\|_2^2$$

holds for all s -sparse vector \mathbf{x} , i.e. \mathbf{x} which has at most s nonzero elements. Theorem 9.2 from [11] with $s = 2$ shows that the RIP constant δ_2 of $\mathbf{W} \in \mathbb{R}^{d \times N}$ is less than t with probability at least $1 - \delta$ if

$$d \geq Ct^{-2} \log\left(\frac{N}{\delta}\right), \tag{4.3}$$

for some $C > 0$ which depends only on the subgaussian parameters. Note that the RIP constant δ_2 can also be defined as

$$\delta_2 := \max_{j,k \in [N], j \neq k} \|\mathbf{W}_{\{j,k\}}^* \mathbf{W}_{\{j,k\}} - \mathbf{I}_2\|_2,$$

where $\mathbf{W}_{\{j,k\}}$ is the submatrix of \mathbf{W} consisting of j -th and k -th column. For any $j, k \in [N], j \neq k$, the matrix $\mathbf{W}_{\{j,k\}}^* \mathbf{W}_{\{j,k\}} - \mathbf{I}_2$ takes the form

$$\mathbf{W}_{\{j,k\}}^* \mathbf{W}_{\{j,k\}} - \mathbf{I}_2 = \frac{1}{d} \begin{bmatrix} \|\boldsymbol{\omega}_j\|_2^2 - d & \langle \boldsymbol{\omega}_k, \boldsymbol{\omega}_j \rangle \\ \langle \boldsymbol{\omega}_j, \boldsymbol{\omega}_k \rangle & \|\boldsymbol{\omega}_k\|_2^2 - d \end{bmatrix}.$$

This is a symmetric matrix and its eigenvalues are

$$\lambda^\pm = \frac{(\|\boldsymbol{\omega}_j\|_2^2 + \|\boldsymbol{\omega}_k\|_2^2 - 2d) \pm \sqrt{(\|\boldsymbol{\omega}_j\|_2^2 - \|\boldsymbol{\omega}_k\|_2^2)^2 + 4|\langle \boldsymbol{\omega}_j, \boldsymbol{\omega}_k \rangle|^2}}{2d}.$$

Therefore, $\delta_2 \leq t$ implies that both eigenvalues are in $[-t, t]$ and consequently,

$$\frac{(\|\boldsymbol{\omega}_j\|_2^2 + \|\boldsymbol{\omega}_k\|_2^2 - 2d) - 2|\langle \boldsymbol{\omega}_j, \boldsymbol{\omega}_k \rangle|}{2d} \geq \lambda^- \geq -t$$

Thus, $\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2^2 = \|\boldsymbol{\omega}_j\|_2^2 + \|\boldsymbol{\omega}_k\|_2^2 - 2\langle \boldsymbol{\omega}_j, \boldsymbol{\omega}_k \rangle \geq (2 - 2t)d$ for all $j, k \in [N], j \neq k$ with probability at least $1 - \delta$ if condition (4.3) is satisfied. \square

Theorem 4.2 holds for arbitrary mean-zero subgaussian distribution. In particular, if $\boldsymbol{\omega}_k$ are sampled from a mean-zero Gaussian distribution, they will be separated with high probability.

Corollary 4.3 (Separation of Gaussian Weights). *Suppose that $\{\boldsymbol{\omega}_k\}_{k \in [N]} \subset \mathbb{R}^d$ are sampled from $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_d)$. If the dimension d satisfies*

$$d \geq Ct^{-2} \log \left(\frac{N}{\delta} \right),$$

for $\delta, t \in (0, 1)$, where $C > 0$ is a universal constant. Then we have

$$\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2^2 \geq (2 - 2t)\sigma^2 d \quad \text{for all } j, k \in [N], j \neq k$$

with probability at least $1 - \delta$.

4.3 Concentration of Random Feature Matrices with Subgaussian Weights

In the proof of Theorem 4.1, we use the condition $\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2^2 \geq Rd$ to obtain an estimate of $\|\mathbb{E}_{\mathbf{x}}[\mathbf{X}_\ell \mathbf{X}_\ell^*]\|_2$. Note that this is a self-adjoint matrix with ones along the diagonal and its off diagonal terms are relatively small when $\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2$ are sufficiently large. Therefore, if $\{\boldsymbol{\omega}_k\}_{k \in [N]}$ are sampled randomly such that each component of $\boldsymbol{\omega}_k$ is subgaussian with variance σ^2 , then we would expect $\mathbb{E}_{\mathbf{x}}[\mathbf{A}^* \mathbf{A}/m]$ to be close to identity \mathbf{I}_N .

Theorem 4.4 (Concentration with Subgaussian Weights). *Suppose that $\{\mathbf{x}_j\}_{j \in [m]} \subset \mathbb{R}^d$ are data points sampled from $\mathcal{N}(\mathbf{0}, \frac{\gamma^2}{d} \mathbf{I}_d)$ and $\{\boldsymbol{\omega}_k\}_{k \in [N]} \subset \mathbb{R}^d$ are feature weights such that components of $\boldsymbol{\omega}_k$ are independent mean-zero subgaussian random variables with variance σ^2 and the same*

subgaussian parameters β, κ . Then there exists a constant $C > 0$ (depending only on subgaussian parameters) such that if the following conditions hold:

$$d \geq C \log \left(\frac{N}{\delta} \right) \quad (4.4)$$

$$\gamma^2 \sigma^2 \geq 4 \log \left(\frac{2N}{\eta} \right) \quad (4.5)$$

for some $\delta, \eta \in (0, 1)$, we have

$$\left\| \mathbb{E}_{\mathbf{x}} \left[\frac{1}{m} \mathbf{A}^* \mathbf{A} \right] - \mathbf{I}_N \right\|_2 \leq \eta$$

with probability at least $1 - 2\delta$. Moreover, if $\eta \geq 2\delta$ (which holds for practical η and δ), then we simultaneously have

$$\left\| \mathbb{E}_{\mathbf{x}} \left[\frac{1}{m} \mathbf{A}^* \mathbf{A} \right] - \mathbb{E}_{\mathbf{x}, \omega} \left[\frac{1}{m} \mathbf{A}^* \mathbf{A} \right] \right\|_2 \leq \eta.$$

Proof. Denote by \mathbf{B} the matrix $m^{-1} \mathbb{E}_{\mathbf{x}}[\mathbf{A}^* \mathbf{A}] - \mathbf{I}_N$. Then \mathbf{B} is symmetric and $(\mathbf{B})_{j,j} = 0$, $(\mathbf{B})_{j,k} = \exp(-\gamma^2 \|\omega_j - \omega_k\|_2^2 / (2d))$ for all $j, k \in [N]$ with $j \neq k$. By Theorem 4.2 with $t = 3/4$, for all $j, k \in [N]$ with $j \neq k$ we have

$$\|\omega_j - \omega_k\|_2^2 \geq \frac{\sigma^2}{2} d,$$

with probability at least $1 - \delta$ if the dimension d satisfies

$$d \geq C \log \left(\frac{N}{\delta} \right). \quad (4.6)$$

Thus, each off diagonal elements is bounded (in magnitude) by

$$|(\mathbf{B})_{j,k}| = \exp \left(-\frac{\gamma^2}{2d} \|\omega_j - \omega_k\|_2^2 \right) \leq \exp \left(-\frac{\gamma^2 \sigma^2}{4} \right) \quad \text{for all } j, k \in [N], j \neq k,$$

with probability at least $1 - \delta$. By Gershgorin disk theorem and condition (4.5), the induced ℓ^2 norm of \mathbf{B} is bounded by

$$\|\mathbf{B}\|_2 \leq \max_{j \in [N]} \sum_{k \neq j} |(\mathbf{B})_{j,k}| \leq N \exp \left(-\frac{\gamma^2 \sigma^2}{4} \right) \leq \eta,$$

with probability at least $1 - \delta$.

Next, we denote by \mathbf{C} the matrix $m^{-1} \mathbb{E}_{\mathbf{x}}[\mathbf{A}^* \mathbf{A}] - m^{-1} \mathbb{E}_{\mathbf{x}, \omega}[\mathbf{A}^* \mathbf{A}]$. Then \mathbf{C} is also symmetric and $(\mathbf{C})_{j,j} = 0$, $(\mathbf{C})_{j,k} = \exp(-\gamma^2 \|\omega_j - \omega_k\|_2^2 / (2d)) - \mathbb{E}_{\omega}[\exp(-\gamma^2 \|\omega_j - \omega_k\|_2^2 / (2d))]$ for all $j, k \in [N]$ with $j \neq k$. The previous argument shows that the term $\exp(-\gamma^2 \|\omega_j - \omega_k\|_2^2 / (2d))$ is small. Thus, we only need to estimate the expectation, i.e. $\mathbb{E}_{\omega}[\exp(-\gamma^2 \|\omega_j - \omega_k\|_2^2 / (2d))]$. Since ω_j and ω_k are

independent subgaussian vectors, $\boldsymbol{\omega}_j - \boldsymbol{\omega}_k$ is also a subgaussian vector with new parameters that depend on β and κ . Applying Lemma A.2 yields

$$\mathbb{P}(\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2^2 \leq z\sigma^2 d) \leq \exp\left(-C\left(1 - \frac{\sqrt{z}}{2}\right)^2 d\right)$$

for some constant $C > 0$ which depends on the subgaussian parameters. Then by setting $z = 1/2$ and decomposing the expectation (where χ is the characteristic function), we have

$$\begin{aligned} & \mathbb{E}\left[\exp\left(-\frac{\gamma^2}{2d}\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2^2\right)\right] \\ & \leq \mathbb{E}\left[\exp\left(-\frac{\gamma^2}{2d}\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2^2\right) \chi_{\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2^2 \leq \frac{\sigma^2 d}{2}}\right] + \mathbb{E}\left[\exp\left(-\frac{\gamma^2}{2d}\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2^2\right) \chi_{\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2^2 \geq \frac{\sigma^2 d}{2}}\right] \\ & \leq \mathbb{P}\left(\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2^2 \leq \frac{\sigma^2 d}{2}\right) + \exp\left(-\frac{\gamma^2 \sigma^2}{4}\right) \\ & \leq \exp(-Cd) + \exp\left(-\frac{\gamma^2 \sigma^2}{4}\right), \end{aligned}$$

for some redefined constant $C > 0$. Therefore, if

$$d \geq C \log\left(\frac{2N}{\eta}\right) \tag{4.7}$$

$$\gamma^2 \sigma^2 \geq 4 \log\left(\frac{2N}{\eta}\right), \tag{4.8}$$

then we have the bound

$$\begin{aligned} \|\mathbf{C}\|_2 & \leq \max_{j \in [N]} \sum_{k \neq j} |(\mathbf{C})_{j,k}| \\ & \leq N \max\left\{\exp\left(-\frac{\gamma^2}{2d}\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2^2\right), \mathbb{E}_{\boldsymbol{\omega}}\left[\exp\left(-\frac{\gamma^2}{2d}\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2^2\right)\right]\right\} \leq \eta, \end{aligned}$$

with probability at least $1 - \delta$.

Note that condition (4.7) and condition (4.6) are slightly different. However, by assuming $\eta \geq 2\delta$, we can combine (4.7) and (4.6) to obtain (4.4). \square

Remark 4.5. If $\boldsymbol{\omega}_k \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$, then the expectation $\mathbb{E}_{\boldsymbol{\omega}}[\exp(-\gamma^2 \|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2^2 / (2d))]$ can be computed explicitly

$$\mathbb{E}_{\boldsymbol{\omega}}\left[\exp\left(-\frac{\gamma^2}{2d}\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2^2\right)\right] = \left(\frac{2\gamma^2 \sigma^2}{d} + 1\right)^{-\frac{d}{2}},$$

which is approximately $\exp(-\gamma^2 \sigma^2)$ for large d . Thus, in this case, the product of γ and σ only needs to satisfy

$$\gamma^2 \sigma^2 \geq \log\left(\frac{N}{\eta}\right).$$

4.4 Proof of Main Results

Using the previous results, we can prove Theorem 2.1.

Proof of Theorem 2.1. By Theorem 4.2 with $t = 3/4$, for weights $\{\boldsymbol{\omega}_k\}_{k \in [N]}$ whose components are independent mean-zero subgaussian random variables with variance σ^2 , we have

$$\|\boldsymbol{\omega}_j - \boldsymbol{\omega}_k\|_2 \geq \frac{\sigma^2 d}{2},$$

for all $j, k \in [N]$ with $j \neq k$ and with probability at least $1 - \delta$ if

$$d \geq C_1 \log \left(\frac{N}{\delta} \right),$$

for some $C_1 > 0$ which depends only on subgaussian parameters. Then for $\eta > 0$, by Theorem 4.1, we have

$$\left\| \frac{1}{m} \mathbf{A}^* \mathbf{A} - \mathbb{E}_{\mathbf{x}} \left[\frac{1}{m} \mathbf{A}^* \mathbf{A} \right] \right\|_2 \leq \eta,$$

with probability at least $1 - 2\delta$ if

$$\begin{aligned} m &\geq C_2 \eta^{-2} N \log \left(\frac{2N}{\delta} \right) \\ \gamma^2 \sigma^2 &\geq 4 \log \left(\frac{2N}{\eta} \right), \end{aligned}$$

for some $C_2 > 0$. Lastly, by Theorem 4.4, we also have

$$\left\| \mathbb{E}_{\mathbf{x}} \left[\frac{1}{m} \mathbf{A}^* \mathbf{A} \right] - \mathbf{I}_N \right\|_2 \leq \eta$$

with probability at least $1 - \delta$ if conditions (2.1) and (2.2) are satisfied. Therefore, the difference in (2.4) is bounded by 2η through (2.6). If $\eta \geq 2\delta$, then we also have

$$\left\| \mathbb{E}_{\mathbf{x}} \left[\frac{1}{m} \mathbf{A}^* \mathbf{A} \right] - \mathbb{E}_{\mathbf{x}, \boldsymbol{\omega}} \left[\frac{1}{m} \mathbf{A}^* \mathbf{A} \right] \right\|_2 \leq \eta$$

and the difference in (2.5) is bounded by 2η through (2.7), with probability at $1 - 3\delta$. \square

To prove Theorem 2.2, consider the matrix \mathbf{A}^* which just switches the role of \mathbf{x} and $\boldsymbol{\omega}$ in the theorems. The proof of Theorem 2.3 is the same as Theorem 4.1 with \mathbf{x} and $\boldsymbol{\omega}$ switched and by removing the dimensional scaling in the sampling of $\boldsymbol{\omega}$ and the separation of the data samples (for consistency).

Note that in [6], the union bound was used to estimate $\|\mathbb{E}_{\mathbf{x}}[\mathbf{X}_\ell \mathbf{X}_\ell^*] - \mathbb{E}_{\mathbf{x}, \boldsymbol{\omega}}[\mathbf{X}_\ell \mathbf{X}_\ell^*]\|_2$ which led to a condition for N that depends algebraically on the probability δ . By showing that i.i.d. subgaussian vectors are well-separated in high dimensions, we obtain a uniform bound which does not depend on N for all entries in the matrix. Thus, the probability δ does not restrict N when d is large.

5 Summary

The main results show that the spectrum of an asymmetric rectangular (nonlinear) random matrix whose entries are of the form $\mathbf{A}_{j,k} = \phi(\langle \mathbf{x}_j, \boldsymbol{\omega}_k \rangle)$ concentrates around its expectation and around 1 given particular (finite) complexity scales. We showed that this holds in the setting where both variables are random (i.e. one is normal and the other is subgaussian) and in the setting where one is a random normal variable and the other is well-separated. The conditions in the theorems relax as the dimension of the input data increases, thus showing that high-dimensional random feature matrices are well-conditioned. In addition, in the case of subgaussian weights (or subgaussian data), we do not require that the weights (or data) follow the same distribution, as long as they have the same subgaussian parameters. This generalizes and extends the results beyond that of previous ones. The results are presented in separate parts, since they may be used for the analysis of other random feature models. For example, using these techniques, one may also be able to find the dependency of the restricted isometry property of random feature matrices on the dimension, which is useful for analyzing sparsity-based approaches for random features models [13, 30, 32, 36, 37].

Acknowledgement

Z.C. and H.S. were supported in part by AFOSR MURI FA9550-21-1-0084 and NSF DMS-1752116. R.W. was supported in part by AFOSR MURI FA9550-19-1-0005, NSF DMS 1952735, NSF HDR-1934932, and NSF 2019844.

References

- [1] 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 *International conference on machine learning*, pages 253–262. PMLR, 2017.
- [2] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.
- [3] Mikhail Belkin, Daniel Hsu, and Ji Xu. Two models of double descent for weak features. *SIAM Journal on Mathematics of Data Science*, 2(4):1167–1180, 2020.
- [4] Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to understand kernel learning. In *International Conference on Machine Learning*, pages 541–549. PMLR, 2018.
- [5] Lucas Benigni and Sandrine Péché. Eigenvalue distribution of nonlinear models of random matrices. *arXiv preprint arXiv:1904.03090*, 2019.
- [6] Zhijun Chen and Hayden Schaeffer. Conditioning of random feature matrices: Double descent and generalization error. *arXiv preprint arXiv:2110.11477*, 2021.

- [7] Xiuyuan Cheng and Amit Singer. The spectrum of random inner-product kernel matrices. *Random Matrices: Theory and Applications*, 2(04):1350010, 2013.
- [8] Nouredine El Karoui. The spectrum of kernel random matrices. *The Annals of Statistics*, 38(1):1–50, 2010.
- [9] Zhou Fan and Andrea Montanari. The spectral norm of random inner-product kernel matrices. *Probability Theory and Related Fields*, 173(1):27–85, 2019.
- [10] Zhou Fan and Zhichao Wang. Spectra of the conjugate kernel and neural tangent kernel for linear-width neural networks. *Advances in neural information processing systems*, 33:7710–7721, 2020.
- [11] Simon Foucart and Holger Rauhut. *A Mathematical Introduction to Compressive Sensing*. Birkhäuser Basel, 2013.
- [12] Anna Gilbert, Albert Gu, Christopher Ré, Atri Rudra, and Mary Wootters. Sparse recovery for orthogonal polynomial transforms. *arXiv preprint arXiv:1907.08362*, 2019.
- [13] Abolfazl Hashemi, Hayden Schaeffer, Robert Shi, Ufuk Topcu, Giang Tran, and Rachel Ward. Generalization bounds for sparse random feature expansions. *arXiv preprint arXiv:2103.03191*, 2021.
- [14] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. *arXiv preprint arXiv:1903.08560*, 2019.
- [15] Zhu Li, Jean-Francois Ton, Dino Oglic, and Dino Sejdinovic. Towards a unified analysis of random fourier features. In *International conference on machine learning*, pages 3905–3914. PMLR, 2019.
- [16] Zhenyu Liao, Romain Couillet, and Michael W Mahoney. A random matrix analysis of random fourier features: beyond the gaussian kernel, a precise phase transition, and the corresponding double descent. *Advances in Neural Information Processing Systems*, 33:13939–13950, 2020.
- [17] Fanghui Liu, Xiaolin Huang, Yudong Chen, and Johan AK Suykens. Random features for kernel approximation: A survey on algorithms, theory, and beyond. *arXiv preprint arXiv:2004.11154*, 2020.
- [18] Cosme Louart, Zhenyu Liao, and Romain Couillet. A random matrix approach to neural networks. *The Annals of Applied Probability*, 28(2):1190–1248, 2018.
- [19] Song Mei and Andrea Montanari. The generalization error of random features regression: Precise asymptotics and the double descent curve. *Communications on Pure and Applied Mathematics*, 2019.
- [20] Kevin P Murphy. *Machine learning: a probabilistic perspective*. MIT press, 2012.
- [21] Roberto Imbuzeiro Oliveira. Concentration of the adjacency matrix and of the laplacian in random graphs with independent edges. *arXiv preprint arXiv:0911.0600*, 2009.

- [22] Ayça Özçelikkale. Sparse recovery with non-linear fourier features. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5715–5719. IEEE, 2020.
- [23] L Pastur and V Slavin. On random matrices arising in deep neural networks: General iid case. *arXiv preprint arXiv:2011.11439*, 2020.
- [24] Leonid Pastur. On random matrices arising in deep neural networks. gaussian case. *arXiv preprint arXiv:2001.06188*, 2020.
- [25] Jeffrey Pennington and Pratik Worah. Nonlinear random matrix theory for deep learning. *Advances in neural information processing systems*, 30, 2017.
- [26] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. *Advances in neural information processing systems*, 20, 2007.
- [27] Ali Rahimi and Benjamin Recht. Uniform approximation of functions with random bases. In *2008 46th Annual Allerton Conference on Communication, Control, and Computing*, pages 555–561. IEEE, 2008.
- [28] Ali Rahimi and Benjamin Recht. Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning. *Advances in neural information processing systems*, 21, 2008.
- [29] Holger Rauhut and Rachel Ward. Interpolation via weighted ℓ_1 minimization. *Applied and Computational Harmonic Analysis*, 40(2):321–351, 2016.
- [30] Nicholas Richardson, Hayden Schaeffer, and Giang Tran. SRMD: Sparse random mode decomposition. *arXiv preprint arXiv:2204.06108*, 2022.
- [31] Alessandro Rudi and Lorenzo Rosasco. Generalization properties of learning with random features. *Advances in neural information processing systems*, 30, 2017.
- [32] Esha Saha, Hayden Schaeffer, and Giang Tran. Harfe: Hard-ridge random feature expansion. *arXiv preprint arXiv:2202.02877*, 2022.
- [33] Joel A Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of computational mathematics*, 12(4):389–434, 2012.
- [34] Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2018.
- [35] Zhichao Wang and Yizhe Zhu. Deformed semicircle law and concentration of nonlinear random matrices for ultra-wide neural networks. *arXiv preprint arXiv:2109.09304*, 2021.
- [36] Yuege Xie, Bobby Shi, Hayden Schaeffer, and Rachel Ward. Shrimp: Sparser random feature models via iterative magnitude pruning. *arXiv preprint arXiv:2112.04002*, 2021.
- [37] Ian En-Hsu Yen, Ting-Wei Lin, Shou-De Lin, Pradeep K Ravikumar, and Inderjit S Dhillon. Sparse random feature algorithm as coordinate descent in hilbert space. *Advances in Neural Information Processing Systems*, 27, 2014.

A Useful Lemmata

We recall the matrix Bernstein's inequality from [11, 21, 33], which is used to prove the main results.

Lemma A.1 (Matrix Bernstein's inequality). *Let $\{\mathbf{X}_j\}_{j \in [m]} \subset \mathbb{C}^{N \times N}$ be independent mean-zero self-adjoint random matrices. Assume that*

$$\|\mathbf{X}_j\|_2 \leq K \quad \text{a.s. for all } j \in [m],$$

and set

$$\sigma^2 := \left\| \sum_{j=1}^m \mathbb{E}[\mathbf{X}_j^2] \right\|_2.$$

Then for $t > 0$,

$$\mathbb{P} \left(\left\| \sum_{j=1}^m \mathbf{X}_j \right\|_2 \geq t \right) \leq 2N \exp \left(-\frac{t^2/2}{\sigma^2 + Kt/3} \right).$$

For a random vector with subgaussian components, we have the following concentration inequality for its ℓ^2 norm.

Lemma A.2 (Theorem 3.1.1 in [34]). *Let $\mathbf{X} = (X_1, \dots, X_d) \in \mathbb{R}^d$ be a random vector with independent X_i . Suppose that $\{X_j\}_{j \in [d]}$ are mean-zero subgaussian random variables with variance 1 and the same subgaussian parameters β, κ . Then there exists a constant $C > 0$ which depends on subgaussian parameters such that*

$$\mathbb{P} \left(\left| \|\mathbf{X}\|_2 - \sqrt{d} \right| \geq t \right) \leq 2 \exp(-Ct^2).$$