Towards a Unified Analysis of Random Fourier Features

Zhu Li Zhu.li@stats.ox.ac.uk

Department of Statistics University of Oxford Oxford, OX1 3LB, UK

Jean-Francois Ton JEAN-FRANCOIS.TON@SPC.OX.AC.UK

Department of Statistics University of Oxford Oxford, OX1 3LB, UK

Dino Oglic DINO.OGLIC@KCL.AC.UK

Department of Informatics King's College London London, WC2R 2LS, UK

Dino Sejdinovic DINO.SEJDINOVIC@STATS.OX.AC.UK

Department of Statistics University of Oxford Oxford, OX1 3LB, UK

Editor:

Abstract

Random Fourier features is a widely used, simple, and effective technique for scaling up kernel methods. The existing theoretical analysis of the approach, however, remains focused on specific learning tasks and typically gives pessimistic bounds which are at odds with the empirical results. We tackle these problems and provide the first unified risk analysis of learning with random Fourier features using the squared error and Lipschitz continuous loss functions. In our bounds, the trade-off between the computational cost and the expected risk convergence rate is problem specific and expressed in terms of the regularization parameter and the *number of effective degrees of freedom*. We study both the standard random Fourier features method for which we improve the existing bounds on the number of features required to guarantee the corresponding minimax risk convergence rate of kernel ridge regression, as well as a data-dependent modification which samples features proportional to *ridge leverage scores* and further reduces the required number of features. As ridge leverage scores are expensive to compute, we devise a simple approximation scheme which provably reduces the computational cost without loss of statistical efficiency.

Keywords: Kernel methods, random Fourier features, stationary kernels, kernel ridge regression, Lipschitz continuous loss, support vector machines, logistic regression, ridge leverage scores.

1. Introduction

Kernel methods are one of the pillars of machine learning (Schölkopf and Smola, 2001; Schölkopf et al., 2004), as they give us a flexible framework to model complex functional relationships in a principled way and also come with well-established statistical properties and theoretical guarantees (Caponnetto and De Vito, 2007; Steinwart and Christmann, 2008). The key ingredient, known as

kernel trick, allows implicit computation of an inner product between rich feature representations of data through the kernel evaluation $k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{H}}$, while the actual feature mapping $\varphi: \mathcal{X} \to \mathcal{H}$ between a data domain \mathcal{X} and some high and often infinite dimensional Hilbert space \mathcal{H} is never computed. However, such convenience comes at a price: due to operating on all pairs of observations, kernel methods inherently require computation and storage which is at least quadratic in the number of observations, and hence often prohibitive for large datasets. In particular, the kernel matrix has to be computed, stored, and often inverted. As a result, a flurry of research into scalable kernel methods and the analysis of their performance emerged (Rahimi and Recht, 2007; Mahoney and Drineas, 2009; Bach, 2013; Alaoui and Mahoney, 2015; Rudi et al., 2015; Rudi and Rosasco, 2017; Rudi et al., 2017; Zhang et al., 2015). Among the most popular frameworks for fast approximations to kernel methods are random Fourier features (RFF) due to Rahimi and Recht (2007). The idea of random Fourier features is to construct an explicit feature map which is of a dimension much lower than the number of observations, but with the resulting inner product which approximates the desired kernel function k(x,y). In particular, random Fourier features rely on Bochner's theorem (Bochner, 1932; Rudin, 2017) which tells us that any bounded, continuous and shift-invariant kernel is a Fourier transform of a bounded positive measure, called spectral measure. The feature map is then constructed using samples drawn from the spectral measure. Essentially, any kernel method can then be adjusted to operate on these explicit feature maps (i.e., primal representations), greatly reducing the computational and storage costs, while in practice mimicking performance of the original kernel method.

Despite their empirical success, the theoretical understanding of statistical properties of random Fourier features is incomplete, and the question of how many features are needed, in order to obtain a method with performance provably comparable to the original one, remains without a definitive answer. Currently, there are two main lines of research addressing this question. The first line considers the approximation error of the kernel matrix itself (e.g., see Rahimi and Recht, 2007; Sriperumbudur and Szabó, 2015; Sutherland and Schneider, 2015, and references therein) and bases performance guarantees on the accuracy of this approximation. However, all of these works require $\Omega(n)$ features (n being the number of observations), which translates to no computational savings at all and is at odds with empirical findings. Realizing that the approximation of kernel matrices is just a means to an end, the second line of research aims at directly studying the risk and generalization properties of random Fourier features in various supervised learning scenarios. Arguably, first such result is already in Rahimi and Recht (2009), where supervised learning with Lipschitz continuous loss functions is studied. However, the bounds therein still require a pessimistic $\Omega(n)$ number of features and due to the Lipschitz continuity requirement, the analysis does not apply to kernel ridge regression (KRR), one of the most commonly used kernel methods. In Bach (2017b), the generalization properties are studied from a function approximation perspective, showing for the first time that fewer features could preserve the statistical properties of the original method, but in the case where a certain data-dependent sampling distribution is used instead of the spectral measure. These results also do not apply to kernel ridge regression and the mentioned sampling distribution is typically itself intractable. Avron et al. (2017) study the empirical risk of kernel ridge regression and show that it is possible to use o(n) features and have the empirical risk of the linear ridge regression estimator based on random Fourier features close to the empirical risk of the original kernel estimator, also relying on a modification to the sampling distribution. However, this result is for the empirical risk only, does not provide any expected risk convergence rates, and a

tractable method to sample from a modified distribution is proposed for the Gaussian kernel only. A highly refined analysis of kernel ridge regression is given by Rudi and Rosasco (2017), where it is shown that $\Omega(\sqrt{n}\log n)$ features suffices for an optimal $O(1/\sqrt{n})$ learning error in a minimax sense (Caponnetto and De Vito, 2007). Moreover, the number of features can be reduced even further if a data-dependent sampling distribution is employed. While these are groundbreaking results, guaranteeing computational savings without any loss of statistical efficiency, they require some technical assumptions that are difficult to verify. Moreover, to what extent the bounds can be improved by utilizing data-dependent distributions still remains unclear. Finally, it does not seem straightforward to generalize the approach of Rudi and Rosasco (2017) to kernel support vector machines (SVM) and/or kernel logistic regression (KLR). Recently, Sun et al. (2018) have provided novel bounds for random Fourier features in the SVM setting, assuming the Massart's low noise condition and that the target hypothesis lies in the corresponding reproducing kernel Hilbert space. The bounds, however, require the sample complexity and the number of features to be exponential in the dimension of the instance space and this can be problematic for high dimensional instance spaces. The theoretical results are also restricted to the hinge loss (without means to generalize to other loss functions) and require optimized features.

In this paper, we address the gaps mentioned above by making the following contributions:

- We devise a simple framework for the unified analysis of generalization properties of random Fourier features, which applies to kernel ridge regression, as well as to kernel support vector machines and logistic regression.
- For the plain random Fourier features sampling scheme, we provide, to the best of our knowledge, the sharpest results on the number of features required. In particular, we show that already with $\Omega(\sqrt{n}\log d_{\mathbf{K}}^{\lambda})$ features, we incur no loss of learning accuracy in kernel ridge regression, where $d_{\mathbf{K}}^{\lambda}$ corresponds to the notion of the number of effective degrees of freedom (Bach, 2013) with $d_{\mathbf{K}}^{\lambda} \ll n$ and $\lambda \coloneqq \lambda(n)$ is the regularization parameter. In addition, $\Omega(1/\lambda)$ features is sufficient to ensure $O(\sqrt{\lambda})$ expected risk rate in kernel support vector machines and kernel logistic regression.
- In the case of a modified data-dependent sampling distribution, the so called *empirical ridge* leverage score distribution, we demonstrate that $\Omega(d_{\mathbf{K}}^{\lambda})$ features suffice for the learning risk to converge at $O(\lambda)$ rate in kernel ridge regression $(O(\sqrt{\lambda}))$ in kernel support vector machines and kernel logistic regression).
- In our refined analysis of kernel ridge regression, we show that the excess risk convergence rate of the estimator based on random Fourier features can (depending on the decay rate of the spectrum of kernel function) be upper bounded by $O(\frac{\log n}{n})$ or even $O(\frac{1}{n})$, which implies much faster convergence than $O(\frac{1}{\sqrt{n}})$ rate featuring in most of previous bounds.
- Similarly, our refined analysis for Lipschitz continuous loss demonstrates that under a realizable case (defined subsequently) one could achieve $O(\frac{\log n}{\sqrt{n}})$ excess risk convergence rate with only $O(\sqrt{n})$ features. To the best of our knowledge, this is the first result offering non-trivial computational savings for approximations in problems with Lipschitz loss functions.
- Finally, as the empirical ridge leverage scores distribution is typically costly to compute, we
 give a fast algorithm to generate samples from the approximated empirical leverage distribution.

Utilizing these samples one can significantly reduce the computation time during the in sample prediction and testing stages, O(n) and $O(\log n \log \log n)$, respectively. We also include a proof that gives a trade-off between the computational cost and the expected risk of the algorithm, showing that the statistical efficiency can be preserved while provably reducing the required computational cost.

2. Background

In this section, we provide some notation and preliminary results that will be used throughout the paper. Henceforth, we denote the Euclidean norm of a vector $a \in \mathbb{R}^n$ with $\|a\|_2$ and the operator norm of a matrix $A \in \mathbb{R}^{n_1 \times n_2}$ with $\|A\|_2$. Furthermore, we denote with $\|A\|_F$ the Frobenious norm of a matrix or operator A. Let \mathcal{H} be a Hilbert space with $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ as its inner product and $\|\cdot\|_{\mathcal{H}}$ as its norm. We use $\mathrm{Tr}(\cdot)$ to denote the trace of an operator or a matrix. Given a measure $d\rho$, we use $L_2(d\rho)$ to denote the space of square-integrable functions with respect to $d\rho$.

2.1 Random Fourier Features

Random Fourier features is a widely used, simple, and effective technique for scaling up kernel methods. The underlying principle of the approach is a consequence of Bochner's theorem (Bochner, 1932), which states that any bounded, continuous and shift-invariant kernel is a Fourier transform of a bounded positive measure. This measure can be transformed/normalized into a probability measure which is typically called the spectral measure of the kernel. Assuming the spectral measure $d\tau$ has a density function $p(\cdot)$, the corresponding shift-invariant kernel can be written as

$$k(x,y) = \int_{\mathcal{V}} e^{-2\pi i v^T (x-y)} d\tau(v) = \int_{\mathcal{V}} \left(e^{-2\pi i v^T x} \right) \left(e^{-2\pi i v^T y} \right)^* p(v) dv, \tag{1}$$

where c^* denotes the complex conjugate of $c \in \mathbb{C}$. Typically, the kernel is real valued and we can ignore the imaginary part in this equation (e.g., see Rahimi and Recht, 2007). The principle can be further generalized by considering the class of kernel functions which can be decomposed as

$$k(x,y) = \int_{\mathcal{V}} z(v,x)z(v,y)p(v)dv,$$
(2)

where $z \colon \mathcal{V} \times \mathcal{X} \to \mathbb{R}$ is a continuous and bounded function with respect to v and x. The main idea behind random Fourier features is to approximate the kernel function by its Monte-Carlo estimate

$$\tilde{k}(x,y) = \frac{1}{s} \sum_{i=1}^{s} z(v_i, x) z(v_i, y),$$
(3)

with reproducing kernel Hilbert space $\tilde{\mathcal{H}}$ (not necessarily contained in the reproducing kernel Hilbert space \mathcal{H} corresponding to the kernel function k) and $\{v_i\}_{i=1}^s$ sampled independently from the spectral measure. In Bach (2017a, Appendix A), it has been established that a function $f \in \mathcal{H}$ can be expressed as 1 :

$$f(x) = \int_{\mathcal{V}} g(v)z(v,x)p(v)dv \qquad (\forall x \in \mathcal{X})$$
 (4)

^{1.} It is not necessarily true that for any $q \in L_2(d\tau)$, there exists a corresponding $f \in \mathcal{H}$.

where $g \in L_2(d\tau)$ is a real-valued function such that $\|g\|_{L_2(d\tau)}^2 < \infty$ and $\|f\|_{\mathcal{H}}$ is equal to the minimum of $\|g\|_{L_2(d\tau)}$, over all possible decompositions of f. Thus, one can take an independent sample $\{v_i\}_{i=1}^s \sim p(v)$ (we refer to this sampling scheme as *plain RFF*) and approximate a function $f \in \mathcal{H}$ at a point $x_j \in \mathcal{X}$ by

$$\tilde{f}(x_j) = \sum_{i=1}^s \alpha_i z(v_i, x_j) := \mathbf{z}_{x_j}(\mathbf{v})^T \alpha \quad \text{with} \quad \alpha \in \mathbb{R}^s.$$

In standard estimation problems, it is typically the case that for a given set of instances $\{x_i\}_{i=1}^n$ one approximates $\mathbf{f}_x = [f(x_1), \cdots, f(x_n)]^T$ by

$$\tilde{\mathbf{f}}_x = [\mathbf{z}_{x_1}(\mathbf{v})^T \alpha, \cdots, \mathbf{z}_{x_n}(\mathbf{v})^T \alpha]^T \coloneqq \mathbf{Z}\alpha,$$

where $\mathbf{Z} \in \mathbb{R}^{n \times s}$ with $\mathbf{z}_{x_i}(\mathbf{v})^T$ as its jth row.

As the latter approximation is simply a Monte Carlo estimate, one could also pick an importance weighted probability density function $q(\cdot)$ and sample features $\{v_i\}_{i=1}^s$ from q (we refer to this sampling scheme as weighted RFF). Then, the function value $f(x_i)$ can be approximated by

$$\tilde{f}_q(x_j) = \sum_{i=1}^s \beta_i z_q(v_i, x_j) \coloneqq \mathbf{z}_{q, x_j}(\mathbf{v})^T \beta,$$

with $z_q(v_i, x_j) = \sqrt{p(v_i)/q(v_i)}z(v_i, x_j)$ and $\mathbf{z}_{q,x_j}(\mathbf{v}) = [z_q(v_1, x_j), \cdots, z_q(v_s, x_j)]^T$. Hence, a Monte-Carlo estimate of \mathbf{f}_x can be written in a matrix form as $\tilde{\mathbf{f}}_{q,x} = \mathbf{Z}_q \beta$, where $\mathbf{Z}_q \in \mathbb{R}^{n \times s}$ with $\mathbf{z}_{q,x_j}(\mathbf{v})^T$ as its jth row.

Let $\tilde{\mathbf{K}}$ and $\tilde{\mathbf{K}}_q$ be Gram-matrices with entries $\tilde{\mathbf{K}}_{ij}=\tilde{k}(x_i,x_j)$ and $\tilde{\mathbf{K}}_{q,ij}=\tilde{k}_q(x_i,x_j)$ such that

$$\tilde{\mathbf{K}} = \frac{1}{s} \mathbf{Z} \mathbf{Z}^T \qquad \wedge \qquad \tilde{\mathbf{K}}_q = \frac{1}{s} \mathbf{Z}_q \mathbf{Z}_q^T.$$

If we now denote the jth column of **Z** by $\mathbf{z}_{v_j}(\mathbf{x})$ and the jth column of \mathbf{Z}_q by $\mathbf{z}_{q,v_j}(\mathbf{x})$, then the following equalities can be derived easily from Eq. (3):

$$\mathbb{E}_{v \sim p}(\tilde{\mathbf{K}}) = \mathbf{K} = \mathbb{E}_{v \sim q}(\tilde{\mathbf{K}}_q) \quad \wedge \quad \mathbb{E}_{v \sim p}[\mathbf{z}_v(\mathbf{x})\mathbf{z}_v(\mathbf{x})^T] = \mathbf{K} = \mathbb{E}_{v \sim q}[\mathbf{z}_{q,v}(\mathbf{x})\mathbf{z}_{q,v}(\mathbf{x})^T].$$

Since we now conduct regularized empirical risk minimization in $\tilde{\mathcal{H}}$, we would like to find out the norm of \tilde{f} , the next proposition gives an upper bound of its norm.

Proposition 1. Assume that the reproducing kernel Hilbert space \mathcal{H} with kernel k admits a decomposition as in Eq. (2) and let $\tilde{\mathcal{H}} := \{\tilde{f} \mid \tilde{f} = \sum_{i=1}^{s} \alpha_i z(v_i, \cdot), \forall \alpha_i \in \mathbb{R}\}$. Then, for all $\tilde{f} \in \tilde{\mathcal{H}}$ it holds that $\|\tilde{f}\|_{\tilde{\mathcal{H}}}^2 \leq s\|\alpha\|_2^2$, where $\tilde{\mathcal{H}}$ is the reproducing kernel Hilbert space with kernel \tilde{k} (see Eq. 3).

Proof. Let us define a space of functions as

$$\mathcal{H}_1 := \{ f \mid f(x) = \alpha z(v, x), \alpha \in \mathbb{R} \}.$$

We now show that \mathcal{H}_1 is a reproducing kernel Hilbert space with kernel defined as $k_1(x,y) = (1/s)z(v,x)z(v,y)$, where s is a constant. Define a map $M: \mathbb{R} \to \mathcal{H}_1$ such that $M\alpha =$

 $\alpha z(v,\cdot), \forall \alpha \in \mathbb{R}$. The map M is a bijection, i.e. for any $f \in \mathcal{H}_1$ there exists a unique $\alpha_f \in \mathbb{R}$ such that $M^{-1}f = \alpha_f$. Now, we define an inner product on \mathcal{H}_1 as

$$\langle f, g \rangle_{\mathcal{H}_1} = \langle \sqrt{s} M^{-1} f, \sqrt{s} M^{-1} g \rangle_{\mathbb{R}} = s \alpha_f \alpha_g.$$

It is easy to show that this is a well defined inner product and, thus, \mathcal{H}_1 is a Hilbert space.

For any instance y, $k_1(\cdot, y) = (1/s)z(v, \cdot)z(v, y) \in \mathcal{H}_1$, since $(1/s)z(v, x) \in \mathbb{R}$ by definition. Take any $f \in \mathcal{H}_1$ and observe that

$$\langle f, k_1(\cdot, y) \rangle_{\mathcal{H}_1} = \langle \sqrt{s} M^{-1} f, \sqrt{s} M^{-1} k_1(\cdot, y) \rangle_{\mathbb{R}}$$

= $s \langle \alpha_f, 1/sz(v, y) \rangle_{\mathbb{R}}$
= $\alpha_f z(v, y) = f(y)$.

Hence, we have demonstrated the reproducing property for \mathcal{H}_1 and $\|f\|_{\mathcal{H}_1} = s\alpha_f^2$.

Now, suppose we have a sample of features $\{v_i\}_{i=1}^s$. For each v_i , we define the reproducing kernel Hilbert space

$$\mathcal{H}_i := \{ f \mid f(x) = \alpha z(v_i, x), \alpha \in \mathbb{R} \}$$

with the kernel $k_i(x,y) = (1/s)z(v_i,x)z(v_i,y)$. Denoting with

$$\tilde{\mathcal{H}} = \bigoplus_{i=1}^{s} \mathcal{H}_i = \{\tilde{f} : \tilde{f} = \sum_{i=1}^{s} f_i, f_i \in \mathcal{H}_i\}$$

and using the fact that the direct sum of reproducing kernel Hilbert spaces is another reproducing kernel Hilbert space (Berlinet and Thomas-Agnan, 2011), we have that $\tilde{k}(x,y) = \sum_{i=1}^s k_i(x,y) = (1/s) \sum_{i=1}^s z(v_i,x) z(v_i,y)$ is the kernel of $\tilde{\mathcal{H}}$ and that the norm of $\tilde{f} \in \tilde{\mathcal{H}}$ is defined as

$$\min_{f_i \in \mathcal{H}_i \mid f = \sum_{i=1}^s f_i} \sum_{i=1}^s \|f_i\|_{\mathcal{H}_i} =$$

$$\min_{\alpha_i \in \mathbb{R} \mid f_i = \alpha_i z(v_i, \cdot)} \sum_{i=1}^s s \alpha_i^2 = \min_{\alpha_i \in \mathbb{R} \mid f_i = \alpha_i z(v_i, \cdot)} s \|\alpha\|_2^2.$$

Hence, we have that $\|\tilde{f}\|_{\tilde{\mathcal{H}}} \leq s \|\alpha\|_2^2$.

An importance weighted density function based on the notion of *ridge leverage scores* is introduced in Alaoui and Mahoney (2015) for landmark selection in the Nyström method (Nyström, 1930; Smola and Schölkopf, 2000; Williams and Seeger, 2001). For landmarks selected using that sampling strategy, Alaoui and Mahoney (2015) establish a sharp convergence rate of the low-rank estimator based on the Nyström method. This result motivates the pursuit of a similar notion for random Fourier features. Indeed, Bach (2017b) propose a leverage score function based on an integral operator defined using the kernel function and the marginal distribution of a data-generating process. Building on this work, Avron et al. (2017) propose the ridge leverage function with respect to a fixed input dataset, i.e.,

$$l_{\lambda}(v) = p(v)\mathbf{z}_{v}(\mathbf{x})^{T}(\mathbf{K} + n\lambda \mathbf{I})^{-1}\mathbf{z}_{v}(\mathbf{x}).$$
(5)

From our assumption on the decomposition of a kernel function, it follows that there exists a constant z_0 such that $|z(v,x)| \leq z_0$ (for all v and x) and $\mathbf{z}_v(\mathbf{x})^T \mathbf{z}_v(\mathbf{x}) \leq nz_0^2$. We can now deduce the following inequality using a result from Avron et al. (2017, Proposition 4):

$$l_{\lambda}(v) \leq p(v) \frac{z_0^2}{\lambda}$$
 with $\int_{\mathcal{V}} l_{\lambda}(v) dv = \text{Tr} \big[\mathbf{K} (\mathbf{K} + n\lambda \mathbf{I})^{-1} \big] := d_{\mathbf{K}}^{\lambda}.$

The quantity $d_{\mathbf{K}}^{\lambda}$ is known for implicitly determining the number of independent parameters in a learning problem and, thus, it is called the *effective dimension of the problem* (Caponnetto and De Vito, 2007) or the *number of effective degrees of freedom* (Bach, 2013; Hastie, 2017).

We can now observe that $q^*(v) = l_{\lambda}(v)/d_{\mathbf{K}}^{\lambda}$ is a probability density function. In Avron et al. (2017), it has been established that sampling according to $q^*(v)$ requires fewer Fourier features compared to the standard spectral measure sampling. We refer to $q^*(v)$ as the *empirical ridge leverage score distribution* and, in the remainder of the manuscript, refer to this sampling strategy as leverage weighted RFF.

2.2 Rademacher Complexity

To characterize the stability of a learning algorithm, we need to take into account the complexity of the space of functions. Below, we first introduce a particular measure of the complexity over function spaces known as *Rademacher complexity*, we then present two lemmas that demonstrate how Rademacher complexity of RKHS can be linked to kernel and how the excess risk can be computed through Rademacher complexity.

Definition 1. Let P_x be a probability distribution on a set \mathcal{X} and suppose that $\{x_1 \cdots, x_n\}$ are independent samples selected according to P_x . Let \mathcal{H} be a class of functions mapping \mathcal{X} to \mathbb{R} . Then, the random variable known as the empirical Rademacher complexity is defined as

$$\hat{R}_n(\mathcal{H}) = \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{H}} \left| \frac{2}{n} \sum_{i=1}^n \sigma_i f(x_i) \right| \mid x_1, \cdots, x_n \right]$$

where $\sigma_1, \dots, \sigma_n$ are independent uniform $\{\pm 1\}$ -valued random variables. The corresponding Rademacher complexity is then defined as the expectation of the empirical Rademacher complexity, i.e.,

$$R_n(\mathcal{H}) = \mathbb{E}\Big[\hat{R}_n(\mathcal{H})\Big].$$

The following lemma provides the Rademacher complexity for a certain RKHS with kernel k.

Lemma 1. (Bartlett and Mendelson, 2002) Let \mathcal{H} be a reproducing kernel Hilbert space of functions mapping from \mathcal{X} to \mathbb{R} that corresponds to a positive definite kernel k. Let \mathcal{H}_0 be the unit ball of \mathcal{H} , centered at the origin. Then, we have that $R_n(\mathcal{H}_0) \leq (1/n)\mathbb{E}_X \sqrt{Tr(\mathbf{K})}$, where \mathbf{K} is the Gram matrix for kernel k over an independent and identically distributed sample $X = \{x_1, \dots, x_n\}$.

Lemma 2 states that the expected excess risk convergence rate of a particular estimator in \mathcal{H} not only depends on the number of data points, but also on the complexity of \mathcal{H} .

Lemma 2. (Bartlett and Mendelson, 2002, Theorem 8) Let $\{x_i, y_i\}_{i=1}^n$ be an independent and identically distributed sample from a probability measure P defined on $\mathcal{X} \times \mathcal{Y}$ and let \mathcal{H} be the space of functions mapping from \mathcal{X} to \mathcal{A} . Denote a loss function with $l: \mathcal{Y} \times \mathcal{A} \to [0,1]$ and define the expected risk function for all $f \in \mathcal{H}$ to be $\mathcal{E}(f) = \mathbb{E}_P(l(y,f(x)))$, together with the corresponding empirical risk function $\hat{\mathcal{E}}(f) = (1/n) \sum_{i=1}^n l(y_i,f(x_i))$. Then, for a sample size n, for all $f \in \mathcal{H}$ and $\delta \in (0,1)$, with probability $1-\delta$, we have that

$$\mathcal{E}(f) \le \hat{\mathcal{E}}(f) + R_n(\tilde{l} \circ \mathcal{H}) + \sqrt{\frac{8\log(2/\delta)}{n}}$$

where $\tilde{l} \circ \mathcal{H} = \{(x, y) \to l(y, f(x)) - l(y, 0) \mid f \in \mathcal{H}\}.$

3. Theoretical Analysis

In this section, we provide a unified analysis of the generalization properties of learning with random Fourier features. We start with a bound for learning with the mean squared error loss function and then extend our results to problems with Lipschitz continuous loss functions. Before presenting the results, we briefly review the standard problem setting for supervised learning with kernel methods.

Let $\mathcal X$ be an instance space, $\mathcal Y$ a label space, and $\rho(x,y)=\rho_{\mathcal X}(x)\rho(y\mid x)$ a probability measure on $\mathcal X\times\mathcal Y$ defining the relationship between an instance $x\in\mathcal X$ and a label $y\in\mathcal Y$. A training sample is a set of examples $\{(x_i,y_i)\}_{i=1}^n$ sampled independently from the distribution ρ , known only through the sample. The distribution $\rho_{\mathcal X}$ is called the marginal distribution of a data-generating process. The goal of a supervised learning task defined with a kernel function k (and the associated reproducing kernel Hilbert space $\mathcal H$) is to find a hypothesis $f: \mathcal X\to\mathcal Y$ such that $f\in\mathcal H$ and f(x) is a good estimate of the label $g\in\mathcal Y$ corresponding to a previously unseen instance $g\in\mathcal X$. While in regression tasks $g\in\mathcal X$ in classification tasks it is typically the case that $g\in\mathcal X$ are sult of the representer theorem an empirical risk minimization problem in this setting can be expressed as (Scholkopf and Smola, 2001)

$$\hat{f}^{\lambda} := \underset{f \in \mathcal{H}}{\operatorname{arg\,min}} \ \frac{1}{n} \sum_{i=1}^{n} l(y_i, (\mathbf{K}\alpha)_i) + \lambda \alpha^T \mathbf{K}\alpha, \tag{6}$$

where $f = \sum_{i=1}^{n} \alpha_i k(x_i, \cdot)$ with $\alpha \in \mathbb{R}^n$, $l : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_+$ is a loss function, \mathbf{K} is the kernel matrix, and λ is the regularization parameter. The hypothesis \hat{f}^{λ} is an empirical estimator and its ability to describe ρ is measured by the expected risk (Caponnetto and De Vito, 2007)

$$\mathcal{E}(\hat{f}^{\lambda}) = \int_{\mathcal{X} \times \mathcal{Y}} l(y, \hat{f}^{\lambda}(x)) d\rho(x, y).$$

Similar to Rudi and Rosasco (2017) and Caponnetto and De Vito (2007), we have assumed 3 the existence of $f_{\mathcal{H}} \in \mathcal{H}$ such that $f_{\mathcal{H}} = \inf_{f \in \mathcal{H}} \mathcal{E}(f)$. The assumption implies that there exists some

^{2.} Throughout the paper, we assume (without loss of generality) that our hypothesis space is the unit ball in a reproducing kernel Hilbert space \mathcal{H} , i.e., $\|f\|_{\mathcal{H}} \leq 1$. This is a standard assumption, characteristic to the analysis of random Fourier features (e.g., see Rudi and Rosasco, 2017)

^{3.} The existence of $f_{\mathcal{H}}$ depends on the complexity of \mathcal{H} which is related to the conditional distribution $\rho(y|x)$ and the marginal distribution $\rho_{\mathcal{X}}$. For more details, please see Caponnetto and De Vito (2007) and Rudi and Rosasco (2017).

ball of radius R > 0 containing $f_{\mathcal{H}}$ in its interior. Our theoretical results do not require prior knowledge of this constant and hold uniformly over all finite radii. To simplify our derivations and constant terms in our bounds, we have (without loss of generality) assumed that R = 1.

3.1 Learning with the Squared Error Loss

In this section, we consider learning with the squared error loss, i.e., $l(y,f(x))=(y-f(x))^2$. For this particular loss function, the optimization problem from Eq. (6) is known as *kernel ridge regression*. The problem can be reduced to solving a linear system $(\mathbf{K}+n\lambda\mathbf{I})\alpha=Y$, with $Y=[y_1,\cdots,y_n]^T$. Typically, an approximation of the kernel function based on random Fourier features is employed in order to effectively reduce the computational cost and scale kernel ridge regression to problems with millions of examples. More specifically, for a vector of observed labels Y the goal is to find a hypothesis $\tilde{\mathbf{f}}_x=\mathbf{Z}_q\beta$ that minimizes $\|Y-\tilde{\mathbf{f}}_x\|_2^2$ while having good generalization properties. In order to achieve this, one needs to control the complexity of hypotheses defined by random Fourier features and avoid over-fitting. According to Proposition 1, $\|\tilde{f}\|_{\tilde{\mathcal{H}}}^2$ can be upper bounded by $s\|\beta\|_2^2$, where s is the number of sampled features. Hence, the learning problem with random Fourier features and the squared error loss can be cast as

$$\beta_{\lambda} := \underset{\beta \in \mathbb{R}^s}{\operatorname{arg\,min}} \quad \frac{1}{n} \|Y - \mathbf{Z}_q \beta\|_2^2 + \lambda s \|\beta\|_2^2. \tag{7}$$

This is a linear ridge regression problem in the space of Fourier features and the optimal hypothesis is given by $f_{\beta}^{\lambda} = \mathbf{Z}_q \beta_{\lambda}$, where $\beta_{\lambda} = (\mathbf{Z}_q^T \mathbf{Z}_q + n\lambda \mathbf{I})^{-1} \mathbf{Z}_q^T Y$. Since $\mathbf{Z}_q \in \mathbb{R}^{n \times s}$, the computational and space complexities are $O(s^3 + ns^2)$ and O(ns). Thus, significant savings can be achieved using estimators with $s \ll n$. To assess the effectiveness of such estimators, it is important to understand the relationship between the expected risk and the choice of s.

3.1.1 WORST CASE ANALYSIS

In this section, we assume that the unit ball of the reproducing kernel Hilbert space contains the hypothesis $f_{\mathcal{H}}$ and provide a bound on the required number of random Fourier features with respect to the worst case minimax rate of the corresponding kernel ridge regression problem. The following theorem gives a general result while taking into account both the number of features s and a sampling strategy for selecting them.

Theorem 1. Assume a kernel function k has a decomposition as in Eq. (2) and let $|y| \leq y_0$ be bounded with $y_0 > 0$. Denote with $\lambda_1 \geq \cdots \geq \lambda_n$ the eigenvalues of the kernel matrix \mathbf{K} and assume the regularization parameter satisfies $0 \leq n\lambda \leq \lambda_1$. Let $\tilde{l}: \mathcal{V} \to \mathbb{R}$ be a measurable function such that $\tilde{l}(v) \geq l_{\lambda}(v)$ ($\forall v \in \mathcal{V}$) and $d_{\tilde{l}} = \int_{\mathcal{V}} \tilde{l}(v) dv < \infty$. Suppose $\{v_i\}_{i=1}^s$ are sampled independently from the probability density function $q(v) = \tilde{l}(v)/d_{\tilde{l}}$. If the unit ball of \mathcal{H} contains the optimal hypothesis $f_{\mathcal{H}}$ and

$$s \geq 5d_{\tilde{l}}\log\frac{16d_{\mathbf{K}}^{\lambda}}{\delta},$$

then for all $\delta \in (0,1)$, with probability $1-\delta$, the excess risk of f^{λ}_{β} can be upper bounded as

$$\mathcal{E}(f_{\beta}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}}) \leq 2\lambda + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right) + \mathcal{E}(\hat{f}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}}). \tag{8}$$

Theorem 1 expresses the trade-off between the computational and statistical efficiency through the regularization parameter λ , the effective dimension of the problem $d_{\mathbf{K}}^{\lambda}$, and the normalization constant of the sampling distribution $d_{\tilde{l}}$. The regularization parameter can be considered as some function of the number of training examples (Caponnetto and De Vito, 2007; Rudi and Rosasco, 2017) and we use its decay rate as the sample size increases to quantify the complexity of the target regression function $f_{\rho}(x) = \int y d\rho(y \mid x)$. In particular, Caponnetto and De Vito (2007) have shown that the minimax risk convergence rate for kernel ridge regression is $O(\frac{1}{\sqrt{n}})$. Setting $\lambda \propto \frac{1}{\sqrt{n}}$, we observe that the estimator f_{β}^{λ} attains the worst case minimax rate of kernel ridge regression.

To prove Theorem 1, we need Theorem 2 and Lemma 3 to analyse the learning risk. In Theorem 2, we give a general result that provides an upper bound on the approximation error between any function $f \in \mathcal{H}$ and its estimator based on random Fourier features. As discussed in Section 2, we would like to approximate a function $f \in \mathcal{H}$ at observation points with preferably as small function norm as possible. The estimation of \mathbf{f}_x can be formulated as the following optimization problem:

$$\min_{\beta \in \mathbb{R}^s} \frac{1}{n} \|\mathbf{f}_x - \mathbf{Z}_q \beta\|_2^2 + \lambda s \|\beta\|_2^2.$$

Below we provide the desired upper bound on the approximation error of the estimator based on random Fourier features (proof presented in Appendix B).

Theorem 2. Let $\lambda_1 \geq \cdots \geq \lambda_n$ be the eigenvalues of the kernel matrix **K** and assume that the regularization parameter satisfies $0 \leq n\lambda \leq \lambda_1$. Let $\tilde{l}: \mathcal{V} \to \mathbb{R}$ be a measurable function such that $\tilde{l}(v) \geq l_{\lambda}(v) \ (\forall v \in \mathcal{V})$ and

$$d_{\tilde{l}} = \int_{\mathcal{V}} \tilde{l}(v) dv < \infty.$$

Suppose $\{v_i\}_{i=1}^s$ are sampled independently according to probability density function $q(v) = \frac{\tilde{l}(v)}{d_{\tilde{l}}}$. If

$$s \ge 5d_{\tilde{l}} \log \frac{16d_{\mathbf{K}}^{\lambda}}{\delta},$$

then for all $\delta \in (0,1)$ and $||f||_{\mathcal{H}} \leq 1$, with probability greater than $1-\delta$, we have that it holds

$$\sup_{\|f\|_{\mathcal{H}} \le 1} \inf_{\sqrt{s} \|\beta\|_2 \le \sqrt{2}} \frac{1}{n} \|\mathbf{f}_x - \mathbf{Z}_q \beta\|_2^2 \le 2\lambda.$$
 (9)

The next lemma, following Rudi and Rosasco (2017), is important in demonstrating the risk convergence rate as it illustrates the relationship between Y, \hat{f}^{λ} and f^{λ}_{β} , the proof is in Appendix C.

Lemma 3. Assuming that the conditions of Theorem 1 hold, let \hat{f}^{λ} and f^{λ}_{β} be the empirical estimators from problems (6) and (7), respectively. In addition, suppose that $\{v_i\}_{i=1}^s$ are independent samples selected according to a probability measure τ_q with probability density function q(v) such that p(v)/q(v) > 0 almost surely. Then, we have

$$\langle Y - \hat{f}^{\lambda}, f^{\lambda}_{\beta} - \hat{f}^{\lambda} \rangle = 0.$$

Equipped with Theorem 2 and Lemma 3, we are now ready to prove Theorem 1.

Proof. The proof relies on the decomposition of the expected risk of $\mathcal{E}(f_{\beta}^{\lambda})$ as follows

$$\mathcal{E}(f_{\beta}^{\lambda}) = \mathcal{E}(f_{\beta}^{\lambda}) - \hat{\mathcal{E}}(f_{\beta}^{\lambda}) \tag{10}$$

$$+\hat{\mathcal{E}}(f_{\beta}^{\lambda}) - \hat{\mathcal{E}}(\hat{f}^{\lambda}) \tag{11}$$

$$+\hat{\mathcal{E}}(\hat{f}^{\lambda}) - \mathcal{E}(\hat{f}^{\lambda}) \tag{12}$$

$$+\mathcal{E}(\hat{f}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}}) + \mathcal{E}(f_{\mathcal{H}}). \tag{13}$$

For (10), the bound is based on the Rademacher complexity of the reproducing kernel Hilbert space $\tilde{\mathcal{H}}$, where $\tilde{\mathcal{H}}$ corresponds to the approximated kernel \tilde{k} . We can upper bound the Rademacher complexity of this hypothesis space with Lemma 1. As l(y,f(x)) is the squared error loss function with y and f(x) bounded, we have that l is a Lipschitz continuous function with some constant L>0. Hence,

$$(10) \leq R_{n}(\tilde{l} \circ \tilde{\mathcal{H}}) + \sqrt{\frac{8 \log(2/\delta)}{n}}$$

$$\leq \sqrt{2}L \frac{1}{n} \mathbb{E}_{X} \sqrt{\text{Tr}(\tilde{\mathbf{K}})} + \sqrt{\frac{8 \log(2/\delta)}{n}}$$

$$\leq \sqrt{2}L \frac{1}{n} \sqrt{\mathbb{E}_{X} \text{Tr}(\tilde{\mathbf{K}})} + \sqrt{\frac{8 \log(2/\delta)}{n}}$$

$$\leq \sqrt{2}L \frac{1}{n} \sqrt{nz_{0}^{2}} + \sqrt{\frac{8 \log(2/\delta)}{n}}$$

$$\leq \frac{\sqrt{2}Lz_{0}}{\sqrt{n}} + \sqrt{\frac{8 \log(2/\delta)}{n}} \in O(\frac{1}{\sqrt{n}}), \tag{14}$$

where in the last inequality we applied Lemma 2 to $\tilde{\mathcal{H}}$, which is a reproducing kernel Hilbert space with radius $\sqrt{2}$. For (12), a similar reasoning can be applied to the unit ball in the reproducing kernel Hilbert space \mathcal{H} .

For (11), we observe that

$$\begin{split} \hat{\mathcal{E}}(f_{\beta}^{\lambda}) - \hat{\mathcal{E}}(\hat{f}^{\lambda}) &= \frac{1}{n} \|Y - f_{\beta}^{\lambda}\|_{2}^{2} - \frac{1}{n} \|Y - \hat{f}^{\lambda}\|_{2}^{2} \\ &= \frac{1}{n} \inf_{\|f_{\beta}\|} \|Y - f_{\beta}\|_{2}^{2} - \frac{1}{n} \|Y - \hat{f}^{\lambda}\|_{2}^{2} \\ &= \frac{1}{n} \inf_{\|f_{\beta}\|} \left(\|Y - \hat{f}^{\lambda}\|_{2}^{2} + \|\hat{f}^{\lambda} - f_{\beta}\|_{2}^{2} \right. \\ &\quad + 2\langle Y - \hat{f}^{\lambda}, \hat{f}^{\lambda} - f_{\beta} \rangle \right) - \frac{1}{n} \|Y - \hat{f}^{\lambda}\|_{2}^{2} \\ &\leq \frac{1}{n} \inf_{\|f_{\beta}\|} \|\hat{f}^{\lambda} - f_{\beta}\|_{2}^{2} \\ &\quad + \frac{2}{n} \inf_{\|f_{\beta}\|} \langle Y - \hat{f}^{\lambda}, \hat{f}^{\lambda} - f_{\beta} \rangle \\ &\leq \frac{1}{n} \inf_{\|f_{\beta}\|} \|\hat{f}^{\lambda} - f_{\beta}\|_{2}^{2} + \frac{2}{n} \langle Y - \hat{f}^{\lambda}, \hat{f}^{\lambda} - f_{\beta}^{\lambda} \rangle \\ &= \frac{1}{n} \inf_{\|f_{\beta}\|} \|\hat{f}^{\lambda} - f_{\beta}\|_{2}^{2} \\ &\leq \sup_{\|f\|} \|f_{\beta}\| \frac{1}{n} \|f - f_{\beta}\|_{2}^{2} \\ &\leq 2\lambda, \end{split}$$

where in the last step we employ Theorem 2. Combining the three results, we derive

$$\mathcal{E}(f_{\beta}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}}) \le 2\lambda + O(\frac{1}{\sqrt{n}}) + \mathcal{E}(\hat{f}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}}). \tag{15}$$

As a consequence of Theorem 1, we have the following bounds on the number of required features for the two strategies: *leverage weighted* RFF (Corollary 1) and *plain* RFF (Corollary 2).

Corollary 1. If the probability density function from Theorem 1 is the empirical ridge leverage score distribution $q^*(v)$, then the upper bound on the risk from Eq. (8) holds for all $s \geq 5d_{\mathbf{K}}^{\lambda} \log \frac{16d_{\mathbf{K}}^{\lambda}}{2}$.

Proof. For Corollary 1, we set $\tilde{l}(v) = l_{\lambda}(v)$ and deduce

$$d_{\tilde{l}} = \int_{\mathcal{V}} l_{\lambda}(v) dv = d_{\mathbf{K}}^{\lambda}.$$

Theorem 1 and Corollary 1 have several implications on the choice of λ and s. First, we could pick $\lambda \in O(n^{-1/2})$ that implies the worst case minimax rate for kernel ridge regression (Caponnetto and De Vito, 2007; Rudi and Rosasco, 2017; Bartlett et al., 2005) and observe that in this case s is proportional to $d_{\mathbf{K}}^{\lambda} \log d_{\mathbf{K}}^{\lambda}$. As $d_{\mathbf{K}}^{\lambda}$ is determined by the learning problem (i.e., the marginal distribution $\rho_{\mathcal{X}}$), we can consider several different cases. In the best case (e.g., the Gaussian kernel with a sub-Gaussian marginal distribution $\rho_{\mathcal{X}}$), the eigenvalues of \mathbf{K} exhibit a geometric/exponential

decay, i.e., $\lambda_i \propto R_0 r^i$ (R_0 is some constant). From Bach (2017b), we know that $d_{\mathbf{K}}^{\lambda} \leq \log(R_0/\lambda)$, implying $s \geq \log^2 n$. Hence, significant savings can be obtained with $O(n\log^4 n + \log^6 n)$ computational and $O(n\log^2 n)$ storage complexities of linear ridge regression over random Fourier features, as opposed to $O(n^3)$ and $O(n^2)$ costs (respectively) in the kernel ridge regression setting.

In the case of a slower decay (e.g., \mathcal{H} is a Sobolev space of order $t \geq 1$) with $\lambda_i \propto R_0 i^{-2t}$, we have $d_{\mathbf{K}}^{\lambda} \leq (R_0/\lambda)^{1/(2t)}$ and $s \geq n^{1/(4t)} \log n$. Hence, even in this case a substantial computational saving can be achieved. Furthermore, in the worst case with λ_i close to $R_0 i^{-1}$, our bound implies that $s \geq \sqrt{n} \log n$ features is sufficient, recovering the result from Rudi and Rosasco (2017).

Corollary 2. If the probability density function from Theorem 1 is the spectral measure p(v) from Eq. (2), then the upper bound on the risk from Eq. (8) holds for all $s \geq 5\frac{z_0^2}{\lambda}\log\frac{16d_{\mathbf{K}}^{\lambda}}{\delta}$.

Proof. For Corollary 2, we set $\tilde{l}(v) = p(v) \frac{z_0^2}{\lambda}$ and derive

$$d_{\tilde{l}} = \int_{\mathcal{V}} p(v) \frac{z_0^2}{\lambda} dv = \frac{z_0^2}{\lambda}.$$

Corollary 2 addresses plain random Fourier features and states that if s is chosen to be greater than $\sqrt{n}\log d_{\mathbf{K}}^{\lambda}$ and $\lambda \propto \frac{1}{\sqrt{n}}$ then the minimax risk convergence rate is guaranteed. When the eigenvalues have an exponential decay, we obtain the same convergence rate with only $s \geq \sqrt{n}\log\log n$ features, which is an improvement compared to a result by Rudi and Rosasco (2017) where $s \geq \sqrt{n}\log n$. For the other two cases, we derive $s \geq \sqrt{n}\log n$ and recover the results from Rudi and Rosasco (2017).

3.1.2 REFINED ANALYSIS

In this section, we provide a more refined analysis with expected risk convergence rates faster than $\mathcal{O}(\frac{1}{\sqrt{n}})$, depending on the spectrum decay of the kernel function and/or the complexity of the target regression function.

Theorem 3. Suppose that the conditions from Theorem 1 apply and let

$$s \geq 5d_{\tilde{l}}\log\frac{16d_{\mathbf{K}}^{\lambda}}{\delta}.$$

Then, for all D > 1 and $\delta \in (0,1)$, with probability $1 - \delta$, the excess risk of f_{β}^{λ} can be upper bounded as

$$\mathcal{E}(f_{\beta}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}}) \leq 2\hat{r}_{\mathcal{H}}^* + 2\lambda \frac{D}{(D-1)} + \mathcal{O}\left(\frac{1}{n}\right) + \mathcal{E}(\hat{f}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}}). \tag{16}$$

Furthermore, denoting the eigenvalues of the normalized kernel matrix $(1/n)\mathbf{K}$ with $\{\hat{\lambda}_i\}_{i=1}^n$, we have that

$$\hat{r}_{\mathcal{H}}^* \le \min_{0 \le h \le n} \left(\frac{h}{n} * \frac{e_7}{n^2 \lambda^2} + \sqrt{\frac{1}{n} \sum_{i > h} \hat{\lambda}_i} \right), \tag{17}$$

where $e_7 > 0$ is a constant and $\hat{\lambda}_1 \ge \cdots \ge \hat{\lambda}_n$.

Theorem 3 covers a wide range of cases and can provide sharper risk convergence rates. In particular, note that $\hat{r}_{\mathcal{H}}^*$ is of order $O(1/\sqrt{n})$, which happens when the spectrum decays approximately as 1/n and h=0. In this case, the excess risk converges with the rate $O(1/\sqrt{n})$, which corresponds to the considered worst case minimax rate. On the other hand, if the eigenvalues decay exponentially, then setting $h=\lceil \log n \rceil$ implies that $\hat{r}_{\mathcal{H}}^* \leq O(\log n/n)$. Furthermore, setting $\lambda \propto \frac{\log n}{n}$, we can show that the excess risk converges at a much faster rate of $O(\log n/n)$. In the best case, when the kernel function has only finitely many positive eigenvalues, we have that $\hat{r}_{\mathcal{H}}^* \leq O(1/n)$ by letting h be any fixed value larger than the number of positive eigenvalues. In this case, we obtain the fastest rate of O(1/n) for the regularization parameter $\lambda \propto \frac{1}{n}$.

To prove Theorem 3, we rely on the notion of local Rademacher complexity and adjust our notation so that it is easier to cross-reference relevant auxiliary claims from Bartlett et al. (2005). Suppose P is a probability measure on $\mathcal{X} \times \mathcal{Y}$ and let $\{x_i, y_i\}_{i=1}^n$ be an independent sample from P. For any reproducing kernel Hilbert space \mathcal{H} and a loss function l, we define the transformed function class as $l_{\mathcal{H}} \coloneqq \{l(f(x), y) \mid f \in \mathcal{H}\}$. We also abbreviate the notation and denote with $l_f = l(f(x), y)$ and $\mathbb{E}_n(f) = 1/n \sum_{i=1}^n f(x_i)$. For the reproducing kernel Hilbert space \mathcal{H} , we denote the solution of the kernel ridge regression problem by \hat{f} .

In general, the reason Theorem 1 is not sharp is because when analysing Eqs.(10 and 12), we used the global Rademacher complexity of the whole RKHS. However, as pointed out by Bartlett et al. (2005), regularised ERM is likely to return a function that are around $f_{\mathcal{H}}$. Hence, instead of estimating the global function space complexity, we could just analyse the local space around $f_{\mathcal{H}}$. Specifically, we would like to apply Theorem 4 to Eqs.(10 and 12) to compute the local Rademacher complexity. In order to do so, we need to find a proper *sub-root* function (see Appendix D for its definition and property). Below, we first state the theorem for local Rademacher complexity. We then present Lemma 4 and 11 in order to prove Theorem 5 to find the proper sub-root function. By combining Theorem 4 and 5 and apply them to Eqs.(10 and 12), we obtain Theorem 3.

Theorem 4. (Bartlett et al., 2005, Theorem 4.1) Let \mathcal{H} be a class of functions with ranges in [-1,1] and assume that there is some constant B_0 such that for all $f \in \mathcal{H}$, $\mathbb{E}(f^2) \leq B_0\mathbb{E}(f)$. Let $\hat{\psi}_n$ be a sub-root function and let \hat{r}^* be the fixed point of $\hat{\psi}_n$, i.e., $\hat{\psi}_n(\hat{r}^*) = \hat{r}^*$. Fix any $\delta > 0$, and assume that for any $r \geq \hat{r}^*$,

$$\hat{\psi}_n(r) \ge e_1 \hat{R}_n \{ f \in \text{star}(\mathcal{H}, 0) \mid \mathbb{E}_n(f^2) \le r \} + \frac{e_2 \delta}{n}$$

and

$$\operatorname{star}(\mathcal{H}, f_0) = \{ f_0 + \alpha(f - f_0) \mid f \in \mathcal{H} \land \alpha \in [0, 1] \}.$$

Then, for all $f \in \mathcal{H}$ and D > 1, with probability greater than $1 - 3e^{-\delta}$,

$$\mathbb{E}(f) \le \frac{D}{D-1} \mathbb{E}_n(f) + \frac{6D}{B} \hat{r}^* + \frac{e_3 \delta}{n}$$

where e_1 , e_2 and e_3 are some constants.

As stated before, in order to obtain a sharper rate, we need to compute the local Rademacher complexity. The following lemma allows us to compute it through the eigenvalues of the Gram matrix.

Lemma 4. (Bartlett et al., 2005, Lemma 6.6) Let k be a positive definite kernel function with reproducing kernel Hilbert space \mathcal{H} and let $\hat{\lambda}_1 \geq \cdots \geq \hat{\lambda}_n$ be the eigenvalues of the normalized Gram-matrix $(1/n)\mathbf{K}$. Then, for all r > 0

$$\hat{R}_n\{f \in \mathcal{H} \mid \mathbb{E}_n(f^2) \le r\} \le \left(\frac{2}{n} \sum_{i=1}^n \min\{r, \hat{\lambda}_i\}\right)^{1/2}$$

Now we would like to apply Theorem 4, our task is to find a proper sub-root function. To this end, we let l be the squared error loss function and observe that for all $f \in \mathcal{H}$ it holds that

$$\mathbb{E}_{n}(l_{f}^{2}) \geq (\mathbb{E}_{n}(l_{f}))^{2} \quad (x^{2} \text{ is convex})$$

$$\geq (\mathbb{E}_{n}l_{f})^{2} - (\mathbb{E}_{n}l_{\hat{f}})^{2}$$

$$= (\mathbb{E}_{n}l_{f} + \mathbb{E}_{n}l_{\hat{f}})(\mathbb{E}_{n}l_{f} - \mathbb{E}_{n}l_{\hat{f}})$$

$$\geq 2\mathbb{E}_{n}l_{\hat{f}} \, \mathbb{E}_{n}(l_{f} - l_{\hat{f}})$$

$$\geq \frac{2}{B}\mathbb{E}_{n}l_{\hat{f}}\mathbb{E}_{n}(f - \hat{f})^{2}. \text{ (Lemma 11 in Appendix D)}$$
(18)

The third inequality holds because \hat{f} achieves the minimal empirical risk. The last inequality is a consequence of Lemma 11 applied to the empirical probability distribution P_n . Hence, to obtain a lower bound on $\mathbb{E}_n l_f^2$ expressed solely in terms of $\mathbb{E}_n (f - \hat{f})^2$, we need to find a lower bound of $\mathbb{E}_n l_{\hat{f}}$. First, observe that it holds

$$\mathbb{E}_n l_{\hat{f}} = \frac{1}{n} \|Y - \mathbf{K} (\mathbf{K} + n\lambda I)^{-1} Y\|^2.$$

Then, using this expression we derive

$$\mathbb{E}_{n}l_{\hat{f}} = \frac{1}{n} \|Y - \mathbf{K}(\mathbf{K} + n\lambda I)^{-1}Y\|^{2}$$

$$= n\lambda^{2}Y^{T}(\mathbf{K} + n\lambda I)^{-2}Y$$

$$\geq \frac{n\lambda^{2}}{(\lambda_{1} + n\lambda)^{2}}Y^{T}Y$$

$$= \left(\frac{n\lambda}{\lambda_{1} + n\lambda}\right)^{2} \frac{1}{n} \sum_{i=1}^{n} y_{i}^{2}$$

$$\geq \left(\frac{n\lambda}{\lambda_{1} + n\lambda}\right)^{2} \sigma_{y}^{2} \quad (\text{with } \frac{1}{n} \sum_{i=1}^{n} y_{i}^{2} \geq \sigma_{y}^{2})$$

$$= \sigma_{y}^{2} \left(\frac{1}{1 + \frac{\lambda_{1}}{n\lambda}}\right)^{2}$$

$$\geq \sigma_{y}^{2} \left(\frac{1}{\frac{\lambda_{1}}{n\lambda} + \frac{\lambda_{1}}{n\lambda}}\right)^{2}$$

$$= \frac{\sigma_{y}^{2}}{4} \left(\frac{n\lambda}{\lambda_{1}}\right)^{2} = e_{4}(n\lambda)^{2}, \tag{19}$$

where $e_4 = (\frac{\sigma_y}{2\lambda_1})^2$ is a constant.

The last equality follows because λ_1 is independent of n and λ , as well as bounded. Hence, Eq.(18) becomes

$$\mathbb{E}_n l_f^2 \ge \frac{2e_4(n\lambda)^2}{B} \mathbb{E}_n (f - \hat{f})^2 =: e_5(n\lambda)^2 \mathbb{E}_n (f - \hat{f})^2.$$

As a result of this, we have the following inequality for the two function classes

$$\{l_f \in l_{\mathcal{H}} \mid \mathbb{E}_n l_f^2 \leq r\} \subseteq \{l_f \in l_{\mathcal{H}} \mid \mathbb{E}_n (f - \hat{f})^2 \leq \frac{r}{e_5(n\lambda)^2}\}.$$

Recall that for a function class \mathcal{H} , we denote its empirical Rademacher complexity by $\hat{R}_n(\mathcal{H})$. Then, we have the following inequality

$$\hat{R}_{n}\{l_{f} \in l_{\mathcal{H}} \mid \mathbb{E}_{n}l_{f}^{2} \leq r\} \leq
\hat{R}_{n}\{l_{f} \in l_{\mathcal{H}} \mid \mathbb{E}_{n}(f - \hat{f})^{2} \leq \frac{r}{e_{5}n^{2}\lambda^{2}}\} =
\hat{R}_{n}\{l_{f} - l_{\hat{f}} \mid \mathbb{E}_{n}(f - \hat{f})^{2} \leq \frac{r}{e_{5}n^{2}\lambda^{2}} \wedge l_{f} \in l_{\mathcal{H}}\} \leq
L\hat{R}_{n}\{f - \hat{f} \mid \mathbb{E}_{n}(f - \hat{f})^{2} \leq \frac{r}{e_{5}n^{2}\lambda^{2}} \wedge f \in \mathcal{H}\} \leq
L\hat{R}_{n}\{f - g \mid \mathbb{E}_{n}(f - g)^{2} \leq \frac{r}{e_{5}n^{2}\lambda^{2}} \wedge f, g \in \mathcal{H}\} \leq
2L\hat{R}_{n}\{f \in \mathcal{H} \mid \mathbb{E}_{n}f^{2} \leq \frac{1}{4e_{5}} \frac{r}{n^{2}\lambda^{2}}\} =
2L\hat{R}_{n}\{f \in \mathcal{H} \mid \mathbb{E}_{n}f^{2} \leq \frac{e_{6}r}{n^{2}\lambda^{2}}\},$$

where the last inequality is due to Bartlett et al. (2005, Corollary 6.7). Now, combining Lemma 4 and Eq. (20) we can derive the following theorem which gives us the proper sub-root function $\hat{\psi}_n$. The theorem is proved in Appendix E.

Theorem 5. Assume $\{x_i, y_i\}_{i=1}^n$ is an independent sample from a probability measure P defined on $\mathcal{X} \times \mathcal{Y}$, with $\mathcal{Y} \in [-1,1]$. Let k be a positive definite kernel with the reproducing kernel Hilbert space \mathcal{H} and let $\hat{\lambda}_1 \geq \cdots, \geq \hat{\lambda}_n$ be the eigenvalues of the normalized kernel Gram-matrix. Denote the squared error loss function by $l(f(x), y) = (f(x) - y)^2$ and fix $\delta > 0$. If

$$\hat{\psi}_n(r) = 2Le_1\left(\frac{2}{n}\sum_{i=1}^n \min\{r, \hat{\lambda}_i\}\right)^{1/2} + \frac{e_3\delta}{n},$$

then for all $l_f \in l_H$ and D > 1, with probability $1 - 3e^{-\delta}$,

$$\mathcal{E}(f) = \mathbb{E}(l_f) \le \frac{D}{D-1} \mathbb{E}_n l_f + \frac{6D}{B} \hat{r}^* + \frac{e_3 \delta}{n}.$$

Moreover, the fixed point \hat{r}^* defined with $\hat{r}^* = \hat{\psi}_n(\hat{r})^*$ can be upper bounded by

$$\hat{r}^* \le \min_{0 \le h \le n} \left(\frac{h}{n} * \frac{e_7}{n^2 \lambda^2} + \sqrt{\frac{1}{n} \sum_{i > h} \hat{\lambda}_i} \right),$$

where e_7 is a constant, and λ is the regularization parameter used in kernel ridge regression.

We are now ready to deliver the proof of Theorem 3.

Proof. We decompose $\mathcal{E}(f_{\beta}^{\lambda})$ with D>1 as follows

$$\begin{split} \mathcal{E}(f_{\beta}^{\lambda}) &= \mathcal{E}(f_{\beta}^{\lambda}) - \frac{D}{D-1}\hat{\mathcal{E}}(f_{\beta}^{\lambda}) \\ &+ \frac{D}{D-1}\hat{\mathcal{E}}(f_{\beta}^{\lambda}) - \frac{D}{D-1}\hat{\mathcal{E}}(\hat{f}^{\lambda}) \\ &+ \frac{D}{D-1}\hat{\mathcal{E}}(\hat{f}^{\lambda}) - \mathcal{E}(\hat{f}^{\lambda}) \\ &+ \mathcal{E}(\hat{f}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}}) \\ &+ \mathcal{E}(f_{\mathcal{H}}). \end{split}$$

Hence,

$$\mathcal{E}(f_{\beta}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}}) \le \left| \mathcal{E}(f_{\beta}^{\lambda}) - \frac{D}{D-1} \hat{\mathcal{E}}(f_{\beta}^{\lambda}) \right|$$
 (21)

$$+\frac{D}{D-1}(\hat{\mathcal{E}}(f_{\beta}^{\lambda}) - \hat{\mathcal{E}}(\hat{f}^{\lambda})) \tag{22}$$

$$+ \left| \frac{D}{D-1} \hat{\mathcal{E}}(\hat{f}^{\lambda}) - \mathcal{E}(\hat{f}^{\lambda}) \right| \tag{23}$$

$$+\mathcal{E}(\hat{f}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}}).$$
 (24)

We have already demonstrated that

Eq. (22)
$$\leq 2 \frac{D}{D-1} \lambda$$
.

For Eqs. (21) and (23) we apply Theorem 5. However, note that f_{β}^{λ} and \hat{f}^{λ} belong to different reproducing kernel Hilbert spaces. As a result, we have

Eq. (21)
$$\leq \hat{r}^*_{\tilde{\mathcal{H}}} + O(1/n)$$

Eq. (23)
$$\leq \hat{r}_{\mathcal{H}}^* + O(1/n)$$

Now, combining these inequalities together we deduce

$$\mathcal{E}(f_{\beta}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}}) \leq \hat{r}_{\tilde{\mathcal{H}}}^* + \hat{r}_{\mathcal{H}}^* + 2\frac{D}{D-1}\lambda + O(1/n)$$

$$+ \mathcal{E}(\hat{f}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}})$$

$$\leq 2\hat{r}_{\mathcal{H}}^* + 2\frac{D}{D-1}\lambda + O(1/n)$$

$$+ \mathcal{E}(\hat{f}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}}).$$

The last inequality holds because the eigenvalues of the Gram-matrix for the reproduing kernel Hilbert space $\tilde{\mathcal{H}}$ decay faster than the eigenvalues of \mathcal{H} . As a result of this, we have that $\hat{r}^*_{\tilde{\mathcal{H}}} \leq \hat{r}^*_{\mathcal{H}}$.

Now, Theorem 5 implies that

$$\hat{r}_{\mathcal{H}}^* \le \min_{0 \le h \le n} \left(\frac{h}{n} * \frac{e_7}{n^2 \lambda^2} + \sqrt{\frac{1}{n} \sum_{i > h} \hat{\lambda}_i} \right). \tag{25}$$

There are two cases worth discussing here. On the one hand, if the eigenvalues of K decay exponentially, we have

$$\hat{r}_{\mathcal{H}}^* \le O\left(\frac{\log n}{n}\right)$$

by substituting $h = \lceil \log n \rceil$. Now, according to Caponnetto and De Vito (2007)

$$\mathcal{E}(\hat{f}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}}) \in O\left(\frac{\log n}{n}\right),$$

and, thus, if we set $\lambda \propto \frac{\log n}{n}$ then the expected risk rate can be upper bounded by

$$\mathcal{E}(f_{\beta}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}}) \in O\left(\frac{\log n}{n}\right).$$

On the other hand, if K has finitely many non-zero eigenvalues (t), we then have that

$$\hat{r}_{\mathcal{H}}^* \in O\left(\frac{1}{n}\right),$$

by substituting $h \geq t$. Moreover, in this case, $\mathcal{E}(\hat{f}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}}) \in O(\frac{1}{n})$ and setting $\lambda \propto \frac{1}{n}$, we deduce that

$$\mathcal{E}(f_{\beta}^{\lambda}) - \mathcal{E}(f_{\mathcal{H}}) \le O\left(\frac{1}{n}\right).$$

3.2 Learning with a Lipschitz Continuous Loss

We next consider kernel methods with Lipschitz continuous loss, examples of which include kernel support vector machines and kernel logistic regression. Similar to the squared error loss case, we approximate y_i with $g_{\beta}(x_i) = \mathbf{z}_{q,x_i}(\mathbf{v})^T \beta$ and formulate the following learning problem

$$g_{\beta}^{\lambda} = \underset{\beta \in \mathbb{R}^s}{\min} \ \frac{1}{n} \sum_{i=1}^n \mathbf{L}(y_i, \mathbf{z}_{q, x_i}(\mathbf{v})^T \beta) + \lambda s \|\beta\|_2^2.$$

3.2.1 WORST CASE ANALYSIS

The following theorem describes the trade-off between the selected number of features s and the expected risk of the estimator, providing an insight into the choice of s for Lipschitz continuous loss functions.

Theorem 6. Suppose that all the assumptions from Theorem 1 apply to the setting with a Lipschitz continuous loss. If

$$s \geq 5d_{\tilde{l}}\log\frac{(16d_{\mathbf{K}}^{\lambda})}{\delta},$$

then for all $\delta \in (0,1)$, with probability $1-\delta$, the expected risk of g^{λ}_{β} can be upper bounded as

$$\mathcal{E}(g_{\beta}^{\lambda}) \leq \mathcal{E}(g_{\mathcal{H}}) + \sqrt{2\lambda} + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right).$$
 (26)

This theorem, similar to Theorem 1, describes the relationship between s and $\mathcal{E}(g_{\beta}^{\lambda})$ in the Lipschitz continuous loss case. However, a key difference here is that the expected risk can only be upper bounded by $\sqrt{\lambda}$, requiring $\lambda \propto \frac{1}{n}$ in order to preserve the convergence properties of the risk.

Proof. The proof is similar to Theorem 1. In particular, we decompose the expected learning risk as

$$\mathcal{E}(g_{\beta}^{\lambda}) = \mathcal{E}(g_{\beta}^{\lambda}) - \hat{\mathcal{E}}(g_{\beta}^{\lambda}) \tag{27}$$

$$+\hat{\mathcal{E}}(g_{\beta}^{\lambda}) - \hat{\mathcal{E}}(g_{\mathcal{H}}) \tag{28}$$

$$+\hat{\mathcal{E}}(g_{\mathcal{H}}) - \mathcal{E}(g_{\mathcal{H}}) + \mathcal{E}(g_{\mathcal{H}}). \tag{29}$$

Now, (27) and (29) can be upper bounded similar to Theorem 1, through the Rademacher complexity bound from Lemma 2. For (28), we have

$$\begin{split} \hat{\mathcal{E}}(g_{\beta}^{\lambda}) - \hat{\mathcal{E}}(g_{\mathcal{H}}) &= \\ \frac{1}{n} \sum_{i=1}^{n} l(y_{i}, g_{\beta}^{\lambda}(x_{i})) - \frac{1}{n} \sum_{i=1}^{n} l(y_{i}, g_{\mathcal{H}}(x_{i})) &= \\ \frac{1}{n} \inf_{\|g_{\beta}\|} \sum_{i=1}^{n} l(y_{i}, g_{\beta}(x_{i})) - \frac{1}{n} \sum_{i=1}^{n} l(y_{i}, g_{\mathcal{H}}(x_{i})) \\ &\leq \inf_{\|g_{\beta}\|} \frac{1}{n} \sum_{i=1}^{n} |g_{\beta}(x_{i}) - g_{\mathcal{H}}(x_{i})| \\ &\leq \inf_{\|g_{\beta}\|} \sqrt{\frac{1}{n} \sum_{i=1}^{n} |g_{\beta}(x_{i}) - g_{\mathcal{H}}(x_{i})|^{2}} \\ &\leq \sup_{\|g\|} \inf_{\|g_{\beta}\|} \sqrt{\frac{1}{n} \|g - g_{\beta}\|_{2}^{2}} \\ &\leq \sqrt{2\lambda}. \end{split}$$

Corollaries 3 and 4 provide bounds for the cases of leverage weighted and plain RFF, respectively. The proofs are similar to the proofs of Corollaries 1 and 2.

Corollary 3. If the probability density function from Theorem 6 is the empirical ridge leverage score distribution $q^*(v)$, then the upper bound on the risk from Eq. (26) holds for all $s \geq 5d_{\mathbf{K}}^{\lambda} \log \frac{16d_{\mathbf{K}}^{\lambda}}{\delta}$.

In the three considered cases for the effective dimension of the problem $d_{\mathbf{K}}^{\lambda}$, Corollary 3 states that the statistical efficiency is preserved if the leverage weighted RFF strategy is used with $s \geq \log^2 n$, $s \geq n^{1/(2t)} \log n$, and $s \geq n \log n$, respectively. Again, significant computational savings can be achieved if the eigenvalues of the kernel matrix \mathbf{K} have either a geometric/exponential or a polynomial decay.

Corollary 4. If the probability density function from Theorem 6 is the spectral measure p(v) from Eq. (2), then the upper bound on the risk from Eq. (26) holds for all $s \geq 5\frac{z_0^2}{\lambda}\log\frac{(16d_{\mathbf{K}}^{\lambda})}{\delta}$.

Corollary 4 states that $n \log n$ features are required to attain $O(n^{-1/2})$ convergence rate of the expected risk with plain RFF, recovering results from Rahimi and Recht (2009). Similar to the analysis in the squared error loss case, Theorem 6 together with Corollaries 3 and 4 allows theoretically motivated trade-offs between the statistical and computational efficiency of the estimator g_{β}^{λ} .

3.2.2 REFINED ANALYSIS

It is generally difficult to achieve a better trade-off between computation cost and prediction accuracy in the Lipschitz continuous loss case. In order to do this, we need some further assumptions. Fortunately, for a loss function l, by assuming that the Bayes classifier function g_l^* is contained in the RKHS $\tilde{\mathcal{H}}$, we can do that.

Theorem 7. Suppose that all the assumptions from Theorem 1 apply to the setting with a Lipschitz continuous loss. In addition, we assume that $\lambda > 1/2$ and $g_l^* \in \tilde{\mathcal{H}}$ with $\mathcal{E}(g_l^*) = R^*$, we have that, for $\delta \in (0,1)$ with probability greater than $1-2e^{-\delta}$,

$$\mathcal{E}(g_{\beta}^{\lambda}) - \mathcal{E}(g_l^*) \le C_3 \frac{s}{n} \log \frac{1}{\lambda} + \frac{2c_1 \delta}{n},$$

where C_3 is a constant with respect to s, n and λ .

In realizable case, i.e., the Bayes classifier belongs to the RKHS spanned by the features, Theorem 7 describes a refined relationship between the learning risk and the number of features. In addition, it also implicitly states how the complexity of g_l^* can affects the learning risk convergence rate. Basically, if choosing $s = O(\sqrt{n})$ is sufficient to make the RKHS $\tilde{\mathcal{H}}$ large enough to include g_l^* , and we let $\lambda = O(1/\sqrt{n})$, we can then achieve the learning rate of $O(\frac{\log n}{\sqrt{n}})$ with only $O(\sqrt{n})$ features. To our knowledge, this is by far the first result that shows we can obtain computational savings in Lipschitz continuous loss case. Furthermore, if g_l^* has low complexity in the sense that with only finitely many features $c_s, c_s < \infty$, $\tilde{\mathcal{H}}$ can include g_l^* , then we can achieve O(1/n) convergence rate with only c_s features. That being said, however, we are in the realizable case which is somewhat limiting. Also, we lack of a specific way to describe the exact complexity of g_l^* in terms of number of features. Hence, how we can extend our analysis to the unrealizable case and how to analyze the complexity of g_l^* would be an interesting future direction.

To prove Theorem 7, we need the help of the following results on the analysis of the local Rademacher complexities from Bartlett et al. (2005).

Theorem 8. (Bartlett et al., 2005, Theorem 3.3) Let $\mathcal{F} := \{f : f(x) \in [a,b]\}$ be a class of functions and B be a constant. Assume that there is a functional $T : \mathcal{F} \to \mathbb{R}^+$ such that $Var(f) \le T(f) \le T(f)$

 $B\mathbb{E}(f)$ for all $f \in \mathcal{F}$ and for $\alpha \in [0,1]$, $T(\alpha f) \leq \alpha^2 T(f)$. Suppose ψ is a sub-root function with fixed point r^* and it satisfies for all $r \geq r^*$,

$$\psi(r) \ge BR_n\{f \in star(\mathcal{F}, 0) : T(f) \le r\},\$$

then we have $\forall f \in \mathcal{F}$, there is a constant D > 1 such that

$$\mathbb{E}(f) \le \frac{D}{D-1} \mathbb{E}_n(f) + \frac{6D}{B} r^* + \frac{c_1 \delta}{n},$$

with probability greater than $1 - e^{-\delta}$.

We would like to apply Theorem 8 to the decomposition of the learning risk in the Lipschitiz continuous loss case, namely Eqs. (27, 28, 29). In order to do that, we need three steps. The first step is to find a proper T functional, and the second one is to find the sub-root function ψ associated with T. Our final job is to find the unique fixed point of ψ . Hence, the following is devoted to solve these three problems.

To this end, we first notice that for all $l_{\tilde{f}} \in l_{\tilde{\mathcal{H}}}$, we have that $l_{\tilde{f}} \in [0, l_b]$ for some constant l_b as the loss is upper bounded. In addition, since we assume $z(v, x) < \infty$, we have for all $\tilde{f} \in \tilde{\mathcal{H}}$, $\|\tilde{f}\|_{\tilde{\mathcal{H}}} < \infty$. We let

$$\forall \tilde{f} \in \tilde{\mathcal{H}}, \ T(l_{\tilde{f}}) = \mathbb{E}(l_{\tilde{f}}^2) + \lambda \|\tilde{f}\|_{\tilde{\mathcal{H}}}^2 \mathbb{E}(l_{f^*}),$$

where $\lambda \in [0, \infty)$ is the regularization parameter and f^* is the Bayes classifier. Now immediately, we have $\operatorname{Var}(l_{\tilde{f}}) \leq T(l_{\tilde{f}})$. Also we have

$$T(l_{\tilde{f}}) = \mathbb{E}(l_{\tilde{f}}^2) + \lambda \|\tilde{f}\|_{\tilde{\mathcal{H}}}^2 \mathbb{E}(l_{f^*})$$

$$\leq l_b \mathbb{E}(l_{\tilde{f}}) + \lambda \|\tilde{f}\|_{\tilde{\mathcal{H}}}^2 \mathbb{E}(l_{\tilde{f}})$$

$$\leq (l_b + \lambda \|\tilde{f}\|_{\tilde{\mathcal{H}}}^2) \mathbb{E}(l_{\tilde{f}})$$

$$\leq B \mathbb{E}(l_{\tilde{f}}),$$

where $B<\infty$ is some constant. Also, it is easy to verify for $\alpha\in[0,1]$, $T(\alpha l_{\tilde{t}})\leq\alpha^2T(l_{\tilde{t}})$.

Now that we have a proper T functional, our next job is to find the sub-root function $\psi.$ We notice that

$$\begin{split} T(l_{\tilde{f}}) &= \mathbb{E}l_{\tilde{f}}^2 + \lambda \|\tilde{f}\|_{\tilde{\mathcal{H}}}^2 \mathbb{E}l_{f^*} \\ &\geq (\mathbb{E}l_{\tilde{f}})^2 + \lambda \|\tilde{f}\|_{\tilde{\mathcal{H}}}^2 \mathbb{E}l_{f^*} \\ &\geq (\mathbb{E}l_{\tilde{f}})^2 - (\mathbb{E}l_{f^*})^2 + \lambda \|\tilde{f}\|_{\tilde{\mathcal{H}}}^2 \mathbb{E}l_{\tilde{f}^*} \\ &= (\mathbb{E}l_{\tilde{f}} + \mathbb{E}l_{f^*})(\mathbb{E}l_{\tilde{f}} - \mathbb{E}l_{f^*}) + \lambda \|\tilde{f}\|_{\tilde{\mathcal{H}}}^2 \mathbb{E}l_{f^*} \\ &\geq 2\mathbb{E}l_{f^*}(\mathbb{E}l_{\tilde{f}} - \mathbb{E}l_{f^*}) + \lambda \|\tilde{f}\|_{\tilde{\mathcal{H}}}^2 \mathbb{E}l_{f^*} \\ &= 2\mathbb{E}l_{f^*}(\mathbb{E}l_{\tilde{f}} - \mathbb{E}l_{f^*} + \frac{\lambda}{2}\|\tilde{f}\|_{\tilde{\mathcal{H}}}) \end{split}$$

Since $\mathbb{E}l_{f^*} = \mathcal{E}(f^*) = R^*$, we have that

$$\{l_{\tilde{f}}: T(l_{\tilde{f}}) \le r\} \subseteq \{l_{\tilde{f}}: \mathbb{E}l_{\tilde{f}} - \mathbb{E}l_{f^*} + \frac{\lambda}{2} \|\tilde{f}\|_{\tilde{\mathcal{H}}} \le \frac{r}{2R^*}\}.$$

As a result,

$$R_{n}\{l_{\tilde{f}}: T(l_{\tilde{f}}) \leq r\} \leq R_{n}\{l_{\tilde{f}}: \mathbb{E}l_{\tilde{f}} - \mathbb{E}l_{f^{*}} + \frac{\lambda}{2} \|\tilde{f}\|_{\tilde{\mathcal{H}}} \leq \frac{r}{2R^{*}}\}$$

$$= R_{n}\{l_{\tilde{f}} - l_{f^{*}}: \mathbb{E}l_{\tilde{f}} - \mathbb{E}l_{f^{*}} + \frac{\lambda}{2} \|\tilde{f}\|_{\tilde{\mathcal{H}}} \leq \frac{r}{2R^{*}}\}$$

Now if we can upper bound the Rademacher complexity and the bound happens to be a sub-root function, then we could apply Theorem 8 directly. With this aim, we now define two function spaces as follows:

$$\tilde{H}_r := \{ \tilde{f} \in \tilde{\mathcal{H}} : \mathbb{E}l_{\tilde{f}} - \mathbb{E}l_{f^*} + \frac{\lambda}{2} \|\tilde{f}\|_{\tilde{\mathcal{H}}} \le \frac{r}{2R^*} \},$$

$$l_{\tilde{\mathcal{H}}_r} := \{ l_{\tilde{f}} - l_{f^*} : \tilde{f} \in \tilde{H}_r \}.$$

Our job is to find the upper bound of $R_n(l_{\tilde{\mathcal{H}}_r})$. To do this, we utilize the theories from entropy number.

Definition 2. (Steinwart and Christmann, 2008) Let $n \geq 1$ be an integer and $(E, \|\cdot\|)$ be a semi-normed space. We define the entropy number for E as

$$e_n(E, \|\cdot\|) := \inf\{\epsilon > 0 : \exists a_1, \dots, a_{2^n - 1} \in E, s.t.E \subset \bigcup_{i = 1}^{2^n - 1} B(a_i, \epsilon)\},\$$

where $B(a, \epsilon)$ is the ball with center a and radius ϵ under the norm $\|\cdot\|$.

Also, given a function space \mathcal{F} and a sequence of numbers $D := \{x_1, \dots, x_n\}$, we define the semi-norm as $\|\cdot\|_D = (\frac{1}{n}\sum_i f^2(x_i))^{\frac{1}{2}}, \forall f \in \mathcal{F}$. Equipped with entropy number, the next two lemmas give us the control on the empirical Rademacher complexity and the expected Rademacher complexity. The proofs (in Appendix F) are similar to Sun et al. (2018) with slight differences on the conditions of each function space.

Lemma 5. Assume $\lambda < \frac{1}{2}$, we have that

$$\hat{R}_n(l_{\tilde{\mathcal{H}}_r}) \le L\sqrt{\frac{s\log 16\log_2(1/\lambda)}{n}} \left(\rho + 9c_2\sqrt{r}\right),$$

where $\rho = \sup_{l_{\tilde{H}_n}} \|l_{\tilde{f}}\|_D$ and c_1 is some constant.

The next lemma gives us an upper bound on the expected Rademacher complexity.

Lemma 6. Assume $\lambda < 1/2$, we have

$$R_n(l_{\tilde{\mathcal{H}}_r}) \le C_1 \sqrt{\frac{s}{n} \log_2 \frac{1}{\lambda}} \sqrt{r} + C_2 \frac{s}{n} \log_2 \frac{1}{\lambda}.$$

Armed with the expected Rademacher complexity, we are now ready to prove Theorem 7.

Proof. We decompose the learning risk as follows:

$$\mathcal{E}(g_{\beta}^{\lambda}) = \mathcal{E}(g_{\beta}^{\lambda}) - \frac{D}{D-1}\hat{\mathcal{E}}(g_{\beta}^{\lambda}) + \frac{D}{D-1}\hat{\mathcal{E}}(g_{\beta}^{\lambda}) - \frac{D}{D-1}\hat{\mathcal{E}}(g^{*}) + \frac{D}{D-1}\hat{\mathcal{E}}(g^{*}) - \mathcal{E}(g^{*}) + \mathcal{E}(g^{*}).$$

$$(30)$$

Since we assume $g_{\mathcal{H}} \in \tilde{\mathcal{H}}$, we have Eq. (30) ≤ 0 . Hence, the learing risk is now:

$$\mathcal{E}(g_{\beta}^{\lambda}) \leq \mathcal{E}(g_{\beta}^{\lambda}) - \frac{D}{D-1}\hat{\mathcal{E}}(g_{\beta}^{\lambda})$$

$$+ \frac{D}{D-1}\hat{\mathcal{E}}(g^{*}) - \mathcal{E}(g^{*})$$

$$+ \mathcal{E}(g^{*}).$$

$$\leq |\mathcal{E}(g_{\beta}^{\lambda}) - \frac{D}{D-1}\hat{\mathcal{E}}(g_{\beta}^{\lambda})|$$

$$+ |\mathcal{E}(g^{*}) - \frac{D}{D-1}\hat{\mathcal{E}}(g^{*})|$$

$$+ \mathcal{E}(g^{*}).$$

$$(32)$$

$$+ \mathcal{E}(g^{*}).$$

We now would like to apply Theorem 8 to Eqs. (32,33). To this end, we have already discussed that $T(l_{\tilde{q}})$ is a proper T functional, moreover, if we let

$$\psi(r) = BC_1 \sqrt{\frac{s}{n} \log \frac{1}{\lambda}} \sqrt{r} + BC_2 \frac{s}{n} \log \frac{1}{\lambda},$$

we can easily check that $\psi(r)$ is a sub-root function and $\psi(r) \ge BR_n\{l_{\tilde{g}} : T(l_{\tilde{g}}) \le r\}$. Applying Theorem 8 to Eqs. (32,33), we have

$$\mathcal{E}(g_{\beta}^{\lambda}) - \mathcal{E}(g_{\mathcal{H}}) \le \frac{12D}{B}r^* + \frac{2c_1\delta}{n},$$

with probability greater than $1-2e^{-\delta}$. r^* is the unique fixed point of $\psi(r)$. By letting $\psi(r)=r$ and solve the equation, we can easily show that

$$r^* \le B^2 C_1^2 \frac{s}{n} \log \frac{1}{\lambda} + 2BC_2 \frac{s}{n} \log \frac{1}{\lambda}.$$

Hence, the excess risk can be upper bounded as

$$\mathcal{E}(g_{\beta}^{\lambda}) - \mathcal{E}(g^*) \le C_3 \frac{s}{n} \log \frac{1}{\lambda} + \frac{2c_1 \delta}{n}.$$

3.3 A Fast Approximation of Leverage Weighted RFF

As discussed in Sections 3.1 and 3.2, sampling according to the empirical ridge leverage score distribution (i.e., leverage weighted RFF) could speed up kernel methods. However, computing ridge leverage scores is as costly as inverting the Gram matrix. To address this computational shortcoming, we propose a simple algorithm to approximate the empirical ridge leverage score distribution and the leverage weights. In particular, we propose to first sample a pool of s features from the spectral measure $p(\cdot)$ and form the feature matrix $\mathbf{Z}_s \in \mathbb{R}^{n \times s}$ (Algorithm 1, lines 1-2). Then, the algorithm associates an approximate empirical ridge leverage score to each feature (Algorithm 1, lines 3-4) and samples a set of $l \ll s$ features from the pool proportional to the computed scores (Algorithm 1, line 5). The output of the algorithm can be compactly represented via the feature matrix $\mathbf{Z}_l \in \mathbb{R}^{n \times l}$ such

Algorithm 1 APPROXIMATE LEVERAGE WEIGHTED RFF

Input: sample of examples $\{(x_i, y_i)\}_{i=1}^n$, shift-invariant kernel function k, and regularization parameter λ

Output: set of features $\{(v_1, p_1), \cdots, (v_l, p_l)\}$ with l and each p_i computed as in lines 3-4

- 1: sample s features $\{v_1, \ldots, v_s\}$ from p(v)
- 2: create a feature matrix \mathbf{Z}_s such that the *i*th row of \mathbf{Z}_s is

$$[z(v_1,x_i),\cdots,z(v_s,x_i)]^T$$

3: associate with each feature v_i a real number p_i such that p_i is equal to the ith diagonal element of the matrix

$$\mathbf{Z}_{s}^{T}\mathbf{Z}_{s}((1/s)\mathbf{Z}_{s}^{T}\mathbf{Z}_{s}+n\lambda I)^{-1}$$

- 4: $l \leftarrow \sum_{i=1}^{s} p_i$ and $M \leftarrow \{(v_i, p_i/l)\}_{i=1}^{s}$ 5: sample l features from M using the multinomial distribution given by the vector $(p_1/l,\cdots,p_s/l)$

that the *i*th row of
$$\mathbf{Z}_l$$
 is given by $\mathbf{z}_{x_i}(\mathbf{v}) = \left[\sqrt{\frac{l}{p_1}}z(v_1, x_i), \cdots, \sqrt{\frac{l}{p_l}}z(v_l, x_i)\right]^T$.

The computational cost of Algorithm 1 is dominated by the operations in step 3. As $\mathbf{Z}_s \in \mathbb{R}^{n \times s}$, the multiplication of matrices $\mathbf{Z}_s^T \mathbf{Z}_s$ costs $O(ns^2)$ and inverting $\mathbf{Z}_s^T \mathbf{Z}_s + n\lambda I$ costs only $O(s^3)$. Hence, for $s \ll n$, the overall runtime is only $O(ns^2 + s^3)$. Moreover, $\mathbf{Z}_s^T \mathbf{Z}_s = \sum_{i=1}^n \mathbf{z}_{x_i}(\mathbf{v}) \mathbf{z}_{x_i}(\mathbf{v})^T$ and it is possible to store only the rank-one matrix $\mathbf{z}_{x_i}(\mathbf{v})\mathbf{z}_{x_i}(\mathbf{v})^T$ into the memory. Thus, the algorithm only requires to store an $s \times s$ matrix and can avoid storing \mathbf{Z}_s , which would incur a cost of $O(n \times s)$.

The following theorem gives the convergence rate for the expected risk of Algorithm 1 in the kernel ridge regression setting.

Theorem 9. Suppose the conditions from Theorem 1 apply to the regression problem defined with a shift-invariant kernel k, a sample of examples $\{(x_i, y_i)\}_{i=1}^n$, and a regularization parameter λ . Let s be the number of random Fourier features in the pool of features from Algorithm 1, sampled using the spectral measure $p(\cdot)$ from Eq. (2) and the regularization parameter λ . Denote with $f_1^{\lambda^*}$ the ridge regression estimator obtained using a regularization parameter λ^* and a set of random Fourier features $\{v_i\}_{i=1}^l$ returned by Algorithm 1. If

$$s \geq \frac{7z_0^2}{\lambda} \log \frac{(16d_{\mathbf{K}}^{\lambda})}{\delta}$$
 and $l \geq 5d_{\mathbf{K}}^{\lambda^*} \log \frac{(16d_{\mathbf{K}}^{\lambda^*})}{\delta}$,

then for all $\delta \in (0,1)$, with probability $1-\delta$, the expected risk of $\tilde{f}_l^{\lambda^*}$ can be upper bounded as

$$\mathcal{E}(\tilde{f}_l^{\lambda^*}) \leq \mathcal{E}(f_{\mathcal{H}}) + 2\lambda + 2\lambda^* + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right).$$

Moreover, this upper bound holds for $l \in \Omega(\frac{s}{n\lambda})$.

Theorem 9 bounds the expected risk of the ridge regression estimator over random features generated by Algorithm 1. We can now observe that using the minimax choice of the regularization parameter for kernel ridge regression $\lambda, \lambda^* \propto n^{-1/2}$, the number of features that Algorithm 1 needs to sample from the spectral measure of the kernel k is $s \in \Omega(\sqrt{n}\log n)$. Then, the ridge regression estimator $\tilde{f}_l^{\lambda^*}$ converges with the minimax rate to the hypothesis $f_{\mathcal{H}} \in \mathcal{H}$ for $l \in \Omega(\log n \cdot \log \log n)$. This is a significant improvement compared to the spectral measure sampling (plain RFF), which requires $\Omega(n^{3/2})$ features for in-sample training and $\Omega(\sqrt{n}\log n)$ for out-of-sample test predictions.

Proof. Suppose the examples $\{x_i, y_i\}_{i=1}^n$ are independent and identically distributed and that the kernel k can be decomposed as in Eq. (2). Let $\{v_i\}_{i=1}^s$ be an independent sample selected according to p(v). Then, using these s features we can approximate the kernel as

$$\tilde{k}(x,y) = \frac{1}{s} \sum_{i=1}^{s} z(v_i, x) z(v_i, y)$$

$$= \int_{V} z(v, x) z(v, y) d\hat{P}(v),$$
(34)

where \hat{P} is the empirical measure on $\{v_i\}_{i=1}^s$. Denote the reproducing kernel Hilbert space associated with kernel \tilde{k} by $\tilde{\mathcal{H}}$ and suppose that kernel ridge regression was performed with the approximate kernel \tilde{k} . From Theorem 1 and Corollary 2, it follows that if

$$s \ge \frac{7z_0^2}{\lambda} \log \frac{16d_{\mathbf{K}}^{\lambda}}{\delta},$$

then for all $\delta \in (0,1)$, with probability $1-\delta$, the risk convergence rate of the kernel ridge regression estimator based on random Fourier features can be upper bounded by

$$\mathcal{E}(f_{\alpha}^{\lambda}) \le 2\lambda + O\left(\frac{1}{\sqrt{n}}\right) + \mathcal{E}(f_{\mathcal{H}}).$$
 (35)

Let $f_{\tilde{\mathcal{H}}}$ be the function in the reproducing kernel Hilbert space $\tilde{\mathcal{H}}$ achieving the minimal risk, i.e., $\mathcal{E}(f_{\tilde{\mathcal{H}}}) = \inf_{f \in \tilde{\mathcal{H}}} \mathcal{E}(f)$. We now treat \tilde{k} as the actual kernel that can be decomposed via the expectation with respect to the empirical measure in Eq. (34) and re-sample features from the set $\{v_i\}_{i=1}^s$, but this time the sampling is performed using the optimal ridge leverage scores. As \tilde{k} is the actual kernel, it follows from Eq. (5) that the leverage function in this case can be defined by

$$l_{\lambda}(v) = p(v)\mathbf{z}_{v}(\mathbf{x})^{T}(\tilde{\mathbf{K}} + n\lambda I)^{-1}\mathbf{z}_{v}(\mathbf{x}).$$

Now, observe that

$$l_{\lambda}(v_i) = p(v_i)[\mathbf{Z}_s^T(\tilde{\mathbf{K}} + n\lambda I)^{-1}\mathbf{Z}_s]_{ii}$$

where $[A]_{ii}$ denotes the *i*th diagonal element of matrix A. As $\tilde{\mathbf{K}} = (1/s)\mathbf{Z}_s\mathbf{Z}_s^T$, then the Woodbury inversion lemma implies that

$$l_{\lambda}(v_i) = p(v_i) [\mathbf{Z}_s^T \mathbf{Z}_s (\frac{1}{s} \mathbf{Z}_s^T \mathbf{Z}_s + n\lambda I)^{-1}]_{ii}.$$

If we let $l_{\lambda}(v_i) = p_i$, then the optimal distribution for $\{v_i\}_{i=1}^s$ is multinomial with individual probabilities $q(v_i) = p_i/(\sum_{j=1}^s p_j)$. Hence, we can re-sample l features according to q(v) and

perform linear ridge regression using the sampled leverage weighted features. Denoting this estimator with $\tilde{f}_l^{\lambda^*}$ and the corresponding number of degrees of freedom with $d_{\tilde{\mathbf{K}}}^{\lambda} = \text{Tr}\tilde{\mathbf{K}}(\tilde{\mathbf{K}} + n\lambda)^{-1}$, we deduce (using Theorem 1 and Corollary 1)

$$\mathcal{E}(\tilde{f}_l^{\lambda^*}) \le 2\lambda^* + O\left(\frac{1}{\sqrt{n}}\right) + \mathcal{E}(f_{\tilde{\mathcal{H}}}),\tag{36}$$

with the number of features $l \propto d_{ ilde{\mathbf{K}}}^{\lambda}$.

As $f_{\tilde{\mathcal{H}}}$ is the function achieving the minimal risk over $\tilde{\mathcal{H}}$, we can conclude that $\mathcal{E}(f_{\tilde{\mathcal{H}}}) \leq \mathcal{E}(f_{\alpha}^{\lambda})$. Now, combining Eq. (35) and (36), we obtain the final bound on $\mathcal{E}(\tilde{f}_l^{\lambda^*})$.

Theorem 10 provides a generalization convergence analysis to kernel support vector machines and logistic regression. Compared to previous result, the convergence rate of the expected risk, however, is at a slightly slower $O(\sqrt{\lambda} + \sqrt{\lambda^*})$ rate due to the difference in the employed loss function (see Section 3.2).

Theorem 10. Suppose the conditions from Theorem 6 apply to a learning problem with Lipschitz continuous loss, a shift-invariant kernel k, a sample of examples $\{(x_i, y_i)\}_{i=1}^n$, and a regularization parameter λ . Let s be the number of random Fourier features in the pool of features from Algorithm 1, sampled using the spectral measure $p(\cdot)$ from Eq. (2) and the regularization parameter λ . Denote with $\tilde{g}_l^{\lambda^*}$ the estimator obtained using a regularization parameter λ^* and a set of random Fourier features $\{v_i\}_{i=1}^l$ returned by Algorithm 1. If

$$s \geq \frac{5z_0^2}{\lambda} \log \frac{(16d_{\mathbf{K}}^{\lambda})}{\delta} \quad and \quad l \geq 5d_{\mathbf{K}}^{\lambda^*} \log \frac{(16d_{\mathbf{K}}^{\lambda^*})}{\delta},$$

then for all $\delta \in (0,1)$, with probability $1-\delta$, the expected risk of $\tilde{g}_l^{\lambda^*}$ can be upper bounded as

$$\mathcal{E}(\tilde{g}_l^{\lambda^*}) \ \leq \ \mathcal{E}(g_{\mathcal{H}}) + \sqrt{2\lambda} + \sqrt{2\lambda^*} + \mathcal{O}\left(rac{1}{\sqrt{n}}
ight).$$

We conclude by pointing out that the proposed algorithm provides an interesting new trade-off between the computational cost and prediction accuracy. In particular, one can pay an upfront cost (same as plain RFF) to compute the leverage scores, re-sample significantly fewer features and employ them in the training, cross-validation, and prediction stages. This can reduce the computational cost for predictions at test points from $O(\sqrt{n}\log n)$ to $O(\log n \cdot \log\log n)$. Moreover, in the case where the amount of features with approximated leverage scores utilized is the same as in plain RFF, the prediction accuracy would be significantly improved as demonstrated in our experiments.

4. Numerical Experiments

In this section, we report the results of our numerical experiments (on both simulated and real-world datasets) aimed at validating our theoretical results and demonstrating the utility of Algorithm 1. We first verify our results through a simulation experiment. Specifically, we consider a spline kernel of order r where $k_{2r}(x,y)=1+\sum_{i=1}^{\infty}\frac{1}{m^{2r}}\cos 2\pi m(x-y)$ (also considered by Bach, 2017b; Rudi and Rosasco, 2017). If the marginal distribution of X is uniform on [0,1], we can show that $k_{2r}(x,y)=\int_0^1 z(v,x)z(v,y)q^*(v)dv$, where $z(v,x)=k_r(v,x)$ and $q^*(v)$ is also uniform on [0,1].

We let y be a Gaussian random variable with mean $f(x) = k_t(x, x_0)$ (for some $x_0 \in [0, 1]$) and variance σ^2 . We sample features according to $q^*(v)$ to estimate f and compute the excess risk. By Theorem 1 and Corollary 1, if the number of features is proportional to $d_{\mathbf{K}}^{\lambda}$ and $\lambda \propto n^{-1/2}$, we should expect the excess risk converging at $O(n^{-1/2})$, or at $O(n^{-1/3})$ if $\lambda \propto n^{-1/3}$. Figure 1 demonstrates this with different values of r and t.

Next, we make a comparison between the performances of leverage weighted (computed according to Algorithm 1) and plain RFF on real-world datasets. We use four datasets from Chang and Lin (2011) and Dheeru and Karra Taniskidou (2017) for this purpose, including two for regression and two for classification: CPU, KINEMATICS, COD-RNA and COVTYPE. Except KINEMATICS, the other three datasets were used in Yang et al. (2012) to investigate the difference between the Nyström method and plain RFF. We use the ridge regression and SVM package from Pedregosa et al. (2011) as a solver to perform our experiments. We evaluate the regression tasks using the root mean squared error and the classification ones using the average percentage of misclassified examples. The Gaussian/RBF kernel is used for all the datasets with hyper-parameter tuning via 5-fold inner cross validation. We have repeated all the experiments 10 times and reported the average test error for each dataset. Figure 2 compares the performances of leverage weighted and plain RFF. In regression tasks, we observe that the upper bound of the confidence interval for the root mean squared error corresponding to leverage weighted RFF is below the lower bound of the confidence interval for the error corresponding to plain RFF. Similarly, the lower bound of the confidence interval for the classification accuracy of leverage weighted RFF is (most of the time) higher than the upper bound on the confidence interval for plain RFF. This indicates that leverage weighted RFFs perform statistically significantly better than plain RFFs in terms of the learning accuracy and/or prediction error.

5. Discussion

We have investigated the generalization properties of learning with random Fourier features in the context of different kernel methods: kernel ridge regression, support vector machines, and kernel logistic regression. In particular, we have given generic bounds on the number of features required for consistency of learning with two sampling strategies: *leverage weighted* and *plain random Fourier features*. The derived convergence rates account for the complexity of the target hypothesis and the structure of the reproducing kernel Hilbert space with respect to the marginal distribution of a datagenerating process. In addition to this, we have also proposed an algorithm for fast approximation of empirical ridge leverage scores and demonstrated its superiority in both theoretical and empirical analyses.

For kernel ridge regression, Avron et al. (2017) and Rudi and Rosasco (2017) have extensively analyzed the performance of learning with random Fourier features. In particular, Avron et al. (2017) have shown that o(n) features are enough to guarantee a good estimator in terms of its *empirical risk*. The authors of that work have also proposed a modified data-dependent sampling distribution and demonstrated that a further reduction in the number of random Fourier features is possible for leverage weighted sampling. However, their results do not provide a convergence rate for the *expected risk* of the estimator which could still potentially imply that computational savings come at the expense of statistical efficiency. Furthermore, the modified sampling distribution can only be used in the 1D Gaussian kernel case. While Avron et al. (2017) focus on bounding the empirical risk

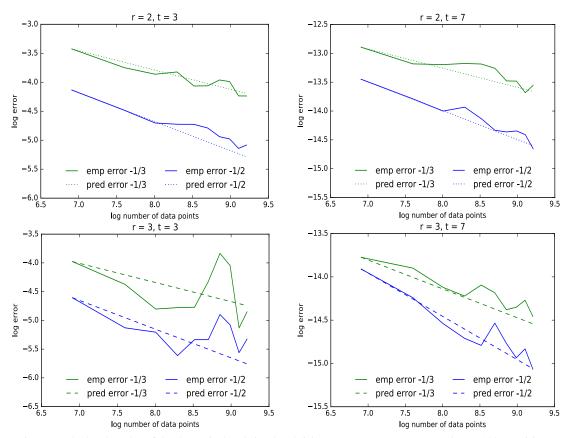


Figure 1: The log-log plot of the theoretical and simulated risk convergence rates, averaged over 100 repetitions.

of an estimator, Rudi and Rosasco (2017) give a comprehensive study of the generalization properties of random Fourier features for kernel ridge regression by bounding the expected risk of an estimator. The latter work for the first time shows that $\Omega(\sqrt{n}\log n)$ features are sufficient to guarantee the (kernel ridge regression) minimax rate and observes that further improvements to this result are possible by relying on a data-dependent sampling strategy. However, such a distribution is defined in a complicated way and it is not clear how one could devise a practical algorithm by sampling from it. While in our analysis of learning with random Fourier features we also bound the expected risk of an estimator, the analysis is not restricted to kernel ridge regression and covers other kernel methods such as support vector machines and kernel logistic regression. In addition to this, our derivations are much simpler compared to Rudi and Rosasco (2017) and provide sharper bounds in some cases. More specifically, we have demonstrated that $\Omega(\sqrt{n}\log\log n)$ features are sufficient to attain the minimax rate in the case where eigenvalues of the Gram matrix have a geometric/exponential decay. In other cases, we have recovered the results from Rudi and Rosasco (2017). Another important difference with respect to this work is that we consider a data-dependent sampling distribution based on empirical ridge leverage scores, showing that it can further reduce the number of features and in this way provide a more effective estimator.

In addition to the squared error loss, we also investigate the properties of learning with random Fourier features using the Lipschitz continuous loss functions. Both Rahimi and Recht (2009) and

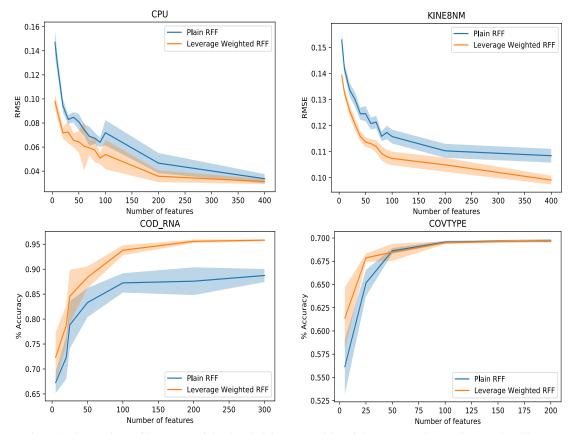


Figure 2: Comparison of leverage weighted and plain RFFs, with weights computed according to Algorithm 1.

Bach (2017b) have studied this problem setting and obtained that $\Omega(n)$ features are needed to ensure $O(1/\sqrt{n})$ expected risk convergence rate. Moreover, Bach (2017b) has defined an optimal sampling distribution by referring to the leverage score function based on the integral operator and shown that the number of features can be significantly reduced when the eigenvalues of a Gram matrix exhibit a fast decay. The $\Omega(n)$ requirement on the number of features is too restrictive and precludes any computational savings. Also, the optimal sampling distribution is typically intractable. In our analysis, through assuming the realizable case, we have demonstrated that for the first time, $O(\sqrt{n})$ features are possible to guarantee $O(\frac{1}{\sqrt{n}})$ risk convergence rate. In extreme cases, where the complexity of target function is small, constant features is enough to guarantee fast risk convergence. Moreover, We also provide a much simpler form of the empirical leverage score distribution and demonstrate that the number of features can be significantly smaller than n, without incurring any loss of statistical efficiency.

Having given risk convergence rates for learning with random Fourier features, we provide a fast and practical algorithm for sampling them in a data-dependent way, such that they approximate the ridge leverage score distribution. In the kernel ridge regression setting, our theoretical analysis demonstrates that compared to spectral measure sampling significant computational savings can be achieved while preserving the statistical properties of the estimator. We further test our findings on several different real-world datasets and verify this empirically. An interesting extension of our empirical analysis would be a thorough and comprehensive comparison of the proposed leverage

weighted sampling scheme to other recently proposed data-dependent strategies for selecting good features (e.g., Rudi et al., 2018; Zhang et al., 2018), as well as a comparison to the Nyström method.

Acknowledgments: We thank Fadhel Ayed, Qinyi Zhang and Anthony Caterini for fruitful discussion on some of the results as well as for proofreading of this paper. This work was supported by the EPSRC and MRC through the OxWaSP CDT programme (EP/L016710/1). Dino Oglic was supported in part by EPSRC grant EP/R012067/1. Zhu Li was supported in part by Huawei UK.

References

- Ahmed Alaoui and Michael W Mahoney. Fast randomized kernel ridge regression with statistical guarantees. In *Advances in Neural Information Processing Systems*, pages 775–783, 2015.
- 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, 2017.
- Francis Bach. Sharp analysis of low-rank kernel matrix approximations. In *Conference on Learning Theory*, pages 185–209, 2013.
- Francis Bach. Breaking the curse of dimensionality with convex neural networks. *Journal of Machine Learning Research*, 18(19):1–53, 2017a.
- Francis Bach. On the equivalence between kernel quadrature rules and random feature expansions. *Journal of Machine Learning Research*, 18(21):1–38, 2017b.
- Peter L Bartlett and Shahar Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- Peter L Bartlett, Olivier Bousquet, Shahar Mendelson, et al. Local rademacher complexities. *The Annals of Statistics*, 33(4):1497–1537, 2005.
- Alain Berlinet and Christine Thomas-Agnan. *Reproducing kernel Hilbert spaces in probability and statistics*. Springer Science & Business Media, 2011.
- Salomon Bochner. Vorlesungen über Fouriersche Integrale. In Akademische Verlagsgesellschaft, 1932.
- Andrea Caponnetto and Ernesto De Vito. Optimal rates for the regularized least-squares algorithm. *Foundations of Computational Mathematics*, 7(3):331–368, 2007.
- Chih-Chung Chang and Chih-Jen Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at http://www.csie.ntu.edu.tw/~cjlin/libsvm.
- Dua Dheeru and Efi Karra Taniskidou. UCI machine learning repository, 2017. URL http://archive.ics.uci.edu/ml.
- Trevor J Hastie. Generalized additive models. In Statistical models in S, pages 249–307. Routledge, 2017.
- Michael W Mahoney and Petros Drineas. CUR matrix decompositions for improved data analysis. *Proceedings of the National Academy of Sciences*, 106(3):697–702, 2009.
- Evert J. Nyström. Über die praktische Auflösung von Integralgleichungen mit Anwendungen auf Randwertaufgaben. *Acta Mathematica*, 1930.

- Dino Oglic and Thomas Gärtner. Greedy feature construction. In *Advances in Neural Information Processing Systems* 29, pages 3945–3953. Curran Associates, Inc., 2016.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pages 1177–1184, 2007.
- Ali Rahimi and Benjamin Recht. Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning. In *Advances in neural information processing systems*, pages 1313–1320, 2009.
- Alessandro Rudi and Lorenzo Rosasco. Generalization properties of learning with random features. In *Advances in Neural Information Processing Systems*, pages 3218–3228, 2017.
- Alessandro Rudi, Raffaello Camoriano, and Lorenzo Rosasco. Less is more: Nyström computational regularization. In *Advances in Neural Information Processing Systems*, pages 1657–1665, 2015.
- Alessandro Rudi, Luigi Carratino, and Lorenzo Rosasco. Falkon: An optimal large scale kernel method. In *Advances in Neural Information Processing Systems*, pages 3891–3901, 2017.
- Alessandro Rudi, Daniele Calandriello, Luigi Carratino, and Lorenzo Rosasco. On fast leverage score sampling and optimal learning. In *Advances in Neural Information Processing Systems*, pages 5672–5682, 2018.
- Walter Rudin. Fourier analysis on groups. Courier Dover Publications, 2017.
- Bernhard Schölkopf and Alexander J. Smola. *Learning with kernels: Support vector machines, regularization, optimization, and beyond.* MIT Press, 2001.
- Bernhard Scholkopf and Alexander J Smola. *Learning with kernels: support vector machines, regularization, optimization, and beyond.* MIT press, 2001.
- Bernhard Schölkopf, Koji Tsuda, and Jean-Philippe Vert. *Kernel methods in computational biology*. MIT press, 2004.
- Alexander J. Smola and Bernhard Schölkopf. Sparse greedy matrix approximation for machine learning. In *Proceedings of the 17th International Conference on Machine Learning*, 2000.
- Bharath Sriperumbudur and Zoltán Szabó. Optimal rates for random Fourier features. In *Advances in Neural Information Processing Systems*, pages 1144–1152, 2015.
- Ingo Steinwart and Andreas Christmann. Support vector machines. Springer Science & Business Media, 2008.
- Yitong Sun, Anna Gilbert, and Ambuj Tewari. But how does it work in theory? linear svm with random features. In *Advances in Neural Information Processing Systems*, pages 3379–3388, 2018.
- Dougal J Sutherland and Jeff Schneider. On the error of random Fourier features. In *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence*, pages 862–871. AUAI Press, 2015.
- Joel A Tropp. An introduction to matrix concentration inequalities. *Foundations and Trends*® *in Machine Learning*, 8(1-2):1–230, 2015.
- Christopher K. I. Williams and Matthias Seeger. Using the Nyström method to speed up kernel machines. In *Advances in Neural Information Processing Systems* 13. 2001.

LI, TON, OGLIC AND SEJDINOVIC

- Tianbao Yang, Yu-Feng Li, Mehrdad Mahdavi, Rong Jin, and Zhi-Hua Zhou. Nyström method vs random Fourier features: A theoretical and empirical comparison. In *Advances in neural information processing systems*, pages 476–484, 2012.
- Jian Zhang, Avner May, Tri Dao, and Christopher Ré. Low-precision random fourier features for memory-constrained kernel approximation. *arXiv preprint arXiv:1811.00155*, 2018.
- Yuchen Zhang, John Duchi, and Martin Wainwright. Divide and conquer kernel ridge regression: A distributed algorithm with minimax optimal rates. *The Journal of Machine Learning Research*, 16(1):3299–3340, 2015.

Appendix A. Bernstein Inequality

The next lemma is the matrix Bernstein inequality, which is a restatement of Corollary 7.3.3 in Tropp (2015).

Lemma 7. (Bernstein inequality, Tropp, 2015, Corollary 7.3.3) Let \mathbf{R} be a fixed $d_1 \times d_2$ matrix over the set of complex/real numbers. Suppose that $\{\mathbf{R}_1, \cdots, \mathbf{R}_n\}$ is an independent and identically distributed sample of $d_1 \times d_2$ matrices such that

$$\mathbb{E}[\mathbf{R}_i] = \mathbf{R} \quad and \quad \|\mathbf{R}_i\|_2 \le L,$$

where L > 0 is a constant independent of the sample. Furthermore, let $\mathbf{M}_1, \mathbf{M}_2$ be semidefinite upper bounds for the matrix-valued variances

$$\operatorname{Var}_{1}[\mathbf{R}_{i}] \leq \mathbb{E}[\mathbf{R}_{i}\mathbf{R}_{i}^{T}] \leq \mathbf{M}_{1}$$

 $\operatorname{Var}_{2}[\mathbf{R}_{i}] \leq \mathbb{E}[\mathbf{R}_{i}^{T}\mathbf{R}_{i}] \leq \mathbf{M}_{2}.$

Let $m = \max(\|\mathbf{M}_1\|_2, \|\mathbf{M}_2\|_2)$ and $d = \frac{Tr(\mathbf{M}_1) + Tr(\mathbf{M}_2)}{m}$. Then, for $\epsilon \geq \sqrt{m/n} + 2L/3n$, we can bound

$$\bar{\mathbf{R}}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{R}_i$$

around its mean using the concentration inequality

$$P(\|\bar{\mathbf{R}}_n - \mathbf{R}\|_2 \ge \epsilon) \le 4d \exp\left(\frac{-n\epsilon^2/2}{m + 2L\epsilon/3}\right).$$

Appendix B. Proof of Theorem 2

The following two lemmas are required for our proof of Theorem 2, presented subsequently.

Lemma 8. Suppose that the assumptions from Theorem 2 hold and let $\epsilon \geq \sqrt{\frac{m}{s}} + \frac{2L}{3s}$ with constants m and L (see the proof for explicit definition). If the number of features

$$s \ge d_{\tilde{l}}(\frac{1}{\epsilon^2} + \frac{2}{3\epsilon})\log\frac{16d_{\mathbf{K}}^{\lambda}}{\delta},$$

then for all $\delta \in (0,1)$, with probability greater than $1-\delta$,

$$-\epsilon \mathbf{I} \preceq (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} (\tilde{\mathbf{K}} - \mathbf{K}) (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} \preceq \epsilon \mathbf{I}.$$

Proof. Following the derivations in Avron et al. (2017), we utilize the matrix Bernstein concentration inequality to prove the result. More specifically, we observe that

$$(\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} \tilde{\mathbf{K}} (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} =$$

$$\frac{1}{s} \sum_{i=1}^{s} (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} \mathbf{z}_{q,v_i} (\mathbf{x}) \mathbf{z}_{q,v_i} (\mathbf{x})^{T} (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} =$$

$$\frac{1}{s} \sum_{i=1}^{s} \mathbf{R}_{i} =: \bar{\mathbf{R}}_{s},$$

with

$$\mathbf{R}_i = (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} \mathbf{z}_{q,v_i}(\mathbf{x}) \mathbf{z}_{q,v_i}(\mathbf{x})^T (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}}.$$

Now, observe that

$$\mathbf{R} = \mathbb{E}[\mathbf{R}_i] = (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} \mathbf{K} (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}}.$$

The operator norm of \mathbf{R}_i is equal to

$$\|(\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}}\mathbf{z}_{q,v_i}(\mathbf{x})\mathbf{z}_{q,v_i}(\mathbf{x})^T(\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}}\|_2.$$

As $\mathbf{z}_{q,v_i}(\mathbf{x})\mathbf{z}_{q,v_i}(\mathbf{x})^T$ is a rank one matrix, we have that the operator norm of this matrix is equal to its trace, i.e.,

$$\begin{split} &\|\mathbf{R}_i\|_2 = \\ &\operatorname{Tr}((\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} \mathbf{z}_{q,v_i}(\mathbf{x}) \mathbf{z}_{q,v_i}(\mathbf{x})^T (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}}) = \\ &\frac{p(v_i)}{q(v_i)} \operatorname{Tr}((\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} \mathbf{z}_{v_i}(\mathbf{x}) \mathbf{z}_{v_i}(\mathbf{x})^T (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}}) = \\ &\frac{p(v_i)}{q(v_i)} \operatorname{Tr}(\mathbf{z}_{v_i}(\mathbf{x})^T (\mathbf{K} + n\lambda \mathbf{I})^{-1} \mathbf{z}_{v_i}(\mathbf{x})) = \\ &\frac{l_{\lambda}(v_i)}{q(v_i)} =: L_i \quad \text{and} \quad L := \sup_i L_i. \end{split}$$

On the other hand,

$$\begin{aligned} &\mathbf{R}_{i}\mathbf{R}_{i}^{T} = \\ &(\mathbf{K} + n\lambda\mathbf{I})^{-\frac{1}{2}}\mathbf{z}_{q,v_{i}}(\mathbf{x})\mathbf{z}_{q,v_{i}}(\mathbf{x})^{T}(\mathbf{K} + n\lambda\mathbf{I})^{-1}\mathbf{z}_{q,v_{i}}(\mathbf{x}) \\ &\cdot \mathbf{z}_{q,v_{i}}(\mathbf{x})^{T}(\mathbf{K} + n\lambda\mathbf{I})^{-\frac{1}{2}} = \\ &\frac{p(v_{i})l_{\lambda}(v_{i})}{q^{2}(v_{i})}(\mathbf{K} + n\lambda\mathbf{I})^{-\frac{1}{2}}\mathbf{z}_{v_{i}}(\mathbf{x})\mathbf{z}_{v_{i}}(\mathbf{x})^{T}(\mathbf{K} + n\lambda\mathbf{I})^{-\frac{1}{2}} \leq \\ &\frac{\tilde{l}(v_{i})}{q(v_{i})}\frac{p(v_{i})}{q(v_{i})}(\mathbf{K} + n\lambda\mathbf{I})^{-\frac{1}{2}}\mathbf{z}_{v_{i}}(\mathbf{x})\mathbf{z}_{v_{i}}(\mathbf{x})^{T}(\mathbf{K} + n\lambda\mathbf{I})^{-\frac{1}{2}} = \\ &d_{\tilde{l}}\frac{p(v_{i})}{q(v_{i})}(\mathbf{K} + n\lambda\mathbf{I})^{-\frac{1}{2}}\mathbf{z}_{v_{i}}(\mathbf{x})\mathbf{z}_{v_{i}}(\mathbf{x})^{T}(\mathbf{K} + n\lambda\mathbf{I})^{-\frac{1}{2}}. \end{aligned}$$

From the latter inequality, we obtain that

$$\mathbb{E}[\mathbf{R}_i \mathbf{R}_i^T] \leq d_{\tilde{l}} (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} \mathbf{K} (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} =: \mathbf{M}_1.$$

We also have the following two inequalities

$$m = \|\mathbf{M}_1\|_2 = d_{\tilde{l}} \frac{\lambda_1}{\lambda_1 + n\lambda} =: d_{\tilde{l}} d_1$$
$$d = \frac{2 \operatorname{Tr}(\mathbf{M}_1)}{m} = 2 \frac{\lambda_1 + n\lambda}{\lambda_1} d_{\mathbf{K}}^{\lambda} = 2 d_1^{-1} d_{\mathbf{K}}^{\lambda}.$$

We are now ready to apply the matrix Bernstein concentration inequality. More specifically, for $\epsilon \geq \sqrt{m/s} + 2L/3s$ and for all $\delta \in (0,1)$, with probability $1-\delta$, we have that

$$P(\|\bar{\mathbf{R}}_s - \mathbf{R}\|_2 \ge \epsilon) \le 4d \exp\left(\frac{-s\epsilon^2/2}{m + 2L\epsilon/3}\right)$$
$$= 8d_1^{-1}d_{\mathbf{K}}^{\lambda} \exp\left(\frac{-s\epsilon^2/2}{d_{\bar{l}}d_1 + d_{\bar{l}}2\epsilon/3}\right)$$
$$\le 16d_{\mathbf{K}}^{\lambda} \exp\left(\frac{-s\epsilon^2}{d_{\bar{l}}(1 + 2\epsilon/3)}\right) \le \delta.$$

In the third line, we have used the assumption that $n\lambda \leq \lambda_1$ and, consequently, $d_1 \in [1/2, 1)$.

Remark: We note here that the two considered sampling strategies lead to two different results. In particular, if we let $\tilde{l}(v) = l_{\lambda}(v)$ then $q(v) = l_{\lambda}(v)/d_{\mathbf{K}}^{\lambda}$, i.e., we are sampling proportional to the ridge leverage scores. Thus, the leverage weighted random Fourier features sampler requires

$$s \ge d_{\mathbf{K}}^{\lambda} \left(\frac{1}{\epsilon^2} + \frac{2}{3\epsilon}\right) \log \frac{16d_{\mathbf{K}}^{\lambda}}{\delta}.$$
 (37)

Alternatively, we can opt for the plain random Fourier feature sampling strategy by taking $\tilde{l}(v) = z_0^2 p(v)/\lambda$, with $l_{\lambda}(v) \leq z_0^2 p(v)/\lambda$. Then, the plain random Fourier features sampling scheme requires

$$s \ge \frac{z_0^2}{\lambda} \left(\frac{1}{\epsilon^2} + \frac{2}{3\epsilon}\right) \log \frac{16d_{\mathbf{K}}^{\lambda}}{\delta}.$$
 (38)

Thus, the leverage weighted random Fourier features sampling scheme can dramatically change the required number of features, required to achieve a predefined matrix approximation error in the operator norm.

Lemma 9. Let $f \in \mathcal{H}$, where \mathcal{H} is the RKHS associated with a kernel k. Let $x_1, \dots, x_n \in \mathcal{X}$ be a set of instances with $x_i \neq x_j$ for all $i \neq j$. Denote with $\mathbf{f}_x = [f(x_1), \dots, f(x_n)]^T$ and let \mathbf{K} be the Gram-matrix of the kernel k given by the provided set of instances. Then,

$$\mathbf{f}_x^T \mathbf{K}^{-1} \mathbf{f}_x \le 1.$$

Proof. For a vector $\mathbf{a} \in \mathbb{R}^n$ we have that

$$\mathbf{a}^{T} \mathbf{f}_{x} \mathbf{f}_{x}^{T} \mathbf{a} = \left(\mathbf{f}_{x}^{T} \mathbf{a}\right)^{2} = \left(\sum_{i=1}^{n} a_{i} f(x_{i})\right)^{2}$$

$$= \left(\sum_{i=1}^{n} a_{i} \int_{\mathcal{V}} g(v) z(v, x_{i}) d\tau(v)\right)^{2}$$

$$= \left(\int_{\mathcal{V}} g(v) \mathbf{z}_{v}(\mathbf{x})^{T} \mathbf{a} d\tau(v)\right)^{2}$$

$$\leq \int_{\mathcal{V}} g(v)^{2} d\tau(v) \int_{\mathcal{V}} (\mathbf{z}_{v}(\mathbf{x})^{T} \mathbf{a})^{2} d\tau(v)$$

$$= \int_{\mathcal{V}} \mathbf{a}^{T} \mathbf{z}_{v}(\mathbf{x}) \mathbf{z}_{v}(\mathbf{x})^{T} \mathbf{a} d\tau(v)$$

$$= \mathbf{a}^{T} \int_{\mathcal{V}} \mathbf{z}_{v}(\mathbf{x}) \mathbf{z}_{v}(\mathbf{x})^{T} d\tau(v) \mathbf{a}$$

$$= \mathbf{a}^{T} \mathbf{K} \mathbf{a}.$$

The third equality is due to the fact that, for all $f \in \mathcal{H}$, we have that $f(x) = \int_{\mathcal{V}} g(v)z(v,x)p(v)dv$ $(\forall x \in \mathcal{X})$ and

$$||f||_{\mathcal{H}} = \min_{\left\{g \mid f(x) = \int_{\mathcal{V}} g(v)z(v,x)p(v)dv\right\}} ||g||_{L_2(d\tau)}.$$

The first inequality, on the other hand, follows from the Cauchy-Schwarz inequality. The bound implies that $\mathbf{f}_x \mathbf{f}_x^T \leq \mathbf{K}$ and, consequently, we derive $\mathbf{f}_x^T \mathbf{K}^{-1} \mathbf{f}_x \leq 1$.

Now we are ready to prove Theorem 2.

Proof. Our goal is to minimize the following objective:

$$\frac{1}{n} \|\mathbf{f}_x - \mathbf{Z}_q \beta\|_2^2 + s\lambda \|\beta\|_2^2. \tag{39}$$

To find the minimizer, we can directly take the derivative with respect to β and, thus, obtain

$$\beta = \frac{1}{s} (\frac{1}{s} \mathbf{Z}_q^T \mathbf{Z}_q + n\lambda \mathbf{I})^{-1} \mathbf{Z}_q^T \mathbf{f}_x$$
$$= \frac{1}{s} \mathbf{Z}_q^T (\frac{1}{s} \mathbf{Z}_q \mathbf{Z}_q^T + n\lambda \mathbf{I})^{-1} \mathbf{f}_x$$
$$= \frac{1}{s} \mathbf{Z}_q^T (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-1} \mathbf{f}_x,$$

where the second equality follows from the Woodbury inversion lemma.

Substituting β into Eq. (39), we transform the first part as

$$\begin{split} \frac{1}{n} \|\mathbf{f}_x - \mathbf{Z}_q \boldsymbol{\beta}\|_2^2 &= \frac{1}{n} \|\mathbf{f}_x - \frac{1}{s} \mathbf{Z}_q \mathbf{Z}_q^T (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-1} \mathbf{f}_x\|_2^2 \\ &= \frac{1}{n} \|\mathbf{f}_x - \tilde{\mathbf{K}} (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-1} \mathbf{f}_x\|_2^2 \\ &= \frac{1}{n} \|n\lambda (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-1} \mathbf{f}_x\|_2^2 \\ &= n\lambda^2 \mathbf{f}_x^T (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-2} \mathbf{f}_x. \end{split}$$

On the other hand, the second part can be transformed as

$$s\lambda \|\beta\|_{2}^{2} = s\lambda \frac{1}{s^{2}} \mathbf{f}_{x}^{T} (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-1} \mathbf{Z}_{q} \mathbf{Z}_{q}^{T} (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-1} \mathbf{f}_{x}$$

$$= \lambda \mathbf{f}_{x}^{T} (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-1} \tilde{\mathbf{K}} (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-1} \mathbf{f}_{x}$$

$$= \lambda \mathbf{f}_{x}^{T} (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-1} (\tilde{\mathbf{K}} + n\lambda \mathbf{I}) (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-1} \mathbf{f}_{x}$$

$$- n\lambda^{2} \mathbf{f}_{x}^{T} (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-2} \mathbf{f}_{x}$$

$$= \lambda \mathbf{f}_{x}^{T} (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-1} \mathbf{f}_{x} - n\lambda^{2} \mathbf{f}_{x}^{T} (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-2} \mathbf{f}_{x}.$$

Now, summing up the first and the second part, we deduce

$$\frac{1}{n} \|\mathbf{f}_{x} - \mathbf{Z}_{q}\beta\|_{2}^{2} + s\lambda \|\beta\|_{2}^{2} =
\lambda \mathbf{f}_{x}^{T} (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-1} \mathbf{f}_{x} =
\lambda \mathbf{f}_{x}^{T} (\mathbf{K} + n\lambda \mathbf{I} + \tilde{\mathbf{K}} - \mathbf{K})^{-1} \mathbf{f}_{x} =
\lambda \mathbf{f}_{x}^{T} (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} (\mathbf{I} + (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} (\tilde{\mathbf{K}} - \mathbf{K})
\cdot (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}})^{-1} (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} \mathbf{f}_{x}.$$

From Lemma 8, it follows that when

$$s \ge d_{\tilde{l}}(\frac{1}{\epsilon^2} + \frac{2}{3\epsilon})\log\frac{16d_{\mathbf{K}}^{\lambda}}{\delta}$$

then $(\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} (\tilde{\mathbf{K}} - \mathbf{K}) (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} \succeq -\epsilon \mathbf{I}$.

We can now upper bound the error as (with $\epsilon = 1/2$):

$$\lambda \mathbf{f}_{x}^{T} (\tilde{\mathbf{K}} + n\lambda \mathbf{I})^{-1} \mathbf{f}_{x} \leq \lambda \mathbf{f}_{x}^{T} (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} (1 - \epsilon)^{-1} (\mathbf{K} + n\lambda \mathbf{I})^{-\frac{1}{2}} \mathbf{f}_{x} = (1 - \epsilon)^{-1} \lambda \mathbf{f}_{x}^{T} (\mathbf{K} + n\lambda \mathbf{I})^{-1} \mathbf{f}_{x} \leq (1 - \epsilon)^{-1} \lambda \mathbf{f}_{x}^{T} \mathbf{K}^{-1} \mathbf{f}_{x} \leq 2\lambda,$$

where in the last inequality we have used Lemma 9. Moreover, we have that

$$s\|\beta\|_{2}^{2} = \mathbf{f}_{x}^{T}(\tilde{\mathbf{K}} + n\lambda\mathbf{I})^{-1}\mathbf{f}_{x} - n\lambda\mathbf{f}_{x}^{T}(\tilde{\mathbf{K}} + n\lambda\mathbf{I})^{-2}\mathbf{f}_{x} \le \mathbf{f}_{x}^{T}(\tilde{\mathbf{K}} + n\lambda\mathbf{I})^{-1}\mathbf{f}_{x} \le (1 - \epsilon)^{-1}\mathbf{f}_{x}^{T}\mathbf{K}^{-1}\mathbf{f}_{x} \le 2.$$

Hence, the squared norm of our approximated function is bounded by $\|\tilde{f}\|_{\tilde{\mathcal{H}}}^2 \leq s\|\beta\|_2^2 \leq 2$. As such, problem (39) can now be written as $\min_{\beta} (1/n) \|\mathbf{f}_x - \tilde{\mathbf{f}}_{\beta}\|_2^2$ subject to $\|\tilde{f}\|_{\tilde{\mathcal{H}}}^2 \leq s\|\beta\|_2^2 \leq 2$, which is equivalent to

$$\sup_{\|f\|_{\mathcal{H}} \le 1} \inf_{\sqrt{s} \|\beta\|_2 \le \sqrt{2}} \frac{1}{n} \|\mathbf{f}_x - \mathbf{Z}\beta\|_2^2,$$

and we have shown that this can be upper bounded by 2λ .

Appendix C. Proof of Lemma 3

Proof. The solution of problem (7) can be derived as

$$f_{\beta}^{\lambda} = \tilde{\mathbf{K}}(\tilde{\mathbf{K}} + n\lambda I)^{-1}Y = \frac{1}{s}\mathbf{Z}_{q}\mathbf{Z}_{q}^{T}(\frac{1}{s}\mathbf{Z}_{q}\mathbf{Z}_{q}^{T} + n\lambda I)^{-1}Y.$$

For all $f \in \mathcal{H}$, let $\mathbf{f} = [f(x_1), \dots, f(x_n)]^T$. Define $\mathcal{H}_x := \{\mathbf{f} \mid f \in \mathcal{H}\}$. Then we can see that \mathcal{H}_x is a subspace of \mathbb{R}^n . Since $Y \in \mathbb{R}^n$, we know there exists an orthogonal projection operator P such that for any vector $Z \in \mathbb{R}^n$, PZ is the projection of Z into \mathcal{H}_x . In particular, we have $\hat{f}^{\lambda} = PY$. In addition, let $\alpha \in \mathbb{R}^n$ and observe that $P\mathbf{K}\alpha = \mathbf{K}\alpha$, as $\mathbf{K}\alpha \in \mathcal{H}_x$. As such, we have that $(I - P)\mathbf{K}\alpha = 0$ for all $\alpha \in \mathbb{R}^n$, implying that $(I - P)\mathbf{K} = 0$. Hence, we have

$$\langle Y - \hat{f}^{\lambda}, f_{\beta}^{\lambda} - \hat{f}^{\lambda} \rangle =$$

$$\langle Y - PY, \tilde{\mathbf{K}}(\tilde{\mathbf{K}} + n\lambda I)^{-1}Y - PY \rangle =$$

$$\langle (I - P)Y, \tilde{\mathbf{K}}(\tilde{\mathbf{K}} + n\lambda I)^{-1}Y \rangle - \langle (I - P)Y, PY \rangle =$$

$$Y^{T}(I - P)\tilde{\mathbf{K}}(\tilde{\mathbf{K}} + n\lambda I)^{-1}Y - Y^{T}(I - P)PY =$$

$$\frac{1}{s}Y^{T}(I - P)\mathbf{Z}_{q}\mathbf{Z}_{q}^{T}(\mathbf{Z}_{q}\mathbf{Z}_{q}^{T} + n\lambda I)^{-1}Y.$$
(40)

The last equality follows from $(I - P)P = P - P^2 = 0$.

We know that the kernel function admits a decomposition as in Eq. (2). Hence, we can express ${\bf K}$ as

$$\mathbf{K} = \int_{\mathcal{V}} \mathbf{z}_v(\mathbf{x}) \mathbf{z}_v(\mathbf{x})^T d\tau(v),$$

where $\mathbf{z}_v(\mathbf{x}) = [z(v, x_1), \cdots, z(v, x_n)]^T$.

Note that we have $(I - P)\mathbf{K} = 0$, which further implies that $(I - P)\mathbf{K}(I - P) = 0$. As a result, we have the following:

$$0 = \operatorname{Tr}[(I - P)\mathbf{K}(I - P)]$$

$$= \operatorname{Tr}\left[(I - P)\int_{\mathcal{V}} \mathbf{z}_{v}(\mathbf{x})\mathbf{z}_{v}(\mathbf{x})^{T}d\tau(v)(I - P)\right]$$

$$= \operatorname{Tr}\left[\int_{\mathcal{V}} (I - P)\mathbf{z}_{v}(\mathbf{x})\mathbf{z}_{v}(\mathbf{x})^{T}(I - P)d\tau(v)\right]$$

$$= \int_{\mathcal{V}} \operatorname{Tr}[(I - P)\mathbf{z}_{v}(\mathbf{x})\mathbf{z}_{v}(\mathbf{x})^{T}(I - P)]d\tau(v)$$

$$= \int_{\mathcal{V}} \|(I - P)\mathbf{z}_{v}(\mathbf{x})\|_{2}^{2}d\tau(v)$$

$$= \int_{\mathcal{V}} \|(I - P)\mathbf{z}_{q,v}(\mathbf{x})\|_{2}^{2}\frac{p(v)}{q(v)}d\tau_{q}(v), \tag{41}$$

where $\mathbf{z}_{q,v}(\mathbf{x}) = \sqrt{p(v)/q(v)}\mathbf{z}_v(\mathbf{x})$. Hence, we have $\|(I-P)\mathbf{z}_{q,v}(\mathbf{x})\|_2^2 = 0$ almost surely (a.s.) with respect to measure $d\tau_q$, which further shows that $(I-P)\mathbf{z}_{q,v}(\mathbf{x}) = \mathbf{0}$ a.s. For any $\alpha \in \mathbb{R}^s$, we have:

$$\alpha^T Y^T (I - P) \mathbf{Z}_q = \sum_{i=1} \alpha_i Y^T (I - P) \mathbf{z}_{q, v_i}(\mathbf{x}) = 0.$$

We now let $\alpha = Y^T(I - P)\mathbf{Z}_q$ and obtain

$$||Y^T(I-P)\mathbf{Z}_q||_2^2 = 0.$$

Returning back to Eq. (40), we have that

$$\langle Y - \hat{f}^{\lambda}, f_{\beta}^{\lambda} - \hat{f}^{\lambda} \rangle =$$

$$\frac{1}{s} Y^{T} (I - P) \mathbf{Z}_{q} \mathbf{Z}_{q}^{T} (\mathbf{Z}_{q} \mathbf{Z}_{q}^{T} + n\lambda I)^{-1} Y.$$

Now, observe that

$$|Y^{T}(I-P)\mathbf{Z}_{q}\mathbf{Z}_{q}^{T}(\mathbf{Z}_{q}\mathbf{Z}_{q}^{T}+n\lambda I)^{-1}Y| \leq ||Y^{T}(I-P)\mathbf{Z}_{q}||_{2}^{2}||\mathbf{Z}_{q}^{T}(\mathbf{Z}_{q}\mathbf{Z}_{q}^{T}+n\lambda I)^{-1}Y||_{2}^{2} = 0.$$

Hence, we conclude that $\langle Y - \hat{f}^{\lambda}, f^{\lambda}_{\beta} - \hat{f}^{\lambda} \rangle = 0$.

Appendix D. Sub-root Function & Square Loss

In this section, we first define the subroot function in Definition 3 and state its property in Lemma 10.

Definition 3. Let $\psi:[0,\infty)\to [0,\infty)$ be a function. Then $\psi(r)$ is called sub-root if it is nondecreasing and $\frac{\psi(r)}{r}$ is nonincreasing for r>0.

Lemma 10. (Bartlett et al., 2005, Lemma 3.2) If $\psi(r)$ is a sub-root function, then $\psi(r) = r$ has a unique positive solution. In addition, we have that $\psi(r) \leq r$ for all r > 0 if and only if $r^* \leq r$, where r^* is the solution of the equation.

Lemma 11. (Bartlett et al., 2005, Section 5.2) Let l be the squared error loss function and \mathcal{H} a convex and uniformly bounded hypothesis space. Assume that for every probability distribution P in a class of data-generating distributions, there is an $f^* \in \mathcal{H}$ such that $\mathcal{E}(f^*) = \inf_{f \in \mathcal{H}} \mathcal{E}(f)$. Then, there exists a constant $B \geq 1$ such that for all $f \in \mathcal{H}$ and for every probability distribution P

$$\mathbb{E}(f - f^*)^2 \le B\mathbb{E}(l_f - l_{f^*}) \tag{42}$$

Appendix E. Proof of Theorem 5

Proof. As $f(x), y \in [-1, 1]$, we have that $l_f \in [0, 1]$ and $\mathbb{E}(l_f^2) \leq \mathbb{E}(l_f)$. Hence, we can apply Theorem 4 to function class $l_{\mathcal{H}}$ and obtain that for all $l_f \in l_{\mathcal{H}}$

$$\mathbb{E}(l_f) \le \frac{D}{D-1} \mathbb{E}_n l_f + \frac{6D}{B} \hat{r}^* + \frac{e_3 \delta}{n},$$

as long as there is a sub-root function $\hat{\psi}_n(r)$ such that

$$\hat{\psi}_n(r) \ge e_1 \hat{R}_n \{ f \in star(\mathcal{H}, 0) \mid \mathbb{E}_n f^2 \le r \} + \frac{e_2 \delta}{n}. \tag{43}$$

We have previously demonstrated that

$$e_{1}\hat{R}_{n}\left\{f \in star(\mathcal{H},0) \mid \mathbb{E}_{n}f^{2} \leq r\right\} + \frac{e_{2}\delta}{n}$$

$$\leq 2e_{1}L\hat{R}_{n}\left\{f \in \mathcal{H} \mid \mathbb{E}_{n}f^{2} \leq \frac{e_{6}r}{n^{2}\lambda^{2}}\right\} + \frac{e_{2}\delta}{n}$$

$$\leq 2e_{1}L\left(\frac{2}{n}\sum_{i=1}^{n}\min\left\{\frac{e_{6}r}{n^{2}\lambda^{2}},\hat{\lambda}_{i}\right\}\right)^{1/2} + \frac{e_{2}\delta}{n}$$
(by Lemma 4).

Hence, if choose $\hat{\psi}_n(r)$ to be equal to the right hand side of Eq.(44), then $\hat{\psi}_n(r)$ is a sub-root function that satisfies Eq.(43). Now, the upper bound on the fixed point \hat{r}^* follows from Corollary 6.7 in Bartlett et al. (2005).

Appendix F. Proofs of Lemma 5 & 6

The following is the proof of Lemma 5.

Proof. By Steinwart and Christmann (2008, Theorem 7.13), we have that:

$$\hat{R}_n(l_{\tilde{\mathcal{H}}_r}) \leq \sqrt{\frac{\log 16}{n}} \Big(\sum_{i=1}^{\infty} 2^{i/2} e_{2^i}(l_{\tilde{\mathcal{H}}_r} \cup \{0\}, \|\cdot\|_D) + \sup_{l_{\tilde{\mathcal{H}}_r}} \|l_{\tilde{f}}\|_D \Big).$$

We can see that $e_i(l_{\tilde{\mathcal{H}}_r} \cup \{0\}, \|\cdot\|_D) \leq e_{i-1}(l_{\tilde{\mathcal{H}}_r}, \|\cdot\|_D)$. Hence, we would like to find an upper bound on $e_i(l_{\tilde{\mathcal{H}}_r}, \|\cdot\|_D)$. To this end, since l is L-Lipschitz, we have $e_i(l_{\tilde{\mathcal{H}}_r}, \|\cdot\|_D) \leq Le_i(\tilde{\mathcal{H}}_r, \|\cdot\|_D)$. Notice that $\forall \tilde{f} \in \tilde{\mathcal{H}}_r$, we have $\tilde{f}(x) = \sum_i^s \alpha_i z(v_i, x)$. As a result, given data $\{x_1, \cdots, x_n\}$, under the semi-norm $\|\cdot\|_D$, $\tilde{\mathcal{H}}_r$ is isometric with the space $\tilde{\mathcal{Z}}_r \subseteq \mathbb{R}^n$ spanned by $\{[z(v_1, x_1), \cdots, z(v_1, x_n)]^T, \cdots [z(v_s, x_1), \cdots, z(v_s, x_n)]^T\}_{i=1}^s$. By the definition of $\tilde{\mathcal{H}}_r$, we have $\forall \tilde{f} \in \tilde{\mathcal{H}}_r, Pl_{\tilde{f}} - Pl_{f^*} + \frac{\lambda}{2} \|\tilde{f}\|_{\tilde{\mathcal{H}}} \leq \frac{r}{2R^*}$. This implies that $\|\tilde{f}\|_{\tilde{\mathcal{H}}} \leq (\frac{r}{\lambda R^*})^{1/2}$, and also by reproducing property, $|\tilde{f}(x)| \leq \|\tilde{f}\|_{\tilde{\mathcal{H}}} \|\tilde{k}(x,\cdot)\|_{\tilde{\mathcal{H}}} \leq z_0(\frac{r}{\lambda R^*})^{1/2}$. Denoting the isomorphism from $\tilde{\mathcal{H}}_r$ to $\tilde{\mathcal{Z}}_r$ by I, we immediately have

$$I(\tilde{\mathcal{H}}_r) \subset B_2^{\mathbb{R}^n}(z_0(\frac{r}{\lambda R^*})^{1/2}) \cap U,$$

where $B_2^{\mathbb{R}^n}(z_0(\frac{r}{\lambda R^*})^{1/2})$ is the R^n ball of radius $z_0(\frac{r}{\lambda R^*})^{1/2}$ under $\|\cdot\|_2$. Since U is a s dimensional subspace of R^n , $I(\tilde{\mathcal{H}}_r)$ can be regarded as a R^s ball of radius $z_0(\frac{r}{\lambda R^*})^{1/2}$. As such we can upper bound its entropy number through volume estimation:

$$e_{i}(\tilde{\mathcal{H}}_{r}, \|\cdot\|_{D}) \leq e_{i}(B_{2}^{\mathbb{R}^{s}}(z_{0}(\frac{r}{\lambda R^{*}})^{1/2}), \|\cdot\|_{2})$$

$$\leq 3z_{0}(\frac{r}{\lambda R^{*}})^{1/2}2^{-\frac{i}{s}}$$
(45)

Now it is easy to see $e_0(l_{\tilde{\mathcal{H}}_r}, \|\cdot\|_D) = \sup_{l_{\tilde{\mathcal{H}}_r}} \|l_{\tilde{f}}\|_D$, and because entropy number is decreasing with i, we have:

$$\begin{split} e_i(l_{\tilde{\mathcal{H}}_r},\|\cdot\|_D) &\leq \min \left\{ Le_i(\tilde{\mathcal{H}}_r,\|\cdot\|_D), \ \sup_{l_{\tilde{\mathcal{H}}_r}} \|l_{\tilde{f}}\|_D \right\} \\ &\leq \min \left\{ Le_i(\tilde{\mathcal{H}}_r,\|\cdot\|_D), \ L\sup_{l_{\tilde{\mathcal{H}}_r}} \|l_{\tilde{f}}\|_D \right\} \\ & (\text{w.l.o.g, we assume } L \geq 1). \end{split}$$

Hence,

$$\hat{R}_{n}(l_{\tilde{\mathcal{H}}_{r}}) \leq \sqrt{\frac{\log 16}{n}} \Big(\sum_{i=1}^{\infty} 2^{i/2} e_{2^{i}-1}(l_{\tilde{\mathcal{H}}_{r}} \cup \{0\}, \| \cdot \|_{D}) + \sup_{l_{\tilde{\mathcal{H}}_{r}}} \|l_{\tilde{f}}\|_{D} \Big)
\leq L \sqrt{\frac{\log 16}{n}} \Big(\sum_{i=1}^{\infty} 2^{i/2} e_{2^{i}}(\tilde{\mathcal{H}}_{r}, \| \cdot \|_{D}) + \sup_{l_{\tilde{\mathcal{H}}_{r}}} \|l_{\tilde{f}}\|_{D} \Big)
\leq L \sqrt{\frac{\log 16}{n}} \Big(\sum_{i=0}^{T-1} 2^{i/2} \sup_{l_{\tilde{\mathcal{H}}_{r}}} \|l_{\tilde{f}}\|_{D} + \sum_{T}^{\infty} 2^{i/2} 3z_{0} (\frac{r}{\lambda R^{*}})^{1/2} 2^{-\frac{2^{i}-1}{s}} \Big)
\leq L \sqrt{\frac{\log 16}{n}} \Big(\sum_{i=0}^{T-1} 2^{i/2} \rho + 3z_{0} (\frac{r}{\lambda R^{*}})^{1/2} \sum_{T}^{\infty} 2^{i/2} 2^{-\frac{2^{i}-1}{s}} \Big)$$
(46)

Note the second inequality is because $L \ge 1$, and the third inequality is because we sum the first T-1 term. The reason for summing the first T-1 term is that we control the sequence such that when $i \le T-1$,

$$\rho \le 3z_0 \left(\frac{r}{\lambda R^*}\right)^{1/2} 2^{-\frac{2^i - 1}{s}}.$$

But $3z_0(\frac{r}{\lambda R^*})^{1/2}2^{-\frac{2^i-1}{s}}$ decreases exponentially fast, so after T, we have

$$\rho \ge 3z_0 \left(\frac{r}{\lambda R^*}\right)^{1/2} 2^{-\frac{2^i - 1}{s}}.$$

The first sum in Eq. (46) is $\frac{2^{T/2}-1}{\sqrt{2}-1}$, the second one can be upper bounded by the following integral

$$\int_{T}^{\infty} 2^{\frac{i}{2}} 2^{-\frac{2^{i}-1}{s}} di.$$

Through some algebra we can show that this integral has upper bound

$$\frac{3s}{2^{T/2}}2^{-\frac{2^T}{2s}}$$
.

By taking $T = \log_2(s \log_2(\frac{1}{\lambda}))$, we have

$$\hat{R}_{n}(l_{\tilde{\mathcal{H}}_{r}}) \leq L\sqrt{\frac{\log 16}{n}} \left(\rho\sqrt{s\log_{2}\frac{1}{\lambda}} + 9z_{0}\sqrt{\frac{sr}{R^{*}\log_{2}\frac{1}{\lambda}}}\right)$$

$$\leq L\sqrt{\frac{s\log 16\log_{2}\frac{1}{\lambda}}{n}} \left(\rho + 9z_{0}\sqrt{\frac{r}{R^{*}}}\right).$$

$$:= L\sqrt{\frac{s\log 16\log_{2}\frac{1}{\lambda}}{n}} \left(\rho + 9c_{2}\sqrt{r}\right) \tag{47}$$

The last inequality is because $\lambda \leq 1/2$ implies that $\log_2(1/\lambda) \geq 1$.

The following is the proof of Lemma 6.

Proof. By Lemma 5, we directly take expectation on the upper bound of the empirical Rademacher complexity, we have:

$$R_n(l_{\tilde{\mathcal{H}}_r}) \leq \mathbb{E}\Big\{L\sqrt{\frac{s\log 16\log_2\frac{1}{\lambda}}{n}}\Big(\rho + 9c_2\sqrt{r}\Big)\Big\}.$$

$$= L\sqrt{\frac{s\log 16\log_2\frac{1}{\lambda}}{n}}\Big(\mathbb{E}\rho + 9c_2\sqrt{r}\Big). \tag{48}$$

Since $\rho = \sup_{l_{\tilde{\mathcal{H}}_r}} \|l_{\tilde{f}}\|_D$, we have

$$\begin{split} \mathbb{E}\rho &= \mathbb{E}(\sup_{l_{\tilde{\mathcal{H}}_r}} \|l_{\tilde{f}}\|_D) \\ &\leq \left(\mathbb{E}\sup_{l_{\tilde{\mathcal{H}}_r}} \|l_{\tilde{f}}\|_D^2\right)^{1/2} \text{ (By Jensen's Inequality)} \\ &= \left(\mathbb{E}\sup_{l_{\tilde{\mathcal{H}}_r}} \frac{1}{n} \sum_{i=1} (l_{\tilde{f}}(x_i, y_i) - l_{f^*}(x_i, y_i))^2\right)^{1/2} \\ &:= \left(\mathbb{E}\sup_{l_{\tilde{\mathcal{H}}_r}} \frac{1}{n} \sum_{i=1} h_i^2\right)^{1/2} \\ &\leq \left(c_3 \mathbb{E}\sup_{l_{\tilde{\mathcal{H}}_r}} \frac{1}{n} \sum_{i=1} h_i\right)^{1/2} \end{split}$$

The last inequality is because we assume x_i, y_i are bounded, hence, for Lipschitz continuous loss, $l_{\tilde{f}} - l_{f^*}$ is upper bounded. Notice that:

$$\frac{1}{n} \sum_{i=1}^{n} h_i = \frac{1}{n} \sum_{i=1}^{n} h_i - \mathbb{E}(h) + \mathbb{E}(h)$$

$$\leq |\frac{1}{n} \sum_{i=1}^{n} h_i - \mathbb{E}(h)| + |\mathbb{E}(h)|$$

Hence, we have

$$\mathbb{E} \sup_{l_{\tilde{\mathcal{H}}_r}} \frac{1}{n} \sum_{i=1} h_i \leq \mathbb{E} \sup_{l_{\tilde{\mathcal{H}}_r}} \left| \frac{1}{n} \sum_{i=1} h_i - \mathbb{E}(h) \right| + \mathbb{E} \sup_{l_{\tilde{\mathcal{H}}_r}} \mathbb{E}(h)$$

$$\leq \mathbb{E} \sup_{l_{\tilde{\mathcal{H}}_r}} \left| \frac{1}{n} \sum_{i=1} h_i - \mathbb{E}(h) \right| + r$$

$$\leq 2\mathbb{E} \mathbb{E}_{\epsilon} \sup_{l_{\tilde{\mathcal{H}}_r}} \left| \frac{1}{n} \sum_{i=1} \epsilon_i h_i \right| + r$$

$$= 2R_n(l_{\tilde{\mathcal{H}}_r}) + r$$

As such, we can upper bound Eq. (48) as:

$$R_n(l_{\tilde{\mathcal{H}}_r}) \leq L\sqrt{\frac{s\log 16\log_2\frac{1}{\lambda}}{n}} \left(c_3^{1/2}\sqrt{2R_n(l_{\tilde{\mathcal{H}}_r}) + r} + 9c_2\sqrt{r}\right)$$

$$\leq \sqrt{\frac{s}{n}\log_2\frac{1}{\lambda}} \left(c_4\sqrt{2R_n(l_{\tilde{\mathcal{H}}_r}) + r} + c_5\sqrt{2R_n(l_{\tilde{\mathcal{H}}_r}) + r}\right)$$

$$\leq c_6\sqrt{\frac{s}{n}\log_2\frac{1}{\lambda}}\sqrt{2R_n(l_{\tilde{\mathcal{H}}_r}) + r}$$

$$:= c_7\sqrt{\frac{s}{n}\log\frac{1}{\lambda}}\sqrt{2R_n(l_{\tilde{\mathcal{H}}_r}) + r}$$

With some algebra, we can easily show that

$$R_n(l_{\tilde{\mathcal{H}}_r}) \le C_1 \sqrt{\frac{s}{n} \log \frac{1}{\lambda}} \sqrt{r} + C_2 \frac{s}{n} \log \frac{1}{\lambda}.$$

Appendix G. Code of Algorithm 1

```
def feat_gen(x,n_feat,lns):
    #function to generate the features for gaussian kernel
    : param x: the data
    :param n_feat: number of features we need
    :param lns: the inverse landscale of the gaussian kernel
    :return: a sequence of features ready for KRR
    n, d = np. shape(x)
   w = np.random.multivariate_normal(np.zeros(d),lns*np.eye(d),n_feat)
    return w
def feat_matrix(x,w):
    #funtion to generate the feature matrix Z
    : param x: the data
    :param w: the features
    : return: the feature matrix of size len(n)*len(s)
    s, dim = w.shape
    # perform the product of x and w transpose
    prot mat = np.matmul(x, w.T)
    feat1 = np.cos(prot mat)
    feat2 = np. sin(prot_mat)
    feat_final = np. sqrt(1.0/s)*np. concatenate((feat1, feat2), axis = 1)
    return feat_final
def opm_feat(x,w,lmba):
    ** ** **
    #function to select the optimum features
    :param x: the independent variable
    :param w: the first layer features generated according to spectral density
    :return: optimum features with importance weight
    n_num, dim = x.shape
    s, dim = w.shape
    prot mat = np.matmul(x, w.T)
```

```
feat1 = np. sqrt (1.0/s)*np. cos (prot_mat)
feat2 = np. sqrt (1.0/s)*np. sin (prot_mat)
Z s = feat1 + feat2
ZTZ = np.matmul(Z_s.T, Z_s)
ZTZ_{inv} = np.linalg.inv(ZTZ +n_num*lmba*np.eye(s))
M = np.matmul(ZTZ, ZTZ_inv)
\#M = np.matmul(M0, ZTZ)
1 = np.trace(M)
#print 1
\#n_feat_draw = min(s, max(50, 1))
#print n_feat_draw
#n_feat_draw = int(round(n_feat_draw))
n feat draw = s
#print n_feat_draw
pi_s = np.diag(M)
#print pi_s
qi_s = pi_s/1
is\_wgt = np.sqrt(1/qi\_s)
#print is_wgt
wgt_order = np.argsort(is_wgt)
w_order = wgt_order[(s-n_feat_draw):]
#print len(w_order)
#w_order = np.random.choice(s, n_feat_draw, replace=False, p=qi_s)
#print w order
w_opm = np.zeros((n_feat_draw,dim))
wgh_opm = np.zeros(n_feat_draw)
for ii in np.arange(n_feat_draw):
    order = w_order[ii]
    w_{opm}[ii,:] = w[order,:]
    wgh_opm[ii] = is_wgt[order]
return w_opm, wgh_opm
```