

# MA4199 FYP – Bias Variance Tradeoff

Ng Wei Le

## 1 Kernels

Notation: we use the symbol  $\mathbb{K}$  when it can refer to both  $\mathbb{R}$  or  $\mathbb{C}$ . Also, let  $z^*$  or  $(z)^*$  denote the conjugate of  $z$  for any  $z \in \mathbb{C}$ . The sections covering Kernels and reproducing kernel Hilbert spaces are heavily referenced using Steinwart, Christman [7].

**Definition 1.** For a non-empty set  $X$ , let  $k : X \times X \rightarrow \mathbb{K}$  be known as a kernel if there exists a function  $\phi : X \rightarrow \mathcal{H}$  (known as a feature map of  $k$ ) where  $\mathcal{H}$  is a  $\mathbb{K}$ -Hilbert space (known as a feature space of  $k$ ) such that

$$k(x_1, x_2) = \langle \phi(x_2), \phi(x_1) \rangle_{\mathcal{H}}. \quad (1)$$

**Lemma 1.** For any kernel  $k$  on  $X$ ,  $k(x_1, x_2) = k(x_2, x_1)^*$ .

From the properties of the inner product, we know that  $k(x_1, x_2) = \langle \phi(x_1), \phi(x_2) \rangle^* = k(x_2, x_1)^*$ . Therefore, for kernels on  $\mathbb{R}$ , the symmetric property:  $k(x_1, x_2) = k(x_2, x_1)$  holds.

**Lemma 2.** Let  $k_1, k_2$  be kernels on a non-empty set  $X$ . Then  $k_1 + k_2$  and  $ak_1, a \in \mathbb{R}^+ \cup \{0\}$  are kernels.

Below, we define the Gaussian RBF kernel:

**Definition 2.** Let the complex Gaussian RBF kernel (on  $\mathbb{C}^d$ ) be defined as:

$$k_{\gamma, \mathbb{C}^d}(z, z') := e^{-\gamma^{-2} \sum_{i=1}^d (z_i - z'_i)^*{}^2}.$$

We then define the real Gaussian RBF kernel (or simply the Gaussian RBF kernel for short) acting on  $\mathbb{R}^d$  as:

$$k_{\gamma}(x, x') = e^{-\gamma^{-2} \|x - x'\|_2^2}$$

It can be shown ([7]) that the complex and real Gaussian RBF kernels are kernels.

**Definition 3.** For a non-empty set  $X$ , a function  $k : X \times X \rightarrow \mathbb{R}$  is said to be a positive definite if, for any  $m \in \mathbb{Z}^+ \cup \{0\}$  and all  $x_1, \dots, x_n \in X$ , we have the following matrix (called the Gram matrix) being positive semi-definite:

$$K := (k(x_i, x_j))_{i,j}.$$

Equivalently: for all  $a_1, \dots, a_n \in \mathbb{R}$ , we have:

$$\sum_{j=1}^n \sum_{i=1}^n a_j a_i k(x_j, x_i).$$

**Definition 4.** The positive definite function  $k : X \times X \rightarrow \mathbb{R}$  is said to be symmetric if  $k(x_1, x_2) = k(x_2, x_1)$  for all  $x_1, x_2 \in X$

**Theorem 1.** A real function  $k : X \times X \rightarrow \mathbb{R}$  is a kernel if and only if  $k$  is a positive definite symmetric function (also known as a positive definite kernel).

*Proof.* Suppose  $k$  is a kernel. Then there exists some feature map  $\Phi : X \rightarrow \mathcal{H}$ .

$$\begin{aligned} \sum_{j=1}^n \sum_{i=1}^n a_j a_i k(x_j, x_i) &= \sum_{j=1}^n \sum_{i=1}^n a_j a_i \langle \phi(x_i), \phi(x_j) \rangle_{\mathcal{H}} \\ &= \left\langle \sum_{i=1}^n a_i \phi(x_i), \sum_{j=1}^n a_j \phi(x_j) \right\rangle_{\mathcal{H}} \\ &= \left\| \sum_{i=1}^n a_i \phi(x_i) \right\|^2 \\ &\geq 0. \end{aligned}$$

Also, from Lemma 1, we know that the real kernel  $k$  is symmetric, proving one side of the theorem. To prove the other side:

Given  $k : X \times X \rightarrow \mathbb{R}$  a positive definite symmetric function, we prove that  $\Phi : X \rightarrow H$  where  $x \mapsto k(\cdot, x)$  is a valid feature map for some feature space  $H$ . First, we define

$$\hat{\mathcal{H}} := \sum_{i=1}^n a_i k(\cdot, x_i), n \in \mathbb{Z}^+ \cup \{0\}, a_i \in \mathbb{R} \text{ for all } i, x_i \in X \text{ for all } i.$$

For  $f, g \in \hat{\mathcal{H}}$  where  $f = \sum_{i=1}^n a_i k(\cdot, x_i)$  and  $g = \sum_{j=1}^m b_j k(\cdot, y_j)$ , we define the inner product as such:

$$\begin{aligned} \langle f, g \rangle &:= \sum_{i=1}^n \sum_{j=1}^m a_i b_j k(y_j, x_i) \\ &= \sum_{j=1}^m b_j f(y_j) \\ &= \sum_{i=1}^n a_i g(x_i) \end{aligned} \tag{2}$$

This definition is bilinear and symmetric.

We also have:  $\langle f, f \rangle = \sum_{i=1}^n \sum_{j=1}^n a_i a_j k(x_j, x_i) \geq 0$  since  $k$  is a positive definite function. It can be shown that  $\langle \cdot, \cdot \rangle$  follows Cauhy-Schwarz Inequality ([7]), hence we have:

$$\begin{aligned} |f(x)|^2 &= \left| \sum_{i=1}^n a_i k(\cdot, x_i) \right|^2 \\ &= |\langle f, k(\cdot, x) \rangle|^2 \quad (\because (2) \text{ with } g = \sum_{j=1}^m b_j k(\cdot, y_j) = k(\cdot, x) \text{ with } m = 1, b_1 = 1, y_1 = x) \\ &\leq \langle k(\cdot, x), k(\cdot, x) \rangle \langle f, f \rangle. \end{aligned}$$

Therefore, if  $\langle f, f \rangle = 0$ , then  $f = 0$ , hence showing that  $\langle f, f \rangle > 0$  if and only if  $f \neq 0$ . Hence,  $\langle \cdot, \cdot \rangle$  defines a proper inner product in  $\hat{\mathcal{H}}$ .

Let  $\mathcal{H}$  be the completion of  $\hat{\mathcal{H}}$  and the map  $U : \hat{\mathcal{H}} \rightarrow \mathcal{H}$  be the map where  $\langle Ux, Uy \rangle_{\mathcal{H}} = \langle x, y \rangle_{\hat{\mathcal{H}}}$  for all  $x, y \in \hat{\mathcal{H}}$ . Then we have, for all  $x, x' \in X$ :

$$k(x, x') = \langle k(\cdot, x'), k(\cdot, x) \rangle_{\hat{\mathcal{H}}} = \langle Uk(\cdot, x'), Uk(\cdot, x) \rangle_{\mathcal{H}}$$

. Thus we find a feature map of  $k$ , proving that  $k$  is a kernel. □

## 2 Reproducing Kernel Hilbert Spaces

Initially introduced by Stanislaw Zaremba, reproducing kernel Hilbert spaces have many applications in the fields such as Statistical Learning and complex analysis. An RKHS is a  $\mathbb{K}$ -Hilbert function space where point evaluation is continuous linear functional.

**Definition 5. (RKHS).** Let  $\mathcal{H}$  be a  $\mathbb{K}$ -Hilbert space of functions over a non-empty set  $X$ .  $\mathcal{H}$  is called an RKHS over  $X$  if the Dirac function  $\delta_x : \mathcal{H} \rightarrow \mathbb{K}$  defined as:

$$\delta_x(f) := f(x), \quad x \in X, \quad f \in \mathcal{H}$$

is continuous. Equivalently, there exists  $0 < M_x < \infty$  such that

$$\delta_x(f) \leq M_x \|f\|_{\mathcal{H}}, \quad \text{for all } f \in \mathcal{H}.$$

$\delta_x$  is called a bounded operator on  $\mathcal{H}$ .

This is not easy to put into practice, hence the reproducing kernel is defined.

**Definition 6. (Reproducing Kernel).** For a non-empty set  $X$  and a function  $k : X \times X \rightarrow \mathbb{K}$  where  $k(\cdot, x) \in \mathcal{H}$  for all  $x \in X$  and the following property hold for all  $x \in X$  and  $f \in \mathcal{H}$ :

$$f(x) = \langle f, k(\cdot, x) \rangle \tag{3}$$

The condition in equation (3) is also known as the reproducing property.

**Definition 7. (Canonical Feature Maps).** Let  $\mathcal{H}$  be an RKHS over  $X$  with reproducing kernel  $k$ . Let the function  $\Phi : X \rightarrow \mathcal{H}$  be defined such that for all  $x \in X$ ,

$$\Phi(x) = k(\cdot, x).$$

We call  $\Phi$  the canonical feature map of  $k$ .

**Lemma 3. (A reproducing kernel of an RKHS is a kernel).** Let  $\mathcal{H}$  be an RKHS over  $X$  with reproducing kernel  $k$ . Then  $k$  is a kernel.

*Proof.* We simply proof that  $\Phi$  is a feature map of  $k$ .

$$\begin{aligned}\langle \Phi(x_2), \Phi(x_1) \rangle &= \langle k(\cdot, x_2), k(\cdot, x_1) \rangle \\ &= k(x_1, x_2) \quad (\because \text{Reproducing Property (3)})\end{aligned}$$

So  $\mathcal{H}$  is also a feature space of  $k$ . □

**Lemma 4.** *Let  $\mathcal{H}$  be an  $\mathbb{K}$ -Hilbert functional RKHS over  $X$  with reproducing kernel  $k$ . Then  $H$  is a Reproducing Kernel Hilbert Space.*

*Proof.* Recall the Dirac functional  $\delta_x : H \rightarrow \mathbb{K}$  where:

$$\delta_x(f) = f(x), \quad x \in X, \quad f \in H.$$

Then we have:

$$\begin{aligned}|\delta_x(f)| &= |f(x)| \\ &= |\langle f, k(\cdot, x) \rangle| \quad (\because \text{Reproducing Property (3)}) \\ &\leq \|k(\cdot, x)\|_{\mathcal{H}} \|f\|_{\mathcal{H}} \quad (\because \text{Cauchy-Schwarz Inequality})\end{aligned}$$

This shows that the Dirac functionals are continuous. □

## 2.1 Representer Theorem

Representer Theorem ensures that the *argmin* of an empirical risk expression involving a function over an RKHS can be expressed as a linear combination of kernels applied on the training data points as proven in [6].

**Theorem 2.** *Given a non-empty set  $X$ , training data  $\{(x_1, y_1), \dots, (x_n, y_n)\} \in X \times \mathbb{R}$ , and RKHS  $\mathcal{H}$  be an  $\mathbb{R}$ -Hilbert function space over  $X$  with reproducing kernel  $k : X \times X \rightarrow \mathbb{R}$ . Let  $g$  be a strictly increasing function  $g : [0, \infty] \rightarrow \mathbb{R}$ , and  $l$  be an arbitrary loss function, where  $l : (X \times \mathbb{R}^2)^n \rightarrow \mathbb{R} \cup \{\infty\}$ .*

*We want to minimize the following empirical risk term:*

$$E(f, (x_1, y_1), \dots, (x_n, y_n)) := l((x_1, y_1, f(x_1)), \dots, (x_n, y_n, f(x_n))) + g(\|f\|).$$

*For  $\hat{f} = \operatorname{argmin}_{f \in \mathcal{H}} E(f, (x_1, y_1), \dots, (x_n, y_n))$ ,  $\hat{f}$  can be represented in the form:*

$$\hat{f}(\cdot) = \sum_{i=1}^n a_i k(\cdot, x_i)$$

*with  $a_i \in \mathbb{R}$  for all  $i$ .*

*Proof.* First we let  $\Phi$  be the canonical feature map of  $k$  as defined in 7. Recall: function  $\Phi : X \rightarrow \mathcal{H}$  where

$\Phi(x)(\cdot) = k(\cdot, x)$ . Due to the reproducing property where  $\Phi(x)(x') = \langle \Phi(x), k(\cdot, x') \rangle$ , we have:

$$\begin{aligned}\Phi(x)(x') &= k(x', x) \\ &= \langle \Phi(x), k(\cdot, x') \rangle \\ &= \langle \Phi(x), \Phi(x') \rangle.\end{aligned}$$

So  $\Phi$  is a feature space of  $k$ . Using orthogonal decomposition, we decompose  $f \in \mathcal{H}$  into a component projected onto the span of  $\Phi(x_1), \dots, \Phi(x_n)$ , and the other component orthogonal to this span. We will then prove this orthogonal component is 0 for any  $f$  that reduces the empirical risk term, hence completing the prove.

$$f = \sum_{i=1}^n a_i \Phi(x_i) + \gamma,$$

where  $\gamma \in \mathcal{H}$ ,  $\langle \Phi(x_i), \gamma \rangle = 0$  for all  $i$ .

Next, applying the reproducing property again,

$$\begin{aligned}f(x_j) &= \langle f, k(\cdot, x_j) \rangle \\ &= \langle \sum_{i=1}^n a_i \Phi(x_i) + \gamma, \Phi(x_j) \rangle \\ &= \langle \sum_{i=1}^n a_i \Phi(x_i), \Phi(x_j) \rangle + \langle \gamma, \Phi(x_j) \rangle \\ &= \sum_{i=1}^n a_i \langle \Phi(x_i), \Phi(x_j) \rangle.\end{aligned}$$

Now, consider:

$$\begin{aligned}\|f\|^2 &= \left\| \sum_{i=1}^n a_i \Phi(x_i) + \gamma \right\|^2 \quad (\because \text{orthogonality}) \\ &= \left\| \sum_{i=1}^n a_i \Phi(x_i) \right\|^2 + \|\gamma\|^2 \\ &\geq \left\| \sum_{i=1}^n a_i \Phi(x_i) \right\|^2 \\ \implies g(\|f\|) &\geq g\left(\left\| \sum_{i=1}^n a_i \Phi(x_i) \right\|\right)\end{aligned}$$

Therefore, if we have  $\gamma = 0$ , since  $f(x_i)$  is unaffected by this for all  $i$ ,  $l((x_1, y_1, f(x_1)), \dots, (x_n, y_n, f(x_n)))$  is also unaffected by  $\gamma$ . For the term  $g(\|f\|)$ , it decreases if we have  $\gamma = 0$ . Hence,  $\hat{f} = \operatorname{argmin}_{f \in \mathcal{H}} E(f, (x_1, y_1), \dots, (x_n, y_n))$ ,

$\hat{f}$  must have  $\gamma = 0$ , and

$$\begin{aligned}\hat{f} &= \sum_{i=1}^n a_i \Phi(x_i) \\ &= \sum_{i=1}^n a_i k(\cdot, x_i)\end{aligned}$$

□

### 3 Notations

Let  $\pi_m(\mathbb{R}^d)$  be a multivariate polynomial with  $d$  variables and degree  $\leq m$ , i.e.

$$\pi_m(\mathbb{R}^d) = \{p(x) = \sum_{k \leq m} c_k x^k\}$$

Let  $C^k(X)$  be the set of functions on  $X$  that are  $k$  times continuously differentiable.

For a point  $x \in \mathbb{R}^d$ , it has the components of its coordinates  $\chi_1, \dots, \chi_d$ , whereas we represent  $n$  points in  $\mathbb{R}^d$  as  $x_1, \dots, x_n$ .

We denote  $\mathbb{N}_0$  as the set of non-negative integers. We denote the multi-index vector with its components as  $\alpha = (\alpha_1, \dots, \alpha_d)^T \in \mathbb{N}_0^d$ , and  $|\alpha| := \|\alpha\|_1$ . For  $X \subseteq \mathbb{R}^d$ ,  $f \in C^k(X)$ ,  $|\alpha| \leq k$  and  $x \in \mathbb{R}^d$ , we denote:

$$D^\alpha f := \frac{\partial^{|\alpha|}}{\partial \chi_1^{\alpha_1} \dots \partial \chi_d^{\alpha_d}} f$$

We will define the power function, as defined in 11.2 in [8]:

**Definition 8.** Suppose  $X \in \mathbb{R}^d$  is open, with  $k : X \times X \rightarrow \mathbb{R}$  be a positive definite kernel. For  $\alpha \in \mathbb{N}_0^d$ ,  $\hat{X} = x_1, x_2, \dots, x_n \subseteq X$  the power function  $P_{k, \hat{X}}^{(\alpha)}(x)$  is defined by:

$$\begin{aligned}(P_{k, \hat{X}}^{(\alpha)}(x))^2 &:= D_1^\alpha D_2^\alpha k(x, x) - 2 \sum_{j=1}^n D^\alpha u_j^*(x) D_1^\alpha k(x, x_j) \\ &+ \sum_{i,j=1}^n D^\alpha u_i^*(x) D^\alpha u_j^*(x) k(x_i, x_j).\end{aligned}$$

**Definition 9.** The fill distance (or sometimes referred to as 'fill' for short) for a set of points  $X = \{x_1, \dots, x_N\} \subseteq \Omega$  for a bounded domain  $\Omega$  is defined to be

$$h_{X, \Omega} := \sup_{x \in \Omega} \min_{1 \leq j \leq N} \|x - x_j\|_2$$

.

Theorem 11.22 in [8]:

**Theorem 3.** Let  $\Omega$  be a cube in  $\mathbb{R}^d$  and  $k = \phi(\|\cdot\|_2)$  be a positive definite kernel with  $f = \phi(\cdot)$  satisfying the condition that there exists,  $l_0$  and constant  $M > 0$  such that for all  $r > 0$  and  $l > l_0$ ,  $|f^{(l)}(r)| \leq l!M^l$ . Then there exists a constant  $c > 0$  such that the error between a function  $f \in \mathcal{H}_\infty$  and its interpolant  $s_{f,X}$  for all data points  $X = \{x_1, \dots, x_n\}$  can be bounded by:

$$\|f - s_{f,X}\|_{L_\infty(\Omega)} \leq \exp(-c/h_{X,\Omega})|f|_{\mathcal{H}_\infty}$$

with sufficiently small fill  $h_{X,\Omega}$ .

*Proof.* From the previous theorem, we have:

$$P^2(x) \leq [1 + c_1(2N)]^2 \|f - p\|_{L_\infty(G)}$$

Where  $G$  is on the interval  $[0, 4(c_2(2N))^2 h^2]$ ,  $x \in \Omega$ ,  $p \in \pi_n(\mathbb{R})$ ,  $h = h_{X,\Omega}$ .

From Theorem TODO: we have for sufficiently small fill distance  $h_{X,\Omega} \leq \frac{c_0}{2n}$ , the constants  $c_1, c_2$  can be replaced by:

$$c_1(2n) = \exp(2d\gamma_d(2n+1))$$

$$c_2(2n) = 2c_2n$$

So  $G$  is on the interval  $[0, 16N^2 c_2^2 h^2]$  For  $p$  the Taylor series of  $f$  about 0, and up to the term  $t^N$ , we then have:

$$\begin{aligned} |f(t) - p(t)| &\leq t^{N+1} \frac{|f^{n+1}(t')|}{(n+1)!} \\ &\leq M^{N+1} t^{N+1} \quad (\text{By assumption}) \\ \implies \|f - p\|_{L_\infty(G)} &\leq (M \cdot 16N^2 c_2^2 h^2)^{N+1} \\ &= (C_0 N^2 h^2)^{N+1} \quad \text{for constant } C_0 = M c_2^2 \end{aligned}$$

Also, we have:

$$\begin{aligned} [1 + c_1(2N)]^2 &= [1 + \exp(2d\gamma_d(2N+1))]^2 \\ &\leq [2 \exp(2d\gamma_d(2N+1))]^2 \\ &= 4 \exp(4d\gamma_d(2N+1)) \\ &= \exp(\log 4 + 4d\gamma_d(2N+1)) \\ &\leq \exp(C_1(N+1)) \quad \text{for sufficiently large } C_1. \end{aligned}$$

We then have:

$$\begin{aligned} P_{\Phi,X}^2(x) &\leq [1 + c_1(2N)]^2 \|f - p\|_{L_\infty(G)} \\ &\leq \exp(C_1(N+1)) (C_0 N^2 h^2)^{N+1} \\ &= (C_0 N^2 h^2 \exp(C_1))^{N+1} \\ &= (C_2 N^2 h^2)^{N+1} \quad \text{for constant } C_2 = C_0 \exp(C_1). \end{aligned} \tag{4}$$

For  $C_3 = \min(\frac{c_0}{2}, \frac{1}{\sqrt{eC_2}})$ , and  $N$  such that

$$\frac{C_3}{N+1} \leq h \leq \frac{C_3}{N}$$

, which gives us:

$$\begin{aligned} h &\leq \frac{c_0}{2N}, \\ -(N+1) &\leq -C_3/h, \\ N^2 h^2 &\leq C_3^2 \leq \frac{1}{eC_2} \\ \implies C_2 N^2 h^2 &\leq 1/e. \end{aligned}$$

We then have:

$$P_{\Phi, X}^2(x) \leq e^{-(N+1)} \leq e^{-C_3/h}.$$

Now, using  $C = C_3/2$  and  $|f(x) - s_{f, X}(x)| \leq P_{\Phi, X}(x)|f|_{\mathcal{H}_\infty}$ , we have:

$$|f(x) - s_{f, X}(x)| \leq e^{-C/h_{X, \Omega}} |f|_{\mathcal{H}_\infty},$$

which is what we wanted to prove.  $\square$

**Theorem 4.** *With the same conditions as Theorem 3, except that  $f$  satisfies the stricter condition  $|f^{(l)}(r)| \leq M^l$ , we can get a better error bound of:*

$$\|f - s_{f, X}\|_{L_\infty(\Omega)} \leq \exp\left(\frac{c \log(h_{X, \Omega})}{h_{X, \Omega}}\right) \|f\|_{\mathcal{H}_\infty}.$$

*Proof.* The inequality at 4 becomes:

$$P_{\Phi, X}^2(x) \leq \frac{(C_2 N^2 h^2)^{N+1}}{(N+1)!}.$$

Using Stirling's inequality  $1/n! \leq (e/n)^n$ , we have:

$$P_{\Phi, X}^2(x) \leq (eC_2 N h^2)^{N+1}.$$

Similarly, with  $C_3 = \min(\frac{c_0}{2}, \frac{1}{eC_2})$ , and  $N$  such that

$$\frac{C_3}{N+1} \leq h \leq \frac{C_3}{N},$$

we then get

$$eC_2 N h \leq 1,$$

which gives us:

$$P_{\Phi, X}^2(x) \leq h^{N+1} \leq h^{C_3/h} = e^{C_3 \log h/h}.$$



Following the steps of the previous theorem then gives us our result.  $\square$

## 4 Approximation Theorem

The below theorem gives us some justification as to why the minimum norm interpolating function was chosen, though this only works under noiseless conditions:

**Theorem 5.** *Fix  $h^* \in \mathcal{H}_\infty$ . Let  $(x_1, y_1), \dots, (x_n, y_n)$  be i.i.d. random variables where  $x_i$  drawn randomly from a compact cube  $\Omega \subseteq \mathbb{R}^d$ ,  $y_i = h^*(x_i) \forall i$ . There exists  $A, B > 0$  such that for any interpolating  $h \in \mathcal{H}_\infty$  with high probability*

$$\sup_{x \in \Omega} |h(x) - h^*(x)| < Ae^{-B(n/\log n)^{1/d}} (\|h^*\|_{\mathcal{H}_\infty} + \|h\|_{\mathcal{H}_\infty})$$

With  $h_{X,\Omega}$  as the fill on the order of  $O(n/\log n)^{-1/d}$  (using the theorem S1 in Belkin's paper which wasn't proved). We consider  $f(x) := h(x) - h^*(x)$ . Since  $h$  is interpolating, we have  $f(x_i) = 0$  for all  $x_i$ . We then let  $s_{f,X}$  be the zero function, since it is an interpolant of  $f$ . Thus, we have:  $s_{f,X}$  can be bounded by:

$$\begin{aligned} \|f\|_{L_\infty(\Omega)} &= \sup_{x \in \Omega} |h(x) - h^*(x)| < \exp(-c(n/\log n)^{1/d}) \|f\|_N(\Omega) \\ &\leq \exp(-c(n/\log n)^{1/d}) (\|h^*\|_{\mathcal{H}_\infty} + \|h\|_{\mathcal{H}_\infty}) \end{aligned}$$

Another form we can have is using proposition 14.1 in [8]:

**Proposition 1.** *Let  $\Omega \subseteq \mathbb{R}^d$  be bounded and measurable. Suppose  $X = \{x_1, \dots, x_N\} \subseteq \Omega$  is quasi-uniform with respect to  $c_{qu} > 0$ . Then there exists constants  $c_1, c_2 > 0$  depending only on space dimension  $d$ , on  $\Omega$  and on  $c_{qu}$  such that:*

$$c_1 N^{-1/d} \leq h_{X,\Omega} \leq c_2 N^{-1/d}$$

.

With the definition of quasi-uniformness being:

**Definition 10.** For the separation distance of  $X = \{x_1, \dots, x_N\}$  being defined as  $q_x := \frac{1}{2} \min_{i \neq j} \|x_i - x_j\|_2$ .

We can then use the above proposition with  $n$  replacing  $n/\log n$ .

In either case, by choosing a the smallest norm for  $h$ , we can see that it corresponds to the smallest upper-bound for  $|h(x) - h^*(x)|$ .

## 5 Existing Bounds Provide No Guarantees for Interpolated Kernel Classifiers

Steps are:

- Find lower bound on function norm of t-overfitted classifiers in RKHS corresponding to Gaussian Kernels.

- Show loss for available bounds for kernel methods based on function norm (can perhaps use this to explain approximation theorem as well?)

Interpolation: 0 regression error. Overfitting: 0 classification error. Interpolation implies overfitting.

**Definition 11.** We say  $h \in H$   $t$ -overfits data, if it achieves zero classification loss (overfits) and  $\forall_i y_i h(x_i) > t > 0$ .

The below shows a theorem on how the function norm changes with respect to  $t$ -overfitting.

**Theorem 6.** Let  $(\mathbf{x}_i, y_i)$  be data sampled from  $P$  on  $\Omega \times \{-1, 1\}$  for  $i = 1, \dots, n$ . Assume that  $y$  is not a deterministic function of  $x$  on a subset of non-zero measure. Then, with high probability, any  $h$  that  $t$ -overfits the data, satisfies

$$\|h\|_H > Ae^{Bn^{1/d}}$$

for some constants  $A, B > 0$  depending on  $t$ .

We define the  $\gamma$ -shattering and fat-shattering dimension below:

**Definition 12.** Let  $F$  be a set of functions mapping from a domain  $X$  to  $\mathbb{R}$ . Suppose  $S = \{x_1, x_2, \dots, x_m\} \subseteq X$ . Suppose also that  $\gamma$  is a positive real number. Then  $S$  is  $\gamma$ -shattered by  $F$  if there are real numbers  $r_1, r_2, \dots, r_m$ , such that for each  $b \in \{0, 1\}^m$  there is a function  $f_b$  in  $F$  with

$$f_b(x_i) \geq r_i + \gamma \text{ if } b_i = 1, \text{ and } f_b(x_i) \leq r_i - \gamma \text{ if } b_i = 0, \text{ for } 1 \leq i \leq m.$$

We say  $r = (r_1, r_2, \dots, r_m)$  witnesses the shattering. Suppose that  $F$  is a set of functions from a domain  $X$  to  $\mathbb{R}$  and that  $\gamma > 0$ . Then  $F$  has  $\gamma$ -dimension  $d$  if  $d$  is the maximum cardinality of a subset  $S$  of  $X$  that is  $\gamma$ -shattered by  $F$ . If no such maximum exists, we say that  $F$  has infinite  $\gamma$ -dimension. The  $\gamma$ -dimension of  $F$  is denoted  $\text{fat}_F(\gamma)$ . This defines a function  $\text{fat}_F : \mathbb{R} \rightarrow N \cup \{0, \infty\}$ , which we call the fat-shattering dimension of  $F$ .

*Proof.* Let  $B_R = \{f \in \mathcal{H}, \|f\|_{\mathcal{H}} < R\}$  be a ball of radius  $R$  in RKHS  $\mathcal{H}$ . Suppose the data is  $\gamma$ -overfitted, [3] gives us a high probability of a bound of

$$L(f) < O\left(\frac{\ln(n)^2}{\sqrt{n}} \sqrt{\text{fat}_{B_R}(\gamma/8)}\right)$$

for  $L(f)$  the expected classification error. Also, from [1] we have

$$\text{fat}_{B_R}(\gamma) < O((\log(R/\gamma))^d)$$

. We then have  $B_R$  containing no function that  $\gamma$  overfits the data unless

$$(\log(R/\gamma))^d > O(n) \implies R > c_1 \exp(c_2 (\frac{n}{\ln n})^{1/d})$$

for some positive constants  $c_1, c_2$ . □

Classical bounds for kernel methods ( [2] ) are in the form:

$$|\frac{1}{n} \sum_i l(f(x_i), y_i) - L(f)| \leq C \frac{\|f\|_{\mathcal{H}}^a}{n^b}, \quad C, a, b \geq 0$$

The right side on this will tend to infinity for bigger  $\|f\|_{\mathcal{H}}$ , which is suggested by Theorem 6.

## 6 Random Fourier Features

For a feature map  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$  the kernel trick allows easy computation for positive definite kernel  $k$  where  $k(x, y) = \langle \phi(x), \phi(y) \rangle$ . We want to find a randomized feature map  $z : \mathbb{R}^d \rightarrow \mathbb{R}^{\bar{d}}$  such that

$$k(x, y) = \langle \phi(x), \phi(y) \rangle \approx \langle z^T(x), z(y) \rangle$$

. As suggested by [4], for a shift-invariant kernel  $k: k(x, y) = k(x - y)$ , we consider the mapping  $z(x) = \cos(w^T x + b)$ , where  $w$  is drawn from the probability distribution  $p$ :

$$p(w) = \frac{1}{2\pi} \int k(h) \exp(-iw^T h) dh \quad (5)$$

when we compute the Fourier transform of the kernel  $k$ , and  $b$  is drawn from the uniform distribution on  $[0, 2\pi]$ .

We know that the fourier transform of  $k(\cdot)$  is a probability distribution from Bochner's theorem:

**Theorem 7.** (Bochner [5]). *For a continuous kernel  $k(x - y)$  it is a positive definite kernel if and only if  $k(\cdot)$  is the fourier transform of a non-negative measure.*

We now have:

$$k(x - y) = \int_{\mathbb{R}^d} p(w) \exp(iw^T(x - y)) dw = \mathbb{E}_w[e^{iw^T x} (e^{iw^T y})^*]$$

. Therefore, we can use  $e^{iw^T x} (e^{iw^T y})^*$  as an estimate (unbiased) of  $k(x, y)$ . Let  $\phi_w(x) = e^{iw^T x}$ . We can also use  $z_w(x) = \sqrt{2} \cos(w^T x + b)$  instead of  $\phi_w(x)$ , as suggested by [4].

**Proposition 2.** *For  $z_w(x) = \sqrt{2} \cos(w^T x + b)$ , where  $w$  is drawn from probability distribution  $p$  in (5) and  $b$  drawn from a uniform random variable on  $[0, 2\pi]$ .*

$$E(z_w(x) z_w(y)) = k(x, y)$$

*Proof.*

$$\begin{aligned} z_w(x) &= 2 \frac{\sqrt{2}}{2} \cos(w^T x + b) \\ &= \frac{1}{\sqrt{2}} (e^{i(w^T x + b)} + e^{-i(w^T x + b)}) \\ &= \frac{1}{\sqrt{2}} (\phi_w(x) e^{ib} + \phi_w(x)^* e^{-ib}) \end{aligned}$$

Where  $\phi_w(x) = e^{iw^\top x}$ .

$$\begin{aligned} z_w(x)z_y(y) &= \frac{1}{2}[\phi_w(x)\phi_w(y)e^{i2b} + \phi_w(x)^*\phi_w(y)^*e^{-i2b} + \phi_w(x)\phi_w(y)^* + \phi_w(x)^*\phi_w(y)] \\ \mathbb{E}[z_w(x)z_y(y)] &= \frac{1}{2}\mathbb{E}[\phi_w(x)\phi_w(y)e^{i2b} + \phi_w(x)^*\phi_w(y)^*e^{-i2b}] + \frac{1}{2}\mathbb{E}[\phi_w(x)\phi_w(y)^*] + \frac{1}{2}\mathbb{E}[\phi_w(x)^*\phi_w(y)] \end{aligned}$$

As mentioned earlier in Theorem 7,  $\mathbb{E}_w[\phi_w(x)\phi_w(y)^*] = k(x - y)$ . Also  $\phi_w(x)\phi_w(y)^* = (\phi_w(x)^*\phi_w(y))^*$ .

$$\begin{aligned} \mathbb{E}[z_w(x)z_y(y)] &= \frac{1}{2}\mathbb{E}[\phi_w(x)\phi_w(y)e^{i2b} + \phi_w(x)^*\phi_w(y)^*e^{-i2b}] + \frac{1}{2}k(x - y) + \frac{1}{2}[k(x - y)]^* \\ &= \frac{1}{2}\mathbb{E}[\phi_w(x)\phi_w(y)e^{i2b} + \phi_w(x)^*\phi_w(y)^*e^{-i2b}] + k(x - y) \end{aligned}$$

For real kernel ,  $k(x - y) = (k(x - y))^*$ .

$$\begin{aligned} \mathbb{E}_{w,b}[\phi_w(x)\phi_w(y)e^{i2b}] &= \frac{1}{2\pi} \int_{\mathbb{R}^d} \int_0^{2\pi} p(w)\phi_w(x)\phi_w(y)e^{i2b} db dw \\ &= \frac{1}{2\pi} \int_{\mathbb{R}^d} p(w)\phi_w(x)\phi_w(y) \int_0^{2\pi} e^{i2b} db dw \\ &= 0 \end{aligned}$$

Since  $\int_0^{2\pi} e^{i2b} db = 0$ . Similarly,  $\mathbb{E}_{w,b}[\phi_w(x)^*\phi_w(y)^*e^{-i2b}] = 0$ .

$$\therefore \mathbb{E}[z_w(x)z_y(y)] = k(x - y).$$

□

As suggested by [4], the variance of the estimate is decreased by using  $z$ , a  $D$  dimensional vector by concatenating  $D$  of  $z_w$  and normalizing by a constant  $\sqrt{D}$ . We let:

$$z(x) = \sqrt{\frac{2}{D}}[\cos(w_1^\top x + b_1) \dots \cos(w_D^\top x + b_D)]$$

with randomly drawn  $w_i$  and  $b_i$  as described previously.

**Theorem 8.** For  $N$  the number of random features, and  $x_1, x_2, \dots, x_n$  the data points, when  $N > n$  and as  $N$  increases, the norm of the minimizer tends to the norm of the minimum norm RKHS interpolant.

*Proof.* Let  $f(x)$  be the minimum norm RKHS interpolant function for the datapoints.

$$f(x) = \sum_i \alpha_i k(x_i, x) \approx \sum_i \alpha_i z(x_i)^\top z(x) = \beta^\top z(x) = \hat{f}(x)$$

(the first equality holds due to Representer Theorem) Where  $\beta = \sum_i \alpha_i z(x_i)$ . The norm of the function from the random fourier features approximation is:

$$\|\beta\| = \beta^\top \bar{\beta} = \left(\sum_i \alpha_i z^\top(x_i)\right) \left(\sum_i \bar{\alpha}_i \bar{z}(x_i)\right) = \sum_i \sum_j \alpha_i \bar{\alpha}_j z^\top(x_i) \bar{z}(x_j) \approx \sum_i \sum_j \alpha_i \bar{\alpha}_j k(x_i, x_j) = \|f\|$$



## References

- [1] Mikhail Belkin. Approximation beats concentration? an approximation view on inference with smooth radial kernels. *arxiv:1801.03437*, 2018.
- [2] Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to understand kernel learning. *arxiv:1802.01396*, 2018.
- [3] Balázs Kégl, Tamás Linder, and Gábor Lugosi. Data-dependent margin-based generalization bounds for classification. In *Lecture Notes in Computer Science*, pages 368–384. Springer Berlin Heidelberg, 2001.
- [4] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In J. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems*, volume 20, pages 1177–1184. Curran Associates, Inc., 2008.
- [5] Walter Rudin. *Fourier Analysis on Groups*. John Wiley & Sons, Inc., jan 1990.
- [6] Bernhard Schölkopf, Ralf Herbrich, and Alex J. Smola. A generalized representer theorem. In *Lecture Notes in Computer Science*, pages 416–426. Springer Berlin Heidelberg, 2001.
- [7] Ingo Steinwart and Andreas Christmann. *Support Vector Machines*. Springer New York, 2008.
- [8] Holger Wendland. *Scattered Data Approximation*. Cambridge Monographs on APplied and Computational Mathematics. Cambridge University Press, 2004.

# Appendix