

Lecture 8: Nonparametric regression in RKHS

March 5, 2026

Lecturer: Ben Dai

“There is Nothing More Practical Than A Good Theory.”

— Kurt Lewin

1 Recall

Based on Lectures 1-7, we are able to compute the convergence rate and establish a probabilistic bound for a general learning method/algorithm. For illustration, we turn to investigate the asymptotics of nonparametric regression in a Reproducing Kernel Hilbert Space (RKHS).

2 RKHS

2.1 Why RKHS?

Motivation. In nonparametric regression, we seek a function space \mathcal{F} that is rich enough to approximate complex functions, yet structured enough to allow for stable estimation.

From Hilbert Spaces to RKHS. A general Hilbert space \mathcal{H} provides a powerful framework for function approximation, but it does not necessarily guarantee that convergence in the norm implies pointwise convergence. In nonparametric regression, we require that if our estimator f_n converges to the true function f in the Hilbert space norm, the predicted values $f_n(x)$ also converge to the true values $f(x)$.

Reproducing Kernel Hilbert Spaces (RKHS) are the specific class of Hilbert spaces that satisfy this requirement. By definition, an RKHS ensures that the evaluation functional $\delta_x(f) = f(x)$ is a continuous linear functional. According to the Riesz representation theorem, this continuity implies the existence of a representer $K_x \in \mathcal{H}$ such that $f(x) = \langle f, K_x \rangle_{\mathcal{H}}$. This reproducing property is the bridge between the abstract Hilbert space structure and the concrete evaluation of functions, enabling the use of kernels to define the space and perform computations.

Theorem 2.1 (The Riesz Representation Theorem for Hilbert Spaces). *Let \mathcal{H} be a Hilbert space, and $L: \mathcal{H} \rightarrow \mathbb{R}$ be a bounded linear functional on \mathcal{H} . Then there exists a unique element $K \in \mathcal{H}$ such that for every $h \in \mathcal{H}$, we have $L(h) = \langle h, K \rangle_{\mathcal{H}}$, and $\|L\| = \|K\|_{\mathcal{H}}$.*

Applying the Riesz representation theorem, if the evaluation functional $\delta_{\mathbf{x}}: f \mapsto f(\mathbf{x})$ is a bounded linear functional on \mathcal{H} , there exists a unique representer $K_{\mathbf{x}} \in \mathcal{H}$ such that

$$f(\mathbf{x}) = \delta_{\mathbf{x}}(f) = \langle f, K_{\mathbf{x}} \rangle_{\mathcal{H}},$$

for all $f \in \mathcal{H}$. This result implies that function evaluation at any point \mathbf{x} can be represented as an inner product with a specific function $K_{\mathbf{x}}$ in the Hilbert space.

Definition 2.2 (Reproducing Kernel Hilbert Space (RKHS)). A Hilbert space \mathcal{H} of functions on \mathcal{X} is called a Reproducing Kernel Hilbert Space (RKHS) if the evaluation functional $\delta_{\mathbf{x}} : f \mapsto f(\mathbf{x})$ is a bounded linear functional on \mathcal{H} for every $\mathbf{x} \in \mathcal{X}$.

Theorem 2.3 (Pointwise Convergence in RKHS). If \mathcal{H} is an RKHS, then convergence in the Hilbert space norm implies pointwise convergence. That is, for any sequence $\{h_n\} \subset \mathcal{H}$ and $h \in \mathcal{H}$,

$$\|h_n - h\|_{\mathcal{H}} \rightarrow 0 \implies h_n(x) \rightarrow h(x), \quad \forall x \in \mathcal{X}.$$

2.2 From kernel function to RKHS

Motivation. The definition of an RKHS guarantees the existence of a representer $K_{\mathbf{x}}$ for every evaluation functional $\delta_{\mathbf{x}}$. This naturally leads to the question: can we reverse this process? That is, given a series of representers $K_{\mathbf{x}} (x \in \mathcal{X})$, can we construct a Hilbert space \mathcal{H} such that K acts as its ‘‘core’’? The answer is yes, and this construction is fundamental to the utility of RKHS in machine learning. By starting with a kernel, we can define the inner product structure of the space, which in turn determines the geometry and the properties of the Hilbert space.

By the Riesz representation theorem, the set $\{K_{\mathbf{x}} : \mathbf{x} \in \mathcal{X}\}$ serves as a *natural basis* for the RKHS.

Step 1. Define the inner product between basis elements.

The inner product between these basis elements is determined by the reproducing property:

$$\langle K_{\mathbf{x}'}, K_{\mathbf{x}} \rangle_{\mathcal{H}} = \delta_{\mathbf{x}}(K_{\mathbf{x}'}) = K_{\mathbf{x}'}(\mathbf{x}).$$

Here, $K(\mathbf{x}, \mathbf{x}') = K_{\mathbf{x}'}(\mathbf{x}) = K_{\mathbf{x}}(\mathbf{x}')$ is a **symmetric bivariate function** known as the kernel function. By defining the kernel, we implicitly define both the basis elements $K_{\mathbf{x}}$ and the inner product between any two basis functions $K(\mathbf{x}, \mathbf{x}')$.

Step 1. Define a pre-RKHS as the linear span of kernel functions.

To mimic the construction in the finite-dimensional case, we first define a pre-RKHS \mathcal{H}_0 as the set of all finite linear combinations of kernel functions:

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \mathbf{x}),$$

where $n \in \mathbb{N}$, $\alpha_i \in \mathbb{R}$, and $\mathbf{x}_i \in \mathcal{X}$. We equip this space with an inner product defined by the reproducing property:

$$\langle f, f' \rangle_{\mathcal{H}_0} = \sum_{i=1}^n \sum_{j=1}^m \alpha_i \alpha'_j K(\mathbf{x}_i, \mathbf{x}'_j),$$

where $f(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \mathbf{x})$ and $f'(\mathbf{x}) = \sum_{j=1}^m \alpha'_j K(\mathbf{x}'_j, \mathbf{x})$.

Step 2. Generate the RKHS by taking the completion of the pre-RKHS.

The pre-RKHS \mathcal{H}_0 is generally not complete, meaning it may not contain the limits of all its Cauchy sequences. To ensure that we have a complete Hilbert space, we define \mathcal{H} as the completion of \mathcal{H}_0 with respect to the norm induced by the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}_0}$. This completion includes all limit points of Cauchy sequences in \mathcal{H}_0 , resulting in a space of functions of the form:

$$f(\mathbf{x}) = \sum_{i=1}^{\infty} \alpha_i K(\mathbf{x}_i, \mathbf{x}),$$

where the series converges in the norm of \mathcal{H} . The inner product on \mathcal{H} is defined by the continuous extension of the inner product on \mathcal{H}_0 :

$$\langle f, f' \rangle_{\mathcal{H}} = \lim_{n \rightarrow \infty} \langle f_n, f'_n \rangle_{\mathcal{H}_0},$$

where $\{f_n\}$ and $\{f'_n\}$ are Cauchy sequences in \mathcal{H}_0 converging to f and f' , respectively.

Finally, we verify that \mathcal{H} is indeed an RKHS by confirming that: (i) \mathcal{H} is a Hilbert space; (ii) the evaluation functional δ_x is continuous on \mathcal{H} . A formal proof of this construction can be found in [Sejdinovic and Gretton, 2012].

Remark 2.4 (Positive-definite kernel). One quick observation is that a valid kernel function should be symmetric and positive definite.

- From **inner product**: a kernel function should be symmetric.

$$K(\mathbf{x}, \mathbf{x}') = \langle K_{\mathbf{x}}, K_{\mathbf{x}'} \rangle_{\mathcal{H}} = \langle K_{\mathbf{x}'}, K_{\mathbf{x}} \rangle_{\mathcal{H}} = K(\mathbf{x}', \mathbf{x})$$

- From **norm**: a kernel function should be positive definite.

$$0 \leq \|f\|_{\mathcal{H}}^2 = \left\| \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \cdot) \right\|_{\mathcal{H}}^2 = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \langle K_{\mathbf{x}_i}, K_{\mathbf{x}_j} \rangle_{\mathcal{H}} = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j),$$

which holds for any $f \in \mathcal{H}$ or for any $n \geq 1$, any $(\alpha_1, \dots, \alpha_n) \in \mathbb{R}^n$, any $(\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathcal{X}^n$.

2.3 Definitions and theorems

In this section, we provide the formal definitions and theorems that underpin the construction of an RKHS.

Definition 2.5 (Reproducing kernel). Let \mathcal{H} be a Hilbert space of functions on \mathcal{X} . A function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is called a reproducing kernel of \mathcal{H} if it satisfies the following two conditions:

1. For every $\mathbf{x} \in \mathcal{X}$, the function $K(\cdot, \mathbf{x})$ belongs to \mathcal{H} .
2. (Reproducing property). For every $\mathbf{x} \in \mathcal{X}$ and every $h \in \mathcal{H}$, the inner product satisfies $\langle h, K(\cdot, \mathbf{x}) \rangle_{\mathcal{H}} = h(\mathbf{x})$.

Definition 2.6 (Positive semi-definite kernel). A symmetric function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is called a positive semi-definite kernel if, for any $n \in \mathbb{N}$, any set of points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subset \mathcal{X}$, and any set of coefficients $\{\alpha_1, \dots, \alpha_n\} \subset \mathbb{R}$, the following inequality holds:

$$\sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) \geq 0.$$

The Moore-Aronszajn theorem guarantees the validity of the construction presented in Section 2.2.

Theorem 2.7 (Moore-Aronszajn theorem). *Let $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a symmetric, positive semi-definite kernel. Then there exists a unique Reproducing Kernel Hilbert Space (RKHS) \mathcal{H} such that K is the reproducing kernel of \mathcal{H} .*

In summary, the following concepts are equivalent:

$$\text{Reproducing kernel} \iff \text{Positive semi-definite kernel} \iff \text{RKHS}.$$

2.4 Examples

- Linear kernel.

$$K(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$$

- Gaussian kernel.

$$K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{\sigma^2}\right)$$

- γ -degree polynomial kernel.

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^\gamma$$

Remark 2.8. What is the difference among different kernels? Theoretically, it affects both estimation and approximation errors; see the discussion in Section 4. Practically, it is highly related to the topic of *multiple kernel learning*; see [Gönen and Alpaydm, 2011] and references therein.

3 Ordinal Regression in RKHS

Let $\mathbf{X} \in \mathbb{R}^d$ be a feature vector and $Y \in \{0, \dots, K\}$ be an ordinal outcome. We consider a continuous prediction function $f(\mathbf{X}) \in \mathbb{R}$ and use the mean absolute error (MAE) as the loss function:

$$R(f) = \mathbb{E}\left(l(Y, f(\mathbf{X}))\right) = \mathbb{E}\left|Y - f(\mathbf{X})\right|.$$

We summarize the quantities of interest as follows:

- **Bayes rule:** $f^*(\mathbf{x}) = \text{Median}(Y|\mathbf{X} = \mathbf{x})$ is the global minimizer of $R(f)$.

- **Excess risk:**

$$\mathcal{E}(f) = R(f) - R(f^*) = \mathbb{E} |Y - f(\mathbf{X})|.$$

- **R-ERM:** Given random samples $\{(\mathbf{X}_i, Y_i)\}_{i=1}^n$ and an RKHS \mathcal{H}_K , the estimator is defined as

$$\hat{f}_n = \arg \min_{f \in \mathcal{H}_K} \left(\frac{1}{n} \sum_{i=1}^n |Y_i - f(\mathbf{X}_i)| + \lambda_n \|f\|_{\mathcal{H}_K}^2 \right).$$

- **Asymptotics:** Finally, we investigate the asymptotic behavior of $\mathcal{E}(\hat{f}_n)$.

First, we consider the empirical optimization of **ERM** on RKHS. Indeed, this can be challenging, since $f \in \mathcal{H}_K$, and the RKHS \mathcal{H}_K is an infinity-dimensional function class. Fortunately, we have the Representer Theorem, which implies that ERM reduces to a finite dimensional optimization problem.

Theorem 3.1 (Representer Theorem [Kimeldorf and Wahba, 1970, Wahba, 1990]). *Let $K(\cdot, \cdot)$ be a kernel function and \mathcal{H}_K be its associated RKHS. Given a training sample $(\mathbf{X}_i, Y_i)_{i=1, \dots, n}$, consider the R-ERM:*

$$\hat{f}_n = \arg \min_{f \in \mathcal{H}_K} \mathcal{L}_n(Y_1, \dots, Y_n, f(\mathbf{X}_1), \dots, f(\mathbf{X}_n)) + g(\|f\|_{\mathcal{H}_K}^2). \quad (1)$$

Suppose $g(\cdot)$ is an increasing function \mathcal{L}_n depends on f only through $f(\mathbf{X}_1), \dots, f(\mathbf{X}_n)$. Then, every minimizer of (1) has form:

$$\hat{f}_n(\mathbf{x}) = \sum_{i=1}^n \hat{\alpha}_i K(\mathbf{X}_i, \mathbf{x}),$$

for some $(\hat{\alpha}_i)_{i=1, \dots, n} \in \mathbb{R}^n$.

Remark 3.2. We refer readers to the summary by Wahba and Wang¹. As they note:

“The significance of the representer theorem is that the solution in an infinite-dimensional space falls in a finite-dimensional space. This property makes it possible to compute estimates for general regularization problems in infinite-dimensional spaces.”

According to the Representer Theorem, the R-ERM can be reduced to:

$$\min_{\boldsymbol{\alpha}} \frac{1}{n} \sum_{i=1}^n |Y_i - \sum_{j=1}^n \alpha_j K(\mathbf{X}_j, \mathbf{X}_i)| + \lambda_n \left\| \sum_{j=1}^n \alpha_j K(\mathbf{X}_j, \cdot) \right\|_{\mathcal{H}_K}^2 = \min_{\boldsymbol{\alpha}} \frac{1}{n} \sum_{i=1}^n |Y_i - \boldsymbol{\alpha}^\top \mathbf{K}_i| + \lambda_n \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha},$$

where $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n)^\top$ and $\mathbf{K} = (K_{ij})_{n \times n}$ is a kernel matrix with $K_{ij} = K(\mathbf{X}_i, \mathbf{X}_j)$.

¹<http://pages.stat.wisc.edu/~wahba/ftp1/wahba.wang.2019submit.pdf>

3.1 A first excess risk bound

Next, we turn to bound the excess risk of RKHS regression. For simplicity, we assume $0 \leq \widehat{f}_n(\mathbf{X}_i) \leq K$, otherwise we can consider its truncated estimation. Note that

$$\lambda_n \|\widehat{f}_n\|_{\mathcal{H}_K}^2 \leq \widehat{R}_n(\widehat{f}_n) + \lambda_n \|\widehat{f}_n\|_{\mathcal{H}_K}^2 \leq \widehat{R}_n(0) \leq K.$$

Thus, $\widehat{f}_n \in \mathcal{H}_n = \{f \in \mathcal{H} : \|f\|_{\mathcal{H}_K} \leq K\lambda_n^{-1/2}\}$. Then, by the standard decomposition of excess risk, we have

$$\mathcal{E}(\widehat{f}_n) = R(\widehat{f}_n) - R(f^*) \leq 2 \sup_{f \in \mathcal{H}_n} |\widehat{R}_n(f) - R(f)| + \mathbf{Approx}(\lambda_n).$$

3.1.1 Estimation error in RKHS

Based on Corollary 3.1 in Lecture 6, it suffices to bound $\mathbf{Rad}_n(l \bullet f)$ and $\mathbf{Approx}(\lambda_n)$. We treat them separately. Note that

$$\mathbb{E}_\rho \left(\sup_{f \in \mathcal{H}_n} \left| \frac{1}{n} \sum_{i=1}^n \rho_i |Y_i - f(\mathbf{X}_i)| \right| \right) \leq \mathbb{E}_\rho \|\mathbf{Rad}_n(f)\|_{\mathcal{H}_n}.$$

This bound follows from the Ledoux-Talagrand contraction inequality, as the absolute value function $\phi(z) = |z|$ is 1-Lipschitz. It suffices to bound the Rademacher complexity of \mathcal{H}_K , we have the following lemma.

Lemma 3.3. *Suppose K is a uniformly bounded kernel with $\sup_{\mathbf{x} \in \mathcal{X}} \sqrt{K(\mathbf{x}, \mathbf{x})} \leq K_0 < \infty$, \mathcal{H}_K is its corresponding RKHS, and $\mathcal{H}_K(r) = \{f \in \mathcal{H}_K : \|f\|_{\mathcal{H}_K} \leq r\}$ is a \mathcal{H}_K -ball with a radius r . Then,*

$$\mathbb{E}_\rho \|\mathbf{Rad}_n(f)\|_{\mathcal{H}_K(r)} \leq rK_0 \sqrt{\frac{1}{n}}.$$

Therefore, $\|\mathbf{Rad}_n(l \bullet f)\|_{\mathcal{H}_n} \leq cK_0(n\lambda_n)^{-1/2}$.

3.1.2 Approximation error in RKHS

Then, we turn to bound $\mathbf{Approx}(\lambda_n)$. We present Proposition 8.5 in [Cucker and Zhou, 2007] to illustrate the approximation error for RKHS. Recall $\mathbf{Approx}(\lambda_n)$ in RKHS regression:

$$\mathbf{Approx}(\lambda_n) = \inf_{f \in \mathcal{H}_n} R(f) - R(f^*) + \lambda_n \|f\|_{\mathcal{H}_n}^2,$$

provided that $R(f) = \mathbb{E}(Y - f(\mathbf{X}))^2$.

Theorem 3.4 ([Cucker and Zhou, 2007]). *Let $\mathcal{X} \subset \mathbb{R}^d$ be a compact domain, and K be a reproducing kernel. Suppose there exists $0 < s \leq 1$, such that $f^* \in \text{Range}(L_K^{s/2})$, then*

$$\mathbf{Approx}(\lambda_n) \leq A_0 \lambda_n^s,$$

where $A_0 = \mathcal{E}(L_K^{-s/2} f^*)$.

3.2 Hyperparameter tuning

Taken together, if we further assume that \mathcal{X} is a compact domain almost surely, and

$$\varepsilon_n \geq A_0 \lambda_n^{s/2} + cK_0 (n\lambda_n)^{-1/2} \geq \mathbf{Approx}(\lambda_n) + 8\mathbb{E}\|\mathbf{Rad}_n(l \bullet f)\|_{\mathcal{H}_n},$$

then

$$\mathbb{P}\left(R(\widehat{f}_n) - R^* \geq \varepsilon_n\right) \leq \exp\left(-\frac{n\varepsilon_n^2}{8(U^2 + (1/2 + U/6)\varepsilon_n)}\right).$$

Note that the developed inequality is valid for any λ_n , thus we can tune λ_n to improve the convergence rate:

$$\varepsilon_n^* = \inf_{\lambda_n} A_0 \lambda_n^{s/2} + cK_0 (n\lambda_n)^{-1/2} = O(n^{-s/2(1+s)}),$$

obtained by $\lambda_n = O(n^{-1/(1+s)})$. Therefore, the convergence rate is given as:

$$\mathcal{E}(\widehat{f}_n) = O_P(\varepsilon_n^*) = O_P(n^{-\frac{s}{2(1+s)}}).$$

References

- [Cucker and Zhou, 2007] Cucker, F. and Zhou, D. X. (2007). *Learning theory: an approximation theory viewpoint*, volume 24. Cambridge University Press.
- [Gönen and Alpaydın, 2011] Gönen, M. and Alpaydın, E. (2011). Multiple kernel learning algorithms. *The Journal of Machine Learning Research*, 12:2211–2268.
- [Kimeldorf and Wahba, 1970] Kimeldorf, G. S. and Wahba, G. (1970). A correspondence between bayesian estimation on stochastic processes and smoothing by splines. *The Annals of Mathematical Statistics*, 41(2):495–502.
- [Sejdinovic and Gretton, 2012] Sejdinovic, D. and Gretton, A. (2012). What is an rkhs? *Lecture Notes*.
- [Wahba, 1990] Wahba, G. (1990). *Spline models for observational data*. SIAM.