

Lecture 7: Nonparametric regression in RKHS

kernel method, Representer theorem

Lecturer: Ben Dai

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

— Kurt Lewin

1 Recall

Based on Lectures 1-6, we are able to compute the convergence rate and establish a probabilistic bound for a general learning method/algorithm. For illustrate, we turn to investigate the asymptotics of the nonparameteric regression in Reproducing kernel Hilbert space (RKHS).

2 RKHS

2.1 Why RKHS?

Requirements. At least, we require pointwise convergence, that is,

$$\|f_n - f\|_{\mathcal{F}} \rightarrow 0 \implies f_n(x) \rightarrow f(x), \text{ for any } x \in \mathcal{X}.$$

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

According to Riesz representation theorem, if $L = \delta_x$ is a linear continuous functional on \mathcal{H} , then we have

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

which means that we can represent function evaluation as the inner production on $K_{\mathbf{x}}$.

Definition 2.2 (RKHS). A Hilbert space \mathcal{H} is said to be a Reproducing Kernel Hilbert Space (RKHS) if $\delta_{\mathbf{x}}$ is a linear continuous functional on \mathcal{H} , for any $\mathbf{x} \in \mathcal{X}$.

Theorem 2.3. *Suppose \mathcal{H} is a RKHS, then*

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

2.2 From kernel function to RKHS

The overall idea in this section is to construct a RKHS from a kernel function. Recall the construction in finite-dimension space: (i) basis; (ii) inner production among basis.

From Riesz representation theorem, $\{K_{\mathbf{x}} : \mathbf{x} \in \mathcal{X}\}$ will be a good option as a basis function, and their inner production is given as:

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

Note that $K_{\mathbf{x}'}(\mathbf{x}) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a symmetric bivariate function (so-called kernel function). Once we define a $K(\mathbf{x}, \mathbf{x}') = K_{\mathbf{x}'}(\mathbf{x}) = K_{\mathbf{x}}(\mathbf{x}')$, then we define the basis $K_{\mathbf{x}}$ and the inner production between two basis $K(\mathbf{x}, \mathbf{x}')$.

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

To mimic the construction in finite-dimension case, we first construct a pre-RKHS \mathcal{H}_0 as a set of functions:

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

equipped with an inner production:

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

Step 2. Generate a RKHS by taking “closure” of the pre-RKHS.

Then, we create \mathcal{H} as an “closure” of \mathcal{H}_0 :

$$\mathcal{H} = \overline{\mathcal{H}_0} = \mathcal{H}_0 \cup \{\text{limit points of all } \mathcal{H}_0\text{-Cauchy sequences}\}, \quad f(\mathbf{x}) = \sum_{i=1}^{\infty} \alpha_i K(\mathbf{x}_i, \mathbf{x}) \in \mathcal{H},$$

equipped with the inner production:

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

Then, we turn to verify \mathcal{H} is indeed a RKHS, and find requirements for a valid kernel function $K(\cdot, \cdot)$. Recall the definition of RKHS: (i) \mathcal{H} equipped with $\langle f, f' \rangle_{\mathcal{H}}$ is a Hilbert space; (ii) $\delta_{\mathbf{x}}$ is a continuous functional on \mathcal{H} . One can find the formal proof 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 production**: 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 give formal definitions and theorems used in Section 2.2.

Definition 2.5 (Reproducing kernel). Let \mathcal{H} be a RKHS. A function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is called a reproducing kernel of \mathcal{H} if it satisfies

- For any $\mathbf{x} \in \mathcal{X}$, $K_{\mathbf{x}} = K(\cdot, \mathbf{x}) \in \mathcal{H}$.
- (Reproducing property). For any $\mathbf{x} \in \mathcal{X}$, and any $h \in \mathcal{H}$, $\langle h, K_{\mathbf{x}} \rangle_{\mathcal{H}} = h(\mathbf{x})$.

Definition 2.6 (Positive semi-definite (Mercer) Kernel). A symmetric function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is positive definite kernel, if for any $n \geq 1$, any $(\alpha_1, \dots, \alpha_n) \in \mathbb{R}^n$, any $(\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathcal{X}^n$,

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

Next, we have Moore-Aronszajn theorem to guarantee the legality of the construction in Section 2.2.

Theorem 2.7 (Moore-Aronszajn theorem). *Let $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be positive semi-definite. Then there is unique RKHS \mathcal{H} with reproducing kernel K .*

Finally, we summarize as:

“legal” kernel \iff reproducing kernel \iff positive semi-definite kernel \iff 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’s the difference among different kernels? Theoretically, it effects both estimation/approximation errors, see discussion in Section 4. Practically, it highly related to the topic *multiple kernel learning*, see [Gönen and Alpaydm, 2011] and references therein.

3 Regression in RKHS

We denote vector of features as $\mathbf{X} \in \mathbb{R}^d$, a scalar outcome as $Y \in \mathbb{R}$. Suppose $\mathbf{Z} = (\mathbf{X}, Y)$ satisfy a nonparametric regression model:

$$Y = f^*(\mathbf{X}) + \varepsilon,$$

where ε is a random noise with $\mathbb{E}(\varepsilon) = 0$ and $\varepsilon \perp \mathbf{X}$, and f^* is the true conditional mean function with $\|f^*\|_\infty < \infty$. Our goal is to find a decision function f minimizing the mean squared loss:

$$R(f) = \mathbb{E}\left(l(Y, f(\mathbf{X}))\right) = \mathbb{E}\left((Y - f(\mathbf{X}))^2\right).$$

Let's summarize the quantities of interests.

- **Bayes rule.** $f^*(\mathbf{x}) = \mathbb{E}(Y|\mathbf{X} = \mathbf{x})$ is the global minimizer of $R(f)$.
- **Excess risk.**

$$\mathcal{E}(f) = R(f) - R(f^*) = \mathbb{E}\left((f(\mathbf{X}) - f^*(\mathbf{X}))^2\right) = \|f - f^*\|_{L^2(\mathbb{P}_{\mathbf{X}})}^2$$

- **R-ERM.** Given random samples $(\mathbf{X}_i, Y_i)_{i=1, \dots, n}$, and a RKHS \mathcal{H}_K associated with a kernel K ,

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

- **Asymptotics.** Finally, we aim to investigate the asymptotics 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 $K(\cdot, \cdot)$ is 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. I recommend readers to check the summary¹ by Grace Wahba and Yuedong Wang. Some quotas:

- *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 of general regularization problems in infinite dimensional spaces.*

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

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

where $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n)^\top$ and $\mathbf{K} = (K_{ij})_{n \times n}$ is a kernel matrix with $K_{ij} = K(\mathbf{X}_i, \mathbf{X}_j)$. Note that the optimization in the right-hand side is reduced to a linear regression with a structured penalization, and the solution is given as:

$$\hat{\boldsymbol{\alpha}} = (\mathbf{K} + \lambda_n \mathbf{I})^{-1} \mathbf{Y}_{1:n}, \quad \hat{f}_n(\mathbf{x}) = \sum_{j=1}^n \hat{\alpha}_j K(\mathbf{X}_j, \mathbf{x}).$$

4 A loose excess risk upper bound

Next, we turn to bound the excess risk of RKHS regression. For simplicity, we assume $|\varepsilon_i| \leq U$, and we have $\hat{f}_n(\mathbf{X}_i) \leq U$, unless we can truncate it by U to further minimize R-ERM. Note that $\hat{f}_n \in \mathcal{H}_K$ is a minimizer of:

$$\frac{1}{n} \sum_{i=1}^n \left((f(\mathbf{X}_i) - f^*(\mathbf{X}_i))^2 + 2\varepsilon_i(f(\mathbf{X}_i) - f^*(\mathbf{X}_i)) \right) + \lambda_n \|f\|_{\mathcal{H}_K}^2 := \hat{R}_n(f) + \lambda_n \|f\|_{\mathcal{H}_K}^2.$$

Then, let $\mathcal{H}_n = \{f \in \mathcal{H} : \|f\|_{\mathcal{H}_K} \leq c\lambda_n^{-1/2}\}$,

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

4.1 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

$$\begin{aligned} \mathbf{Rad}_n(l \bullet f) &= \left\| \frac{1}{n} \sum_{i=1}^n \rho_i \left((f(\mathbf{X}_i) - f^*(\mathbf{X}_i))^2 + 2\varepsilon_i(f(\mathbf{X}_i) - f^*(\mathbf{X}_i)) \right) \right\|_{\mathcal{H}_n} \\ &= \left\| \frac{1}{n} \sum_{i=1}^n \rho_i (f(\mathbf{X}_i) - f^*(\mathbf{X}_i))^2 \right\|_{\mathcal{H}_n} + 2 \left\| \frac{1}{n} \sum_{i=1}^n \rho_i \varepsilon_i (f(\mathbf{X}_i) - f^*(\mathbf{X}_i)) \right\|_{\mathcal{H}_n} \\ &\leq 4U \mathbf{Rad}_n(f). \end{aligned}$$

It suffices to bound the Rademacher complexity of \mathcal{H}_K , we have the following lemma.

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

Lemma 4.1. Suppose K is a uniformly bounded kernel with $\sup_{\mathbf{x} \in \mathcal{X}} K(\mathbf{x}, \mathbf{x}) \leq K_0 < \infty$, \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, $\mathbf{Rad}_n(l \bullet f) \leq cUK_0(n\lambda_n)^{-1/2}$.

4.2 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 4.2 ([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 $A_0 = \mathcal{E}(L_K^{-s/2} f^*)$.

4.3 Hyperparameter tuning

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

$$\varepsilon_n \geq A_0 \lambda_n^s + cUK_0(n\lambda_n)^{-1/2} \geq \mathbf{Approx}(\lambda_n) + 8\|\mathbb{E}\mathbf{Rad}_n(l \bullet f)\|_{\mathcal{H}_n},$$

then

$$\mathbb{P}\left(\|f - f^*\|_{L^2(\mathbb{P}_{\mathbf{x}})}^2 \geq \varepsilon_n\right) \leq \exp\left(-\frac{n\varepsilon_n^2}{8(U^2 + (1/2 + U/3)\varepsilon_n)}\right).$$

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

$$\varepsilon_n^* = \inf_{\lambda_n} A_0 \lambda_n^s + cUK_0(n\lambda_n)^{-1/2} = O(n^{-s/(1+2s)}),$$

obtained by $\lambda_n = O(n^{-1/(1+2s)})$. Therefore, the convergence rate is given as:

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

Remark 4.3. Note that the result is a simple but loose bound for the excess risk, we can improve the convergence rate from $n^{-s/(1+2s)}$ to $n^{-2s/(1+2s)}$ in the next few lectures. Yet, we are looking for where there is potential for improvement...

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.