

REPRODUCING KERNEL HILBERT SPACES

RKHSs are spaces of functions that have an interesting mathematical theory, and that have been applied is statistics and more recently machine learning.

We will see that some familiar supervised kernel methods, such as KRR and SVMs, can be cast in terms of RKHSs.

New algorithms can also be derived from this general framework.

[The "Mathematical Background" section can be skipped]

Mathematical Background

Hilbert Spaces

Let \mathcal{H} be a vector space over \mathbb{R} ,
with inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$.

This inner product induces a norm

$$\|v\|_{\mathcal{H}} = \sqrt{\langle v, v \rangle_{\mathcal{H}}}$$

and distance

$$d_{\mathcal{H}}(u, v) = \|u - v\|_{\mathcal{H}}.$$

With this distance we can define convergent
and Cauchy sequences.

Defn | Let $v_1, v_2, \dots \in \mathcal{H}$.

1) We say the sequence converges if $\exists v \in \mathcal{H}$ s.t.

$$\lim_{n \rightarrow \infty} d(v_n, v) = 0$$

2) We say the sequence is Cauchy if $\forall \varepsilon > 0$,

$$\exists N \text{ s.t. } m, n \geq N \Rightarrow d(v_m, v_n) < \varepsilon.$$

3) \mathcal{H} is a Hilbert space if every Cauchy

sequence in \mathcal{H} converges (more generally, a metric
space is complete iff every Cauchy sequence converges)

Loosely speaking: every sequence that looks like it converges actually converges (to a point in \mathcal{H})

Examples

- $\mathcal{H} = \mathbb{R}^d$, $\langle u, v \rangle = \sum_{i=1}^d u^{(i)} v^{(i)}$

is complete

- $\mathbb{Q}^d = \{\text{rational numbers}\}^d$

is not complete

- $L^2(\mathbb{R}^d) = \{f: \mathbb{R}^d \rightarrow \mathbb{R} \mid \int f(x)^2 dx < \infty\}$

$$\langle f, g \rangle = \int f(x)g(x) dx$$

is complete

- $C(\mathbb{R}^d) \subseteq L^2(\mathbb{R}^d)$ is not complete

↖ continuous functions

Linear Functions

A function $L: \mathcal{H} \rightarrow \mathbb{R}$ is linear iff

$$L(au + bv) = aL(u) + bL(v)$$

$$\forall a, b \in \mathbb{R}, u, v \in \mathcal{H}.$$

Theorem 1 (Riesz Representation Theorem)

Suppose \mathcal{H} is a Hilbert space and $L: \mathcal{H} \rightarrow \mathbb{R}$ linear.

Then L is continuous $\Leftrightarrow \exists u \in \mathcal{H}$ such that

$$L(v) = \langle u, v \rangle \quad \forall v \in \mathcal{H}.$$

RKHSs

Now suppose \mathcal{H} is a Hilbert space whose elements are functions $f: \Omega \rightarrow \mathbb{R}$, where $\Omega \subseteq \mathbb{R}^d$.

Def 1 \mathcal{H} is a RKHS $\Leftrightarrow \forall x \in \Omega$, the linear function

$$L_x(f) := f(x)$$

is continuous.

By the Riesz Representation Theorem, \mathcal{H} is a RKHS

$\Leftrightarrow \forall x \in \Omega, \exists \Phi_x \in \mathcal{H}$ s.t.

$$f(x) = \langle \Phi_x, f \rangle.$$

$\forall f \in \mathcal{H}$.

The function $k(x', x) := \Phi_x(x')$ is called the reproducing kernel of \mathcal{H} .

It can easily be seen that k is a PSD kernel.

• symmetry:

$$\begin{aligned}k(x', x) &= \Phi_x(x') && \text{[def of } k\text{]} \\&= \langle \Phi_{x'}, \Phi_x \rangle && \text{[Riesz rep. of } \Phi_x\text{]} \\&= \langle \Phi_x, \Phi_{x'} \rangle && \text{[symmetry of } \langle \cdot, \cdot \rangle\text{]} \\&= \Phi_{x'}(x) && \text{[Riesz rep. of } \Phi_{x'}\text{]} \\&= k(x, x') && \text{[def of } k\text{]}\end{aligned}$$

• PSD: For any $x_1, \dots, x_n \in \Omega$,

$$K = [k(x_i, x_j)] = [\langle \Phi_{x_i}, \Phi_{x_j} \rangle]$$

is a Gram matrix and therefore a PSD matrix.

Conversely, for every PSD kernel k , there is an associated RKHS.

- Let $k(\cdot, x)$ denote the function $x' \mapsto k(x', x)$.
- Let $\mathcal{H}_0 = \{ f = \sum_{i=1}^n \alpha_i k(\cdot, x_i) \mid x_i \in \Omega, \alpha_i \in \mathbb{R}, n \in \mathbb{N} \}$,
- For $f = \sum_{i=1}^n \alpha_i k(\cdot, x_i)$, $f' = \sum_{i=1}^{n'} \alpha'_i k(\cdot, x'_i)$, set

$$\langle f, f' \rangle = \sum_{i=1}^n \sum_{j=1}^{n'} \alpha_i \alpha'_j k(x_i, x'_j).$$

It can be shown that this is an inner product on \mathcal{H}_0 .

- Now let \mathcal{H} be the completion of \mathcal{H}_0 .

Then \mathcal{H} is a RKHS, and its reproducing kernel is k

Note | The above construction shows that for any PSD kernel, \exists a Hilbert space \mathcal{H} and a feature map $x \mapsto \Phi_x$ s.t. $k(x, x') = \langle \Phi_x, \Phi_{x'} \rangle$.

In particular, $\Phi_x = k(\cdot, x)$.

Ref | Steinwart and Christmann, Support Vector Machines, Springer, 2008.

RKHS Overview

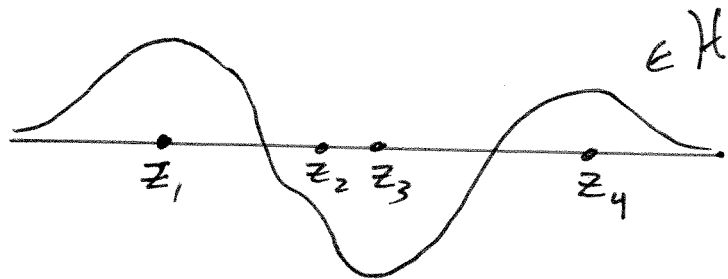
A RKHS is an inner product space of real-valued functions on \mathbb{R}^d , defined in terms of a PSD kernel as follows:

$$\mathcal{H}_0 = \left\{ \sum_{i=1}^m \alpha_i k(\cdot, z_i) \mid \alpha_i \in \mathbb{R}, z_i \in \mathbb{R}^d, m \in \mathbb{N} \right\}$$

$$\begin{aligned} \left\langle \sum_{i=1}^m \alpha_i k(\cdot, z_i), \sum_{j=1}^n \beta_j k(\cdot, z'_j) \right\rangle &:= \sum_{i,j} \alpha_i \beta_j k(z_i, z'_j) \\ &= \alpha^T K \beta, \quad K = [k(z_i, z'_j)] \end{aligned}$$

$\mathcal{H} = \mathcal{H}_0$ plus all its limit points (with finite norm)

Example | Gaussian kernel



Remarks |

• For every $f \in \mathcal{H}$ and $x \in \mathbb{R}^d$, $f(x) = \langle f, k(\cdot, x) \rangle$.
Proof is obvious for $f \in \mathcal{H}_0$. Known as the reproducing property.

• The above construction is how one shows that every PSD kernel is an IP kernel. \mathcal{H} is the feature space and $\Phi(x) := k(\cdot, x)$ is the feature mapping (note $\langle \Phi(x), \Phi(x') \rangle = k(x, x')$ by definition of $\langle \cdot, \cdot \rangle$).

The Representer Theorem

RKHSs are nice because certain optimization problems over infinite-dimensional RKHSs actually reduce to finite dimensional optimization problems.

Theorem | Let \mathcal{H} be a RKHS of functions f defined on a domain $\Omega \subseteq \mathbb{R}^d$. Consider an optimization problem of the form

$$\min_{f \in \mathcal{H}} J(f) \quad (*)$$

where J can be expressed

$$J(f) = L(f(x_1), \dots, f(x_n)) + \Lambda(\|f\|_{\mathcal{H}}^2)$$

\uparrow arbitrary \uparrow nondecreasing

for some $x_1, \dots, x_n \in \Omega$.

Then there exists a minimizer of $(*)$ of the form

$$f^*(x) = \sum_{i=1}^n \alpha_i k(x, x_i)$$

where k is the reproducing kernel of \mathcal{H} and $\alpha_i \in \mathbb{R}$.

Proof Consider the subspace $S \subseteq \mathcal{H}$ given by

$$S = \left\{ \sum_{i=1}^n \alpha_i k(\cdot, x_i) \mid \alpha_i \in \mathbb{R}, i=1, \dots, n \right\}$$

This is a finite dimensional subspace, and therefore

$$\mathcal{H} = S \oplus S^\perp.$$

For any $f \in \mathcal{H}$, let us write $f = f_{\parallel} + f_{\perp}$

where $f_{\parallel} \in S$, $f_{\perp} \in S^\perp$. Since

$$\begin{aligned} f(x_i) &= \langle f, k(x_i, \cdot) \rangle \\ &= \langle f_{\parallel}, k(x_i, \cdot) \rangle + \langle f_{\perp}, k(x_i, \cdot) \rangle \\ &= \langle f_{\parallel}, k(x_i, \cdot) \rangle \\ &= f_{\parallel}(x_i), \end{aligned}$$

the value of $L(f(x_1), \dots, f(x_n))$ depends only on f_{\parallel} .

Let us write $L(f)$ for $L(f(x_1), \dots, f(x_n))$. Now

$$J(f) = L(f) + \Lambda(\|f\|_{\mathcal{H}}^2)$$

$$= L(f_{\parallel}) + \Lambda(\|f\|_{\mathcal{H}}^2)$$

$$\geq L(f_{\parallel}) + \Lambda(\|f_{\parallel}\|_{\mathcal{H}}^2)$$

$$= J(f_{\parallel})$$

Since Λ is nondecreasing and $\|f\|^2 = \|f_{\parallel}\|^2 + \|f_{\perp}\|^2$.

Therefore, if f is a minimizer of J , then so is f_{\perp} .

Since $f_{\perp} \in S$, it has the desired form \square

Note From the proof we see that if λ is

strictly increasing, then $J(f) > J(f_{\perp})$

unless $f_{\perp} = 0$. Therefore all minimizers of J

have the stated form.

Kernel Ridge Regression

Let's apply the preceding with

$$L(f(x_1), \dots, f(x_n)) = \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\Lambda(\|f\|^2) = \lambda \|f\|^2$$

By the Representer Theorem, the solution of

$$\min_{f \in \mathcal{H}} \sum (y_i - f(x_i))^2 + \lambda \|f\|_{\mathcal{H}}^2$$

has the form

$$f = \sum \alpha_i k(\cdot, x_i).$$

Therefore it suffices to solve

$$\min_{\alpha} \sum_{i=1}^n (y_i - \sum_{j=1}^n \alpha_j k(x_i, x_j))^2 + \lambda \left\| \sum_{j=1}^n \alpha_j k(\cdot, x_j) \right\|_{\mathcal{H}}^2$$

Now

$$\begin{aligned} \left\| \sum_{j=1}^n \alpha_j k(\cdot, x_j) \right\|^2 &= \sum_{i,j} \alpha_i \alpha_j k(x_i, x_j) \\ &= \alpha^T K \alpha \end{aligned}$$

and

$$\begin{aligned} \sum_{i=1}^n (y_i - \sum_j \alpha_j k(x_i, x_j))^2 &= \|y - K\alpha\|^2 \\ &= y^T y - 2y^T K\alpha + \alpha^T K^2 \alpha, \end{aligned}$$

So we need to minimize

$$\alpha^T (K^2 + \lambda K) \alpha - 2y^T K \alpha$$

$$\Rightarrow 2(K^2 + \lambda K) \alpha - 2Ky = 0$$

$$\Rightarrow \alpha = (K + \lambda I)^{-1} y \quad [\text{assuming } K \text{ invertible}]$$

This is identical to kernel ridge regression. Easy!

Support Vector Machines

Now consider

$$L(f(x_1), \dots, f(x_n)) = \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i f(x_i))$$

$$\Lambda(\|f\|^2) = \frac{\lambda}{2} \|f\|^2$$

Once again, by the representer theorem, the minimizer has the form

$$f^*(x) = \sum_{i=1}^n r_i k(x, x_i),$$

$r_i \in \mathbb{R}$. This reduces the optimization problem to

$$\min_r \quad \frac{c}{n} \sum_{i=1}^n \max(0, 1 - y_i \cdot \sum_{j=1}^n r_j k(x_i, x_j)) + \frac{1}{2} r^T K r$$

$$[c = \frac{1}{\lambda}]$$

$$\Leftrightarrow \min_{r, \xi} \quad \frac{1}{2} r^T K r + \frac{c}{n} \sum_{i=1}^n \xi_i$$

$$\text{s.t.} \quad y_i \cdot \sum_{j=1}^n r_j k(x_i, x_j) \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

The Lagrangian is

$$L(r, \xi, \alpha, \beta) = \frac{1}{2} r^T K r + \frac{c}{n} \sum \xi_i - \sum_i \alpha_i \left[y_i \left(\sum_j r_j k(x_i, x_j) \right) - 1 + \xi_i \right] - \sum \beta_i \xi_i$$

We can write

$$\sum_i \alpha_i y_i \left(\sum_j r_j k(x_i, x_j) \right) = (\alpha \odot y)^T K r$$

↑ element-wise product

By the KKT conditions,

$$0 = \frac{\partial L}{\partial r} = K r - K (\alpha \odot y)$$

$$\Leftrightarrow r = \alpha \odot y \quad [\text{if } K \text{ invertible}]$$

$$\Leftrightarrow r_i = \alpha_i y_i \quad \forall i$$

and

$$0 = \frac{\partial L}{\partial \xi_i} = \frac{c}{n} - \alpha_i - \beta_i$$

Converting to the dual we have

$$\max_{\alpha} \quad -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(x_i, x_j) + \sum_i \alpha_i$$

$$\text{s.t.} \quad 0 \leq \alpha_i \leq \frac{c}{n}$$

This is the SVM dual. ^(without offset) Furthermore, since $r_i = \alpha_i y_i$, we have

$$f^*(x) = \sum_{i=1}^n \alpha_i y_i k(x, x_i).$$

Thus, we recover the SVM decision function.

Other Instances

- Kernel logistic regression

$$L(f(x_1), \dots, f(x_n)) = \sum_{i=1}^n \log(1 + \exp(-y_i f(x_i)))$$

$$\Lambda(\|f\|) = \frac{\lambda}{2} \|f\|^2$$

Exercise | Derive a Newton-type algorithm to find the optimal α_i .

- Cubic smoothing spline

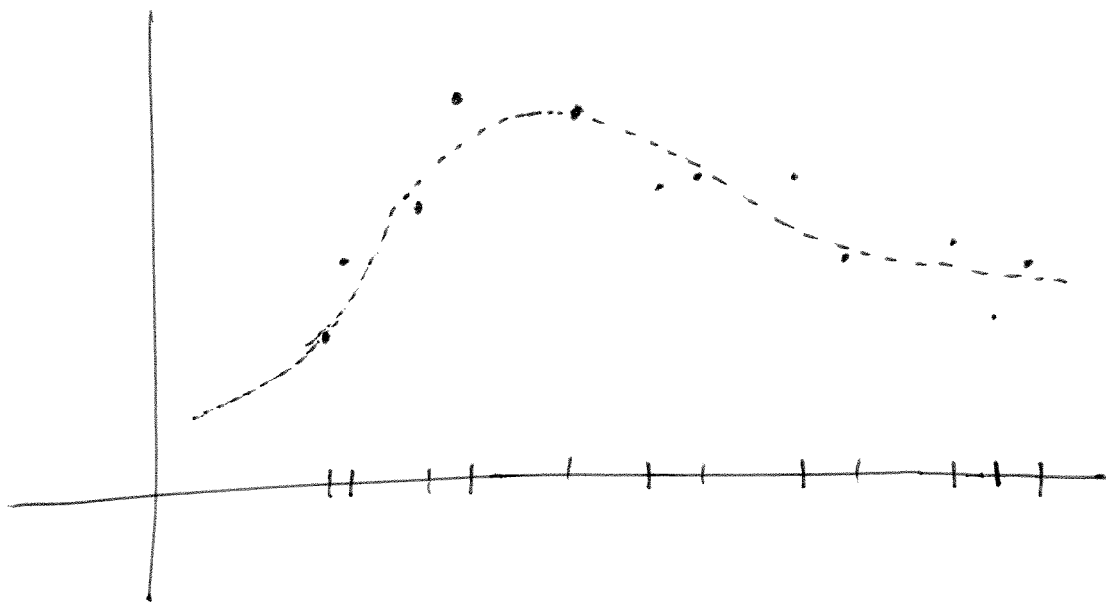
$$\min_{f \in \mathcal{H}} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \int_a^b [f''(x)]^2 dx$$

where $\mathcal{H} = \left\{ f: [a, b] \rightarrow \mathbb{R} \mid \int_a^b [f''(x)]^2 dx < \infty \right\}$.

It can be shown, using a more general version of the representer theorem, that the solution is a cubic spline with knots at the x_i .

↔ piecewise cubic polynomial.

Thus the optimization problem reduces to a finite dimensional one, and the optimal cubic spline results from solving a system of linear equations.



Reference | Nancy Heckman, "The theory and application of penalized methods, or Reproducing Kernel Hilbert Spaces made easy," 2011.

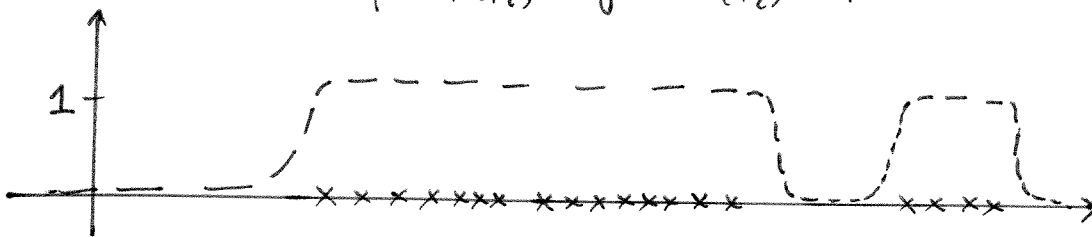
- Robust Kernel Density Estimation
→ Kim and Scott, "On the robustness of kernel density M-estimators," ICML 2011

- one-class SVM for novelty detection
given training data from one class: x_1, \dots, x_n

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \phi(f(x_i)) + \lambda \|f\|_{\mathcal{H}}^2, \quad \lambda > 0$$

where $\phi(t) = \max(0, 1-t)$, so

$$\phi(f(x_i)) = \begin{cases} 0 & \text{if } f(x_i) \geq 1 \\ 1-f(x_i) & \text{if } f(x_i) < 1 \end{cases}$$



Final classifier: $\text{sign}(f(x) - \delta)$

Using representer theorem and KKT conditions (as in the case of the SVM above), one can show the solution is

$f(x) = \sum \alpha_i k(x, x_i)$, where $\alpha = (\alpha_1, \dots, \alpha_n)$ solves

$$\max_{\alpha} -\frac{1}{2} \sum_{i,j} \alpha_i \alpha_j k(x_i, x_j) + \sum \alpha_i$$

$$\text{s.t. } 0 \leq \alpha_i \leq \frac{1}{\lambda n}$$

Note: The OC-SVM is often parameterized in a different way:
see Lee and Scott, "The one-class support vector machine solution path," ICASSP 2007.

- Multitask Learning

→ Evgeniou + Pontil, "Learning multiple tasks with kernel methods," JMLR, 2005.

- Transfer Learning

→ Blanchard, Lee, + Scott, "Generalizing from several related classification tasks to a new unlabeled sample," NIPS 2011.

Summary

- infinite dimensional optimization problem

↓ (representer theorem)

finite dimensional optimization problem

- recover existing methods

- derive new ones