

REPRODUCING KERNEL HILBERT SPACES

Overview

A reproducing kernel Hilbert space (RKHS) is a space of real-valued functions defined in terms of a positive definite kernel. By optimizing over a RKHS, we can derive many ML algorithms - some previously seen and some new.

A Rough Definition

Everything that follows can be made rigorous, but I will only present some rough ideas to enable some intuition.

Let $k: X \times X \rightarrow \mathbb{R}$ where X is the input space. We previously asserted that the following are equivalent :

(a) k is symmetric and positive definite

(b) \exists an inner product space \mathcal{H} and a feature map

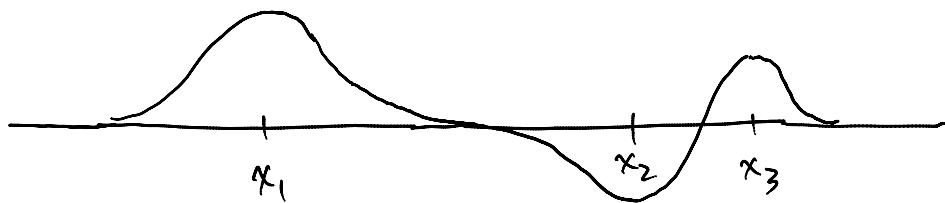
$\Phi: X \rightarrow \mathcal{H}$ such that $k(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\mathcal{H}}$.

The proof of (b) \Rightarrow (a) was easy. To prove (a) \Rightarrow (b), first define

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

where $k(\cdot, x')$ is the function $x \mapsto k(x, x')$. ↑ set of natural numbers

Thus, \mathcal{H}_0 is a space of functions. For example, if k is the Gaussian kernel, then



depicts an element of \mathcal{H}_0 .

Define an inner product on \mathcal{H}_0 by

$$\left\langle \sum_{i=1}^m \alpha_i k(\cdot, x_i), \sum_{j=1}^n \beta_j k(\cdot, x'_j) \right\rangle := \sum_{i=1}^m \sum_{j=1}^n \alpha_i \beta_j k(x_i, x'_j)$$

it can be shown that this is a valid inner product.

Now, consider the feature map $\Phi: X \rightarrow \mathcal{H}_0$ given by

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

By definition of the inner product, observe

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

This establishes (b). The feature map above is known as the canonical feature map. It works for any PD kernel.

The reproducing property states that for any $f \in H_0$, $x \in X$

$$f(x) = \langle f, k(\cdot, x) \rangle.$$

To see this, let $f = \sum_{i=1}^m \alpha_i k(\cdot, x_i)$. Then

$$\begin{aligned}\langle f, k(\cdot, x) \rangle &= \left\langle \sum \alpha_i k(\cdot, x_i), k(\cdot, x) \right\rangle \\ &= \sum \alpha_i k(x_i, x) \\ &= \sum \alpha_i k(x, x_i) \\ &= f(x).\end{aligned}$$

For technical reasons that we won't cover, we need to enlarge H_0 slightly. Define H to be the completion of H_0 , so that H is a Hilbert space (the name

for a complete inner product space). Completion essentially adds to \mathcal{H}_0 those functions that are not in \mathcal{H}_0 , but that can be approximated arbitrarily accurately by elements of \mathcal{H}_0 .

It can be shown that the reproducing property still holds for all $f \in \mathcal{H}$. \mathcal{H} is known as the reproducing kernel Hilbert space associated with the PD kernel k . k is called the reproducing kernel of \mathcal{H} .

Remark] The results and derivations below actually hold if we take $\mathcal{H} = \mathcal{H}_0$. Completeness of \mathcal{H} is needed for more theoretical considerations.

The Representer Theorem

Previously, we have derived kernel methods by optimizing over a class of linear models and then kernelizing. As an alternative, we can optimize over the RKHS directly. Even though an RKHS may be infinite dimensional, by the following result, optimization problems of a certain type reduce to finite-dimensional problems.

reduce to finite-dimensional programs.

Theorem | Let \mathcal{H} be an RKHS consisting of functions defined on X . Consider an optimization problem of the form

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

where

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

for some $x_1, \dots, x_n \in X$, and where Λ is nondecreasing and $\|f\|_{\mathcal{H}} = \langle f, f \rangle_{\mathcal{H}}^{1/2}$. Then there exists a minimizer of the form

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

Furthermore, if Λ is strictly increasing, then every minimizer has the given form.

Remark | The notation $L(f(x_1), \dots, f(x_n))$ indicates that this term does not depend on values of f outside of $\{x_1, \dots, x_n\}$.

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

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

Since S is finite dimensional, the projection theorem holds and

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

$$\text{where } S^\perp = \{g \in \mathcal{H} \mid \langle f, g \rangle = 0 \quad \forall f \in S\}.$$

Thus, for each $f \in \mathcal{H}$, we may write $f = f_{\parallel} + f_{\perp}$
 where $f_{\parallel} \in S$, $f_{\perp} \in S^\perp$ are unique. Since $\forall i=1,\dots,n$

$$f(x_i) = \langle f, k(\cdot, x_i) \rangle$$

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

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

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

$$\begin{aligned} J(f) &= L(f) + \lambda (\|f\|^2) \\ &= L(f_{\parallel}) + \lambda (\|f\|^2) \quad \text{because } \lambda \text{ is} \\ &\geq L(f_{\parallel}) + \lambda (\|f_{\parallel}\|^2) \quad \text{nondecreasing and} \\ &\qquad \|f\|^2 = \|f_{\parallel}\|^2 + \|f_{\perp}\|^2 \\ &= J(f_{\parallel}). \end{aligned}$$

Therefore, if f is a minimizer, then so is f_{\parallel} . Since

$f_{\parallel} \in S$, it has the desired form.

If Λ is strictly increasing, then $J(f_{\parallel}) < J(f)$ unless $f_{\perp} = 0$. Thus, every minimizer has the stated form $\boxed{\square}$

Kernel Ridge Regression

Let's apply the framework with

$$L(f(x_1), \dots, f(x_n)) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$
$$\Lambda(f) = \lambda \|f\|_{\mathcal{H}}^2.$$

By the representer theorem, the solution of

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \|f\|_{\mathcal{H}}^2$$

has the form

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

Therefore it suffices to solve

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

Now

$$\left\| \sum_j \alpha_j k(\cdot, x_j) \right\|^2 = \sum_{i,j} \alpha_i \alpha_j k(x_i, x_j)$$
$$= \cdot^T K \cdot$$

$$= \alpha^\top \alpha$$

where

$$K = [k(x_i, x_j)]_{ij} \quad (n \times n)$$

is the kernel matrix. Also,

$$\begin{aligned} \sum_{i=1}^n (y_i - \sum_j \alpha_j k(x_i, x_j))^2 &= \|y - K\alpha\|_2^2 \\ &= y^\top y - 2y^\top K\alpha + \alpha^\top K^2 \alpha. \end{aligned}$$

So we need to minimize

$$\alpha^\top (K^2 + n\lambda I) \alpha - 2y^\top K\alpha.$$

Equating the gradient to zero we have

$$\begin{aligned} 0 &= 2(K^2 + n\lambda I)\alpha - 2Ky \\ &= 2K[(K + n\lambda I)\alpha - y] \end{aligned}$$

which is solved by

$$\alpha = (K + n\lambda I)^{-1} y.$$

Therefore, we have re-derived kernel ridge regression without offset.

Support Vector Machines

Now consider

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

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

Again by the representer theorem, the minimizer has the form

$$f = \sum_{i=1}^n r_i k(\cdot, x_i), \quad r_i \in \mathbb{R}.$$

Denoting $r = [r_1 \dots r_n]^T$, and substituting $C = \frac{1}{\lambda}$,

the optimization problem reduces to

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

Introducing slack variables $\xi = [\xi_1 \dots \xi_n]^T$,

this is equivalent to

$$\min_{r, \xi} \frac{1}{2} r^T K r + \frac{C}{n} \sum \xi_i$$

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

$$\xi_i \geq 0 \quad \forall i.$$

This is a convex, differentiable optimization problem with affine constraints, so strong duality holds. The problem is also differentiable, and so the KKT conditions hold.

The Lagrangian is

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

Let us write

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

where \odot denotes element-wise product. By the KKT conditions,

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

This is solved by

$$r = \alpha \odot y$$



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

Also,

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

Connecting to the dual we have

$$\begin{aligned} \max_{\alpha} \quad & -\frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j k(x_i, x_j) + \sum_i \alpha_i \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq \frac{c}{n} \quad \forall i. \end{aligned}$$

This is the dual for the SVM without offset. We also recover the SVM decision function:

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

Kernel Logistic Regression

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

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

One - Class SVM

One-class classification refers to binary classification where there are training examples for only one of the classes.

Let x_1, \dots, x_n denote the data from the observed class.

We will consider a classifier of the form

$$x \mapsto \text{sign}(f(x) - 1)$$

where $f \in \mathcal{H}$ and +1 is the observed class. To select f , we would like to trade off between two criteria:

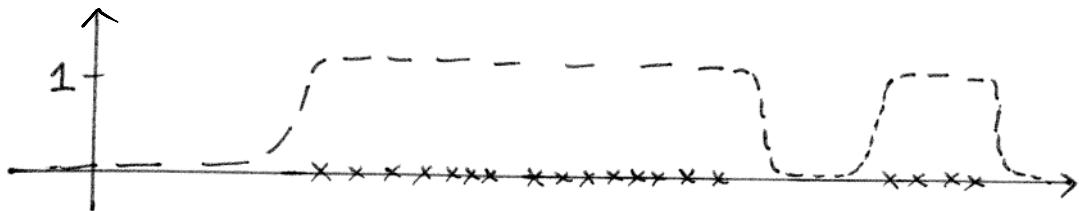
- Correctly classify training data
- Small $\|f\|$ to avoid overfitting.

Consider

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_i \phi(f(x_i)) + \frac{\lambda}{2} \|f\|^2$$

where $\phi(t) = \max(0, 1-t)$ is the hinge loss. Notice

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



Using the representer theorem and KKT conditions as before, one can show that the solution is

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

where $\alpha = [\alpha_1 \dots \alpha_n]^T$ solves

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

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

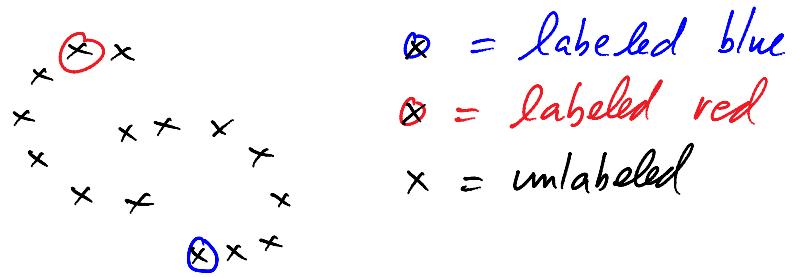
Note] The OC-SVM is often parametrized differently.

See Lee and Scott, "The one-class support vector machine solution path," ICASSP 2007.

Semi-Supervised Learning

Semi-supervised learning refers to a supervised learning problem where, in addition to labeled examples, there are also unlabeled examples. The goal is to leverage the unlabeled examples to improve the performance of a method that uses only the labeled data.

For example, consider a classification problem:



Assuming the clusters are the classes, the unlabeled data help us learn the clusters which allows us to generalize well.

Let $(x_1, y_1), \dots, (x_m, y_m), x_{m+1}, \dots, x_{m+n}$ denote the training data. For regression, consider

$$\min_{f \in \mathcal{H}} \lambda \|f\|_X^2 + \frac{1}{n} \sum_i (y_i - f(x_i))^2 + \gamma \cdot \underbrace{\frac{1}{2} \sum_{i,j} w_{ij} (f(x_i) - f(x_j))^2}_{= \underline{f^T L f}}$$

$$\text{where } \underline{f} = [f(x_{m+1}) \dots f(x_{m+n})]^T$$

where $W = [w_{ij}]_{i,j=1}^n$ is an $n \times n$ weighted adjacency matrix based on the unlabeled data. Intuitively, this term encourages f to take similar values at similar points. The solution can be derived in a manner similar to kernel ridge regression.