

## Examples

### Linear kernel

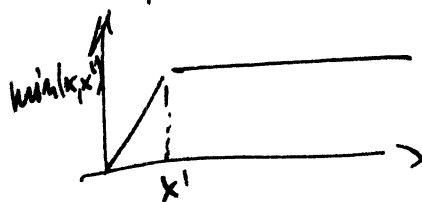
$$\mathcal{H} = \left\{ f : f(x) = \sum_{j=1}^n \alpha_j x_j^T x \text{ for some } n \in \mathbb{N}, \alpha \in \mathbb{R}^n, x_1, \dots, x_n \in \mathbb{R}^p \right\}$$

$$\text{or } \mathcal{H} = \left\{ f : f(x) = \beta^T x \text{ for some } \beta \in \mathbb{R}^p \right\}$$

$$\text{If } f(x) = \beta^T x, \text{ i.e. } f(\cdot) = k(\cdot, \beta), \text{ then } \|f\|_{\mathcal{H}}^2 = \langle k(\cdot, \beta), k(\cdot, \beta) \rangle \\ = k(\beta, \beta) = \|\beta\|_2^2$$

$$\text{Sobolev kernel } k(x, x') = \min(x, x'), \quad x, x' \in [0, 1]$$

$\mathcal{H}$  includes linear combinations of functions  $x \mapsto \min(x, x')$  for  $x' \in [0, 1]$  and their pointwise limits.



Can show that  $\mathcal{H}$  includes all Lipschitz functions on  $[0, 1]$  that are 0 at the origin (Lipschitz means  $|f(x) - f(x')| \leq L|x - x'|$  for some  $L \forall x, x' \in [0, 1]$ ).

$$\text{The norm } \left( \int_0^1 f'(x)^2 dx \right)^{1/2}$$

$$\text{e.g. } \|\min(\cdot, x')\|_{\mathcal{H}}^2 = \min(x', x') = x'$$

$$\text{Also } \int_0^{x'} 1^2 dx = x' \quad \checkmark$$

### 1.4.3 The representer theorem

(KRR)

What optimisation problem is kernel ridge regression solving? An alternative way of writing the usual ridge optimisation is

$$\arg \min_{f \in \mathcal{H}} \left\{ \sum_{i=1}^n (Y_i - \frac{f(x_i)}{x_i^T \beta})^2 + \lambda \frac{\|f\|_{\mathcal{H}}^2}{\|\beta\|_2^2} \right\} \quad \text{where } \mathcal{H} \text{ is the RKHS of the linear kernel.}$$

### Theorem 6 (the representer theorem)

Let  $c$  be an arbitrary loss function and suppose  $J$  is strictly increasing. Let  $\mathcal{H}$  be an RKHS with reproducing kernel  $k$ . Then  $\hat{f} \in \mathcal{H}$  minimises

$$Q_1(f) = c(Y, X, f(x_1), \dots, f(x_n)) + J(\|f\|_{\mathcal{H}}^2) \text{ over } f \text{ iff}$$

$$\hat{f}(\cdot) = \sum_{i=1}^n \hat{\alpha}_i k(\cdot, x_i) \text{ and } \hat{\alpha}_i \in \mathbb{R} \text{ minimises}$$

$$Q_2(\alpha) = c(Y, X, K\alpha) + J(\alpha^T K \alpha) \text{ over } \alpha \in \mathbb{R}^n. \\ \hat{\alpha}^T K \hat{\alpha} = \|\hat{f}\|_{\mathcal{H}}^2$$

Remark Squaring to ridge, (\*) is equivalent to minimizing

$\|Y - K\alpha\|_2^2 + \alpha^T K \alpha$ . Example 1 qn 9 shows that minimizer  $\hat{\alpha}$  satisfies  $K\hat{\alpha} = K(K + \lambda I)^{-1}Y$ . Thus KRR is exactly (\*).

Proof Suppose  $\hat{f}$  minimizes  $Q_1$ . We may write  $\hat{f} = u + v$  where

$u \in V = \text{span}\{k(\cdot, x_1), \dots, k(\cdot, x_n)\}$  and  $v \in V^\perp$ .

Then  $\hat{f}(x_i) = \langle \hat{f}, k(\cdot, x_i) \rangle = \langle u + v, k(\cdot, x_i) \rangle = u(x_i)$

Meanwhile  $J(\|\hat{f}\|_H^2) = J(\|u\|_H^2 + \|v\|_H^2) \geq J(\|u\|_H^2)$  with equality iff  $v = 0$ .

By optimality of  $\hat{f}$ ,  $v = 0$ . So  $\hat{f}(\cdot) = \sum_{i=1}^n \hat{\alpha}_i k(\cdot, x_i)$ , and again by optimality of  $\hat{f}$ ,  $\hat{\alpha}$  must minimize  $Q_2$ .

Now suppose  $\hat{\alpha}$  minimizes  $Q_2$  and let  $\hat{f}(\cdot) = \sum_{i=1}^n \hat{\alpha}_i k(\cdot, x_i)$

if  $\tilde{f} \in \mathcal{H}$  with  $Q_1(\tilde{f}) \leq Q_2(\hat{f})$ , by the argument above, we can write  $\tilde{f} = u + v$  with  $u \in V$ ,  $v \in V^\perp$  and we know  $Q_1(u) \leq Q_1(\tilde{f})$ . But by optimality of  $\hat{f}$ ,

$Q_1(\hat{f}) \leq Q_1(u)$  so  $Q_1(\hat{f}) = Q_1(\tilde{f})$ .  $\square$

The representer theorem gives the form of the entire fitted regression  $\hat{f}^n$  (not just the fitted vals).  
e.g. with KRR, given a new obs  $x$ , our prediction would be

$$\hat{f}(x) = \sum_{i=1}^n \hat{\alpha}_i k(x, x_i)$$

## 1.5 Kernel ridge regression

Consider a model  $Y_i = f^0(x_i) + \varepsilon_i$ ,  $\text{Var}(\varepsilon) = \sigma^2 I$ ,  $\mathbb{E} \varepsilon = 0$ .

Assume  $f^0 \in \mathcal{H}$  where  $\mathcal{H}$  is a RKHS with reproducing kernel  $k$ . Assume  $\|f^0\|_H \leq 1$ .

Let  $K$  be the kernel matrix  $K_{ij} = k(x_i, x_j)$  with eigenvalues  $d_1, d_2, \dots, d_n \geq 0$ .

Define  $\hat{f}_\lambda = \arg\min_{f \in \mathcal{H}} \left\{ \sum_{i=1}^n (Y_i - f(x_i))^2 + \lambda \|f\|_H^2 \right\}$ .

Theorem 7 The mean squared prediction error of  $\hat{f}_\lambda$  (MSPE) is

$$\frac{1}{n} \mathbb{E} \sum_{i=1}^n (f^0(x_i) - \hat{f}_\lambda(x_i))^2 \leq \frac{\sigma^2}{n} \sum_{i=1}^n \frac{d_i^2}{(d_i + \lambda)^2} + \frac{\lambda}{4n} \leq \frac{\sigma^2}{n} \frac{1}{\lambda} \sum_{i=1}^n \min\left(\frac{d_i}{4}, \lambda\right) + \frac{\lambda}{4n}$$