

KERNEL RIDGE REGRESSION

Recall ridge regression: given $(x_i, y_i)_{i=1}^n$, $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$,

solve

$$\min_{\beta_0, \beta} \sum_{i=1}^n (y_i - \beta^T x_i - \beta_0)^2 + \lambda \|\beta\|^2$$

The solution is

$$\hat{\beta} = (A^T A + \lambda I)^{-1} A^T \underline{y},$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}^T \bar{x}$$

$$A = \begin{bmatrix} x_1^{(1)} & \dots & x_1^{(d)} \\ \vdots & & \vdots \\ x_n^{(1)} & \dots & x_n^{(d)} \end{bmatrix}$$

and the final estimator is

$$\hat{f}(x) = \hat{\beta}^T x + \hat{\beta}_0$$

Can these operations be expressed entirely in terms of $\langle x_i, x_j \rangle$ and $\langle x_i, x \rangle$?

Not immediately. Note $A^T A$ is not $[\langle x_i, x_j \rangle]_{i,j=1}^n$

Let's apply the matrix inversion lemma, also known as the Woodbury matrix identity:

$$(P + QRS)^{-1} = P^{-1} - P^{-1}Q(R^{-1} + SP^{-1}Q)^{-1}SP^{-1}$$

where

$$P = \lambda I, \quad Q = A^T, \quad R = I, \quad S = A.$$

We have

$$(\lambda I + A^T A)^{-1} = \frac{1}{\lambda} I - \frac{1}{\lambda} I \cdot A^T (I + \frac{1}{\lambda} A A^T)^{-1} A \cdot \frac{1}{\lambda}$$

$$= \frac{1}{\lambda} \left[I - A^T (\lambda I + A A^T)^{-1} A \right]$$

$$\Rightarrow (A^T A + \lambda I)^{-1} A^T \underline{y} = \frac{1}{\lambda} \left[A^T - A^T (A A^T + \lambda I)^{-1} A A^T \right] \underline{y}$$

$$= \frac{1}{\lambda} \left[A^T - A^T (K + \lambda I)^{-1} K \right] \underline{y}$$

where $K = [\langle x_i, x_j \rangle]_{i,j=1}^n$

What about the remaining A^T ? It is handled when we evaluate the estimate:

$$\begin{aligned}\hat{\beta}^T x &= \frac{1}{\lambda} \underline{y}^T [A - K(K + \lambda I)^{-1} A] x \\ &= \frac{1}{\lambda} \underline{y}^T [I - K(K + \lambda I)^{-1}] \underline{k}(x)\end{aligned}$$

where

$$\underline{k}(x) = \begin{bmatrix} \langle x_1, x \rangle \\ \vdots \\ \langle x_n, x \rangle \end{bmatrix}$$

we can simplify further:

$$I - K(K + \lambda I)^{-1} =$$

(A)

$$\Rightarrow \hat{\beta}^T x = \underline{y}^T (K + \lambda I)^{-1} \underline{k}(x)$$

What about $\hat{\beta}_0 = \bar{y} - \hat{\beta}^T \bar{x}$?

$$\hat{\beta}_0 = \bar{y} - \frac{1}{n} \sum_{i=1}^n \hat{\beta}^T x_i$$

$$= \bar{y} - \frac{1}{n} \sum_{i=1}^n \underline{y}^T (K + \lambda I)^{-1} \underline{k}(x_i)$$

$$= \bar{y} - \underline{y}^T (K + \lambda I)^{-1} \cdot \frac{1}{n} \sum_{i=1}^n \underline{k}(x_i)$$

Note | For some kernels, $k(x, x') = \langle \Phi(x), \Phi(x') \rangle$,

$\Phi(x)$ already contains a constant component,

in which case β_0 is not needed. Examples

include the inhomogeneous polynomial + Gaussian kernels.

Example 1 Gaussian kernel

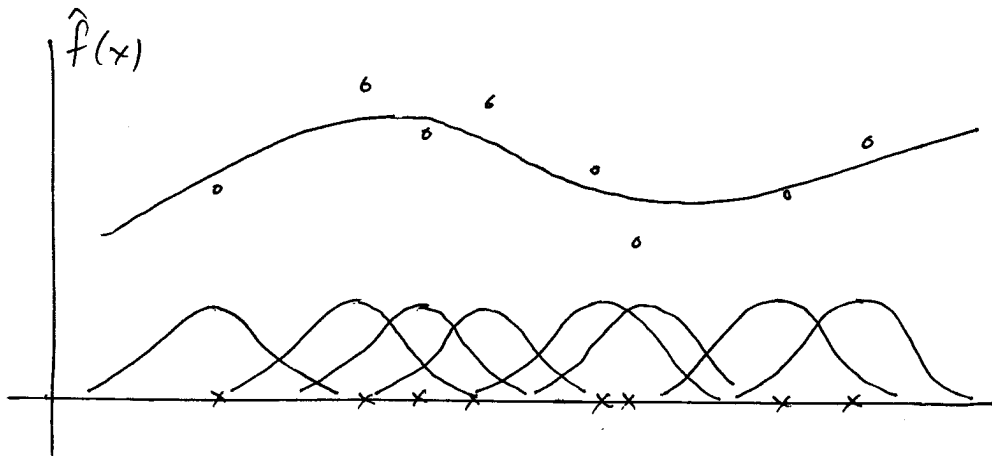
$$k_{\sigma}(x, x_i) = \exp\left\{-\frac{\|x-x_i\|^2}{2\sigma^2}\right\}$$

$$\text{Then } \hat{\beta}^T x = \underline{y}^T (K + \lambda I)^{-1} \underline{k}(x)$$

$$= \underline{\alpha}^T \underline{k}(x)$$

$$= \sum \alpha_i k_{\sigma}(x, x_i)$$

α independent
of x



Key

$$A. \quad I - K(K + \lambda I)^{-1}$$

$$= (K + \lambda I - K)(K + \lambda I)^{-1}$$

$$= \lambda \cdot (K + \lambda I)^{-1}$$