

KERNEL RIDGE

REGRESSION

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

$$(\hat{\beta}, \hat{\beta}_0) \leftarrow \min_{\beta, \beta_0} \sum_{i=1}^n (y_i - \beta^T x_i - \beta_0)^2 + \lambda \|\beta\|^2$$

Solution:

$$\left. \begin{aligned} \hat{\beta} &= (A^T A + \lambda I)^{-1} A^T \tilde{y} \\ \hat{\beta}_0 &= \bar{y} - \hat{\beta}^T \bar{x} \end{aligned} \right\} \Rightarrow \hat{f}(x) = \hat{\beta}^T x + \hat{\beta}_0$$

where

$$A = \begin{bmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_n^T \end{bmatrix}, \quad \begin{aligned} \tilde{x}_i &= x_i - \bar{x} \\ \tilde{y}_i &= y_i - \bar{y} \end{aligned}$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

Can we express RR in terms of inner products

$\langle x_i, x_j \rangle$ and $\langle x_i, x \rangle$?

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

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

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

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

where $K = [\langle \tilde{x}_i, \tilde{x}_j \rangle]_{i,j=1}^n$. Note $\langle \tilde{x}_i, \tilde{x}_j \rangle$

$$= \langle x_i, x_j \rangle - \frac{1}{n} \sum_{r=1}^n \langle x_i, x_r \rangle - \frac{1}{n} \sum_{s=1}^n \langle x_s, x_j \rangle + \frac{1}{n^2} \sum_{r,s=1}^n \langle x_r, x_s \rangle.$$

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

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

where

$$\underline{k}(x) = \begin{bmatrix} \langle x_1, x \rangle \\ \vdots \\ \langle x_n, x \rangle \end{bmatrix} = \begin{bmatrix} \langle x_1, x \rangle - \frac{1}{n} \sum \langle x_r, x \rangle \\ \vdots \\ \langle x_n, x \rangle - \frac{1}{n} \sum \langle x_r, x \rangle \end{bmatrix}$$

We can simplify further:

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

$$\Rightarrow \hat{\beta}^T x =$$

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

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

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

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

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

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 kernels.

Example 1 Gaussian kernel

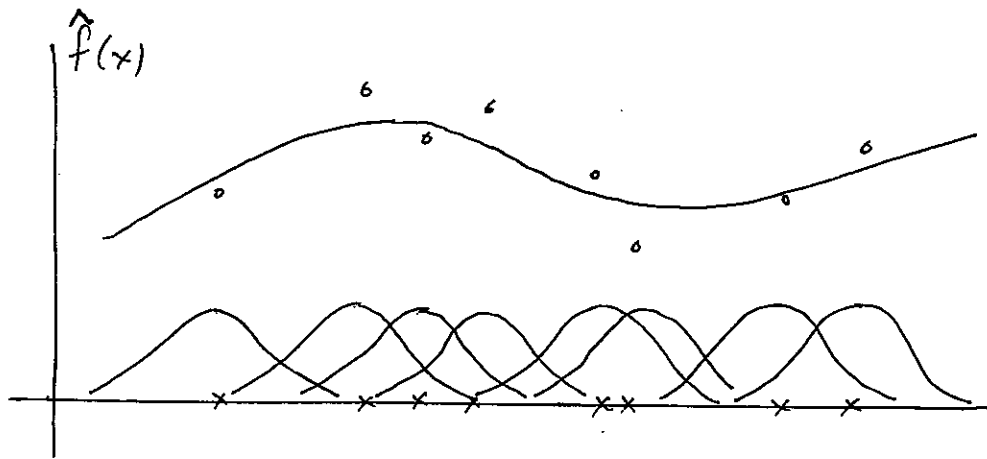
$$k(x, x') = \exp\left\{-\frac{\|x-x'\|^2}{2\sigma^2}\right\}$$

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

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

$$= \sum \alpha_i k(x, x_i)$$

$\underline{\alpha}$ independent
of x



Key

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

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

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

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