

KERNEL RIDGE

REGRESSION

In these notes we will show how to incorporate kernels into ridge regression yielding a non linear regression method.

Ridge Regression

Ridge regression solver

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

The solution is

$$\hat{w} = (A^T A + n\lambda I)^{-1} A^T \hat{y}$$

$$\hat{b} = \bar{y} - \hat{w} \bar{x}$$

where

$$\hat{y} = \begin{bmatrix} \tilde{y}_1 \\ \vdots \\ \tilde{y}_n \end{bmatrix}, \quad A = \begin{bmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_n^T \end{bmatrix}, \quad \begin{aligned} \tilde{x}_i &= x_i - \bar{x} \in \mathbb{R}^d \\ \tilde{y}_i &= y_i - \bar{y} \in \mathbb{R} \\ \bar{x} &= \frac{1}{n} \sum x_i \\ - & 1 \leq i \leq n \end{aligned}$$

$$\begin{bmatrix} \tilde{y}_n \\ \vdots \\ \tilde{y}_1 \end{bmatrix} \quad \begin{bmatrix} \tilde{x}_n^T \\ \vdots \\ \tilde{x}_1^T \end{bmatrix} \quad \hat{\bar{y}} = \frac{1}{n} \sum_{i=1}^n \hat{y}_i$$

The regression function estimate is

$$\hat{f}(x) = \hat{w}^T x + \hat{b} = \bar{y} + \hat{w}(x - \bar{x})$$

To apply kernels, we must show that $\hat{f}(x)$ depends on x, x_1, \dots, x_n only in terms of inner products, such as $\langle x_i, x_j \rangle$ or $\langle x_i, x \rangle$.

Kernelizing Ridge Regression

Note that $A^T A$ is not a Gram matrix. Let's apply the matrix inversion lemma

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

with

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

$\nwarrow \mu = n\lambda$ for brevity

This yields

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

Therefore

$$\begin{aligned}
 (A^T A + \mu I)^{-1} A^T \tilde{y} &= \frac{1}{\mu} [A^T - A^T (A A^T + \mu I)^{-1} A A^T] \tilde{y} \\
 &= \frac{1}{\mu} [A^T - A^T (\bar{R} + \mu I)^{-1} \bar{R}] \tilde{y}
 \end{aligned}$$

where

$$\bar{R} = [\langle \tilde{x}_i, \tilde{x}_j \rangle]_{ij=1}^n. \quad (n \times n)$$

Note

$$\langle \tilde{x}_i, \tilde{x}_j \rangle = \langle x_i - \bar{x}, x_j - \bar{x} \rangle$$

$$= \langle x_i, x_j \rangle - \frac{1}{n} \sum_{r=1}^n \langle x_i, x_r \rangle$$

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

In addition,

$$\begin{aligned}
 \hat{f}(x) &= \bar{y} + \hat{w}^\top (x - \bar{x}) \\
 &= \bar{y} + \frac{1}{\mu} \tilde{y}^\top [A - \tilde{K}(\tilde{K} + \mu I)^{-1} A] (x - \bar{x}) \\
 &= \bar{y} + \frac{1}{\mu} \tilde{y}^\top [\mathbb{I} - \tilde{K}(\tilde{K} + \mu I)^{-1}] \tilde{k}(x)
 \end{aligned}$$

where

$$\tilde{k}(x) = \begin{bmatrix} \langle \tilde{x}_1, x - \bar{x} \rangle \\ \vdots \\ \langle \tilde{x}_n, x - \bar{x} \rangle \end{bmatrix}.$$

The entries of this vector are

$$\begin{aligned}
 \langle \tilde{x}_i, x - \bar{x} \rangle &= \langle x_i - \bar{x}, x - \bar{x} \rangle \\
 &= \langle x_i, x \rangle - \frac{1}{n} \sum_r \langle x_i, x_r \rangle \\
 &\quad - \frac{1}{n} \sum_s \langle x_s, x \rangle + \frac{1}{n^2} \sum_r \sum_s \langle x_r, x_s \rangle.
 \end{aligned}$$

Finally, observe

$$\mathbb{I} - \tilde{K}(\tilde{K} + \mu I)^{-1}$$

$$= [\tilde{K} + \mu I - \tilde{K}] (\tilde{K} + \mu I)^{-1}$$

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

which yields the estimate

$$\hat{f}(x) = \bar{y} + \tilde{\gamma}^\top (\tilde{K} + \mu I)^{-1} \tilde{k}(x).$$

The point of these manipulations is that \hat{f} depends on elements of the feature space only in terms of inner products. If k is an inner product kernel with feature map Φ , and we substitute

$$\langle x, x' \rangle \rightarrow k(x, x')$$

whenever inner products occur, then \hat{f} corresponds to ridge regression after the original feature space is transformed via Φ .

No Offset

For some kernels $k(x, x') = \langle \Phi(x), \Phi(x') \rangle$, $\Phi(x)$ already contains a constant component, in which case the offset/bias b is not necessary. The

inhomogeneous polynomial kernels are an example.

The Gaussian kernel also does not require an offset, even though its feature space does not contain a constant term (Steinwart and Christmann, 2008).

If b is omitted, the KRR estimator is

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

where

$$K = [k(x_i, x_j)]_{i,j=1}^n \quad (n \times n)$$

and

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

Example] Gaussian kernel

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

Then

$$\begin{aligned}\hat{f}(x) &= y^T (K + n\lambda I)^{-1} k(x) \\ &= \alpha^T k(x) \\ &= \sum_{i=1}^n \alpha_i k(x, x_i)\end{aligned}$$

where

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

is independent of x .

