

Lecture 9



• Part 3: kernel regression, Gaussian Process, Support Vector Machine

• Recall Ridge Regression

	$j \dots p$	
i	x_{ij}	y_i
\vdots		
n		
	X $n \times p$	Y $n \times 1$

$$L = \text{Loss}(\beta) = \|Y - X\beta\|^2 + \lambda \|\beta\|^2$$

$$= (Y - X\beta)^T (Y - X\beta) + \lambda \|\beta\|^2$$

$$= Y^T Y - \beta^T X^T Y - Y^T X \beta + \beta^T X^T X \beta + \lambda \|\beta\|^2$$

$$= \text{const} - 2 \langle X^T Y, \beta \rangle + \beta^T X^T X \beta + \lambda \beta^T \beta$$

$$\frac{\partial L}{\partial \beta} = -2 X^T Y - 2 X^T X \beta + 2 \lambda \beta = 0$$

$$(X^T X + \lambda I_p) \beta = X^T Y$$

$$\hat{\beta}_{\text{ridge}} = (X^T X + \lambda I_p)^{-1} X^T Y$$

$$\text{Identity: } (X^T X + \lambda I_p)^{-1} X^T = X^T (X X^T + \lambda I_n)^{-1}$$

$$\text{Proof: } X^T (X X^T + \lambda I_n) = (X^T X + \lambda I_p) X^T$$

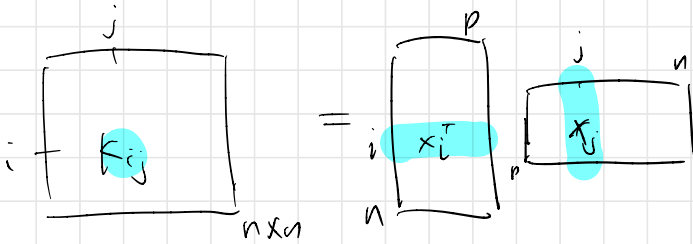
$$\downarrow$$

$$X^T X X^T + \lambda X^T = X^T X X^T + \lambda X^T \quad \square$$

$$\text{Thus } \hat{\beta} = X^T (X X^T + \lambda I_n)^{-1} Y$$

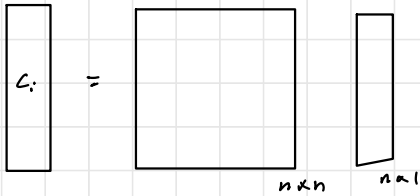
$$\hat{\beta} = X^T (X X^T + \lambda I_n)^{-1} Y$$

Let $K = X X^T$ with $k_{ij} = \langle x_i, x_j \rangle$



$$c = (K + \lambda I_n)^{-1} Y$$

$n \times 1$ $n \times n$ $n \times 1$



So $\hat{\beta} = X^T C$

$$\hat{\beta} = \begin{matrix} p \\ \hline \end{matrix} \begin{matrix} n \\ \hline \end{matrix} \begin{matrix} x_1 & \dots & x_i & \dots & x_n \end{matrix} \begin{matrix} c_1 \\ \vdots \\ c_i \\ \vdots \\ c_n \end{matrix} = \sum_{i=1}^n \begin{matrix} c_i \\ \hline \end{matrix} \begin{matrix} p \times 1 \\ \hline \end{matrix} x_i$$

Representer Form

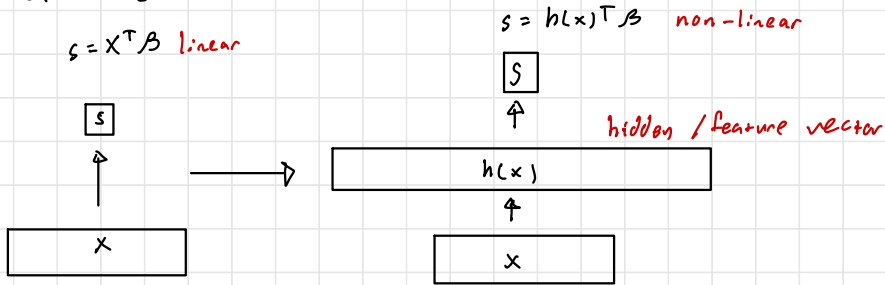
Now consider

	$j \dots p$	
i	x_{ij}	y_i
\vdots		
n		
0	x_0^T	$?$

$$\begin{aligned}
 s &= f(x_0) = x_0^T \hat{\beta} \\
 &= x_0^T \sum_{i=1}^n c_i x_i \\
 &= \sum_{i=1}^n c_i \langle x_i, x_0 \rangle = \sum_{i=1}^n c_i k(x_i, x_0)
 \end{aligned}$$

$k(x, x') = \langle x, x' \rangle$ ← Kernel Function

• Kernel Trick



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

$$k(x, x') = \langle h(x), h(x') \rangle$$

Measures similarity

Implicit

e.g. $\exp(-\gamma |x - x'|^2)$

Gaussian Kernel / Radial Basis Function

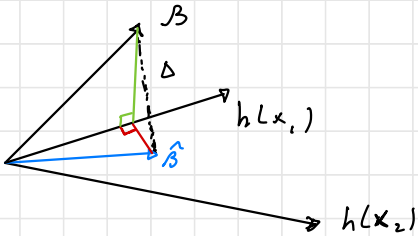
• More systematic :

$$\text{Loss} = \sum_{i=1}^n (y_i - h(x_i)^T \beta)^2 + \lambda |\beta|^2$$

$$\downarrow$$

$$s_i = f(x_i) = h(x_i)^T \beta = \langle h(x_i), \beta \rangle$$

• Representer Theorem



$$h(x_i)^T \beta = h(x_i)^T \hat{\beta}$$

$$|\beta|^2 \geq |\hat{\beta}|^2$$

$$\text{So } \hat{\beta} = \sum_{i=1}^n c_i h(x_i)$$

$$\text{Loss} = \sum_{i=1}^n (y_i - h(x_i)^T \underbrace{\sum_{j=1}^n c_j h(x_j)}_{\hat{\beta}})^2 + \lambda \left\langle \underbrace{\sum_{i=1}^n c_i h(x_i)}_{\hat{\beta}}, \underbrace{\sum_{j=1}^n c_j h(x_j)}_{\hat{\beta}} \right\rangle$$

$$= \sum_{i=1}^n \left(y_i - \sum_{j=1}^n c_j k_{ij} \right)^2 + \lambda \sum_{i,j} c_i c_j k_{ij}$$

$$\text{Loss}(c) = \|y - Kc\|^2 + \lambda c^T Kc$$

$$= (y - Kc)^T (y - Kc) + \lambda c^T Kc$$

$$= y^T y - c^T K y - y^T K c - c^T K^2 c + \lambda c^T Kc$$

$$\frac{\partial L}{\partial c} = -2Ky + 2K^2c + 2\lambda Kc = 0$$

$$K^{-1}(-2Ky + 2K^2c + 2\lambda Kc) = 0$$

$$c = (K + \lambda I)^{-1} y \quad \text{Estimation}$$

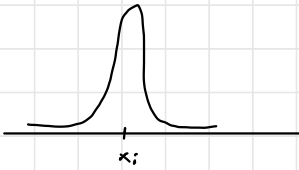
Thus $\hat{\beta} = \sum_{i=1}^n c_i h(x_i)$

Consider $f(x_0) = h(x_0)^T \hat{\beta}_n$
 $= h(x_0)^T \sum_{i=1}^n c_i h(x_i)$
 $= \sum_{i=1}^n c_i k(x_i, x_0)$

Prediction

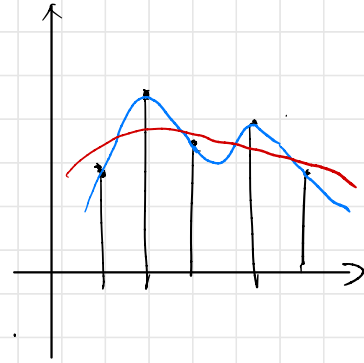
• $f(x) = \sum_{i=1}^n c_i k(x_i, x)$

Suppose $k(x_i, x) = \exp(-\gamma |x - x_i|^2)$



let $\lambda \rightarrow 0, \gamma \rightarrow \infty$

$$k(x_i, x_j) = \begin{cases} 1 & \text{if } x_i = x_j \\ \rightarrow 0 & \text{if } x_i \neq x_j \end{cases}$$



$K \rightarrow I, C \rightarrow Y$

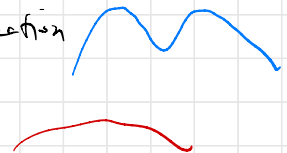
$$f(x) = \sum_{i=1}^n y_i k(x_i, x) = \begin{cases} y_i & \text{if } x = x_i \\ 0 & \text{otherwise} \end{cases} \quad \text{memorization}$$

for $\gamma < \infty, \lambda \rightarrow 0$

$C \rightarrow K^T Y$ so $KC = Y$

$$f(x) = \begin{cases} y_i & \text{if } x = x_i, \text{ memorization} \\ \text{Smooth interpolation} & \end{cases}$$

for $\gamma < \infty, \lambda > 0$, smooth fitting



Theoretical underpinning

- Reproducing Kernel Hilbert Space: (RKHS)

$$\mathcal{F} = \{ f(x) = h(x)^T \beta, \quad |\beta|_{\ell_2}^2 < \infty \}$$

Suppose $f(x) = h(x)^T \beta$
 $g(x) = h(x)^T \alpha$

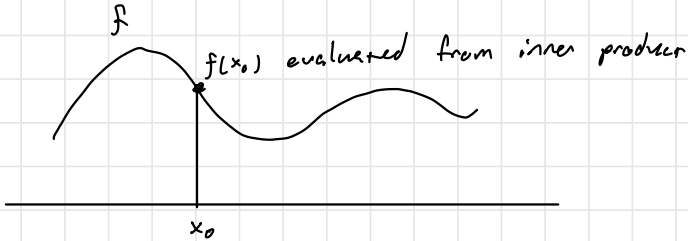
$$\langle f, g \rangle_{\mathcal{F}} = \langle \beta, \alpha \rangle_{\ell_2}$$

$$\|f\|_{\mathcal{F}}^2 = \langle f, f \rangle_{\mathcal{F}} = \|\beta\|_{\ell_2}^2$$

- Reproducing Property:

$$\langle f(x), k(x_0, x) \rangle_{\mathcal{F}} = h(x_0)^T \beta = f(x_0)$$

$$\begin{array}{ccc} \downarrow & & \downarrow \\ h(x)^T \beta & & h^T(x) \underbrace{h(x_0)}_{\alpha} \end{array}$$



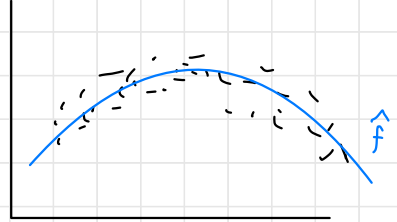
In regular space we need $\langle f, \delta_{x_0} \rangle$

Sweep $h(x)$, β under the rug

non-parametric regression:

$$\text{Loss}(f) = \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \|f\|_{\mathcal{F}}^2$$

\downarrow \downarrow
 $l(y_i, f(x_i))$ \downarrow \downarrow
 smoothness



• Representer:

$$\hat{f}(x) = \sum_{i=1}^n c_i \underbrace{K(x_i, x)}_{f(x)} + \Delta(x)$$

Proof:

$$\langle \Delta, K(x_i, \cdot) \rangle_{\mathcal{F}} = 0 \quad \forall i$$

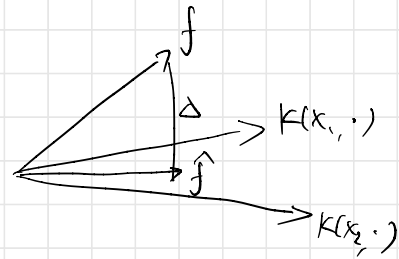
$\Delta(x_i)$

$$\|f\|_{\mathcal{F}}^2 \geq \|\hat{f}\|_{\mathcal{F}}^2$$

$$f(x_j) = \sum_{i=1}^n c_i K(x_i, x_j) + \Delta(x_j)$$

$$\| \langle \Delta(x), K(x_j, x) \rangle_{\mathcal{F}} = 0$$

$$= \hat{f}(x_j)$$



$$\text{Loss}(f) = \text{Loss}(c) = \|y - Kc\|^2 + \lambda c^T K c$$

• Mercer Theorem: Condition for K , so that $K(x, x') = \langle h(x), h(x') \rangle$

• Recall $K_{n \times n} = Q \Lambda Q^T$
symmetric

$K \geq 0$ iff $\lambda_i \geq 0, i=1, \dots, n$
equivalently

$$a^T K a \geq 0 \quad \forall a \neq 0$$

$$\downarrow$$

$$b^T \Lambda b$$

• kernel $k(x, x')$



Mercer Decomposition

$$k(x, x') = \sum_k \lambda_k q_k(x) q_k(x') = \sum_k h_k(x) h_k(x') = \langle h(x), h(x') \rangle$$

$$K = \sum_k \lambda_k q_k q_k^T$$

$$h_k(x) = \sqrt{\lambda_k} q_k(x)$$

$$K \geq 0, \lambda_k \geq 0 \forall k$$

