

today: continue R&HS in all their glory! ☺

$$\langle k(x, \cdot), k(x', \cdot) \rangle_H = k(x, x')$$

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

$$\langle f, k(x, \cdot) \rangle_p \approx f(x)$$

" $\oplus(x)$ "

analog: $\langle \omega, \mathbb{D}(z) \rangle$

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i)) + \alpha \|f\|_H^2$$

is reached for $f^* = \sum_{i=1}^n \alpha_i^* k(x_i, \cdot)$

$(x_i, y_i)_{i=1}^n$
training set

$$\text{Let } f_{\alpha} = \sum_{i=1}^n \alpha_i k(x_i, \cdot) \quad \alpha \in \mathbb{R}^n$$

$$\text{then } \|f_\alpha\|_H^2 = \langle f_\alpha, f_\alpha \rangle_H = \sum_{i,j} \alpha_i \alpha_j k(x_i, x_j) = \alpha^T K \alpha$$

Gram matrix: $(k)_{ij} \triangleq k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle_2$

\hookrightarrow inner product of $\Phi(x_i)$ on data

$$\min_{\alpha \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n \underbrace{\left[y_i (\underbrace{\alpha_j k(x_j, x_i)}_{f_k(x_i)}) + \lambda \alpha^T k \alpha \right]}_{f_k(x_i)}$$

forite dirn. opt.
(thanks to representat's
them)

$$f_x(x) = \sum_j q_j \underbrace{k(x_j, x)}_{\text{重}(x_j), \text{重}(x) > 0}$$

Getting a handle on \mathbb{C}^n : generalize diagonalization of matrices to n -dim

I) start with finite matrices

say X is finite e.g. x_1, \dots, x_n

$$|x|=n$$

$f: X \rightarrow \mathbb{R}$
is fine \Rightarrow just a vector

$$\begin{pmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix} \in \mathbb{R}^n$$

$$(k\alpha)_i = \delta_{ij} (x_j)$$

$$= \sum_j \alpha_j k(x_i, x_j) = (k\alpha)_i$$

(\mathcal{L}_2 visto)

\hookrightarrow vector $(f(x_i))$

$\mathcal{H} = \text{span} \{ k(x_i, \cdot) : i=1, \dots, n \} = \{ k\alpha : \alpha \in \mathbb{R}^n \} \subseteq \mathbb{R}^n$ [“ \mathcal{H} -view”]

Let K be Gram matrix $(K)_{ij} \triangleq k(x_i, x_j)$

If k is a valid kernel
(i.e. defines a inner product)

we can let $\tilde{\Phi} \triangleq \Lambda^{1/2} U^T$

$$\tilde{\Phi} = \begin{pmatrix} -\sqrt{\lambda_1} \psi_1^T & \cdots \\ \vdots & \ddots \\ -\sqrt{\lambda_d} \psi_d^T & \cdots \end{pmatrix}$$

$d = \text{rank}(K) \leq n$

$\tilde{\Phi} : X \rightarrow \mathbb{R}^d$

$$\tilde{\Phi} \triangleq \begin{pmatrix} \tilde{\Phi}(x_1) & \cdots & \tilde{\Phi}(x_n) \end{pmatrix}$$

$\Rightarrow \langle \tilde{\Phi}(x_i), \tilde{\Phi}(x_j) \rangle_{\mathbb{R}^d}$

$= k(x_i, x_j)$

\Rightarrow

note: $K \psi_i = \lambda_i \psi_i$

$$\boxed{\tilde{\Phi}(x) = \begin{pmatrix} \sqrt{\lambda_1} \psi_1(x) \\ \vdots \\ \sqrt{\lambda_d} \psi_d(x) \end{pmatrix} \in \mathbb{R}^d}$$

x-coord. of ψ_i vector

“feature space” pt. of view

$$k(x_i, x_j) = \langle \tilde{\Phi}(x_i), \tilde{\Phi}(x_j) \rangle_{\mathbb{R}^d}$$

$$= \sum_{l=1}^d \lambda_l \psi_l(x_i) \psi_l(x_j)$$

back to \mathcal{H} -view: $\mathcal{H} \subseteq \mathcal{F}_2$

$v \in \mathcal{H} \Rightarrow v = k\alpha$ for some $\alpha \in \mathbb{R}^n$

to get $\|v\|_{\mathcal{H}}$, we compute $\alpha_v = K^+ v$ pseudo-inverse

$$\Rightarrow \|v\|_{\mathcal{H}}^2 = \alpha_v^T K \alpha_v = v^T K^+ K \alpha_v$$

$$= v^T U \underbrace{L^T L}_{\in \mathbb{R}^{d \times d}} \underbrace{U^T}_{\in \mathbb{R}^{n \times n}} \underbrace{(I_d \otimes I_d)}_{\in \mathbb{R}^{n \times n}} U^T$$

$$= v^T \underbrace{L^T L^T}_{\in \mathbb{R}^{d \times d}} \underbrace{U^T}_{\in \mathbb{R}^{n \times n}} U^T$$

$$= v^T \underbrace{(I_d \otimes I_d)}_{\in \mathbb{R}^{n \times n}} U^T$$

$$\boxed{\|v\|_{\mathcal{H}}^2 = \sum_{j=1}^d \frac{\langle v, \psi_j \rangle_{\mathcal{F}_2}^2}{\lambda_j}}$$

representation

$$\text{more generally, } \|v\|_{\mathcal{H}}^2 = \sum_{j=1}^d \underbrace{\beta_j}_{\text{if any } \langle v, \psi_j \rangle \neq 0} \underbrace{\gamma_j}_{\text{for } j > d}$$

i.e. $v = \sum_{j=1}^n \beta_j \psi_j$ i.e. $\beta_j \triangleq \langle v, \psi_j \rangle_{\mathcal{F}_2}$

\hookrightarrow $\|v\|_{\mathcal{H}}^2 = \sum_{j=1}^n \beta_j^2$

$$||v||_{\ell_2^d}^2 = \sum_{j=1}^d \langle v, e_j \rangle_{\ell_2}^2$$

$$\text{Vs. } ||v||_{\ell_2^d}^2 = \sum_{j=1}^d \langle v, e_j \rangle_{\ell_2}^2$$

$$\text{and } \boxed{\frac{\|v_j\|_{\ell_2}^2}{\|v_j\|_{\ell_2}^2} = \frac{1}{\lambda_j}}$$

for $j \leq d$
 $\|v_j\|_{\ell_2}^2 = +\infty$
for $j > d$

so orthonormal basis of \mathcal{H} in ℓ_2^d

$$\{ \sqrt{\lambda_i} e_i \}_{i=1}^d$$

$$\langle v, v \rangle_{\ell_2^d} = \sum_{i=1}^n \beta_i^2$$

$$\langle v, v \rangle_{\mathcal{H}} = \sum_{i=1}^d \frac{\beta_i^2}{\lambda_i}$$

$$\text{thus } ||v||_{\mathcal{H}} \leq 1$$

↳ makes an ellipsoid in ℓ_2^d

i.e. higher coordinates one shrunk more?
(since λ_i 's is smaller as $i \rightarrow$)

10h47

II) generalization to $\text{to-dim } \mathcal{H}$

Suppose X is a compact normed space (e.g. $X = [0, 1]$)

+ Lebesgue measure on it

$$\mathcal{L}_2(X) \triangleq \{ f: X \rightarrow \mathbb{R} \mid \int_X (f(x))^2 dx < \infty \}$$

$$\mathcal{Q}_2 \triangleq \{ (\alpha_i)_{i=1}^\infty \text{ s.t. } \sum_i \alpha_i^2 < \infty \}$$

Let k be a cts. psd kernel fct. (note: it is symmetric)

↳ with respect to standard norm on $X \not\cong \mathbb{R}$

define

$$L_k : \mathcal{L}_2 \rightarrow \mathcal{L}_2 \quad k(\cdot, x)$$

$$\text{s.t. } [L_k f](\cdot) \triangleq \int_X k(x, \cdot) f(x) dx$$

$$Kv = \sum_i K_{i,i} v_i$$

Then can show that

L_k is a "compact self-adjoint positive" operator

and yields an (countable) orthonormal basis (for \mathcal{L}_2)
[Hilbert basis]

of G-functions for L_k $\{ \psi_i \}_{i=1}^\infty$

with non-negative eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq 0$

$$\text{i.e. } L_k \psi_i = \lambda_i \psi_i$$

and we have $\boxed{k(x, z) = \sum_i \lambda_i \psi_i(x) \psi_i(z)}$

$$\langle f, L_k g \rangle_{\mathcal{L}_2}$$

$$= \langle L_k f, g \rangle_{\mathcal{L}_2} \quad \text{if } f, g$$

↓

$$\text{finite version: } \langle v, L_k w \rangle$$

$$\begin{aligned} &= \sqrt{\lambda} k_w \\ &= \sqrt{\lambda} k^T w \\ &= \langle k v, w \rangle \end{aligned}$$

use $k = k^T$

$$U \perp \cup U^T$$

end we have
$$k(x_1, x_2) = \sum_{i=1}^n \lambda_i \psi_i(x) \psi_i(x)$$

Mercer's Thm

" $\cup \cap \cup$ "
 like $k = \sum \lambda_i \psi_i$
 of before

a) Feature space $H \subseteq \ell_2$

view of H

$$\Phi: X \rightarrow \ell_2 \text{ with } (\Phi(x))_i \triangleq \sqrt{\lambda_i} \psi_i(x)$$

$$\text{here: } k(x_1, x_2) = \langle \Phi(x_1), \Phi(x_2) \rangle_{\ell_2}$$

"diagonalized representation"

Here identify $k(x_1, x_2) \in \ell_2$ as $\Phi(x) \in \ell_2$

$$\begin{aligned} & \text{valid element of } \ell_2 \\ & \because \sum_i (\Phi(x)_i)^2 = \sum_i \lambda_i \psi_i^2(x) \\ & = k(x, x) < \infty \end{aligned}$$

(do not know what $H \subseteq \ell_2$ looks like though) $\beta_i \downarrow \|f\|_H^2 < \infty$

b) ℓ_2 view

$$H \subseteq \ell_2 : H = \left\{ f \in \ell_2 : \sum_{i=1}^n \underbrace{\langle f, \psi_i \rangle}_{\beta_i} \psi_i < \infty \right\}$$

$$\langle f, g \rangle_H \triangleq \sum_{i=1}^n \underbrace{\langle f, \psi_i \rangle}_{\beta_i} \underbrace{\langle g, \psi_i \rangle}_{\beta_i} \psi_i$$

→ ellipsoid in ℓ_2

⊗ if K is "universal"

$\Rightarrow H_K$ is dense in ℓ_2 i.e. for any $f \in \ell_2$

\exists a sequence $h_n \in H_K$ s.t. $\|h_n - f\|_{\ell_2} \xrightarrow{n \rightarrow \infty} 0$

note: if $f \notin H \Rightarrow \|h_n\|_H \rightarrow \infty$

* non-parametric learning:

$$\hat{f}_n = \underset{f \in H}{\operatorname{argmin}} \underbrace{\frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i, f(x_i))}_{n \rightarrow \infty \mathbb{E} \mathcal{L}(f)} + \lambda_n \|f\|_H^2$$

$$f^* \triangleq \underset{f \in \ell_2}{\operatorname{argmin}} \mathbb{E} \mathcal{L}(Y, f(X)) \quad \text{perhaps } f^* \notin H$$

but H dense in ℓ_2

→ regularity property \mathcal{L} + decrease λ_n at correct rate

⇒ consistency of \hat{f}_n i.e. $\hat{f}_n \xrightarrow[n \rightarrow \infty]{\ell_2} f^*$

$$\|\hat{f}_n - f^*\|_{\ell_2}$$

e.g. SVM with RBF kernel is "universally consistent"
when $\lambda_n \rightarrow 0$ at correct rate