

today: continue RKHS in all their glory? 

$$\langle k(x, \cdot), k(x', \cdot) \rangle_{\mathcal{H}} = k(x, x')$$

$$\underbrace{k(x)}_{\mathcal{D}(x)} \quad \underbrace{k(x')}_{\mathcal{D}(x')} \quad \langle \mathcal{D}(x), \mathcal{D}(x') \rangle = k(x, x')$$

representer's thm: says that (for \mathcal{H} a RKHS)

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

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

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

training set

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

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

Gram matrix: $(K)_{ij} \triangleq k(x_i, x_j)$ $n \times n$ matrix

↳ inner product of $\mathcal{D}(x_i)$ on data $k_{ij} = \langle \mathcal{D}(x_i), \mathcal{D}(x_j) \rangle$

finite dim.
cpd.
(thanks to
representer's thm)

$$\min_{\alpha \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n \ell(y_i, \underbrace{\sum_{j=1}^n \alpha_j k(x_j, x_i)}_{f_\alpha(x_i)}) + \lambda \alpha^T K \alpha$$

getting a handle on \mathcal{H} : generalize diagonalization of matrices to no-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 finite \Rightarrow just a vector

$$\begin{pmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix} \in \mathbb{R}^n \quad (\text{"f-a view"})$$

$$\mathcal{H} = \text{span} \{ k(x_i, \cdot) : i=1, \dots, n \} = \{ K\alpha : \alpha \in \mathbb{R}^n \} \subseteq \mathbb{R}^n$$

" f -view"

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

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

if K is a valid $\Rightarrow K \geq 0$

(i.e. kernel
i.e. eigenvectors
a inner product)

spectral thm:

we can let $\Phi = L^{1/2} U^T$

$$\Rightarrow K = \Phi^T \Phi$$

$$\Phi = \begin{pmatrix} \sqrt{\lambda_1} \psi_1^T & \\ \vdots & \\ \sqrt{\lambda_d} \psi_d^T & \\ 0 & \\ 0 & \end{pmatrix}$$

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

$$\Phi = \begin{pmatrix} \Phi(x_1) & \cdots & \Phi(x_n) \\ \vdots & & \vdots \\ 0 & & 0 \end{pmatrix}$$

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

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

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

$\text{diag}(\lambda_1, \dots, \lambda_n)$

$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$

and U is orthonormal basis of \mathbb{R}^n

$$\text{i.e. } U = \begin{pmatrix} | \psi_1 & \cdots & | \psi_n \end{pmatrix}$$

$$U^T U = U U^T = I_n \quad \langle \psi_i, \psi_j \rangle_{\mathbb{R}^n} = \delta_{ij}$$

\Rightarrow if we define, *x-coord of ψ_i vector*

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

"feature space" pt. of view

$$\text{note: } K \psi_i = \lambda_i \psi_i$$

back to \mathcal{L}_2 -view: $H \subseteq \mathcal{L}_2$
 $(\subseteq \mathbb{R}^n)$

$v \in H \Rightarrow v = k\alpha$ for some $\alpha \in \mathbb{R}^n$ *Pseudo-inverse*
 to get $\|v\|_H$, we compute $\alpha_v = K^+ v$

$$\begin{aligned} \text{so } \|v\|_H^2 &= \alpha_v^T K \alpha_v = v^T K^+ K K^+ v & K &= U \Lambda U^T \\ &= v^T U \Lambda^T U^T \underbrace{(I_d \ 0)}_{I_n} U^T v & K^+ &= U \Lambda^+ U^T \end{aligned}$$

$$= (U v)^T \Lambda^+ (U v)$$

$$\begin{aligned} K K^+ &= U \Lambda U^T \Lambda^+ U^T \\ &= U \left(\begin{matrix} I_d & 0 \\ 0 & 0 \end{matrix} \right) U^T \end{aligned}$$

$$\boxed{\|v\|_H^2 = \sum_{j=1}^d \frac{\langle v, \psi_j \rangle_{\mathbb{R}^n}^2}{\lambda_j}}$$

B_j: representation

$$\text{vs. } \|v\|_{\mathbb{R}^n}^2 = \sum_{j=1}^n \langle v, \psi_j \rangle_{\mathbb{R}^n}^2$$

$\frac{n}{n \cdot \|v\|^2 - 1}$

$\Rightarrow U v$ is projection of v onto $\{\psi_1, \dots, \psi_n\}$ basis
 i.e. $v = \sum_{j=1}^n \beta_j \psi_j$ i.e. $\beta_j = \langle v, \psi_j \rangle_{\mathbb{R}^n}$

$$\dots \dots \dots \frac{n}{n \cdot \|v\|^2 - 1}$$

" \rightarrow " \leftarrow

and $\boxed{\|v_j\|_H^2 = \frac{1}{\lambda_j}}$

for $j \leq d$

$\|v_j\|_H = +\infty$ "
for $j > d$

so orthonormal basis of H in \mathbb{R}^n
 $\{v_i\}_{i=1}^n$

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

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

thus $\|v\|_H \leq 1$

↳ makes an ellipsoid in \mathbb{R}^n

i.e. higher coordinates are shrunk more
 (since λ_i is smaller as $i \rightarrow \infty$)

14h45

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

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

+ Lebesgue measure on it

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

$$\mathcal{D}_2 \triangleq \{(\alpha_i)_i \text{ s.t. } \sum_i \alpha_i^2 < \infty\}$$

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

↳ with respect to standard norm on $X \subseteq \mathbb{R}$

$$K_V = \sum_i K_{ii} V_i$$

define

$$L_K: \mathcal{L}_2 \rightarrow \mathcal{L}_2$$

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

then can show that

L_K is a "compact self-adjoint positive" operator

$$\langle f, L_K g \rangle_{\mathbb{R}^n}$$

$$\begin{aligned} &= \langle L_K f, g \rangle_{\mathbb{R}^n} \\ &\text{finite version} \end{aligned}$$

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

of e-functions for L_K

$$\sum_i \lambda_i v_i$$

with non-negative e-values $\lambda_1 \geq \lambda_2 \geq \lambda_3 \dots \geq 0$

$$\text{i.e. } L_K v_i = \lambda_i v_i$$

$$\dots \quad \int_{X^n} \dots \int_{X^1} \dots \int_{X^1} \dots$$

"... . T"

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

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

$$\text{thus } \|v\|_H \leq 1$$

↳ makes an ellipsoid in \mathbb{R}^n

i.e. higher coordinates are shrunk more
 (since λ_i is smaller as $i \rightarrow \infty$)

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

and we have

$$k(x, x') = \sum_{i=1}^{\infty} \lambda_i \psi_i(x) \psi_i(x')$$

Mercer's thm.

$$= \langle kx, w \rangle$$

"U U U"

like $k = \Phi^T \Phi$
of before

a) feature space $H \subseteq \ell_2$
new of H

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

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

"diagonalized representation"

here identify $k(x, \cdot) \in \ell_2$ as $\Phi(x) \in \ell_2$

(do not what $H \subseteq \ell_2$ looks like though)

$$\begin{aligned} & \text{[vector element of } \ell_2 \\ & \text{because } \sum_i (\Phi(x)_i)^2 = \sum_{i=1}^{\infty} \lambda_i \psi_i(x)^2 \\ & = k(x, x) < \infty \end{aligned}$$

b) ℓ_2 new:

$$\begin{aligned} H \subseteq \ell_2 : H &= \left\{ f \in \ell_2 : \sum_{i=1}^{\infty} \frac{|\langle f, \psi_i \rangle|^2}{\lambda_i} < \infty \right\} \\ \langle fg \rangle_H &\triangleq \sum_{i=1}^{\infty} \frac{\langle f, \psi_i \rangle g_i \langle g, \psi_i \rangle}{\lambda_i} \end{aligned}$$

ellipses
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

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

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

but H dense in ℓ_2
+ regularity property of ℓ + decrease λ_n at correct rate

but H dense in \mathcal{L}_2
+ regularity property of J + choose λ_n at correct rate
 \Rightarrow consistency of \hat{f}_n i.e. $\lim_{n \rightarrow \infty} \hat{f}_n \xrightarrow{\mathcal{L}_2} f^*$

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