

Lecture 7 - RKHS

Tuesday, January 28, 2020 14:29

Today: continue RKHS in all their glory? :-)

representer's thm: says that (for H a RKHS)

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

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

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

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

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

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

\hookrightarrow inner product of $\varphi(x_i)$ on data $K_{ij} = \langle \varphi(x_i), \varphi(x_j) \rangle$

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

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

getting a handle on H : generalize diagonalization of matrices to ∞ -dim

I) start with finite matrices:

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

$$f: X \rightarrow \mathbb{R}$$

\hookrightarrow finite \Rightarrow just a vector $\begin{pmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix} \in \mathbb{R}^n$ (" \mathbb{R}^n " vector)

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

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

$$H = \overline{\text{span}} \left\{ k(x_i, \cdot); i=1, \dots, n \right\} = \left\{ \sum \alpha_i k(\cdot, x_i) ; \alpha \in \mathbb{R}^n \right\} \subseteq \mathbb{R}^n$$

$\|f\|_H^2 = \sum \alpha_i^2$

constraint.

because $K \geq 0$ spectral thm. \Rightarrow $K = U \Lambda U^T$
 $n \times n$ $\hookrightarrow \text{diag}(\lambda_1, \dots, \lambda_n)$ with $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$

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

we can let $\Phi = \sqrt{\frac{1}{2}} U^T$
 $\Rightarrow K = \Phi^T \Phi$

$\Phi = \begin{pmatrix} \sqrt{\lambda_1} \psi_1^T \\ \vdots \\ \sqrt{\lambda_d} \psi_d^T \\ 0 \end{pmatrix}$
 $d = \text{rank}(K) \leq n$
 $\Phi = \begin{pmatrix} \Phi(x_1) & \dots & \Phi(x_n) \\ \vdots \end{pmatrix}$

$K = \Phi^T \Phi = \sum_{i=1}^n \lambda_i \psi_i \psi_i^T$
 $k(x_i, x_j) = k_{ij} = \sum_{i=1}^n \lambda_i \psi_i(x_i) \psi_i(x_j)$

i.e. $U = \begin{pmatrix} \psi_1 & \dots & \psi_n \end{pmatrix}$

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

If we define x - record of ψ_i , vector
 $\Phi(x) \triangleq \begin{pmatrix} \sqrt{\lambda_1} \psi_1(x) \\ \vdots \\ \sqrt{\lambda_d} \psi_d(x) \end{pmatrix} \in \mathbb{R}^d$
"feature space pt. of view"

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

Note: $K \psi_j = \sum_i \lambda_i \psi_i \underbrace{\psi_i^T \psi_j}_{\delta_{ij}} = \lambda_j \psi_j$
 $\langle \psi_i, \psi_j \rangle_{\mathbb{R}^n} = \delta_{ij}$

back to \mathbb{R}^n -view: $H \subseteq \mathbb{R}^n$; $v \in H \Rightarrow v = K\alpha$ for some $\alpha \in \mathbb{R}^n$
 to get $\|v\|_H$, we compute $\alpha_v = K^+ v$ pseudo-inverse

so $\|v\|_H^2 = \alpha_v^T K \alpha_v = v^T K^+ K v$

$K = U \Lambda U^T$
 $K^+ = U \Lambda^+ U^T$

$= v^T U \Lambda^+ U^T \underbrace{(I_d \otimes \begin{pmatrix} 0 & \\ & I_{n-d} \end{pmatrix})}_{\text{In}} U v$
 $K K^+ = U \Lambda U^T \Lambda^+ U^T$
 $= U \left(\begin{pmatrix} I_d & \\ 0 & 0 \end{pmatrix} \right) U^T$

so $\|v\|_H^2 = \sum_{j=1}^n \underbrace{\langle v, \psi_j \rangle_{\mathbb{R}^n}^2}_{\beta_j \text{ representation}}$

$\rightarrow U^T v$ is projection of v on $\{\psi_i\}_{i=1}^n$ basis
 i.e. $v = \sum_{j=1}^n \beta_j \psi_j$ i.e. $\beta_j = \langle v, \psi_j \rangle_{\mathbb{R}^n}$

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

and $\|\psi_j\|_{\mathbb{R}^n}^2 = \frac{1}{\lambda_j}$

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

$$|\overbrace{\quad}^{x_j^*}|$$

So orthonormal basis of \mathcal{H} in S_2 is $\{\sqrt{\lambda_j} u_j\}_{j=1}^d$
 $\hookrightarrow \|\cdot\|_{l_2}$

$$\begin{aligned}\langle v, v \rangle_{S_2} &= \sum_{i=1}^d \beta_i^2 \\ \langle v, v \rangle_{\mathcal{H}} &= \sum_{i=1}^d \frac{\beta_i^2}{\lambda_i}\end{aligned}$$

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

\hookrightarrow makes an ellipsoid in S_2

|5h3|

Q: higher coordinates are shrunk more?
(since λ_j is smaller as $j \gg$)

II) generalization to n -dim \mathcal{H}

Suppose X is a compact space (e.g. $X = [0, 1]$)
+ Lebesgue measure on it

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

$$l_2 \triangleq \left\{ (\alpha_i)_{i=1}^\infty \text{ s.t. } \sum_i \alpha_i^2 < \infty \right\}$$

Let k be a continuous psd kernel fct. (note it is symmetric)
 \hookrightarrow with respect to standard norm on $X \subseteq \mathbb{R}$

Define $L_k : \mathcal{L}_2 \rightarrow \mathcal{L}_2$
s.t. $[L_k f](\cdot) \triangleq \int_X k(x, \cdot) f(x) dx$

then can show that

L_k is a "compact selfadjoint positive" operator

$$\langle f, L_k g \rangle_{\mathcal{L}_2} = \langle L_k f, g \rangle_{\mathcal{L}_2} \quad \forall f, g \in \mathcal{L}_2$$

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

[Hilbert basis]

of e-functions for L_k $\{u_i\}_{i=1}^\infty$
with non-negative \mathbb{R} -values $\lambda_1 \geq \lambda_2 \geq \dots \geq 0$

$$\begin{aligned}&\text{finite version -} \\ &\langle v, k w \rangle = v^T k w \\ &= (v^T k^T) w \\ &= (Kv)^T w \\ &= \langle Kv, w \rangle\end{aligned}$$

i.e. $L_k u_i = \lambda_i u_i$

and we have

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

Mercator

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

a) feature space $H \subseteq \mathcal{L}_2$
view of \mathcal{H}

$$\Phi : X \rightarrow l_2 \quad \text{with} \quad (\Phi(x))_i \triangleq \sqrt{\lambda_i} u_i(x)$$

a) view of \mathcal{H} $\Phi: \mathcal{X} \rightarrow \ell_2$ with $(\Phi(x))_i \triangleq \sqrt{\lambda_i} \psi_i(x)$

here: $k(x, x') = \langle \Phi(x), \Phi(x') \rangle_{\ell_2}$ [valid element of ℓ_2 because $\sum_i (\Phi(x))_i^2 = \sum_i \lambda_i \psi_i(x)^2 < \infty$]

"diagonalized representation" $\sum_i (\Phi(x))_i^2 = \sum_i \lambda_i \psi_i(x)^2$

here identify $k(x, \cdot) \in \ell_2$ as $\Phi(x) \in \ell_2$ $= k(x, \cdot) \in \mathcal{H}$

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

$$\|f\|_{\mathcal{H}}^2 < \infty$$

b) ℓ_2 view:

$\mathcal{H} \subseteq \ell_2: \mathcal{H} = \left\{ f \in \ell_2 : \sum_{i=1}^n \frac{\langle f, \psi_i \rangle_{\ell_2}^2}{\lambda_i} < \infty \right\} \rightarrow \text{ellipsoid in } \ell_2$

$$\text{and } \langle f, g \rangle_{\mathcal{H}} \triangleq \sum_{i=1}^n \frac{\langle f, \psi_i \rangle_{\ell_2} \langle g, \psi_i \rangle_{\ell_2}}{\lambda_i}$$

④ if \mathcal{K} is "universal"

$\Rightarrow \mathcal{H}_{\mathcal{K}}$ is dense in ℓ_2 i.e. for any $f \in \ell_2$

3 sequence $h_n \in \mathcal{H}_{\mathcal{K}}$ s.t. $\|h_n - f\|_{\ell_2} \xrightarrow{n \rightarrow \infty} 0$

note: if $f \notin \mathcal{H} \Rightarrow \|h_n\|_{\mathcal{H}} \rightarrow \infty$

* non-parametric learning:

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

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

but \mathcal{H} dense in ℓ_2

+ regularity property of \mathcal{L} + correct choice of λ_n

\Rightarrow consistency of \hat{f}_n i.e. $\hat{f}_n \xrightarrow{n \rightarrow \infty} f^*$

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