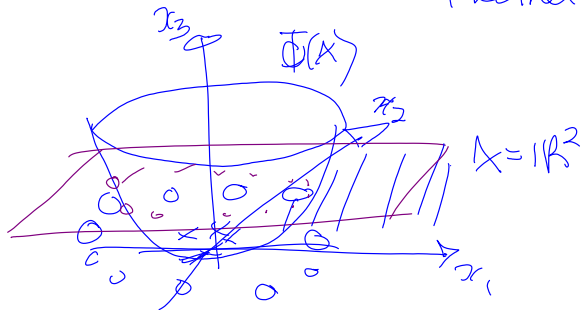


today's more on RKHS

RKHS:

motivation: generalize $\langle w, \phi(x) \rangle$ to higher dimensional space
+ kernel trick $\langle \phi(x), \phi(y) \rangle = k(x, y)$



$$\Phi: X \rightarrow \mathbb{R}^3$$

$$\Phi(x) = \begin{pmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{pmatrix}$$

$$\langle \Phi(x), \Phi(y) \rangle_{\mathbb{R}^3} = (\langle x, y \rangle_{\mathbb{R}^2})^2 = k(x, y)$$

representer's thm says that: for H RKHS produced by a kernel K i.e. $K: X \times X \rightarrow \mathbb{R}$
st. symmetric & psd.

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

is reached for $f^* = \sum_{i=1}^n \alpha_i^* k(x_i, \cdot)$ i.e. $\min_{\alpha \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n f(y_i, \sum_{j=1}^n \alpha_j k(x_j, x_i)) + \lambda \alpha^T K \alpha$

$$\|f_\alpha\|_H^2 = \langle f_\alpha, f_\alpha \rangle_H = \sum_i \sum_j \alpha_i \alpha_j k(x_i, x_j) = \alpha^T K \alpha$$

$$f_\alpha = \sum_{i=1}^n \alpha_i k(x_i, \cdot)$$

$$\langle k(x_i, \cdot), k(x_j, \cdot) \rangle_H = k(x_i, x_j)$$

for any finite $\{x_i\}_{i=1}^m$ and $\alpha \in \mathbb{R}^m$ and any m

$$\alpha^T K \alpha \geq 0$$

\downarrow
 $(K)_{ij} = k(x_i, x_j)$
 "Gram matrix on"
 $\{x_i\}_{i=1}^m$

RKHS:

RKHS:

$$\Phi: X \rightarrow \mathcal{H} \quad \text{s.t.} \quad \langle \Phi(x), \Phi(y) \rangle_{\mathcal{H}} = k(x, y)$$

$$\text{let } \tilde{\mathcal{H}} = \text{span}\{k(x, \cdot) : x \in X\}$$

$$\text{e.g. } f \in \tilde{\mathcal{H}} \Rightarrow f = \sum_i \alpha_i k(x_i, \cdot) \text{ for some finite } \{x_i\}_{i=1}^n \text{ s.t. } \alpha_i \in \mathbb{R}^n$$

"pre-Hilbert" space [inner product space]

$$\text{with } \langle f, g \rangle_{\tilde{\mathcal{H}}} = \sum_{i,j} \alpha_i \alpha_j k(x_i, x_j) \quad \|f\|_{\tilde{\mathcal{H}}} \triangleq \sqrt{\langle f, f \rangle_{\tilde{\mathcal{H}}}}$$

then RKHS \mathcal{H} is = completion($\tilde{\mathcal{H}}$) using $\|\cdot\|_{\tilde{\mathcal{H}}}$ as a norm

i.e. add all limit points of $\tilde{\mathcal{H}}$ -Cauchy sequence to \mathcal{H}

getting a handle on \mathcal{H} : generalizing diagonalization of matrices to ∞ dim

$$\text{say } X \text{ is finite i.e. } \{x_i\}_{i=1}^n; \text{ form Gram matrix } (K)_{i,j} \triangleq k(x_i, x_j)$$

$$K \text{ a valid kernel} \Rightarrow K \succeq 0$$

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

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

$$K \succeq 0 \Rightarrow \text{spectral thm.}$$

$$K = U \Lambda U^T \text{ with } \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \text{ and } \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$$

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

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

$$U^T U = I_n \quad \text{i.e. } \langle \psi_i, \psi_j \rangle = \delta_{ij}$$

$$\text{we can let } \Phi = \Lambda^{1/2} U^T \Rightarrow K = \Phi^T \Phi$$

$$\Phi = \begin{pmatrix} \sqrt{\lambda_1} \psi_1^T \\ \vdots \\ \sqrt{\lambda_n} \psi_n^T \end{pmatrix} \quad \text{rowwise: } \sqrt{\lambda_i} (\psi_i(x_1) \dots \psi_i(x_n))$$

$$\{x_i\}_{i=1}^n$$

by definition, we have that $\langle f, k(x, \cdot) \rangle_{\mathcal{H}} = f(x)$

$$\text{mapping } E_x: \mathcal{H} \rightarrow \mathbb{R}$$

$$E_x(f) = f(x)$$

$$\begin{aligned} & |f(x) - g(x)| \\ &= |\langle f - g, k(x, \cdot) \rangle| \\ &\stackrel{\text{C.S.}}{\leq} \|f - g\|_{\mathcal{H}} \|k(x, \cdot)\|_{\mathcal{H}} \end{aligned}$$

i.e. E_x is Lipschitz continuous with $L = \|k(x, \cdot)\|_{\mathcal{H}}$

suppose $\lambda_{d+1} = 0$

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

"feature space view"

$$K = \Phi \Phi^T = \sum_{i=1}^n \lambda_i \psi_i \psi_i^T$$

$$K_{ij} = \sum_{k=1}^n \lambda_k \psi_k(x_i) \psi_k(x_j)$$

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

$\mathcal{S}_2 \approx \mathbb{R}^n$ when $\mathcal{X} = \text{finite set}$

$$K \psi_j = \sum_{i=1}^n \lambda_i \psi_i \psi_i^T \psi_j = \sum_{i=1}^n \lambda_i \psi_i \langle \psi_i, \psi_j \rangle_{\mathbb{R}^n} = \sum_{i=1}^n \lambda_i \psi_i \langle \psi_i, \psi_j \rangle_{\mathcal{S}_2} \Rightarrow K \psi_j = \lambda_j \psi_j$$

back to \mathcal{S}_2 -view: $H \subseteq \mathcal{S}_2$; $v \in H \Rightarrow v = K \alpha$ for some α
 $= \sum_{i=1}^n \alpha_i K(x_i, \cdot)$

to get $\|v\|_H$, we compute $\alpha = K^+ v$ pseudo-inverse

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

$$\text{so } \|v\|_H^2 = \alpha^T K \alpha = v^T K^+ v$$

$$= (U^T v)^T \Lambda^+ (U^T v) \rightarrow U^T v \text{ is projection of } v \text{ on } \{\psi_j\}_{j=1}^n \text{ basis}$$

so let $v = \sum_j \beta_j \psi_j$ i.e. $\beta_j = \langle v, \psi_j \rangle_{\mathcal{S}_2}$; then

$$\|v\|_H^2 = \sum_{j=1}^n \frac{\langle v, \psi_j \rangle_{\mathcal{S}_2}^2}{\lambda_j}$$

β_j representation

$$\text{vs. } \|v\|_{\mathcal{S}_2}^2 = \sum_{j=1}^n \langle v, \psi_j \rangle_{\mathcal{S}_2}^2$$

$$\text{and } \|\psi_j\|_H^2 = \frac{1}{\lambda_j}$$

so orthonormal basis of H in \mathcal{S}_2 is $\{\frac{\psi_j}{\sqrt{\lambda_j}}\}_{j=1}^d$
 $\hookrightarrow \| \cdot \|_H$

\mathcal{S}_2 (outside space view of H): $H \subseteq \mathbb{R}^n$ (i.e. \mathcal{S}_2) say $v \in H$ ($v \in \text{Im}(K)$) i.e. $v = \sum_{i=1}^d \beta_i \psi_i$

$$\text{then } \langle v, v \rangle_{\mathcal{S}_2} = \sum_{i=1}^d \beta_i^2 \langle \psi_i, \psi_i \rangle_{\mathcal{S}_2} = \sum_{i=1}^d \beta_i^2$$

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

thus $\|v\|_H \leq 1$ makes ellipsoid in \mathcal{S}_2

feature space view : $H \subseteq \mathbb{R}^d$ $\text{span}(\Phi(x))$ $\langle v, v \rangle_H = \langle v, v \rangle_{\mathbb{R}^d}$
 i.e. already diagonalized

(with small coordinates
 for j s.t. λ_j is small)

$$|X| = n$$

$$\mathbb{R}^X \cong \mathbb{R}^n$$

generalization to n -dim space:

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

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

$$\ell_2 \triangleq \{ (\alpha_i)_{i=1}^\infty \text{ s.t. } \sum_{i=1}^\infty \alpha_i^2 < \infty \}$$

Let k be a continuous psd kernel function

↓ with respect to standard norm on $\mathbb{A} \in \mathbb{R}$

$$\langle v, kw \rangle = v^T k w = (v^T k^T) w = \langle v, w \rangle$$

when $k = k^T$

define $L_k: \mathcal{S}_2 \rightarrow \mathcal{S}_2$

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

then can show that

L_k is a "compact self-adjoint positive" operator

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

and yields an orthonormal basis (for \mathcal{S}_2) of e-functions for L_k $\{\psi_i\}_{i=1}^\infty$

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

Mercer's Thm

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

and we have $k(x, y) = \sum_{i=1}^\infty \lambda_i \psi_i(x) \psi_i(y)$

[like the $K = \Phi \Phi^T$ of before?]

feature space $\mathcal{H} \subseteq \ell_2$ $\Phi: X \rightarrow \ell_2$ with $(\Phi(x))_i \triangleq \sqrt{\lambda_i} \psi_i(x)$ $\left(\sum_{i=1}^\infty (\Phi(x))_i^2 = \sum_{i=1}^\infty \lambda_i \psi_i(x)^2 \right)$

of \mathcal{H} .

\mathcal{S}_2
rew

i.e. we identify
 $k(x, \cdot) \in \mathcal{S}_2$ as $\Phi(x) \in \ell_2$

where $k(x, y) = \langle \Phi(x), \Phi(y) \rangle_{\ell_2}$ "diagonalized representation"

$$= k(x, x) \\ < \infty$$

(do not know
what is $\mathcal{H} \subseteq \mathcal{S}_2$ though)

$$\mathcal{H} \subseteq \mathcal{S}_2 : \mathcal{H} = \left\{ f \in \mathcal{S}_2 : \sum_{i=1}^{\infty} \underbrace{\langle f, \psi_i \rangle_{\mathcal{S}_2}}_{\lambda_i} < \infty \right\}$$

$$\text{and } \langle f, g \rangle_{\mathcal{H}} \triangleq \sum_{i=1}^{\infty} \underbrace{\langle f, \psi_i \rangle_{\mathcal{S}_2} \langle g, \psi_i \rangle_{\mathcal{S}_2}}_{\lambda_i}$$

Structured prediction setup:

$$S(\cdot, y) \in \mathcal{H}_y \text{ for each } y$$

before, had $S(x, y; w)$

now, we have $S(x, y; \vec{f})$

$$\vec{f} = (f_{y_1}, \dots, f_{y_K})$$

$\in \mathcal{H}_{y_1}$

i.e. $\vec{f} : \mathcal{X} \rightarrow \mathbb{R}^K$

$$S : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$$

$$\text{consider } \ell(x, y; S) \triangleq \frac{1}{2K} \sum_{\tilde{y}} (S(x, \tilde{y}) + \ell(y, \tilde{y}))^2$$

$K=|\mathcal{Y}|$

can be thought of as generalization
of squared loss for binary classification
to multiclass $\rightarrow (1 - y_i \langle w, \phi(x_i) \rangle)^2$

$$\ell_{pc}(x; S) \triangleq \mathbb{E}_{q(y|x)} \left(\frac{1}{2K} \sum_{\tilde{y}} S(x, \tilde{y})^2 + 2S(x, \tilde{y}) \mathbb{E}_{q(y|x)} [\ell(y, \tilde{y})] + \text{constant} \right)$$

$$= \frac{1}{2K} \sum_{\tilde{y}} (S(x, \tilde{y})^2 + 2S(x, \tilde{y}) \underbrace{\mathbb{E}_{q(y|x)} [\ell(y, \tilde{y})]}_{\ell_{q(x)}(\tilde{y})} + \text{ct.})$$

$$l_{q_x}(\tilde{y})$$

Suppose S is unconstrained;

$$\min_S l_{q_x}(x; s) \rightarrow 2S(\tilde{y}) + 2 l_{q_x}(\tilde{y}) = 0 \quad \forall \tilde{y}$$

$$\Rightarrow S^*(\tilde{y}) = -l_{q_x}(\tilde{y})$$

$$\arg\max_{\tilde{y}} S^*(x, \tilde{y}) = \arg\min_{\tilde{y}} l_{q_x}(\tilde{y}) \quad \text{i.e. you predict optimally}$$

so S here is consistent

$$\hookrightarrow \text{i.e. } S^* \in \arg\min_{\text{over all } S} \mathcal{L}(S) \Rightarrow \text{that } L(h_{S^*}) = \min_{\text{all } h} L(h)$$

$$h_S(x) \triangleq \arg\max_{\tilde{y}} S(x, \tilde{y})$$