

Kernel Machines

Basic idea:

$$(X, Y) \in \mathcal{X} \times \{\pm 1\}$$

$$\left\{ \begin{array}{l} K: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R} \quad (\text{Mercer kernel}) \\ (\mathcal{H}_K, \langle \cdot, \cdot \rangle_K) - \text{RKHS} \end{array} \right.$$

$\mathcal{F} \subset \mathcal{H}_K$

classifiers: $g_f(x) = \text{sgn } f(x)$

$$\mathcal{F}_\lambda := \{f \in \mathcal{H}_K : \|f\|_K^2 \leq \lambda\} \quad (\lambda > 0)$$

Rademacher complexities:

$$x^n = (x_1, x_2, \dots, x_n) \in \mathcal{X}^n$$

$$R_n(\mathcal{F}_\lambda(x^n)) = \frac{1}{n} \mathbb{E}_\varepsilon \left[\sup_{f \in \mathcal{F}_\lambda} \left| \sum_{i=1}^n \varepsilon_i f(x_i) \right| \right]$$

$$f(x_i) = \langle f, K_{x_i} \rangle_K$$

K_{x_i} - rep. of $K(x_i, \cdot)$ in \mathcal{H}_K

$$= \frac{1}{n} \mathbb{E}_\varepsilon \left[\sup_{f \in \mathcal{F}_\lambda} \left| \sum_{i=1}^n \varepsilon_i \langle f, K_{x_i} \rangle_K \right| \right]$$

$$= \frac{1}{n} \mathbb{E}_\varepsilon \left[\sup_{f \in \mathcal{F}_\lambda} \left| \langle f, \underbrace{\sum_{i=1}^n \varepsilon_i K_{x_i}}_{\in \mathcal{H}_K} \rangle_K \right| \right]$$

Cauchy-Schwarz:

$$|\langle f, g \rangle_K| \leq \|f\|_K \|g\|_K$$

for $f, g \in \mathcal{H}_K$

$$\leq \frac{1}{n} \mathbb{E}_\varepsilon \left[\sup_{f \in \mathcal{F}_\lambda} \|f\|_K \left\| \sum_{i=1}^n \varepsilon_i K_{x_i} \right\|_K \right]$$

$$\leq \frac{\lambda}{n} \mathbb{E}_\varepsilon \left\| \sum_{i=1}^n \varepsilon_i K_{x_i} \right\|_K$$

$$\mathbb{E}_{\varepsilon} \left\| \sum_{i=1}^n \varepsilon_i K_{x_i} \right\|_K = \mathbb{E}_{\varepsilon} \sqrt{\left\langle \sum_{i=1}^n \varepsilon_i K_{x_i}, \sum_{i=1}^n \varepsilon_i K_{x_i} \right\rangle_K}$$

$$= \mathbb{E}_{\varepsilon} \sqrt{\sum_{i=1}^n \sum_{j=1}^n \varepsilon_i \varepsilon_j \langle K_{x_i}, K_{x_j} \rangle_K}$$

$$K(x, x') = \langle K_x, K_{x'} \rangle_K$$

$$= \mathbb{E}_{\varepsilon} \sqrt{\sum_{i=1}^n \sum_{j=1}^n \varepsilon_i \varepsilon_j K(x_i, x_j)}$$

$$\leq \sqrt{\sum_{i=1}^n \sum_{j=1}^n K(x_i, x_j) \underbrace{E[\varepsilon_i \varepsilon_j]}_{\mathbb{1}_{\{i=j\}}}}$$

$$= \sqrt{\sum_{i=1}^n K(x_i, x_i)}$$

$$\Rightarrow R_n(\mathcal{F}_\lambda(x^n)) \leq \frac{\lambda}{n} \sqrt{\sum_{i=1}^n K(x_i, x_i)}$$

Remarks:

• Gram matrix: $G = [K(x_i, x_j)]_{i, j=1, \dots, n}$

$$\sum_{i=1}^n K(x_i, x_i) = \text{tr } G$$

$$R_n(\mathcal{F}_\lambda(x^n)) \leq \frac{\lambda}{n} \sqrt{\text{tr } G} \leftarrow \text{data-dependent}$$

- X_1, X_2, \dots, X_n iid elements of \mathcal{X}

$$\begin{aligned} \mathbb{E} R_n(\mathcal{F}_\lambda(x^n)) &\leq \frac{\lambda}{n} \mathbb{E} \sqrt{\sum_{i=1}^n K(X_i, X_i)} \\ &\leq \frac{\lambda}{n} \sqrt{n \mathbb{E}[K(X, X)]} \\ &= \frac{\lambda}{\sigma_n} \sqrt{\mathbb{E}[K(X, X)]} \end{aligned}$$

or if $\sup_{x \in \mathcal{X}} K(x, x) =: C_K < \infty$,

$$\mathbb{E} R_n(\mathcal{F}(x^n)) \leq \frac{\lambda C_K}{\sigma_n}.$$

- the bound is dimension-free: \mathcal{H}_K can be finite- or infinite-dim.

$$K(x, x') = 1 + \langle x, x' \rangle \quad \text{on } \mathbb{R}^d$$

$$\mathcal{F}_\lambda = \left\{ x \mapsto \langle w, x \rangle + b : b^2 + \|w\|^2 \leq \lambda^2 \right\}$$

$$K(x, x') = e^{-\alpha \|x - x'\|^2}$$

\mathcal{H}_K spanned by a countably infinite set of basis functions ψ_1, ψ_2, \dots

Surrogate losses:

$(x_1, y_1), \dots, (x_n, y_n)$ iid in $\mathcal{X} \times \{\pm 1\}$

$$\hat{f}_n \in \mathcal{F}_\lambda$$

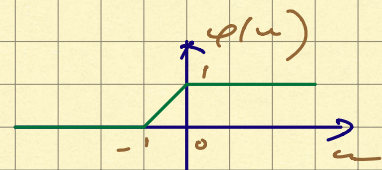
w.p. $\geq 1 - \delta$,

$$L(\text{sgn} \hat{f}_n) \leq A_{\varphi, n}(\hat{f}_n) + \frac{C M_\varphi \lambda}{K} + C \sqrt{\frac{\log(1/\delta)}{n}}$$

- assuming pen. fcn. φ is M_φ -Lipschitz.

e.g. $\varphi(u) = \min\{1, (1+u)_+\}$

or $\varphi(u) \geq \min\{1, (1+u)_+\}$



Kernel Trick / Representer Theorem

• main idea: restrict opt. to fens of the form

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

where $(x_1, y_1), \dots, (x_n, y_n)$ are iid and c_1, \dots, c_n will be tuned based on data.

• orthogonal projection in Hilbert spaces

$(\mathcal{H}, \langle \cdot, \cdot \rangle)$ - Hilbert space

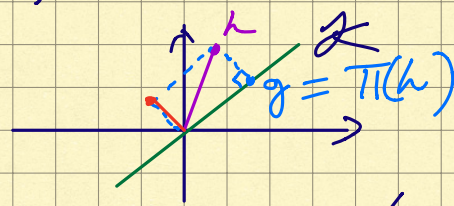
\mathcal{K} : closed subspace of \mathcal{H}

- $h, h' \in \mathcal{K} \Rightarrow \alpha h + \beta h' \in \mathcal{K} \quad \alpha, \beta \in \mathbb{R}$
(\mathcal{K} is a subspace of \mathcal{H})

- (h_n) in \mathcal{K} s.t. $h = \lim_{n \rightarrow \infty} h_n$ exists
then $h \in \mathcal{K}$ (closed subset of \mathcal{H})

Then: for any $h \in \mathcal{H}$, the problem

$$\min_{g \in \mathcal{K}} \|g - h\|^2$$



has a unique solution (the projection of h onto \mathcal{K})

$$\Pi(h) = \operatorname{argmin}_{g \in \mathcal{K}} \|g - h\|^2$$

1) The map $\Pi: \mathcal{H} \rightarrow \mathcal{K}$ is linear:

$$\Pi(\alpha h + \beta h') = \alpha \Pi(h) + \beta \Pi(h')$$

$$2) \Pi^2 = \Pi \circ \Pi = \Pi \quad \Pi(\Pi(h)) = \Pi(h)$$

3) for any $h \in \mathcal{H}$, $g \in \mathcal{K}$,

$$\langle g, h \rangle = \langle \Pi(g), h \rangle = \langle g, \Pi(h) \rangle$$

4) $\mathcal{K}^\perp := \{ h \in \mathcal{H} : \langle g, h \rangle = 0 \ \forall g \in \mathcal{K} \}$
— orthogonal complement of \mathcal{K} in \mathcal{H} —
is also a closed subspace of \mathcal{H} , and any $h \in \mathcal{H}$ can be uniquely represented as

$$h = g + g^\perp$$

where $g = \Pi(h)$ and $g^\perp \in \mathcal{K}^\perp$.

Back to kernels —

$$(X_1, Y_1), \dots, (X_n, Y_n)$$

\mathcal{H}_K RKHS

\mathcal{H}_n : closed subspace of \mathcal{H}_K spanned by
 $K_{X_1}, K_{X_2}, \dots, K_{X_n}$
— random subspace of \mathcal{H}_K

$\Pi_n: \mathcal{H}_K \rightarrow \mathcal{H}_n$ — orthogonal proj.
onto \mathcal{H}_n

Representer theorem If F is a subset of \mathcal{H}_K
s.t. $\Pi_n(F) \subseteq F$, then

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(x_i)) = \min_{\hat{f} \in \mathcal{H}_n} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, \hat{f}(x_i))$$

— minimizer is of the form

$$\hat{f}_n(x) = \sum_{i=1}^n c_i K_{x_i}(x)$$

where $c_1, \dots, c_n \in \mathbb{R}$ depend on $(x_1, Y_1), \dots, (x_n, Y_n)$

Proof idea

$$f(x_i) = \langle f, \underbrace{K_{x_i}}_{\in \mathcal{H}_n} \rangle_K \quad (\text{repr. property})$$

$$= \langle f, \Pi_n(K_{x_i}) \rangle_K$$

$$= \langle \Pi_n(f), K_{x_i} \rangle_K$$

$$= \langle \tilde{f}, K_{x_i} \rangle_K$$

$$\tilde{f} := \Pi_n(f)$$

$$= \tilde{f}(x_i)$$

$$\forall f \in \mathcal{F}, \quad \ell(Y_i, f(x_i)) = \ell(Y_i, \tilde{f}(x_i))$$

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(x_i)) = \min_{\tilde{f} \in \Pi_n(\mathcal{F})} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, \tilde{f}(x_i))$$

$$\hat{f}_n = \underset{\tilde{f} \in \Pi_n(\mathcal{F})}{\text{argmin}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, \tilde{f}(x_i))$$

is an element of \mathcal{H}_n . □

Optimization:

$$\varphi: \mathbb{R} \rightarrow \mathbb{R}$$

$$\min_{f \in \mathcal{F}_\lambda} \frac{1}{n} \sum_{i=1}^n \varphi(-\gamma_i f(x_i))$$

$$\mathcal{F}_\lambda = \left\{ f \in \mathcal{H}_K : \|f\|_K^2 \leq \lambda^2 \right\}$$

$$= \min_{c_1, \dots, c_n} \frac{1}{n} \sum_{i=1}^n \varphi\left(-\gamma_i \sum_{j=1}^n c_j k(x_i, x_j)\right)$$

$$\text{s.t.} \quad \sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) \leq \lambda^2$$

$$= c^T G c \quad \text{where } G = [k(x_i, x_j)]_{i,j}$$

- by rep. thm, look for $\hat{f}_n(x) = \sum_{i=1}^n c_i k_{x_i}(x)$
s.t. $\|\hat{f}_n\|_K^2 \leq \lambda^2$

$$\begin{aligned} \|\hat{f}_n\|_K^2 &= \left\langle \sum_{i=1}^n c_i k_{x_i}, \sum_{i=1}^n c_i k_{x_i} \right\rangle_K \\ &= \sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) \end{aligned}$$

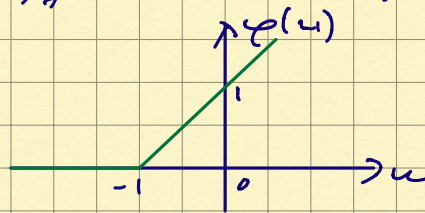
- note: if φ is convex, then

$$\min_{c \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n \varphi\left(-\gamma_i \sum_{j=1}^n c_j k(x_i, x_j)\right)$$

$$\text{s.t.} \quad c^T G c \leq \lambda^2$$

is a convex program w/ quadratic constraints!

- Take $\varphi(u) = (1+u)_+$ [hinge]



$$K(x, x') = \sum_j \psi_j(x) \psi_j(x')$$

where ψ_1, ψ_2, \dots are basis fns

$$f(x) = \arg \left(\sum_j c_j \tilde{\psi}_j(x) \right)$$

$$\text{let } \psi_j(x) = \sqrt{w_j} \tilde{\psi}_j(x) \quad w_j > 0, \sum_j w_j < \infty$$

$$K(x, x') = \sum_j w_j \tilde{\psi}_j(x) \tilde{\psi}_j(x')$$

$$f(x) = \sum_j \frac{c_j}{\sqrt{w_j}} \psi_j(x)$$

$$\|f\|_K^2 = \sum_j \frac{c_j^2}{w_j}$$

$$\min \frac{1}{n} \sum_{i=1}^n \varphi \left(-\gamma_i \sum_j c_j \tilde{\psi}_j(x_i) \right)$$

$$\text{s.t. } \sum_j c_j^2 / w_j \leq \Lambda^2$$

if $\dim \mathcal{H}_K$ is fin-dim, this is more efficient
if $\dim \mathcal{H}_K \gg n$

or if $\tilde{\psi}_1, \tilde{\psi}_2, \dots$ are arranged s.t.

$w_1 > w_2 > \dots$, then can set all but finitely many c_1, \dots, c_K to zero.