

HW2 ✓ 2 different textbooks
 Foundations of Machine Learning Mohri
 2018
 Learning with kernels Smola &
 Schölkopf 2002

Interpret - kernel regression on RKHS

closed form solution
 (finite-dim space)

$$\text{span}\{x_i\}_{i \in [m]} \quad \{(x_i, y_i)\}_{i \in [m]} = S$$

- why does the kernel regression problem on RKHS become finite-dim. instead?

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

\downarrow
RKHS

Regularization term: reduces "capacity" of \mathcal{H} .

- "Smallness" of \mathcal{H} comes from the compactness of Y^*Y (regularization operator)

$$Y^*Y \iff \text{for PD kernels have a one-to-one correspondence}$$

$$\langle Y^*Y f, f \rangle = \|f\|_{\mathcal{H}}^2$$

$$\langle Y f, Y g \rangle = \langle f, g \rangle_{\mathcal{H}}$$

L^2 inner product RKHS product

- Y^*Y PD

$$\equiv T_X^{-1} \quad T_X g(x') = \int g(x) x(x', x) dx$$

- Green's function of Y^*Y

$$\langle Y^*Y f, G(x, \cdot) \rangle = f(x)$$

$$\text{Last time: } G(x, x') = x(x, x')$$

$$\langle Y^*Y f, x(x, \cdot) \rangle = f(x)$$

$$\int (Y^*Y f)(z) x(x, z) dz = f(x)$$

$$\int (Y^*Y)(T_X f)(z) x(x, z) dz = T_X f(x)$$

$$= \int f(z) x(x, z) dz$$

$$\Rightarrow \int ((Y^*Y)(T_X) f(z) - f(z)) x(x, z) dz = 0$$

$$(Y^*Y) T_X = \text{Id} \quad \text{on } \mathcal{H}.$$

$$\|Y^*Y f\|^2 = \langle Y^*Y f, Y f \rangle = \|f\|_{\mathcal{H}}^2$$

space of functions on which $\|f\|_{\mathcal{H}}^2$ is minimized is "small"

e.g. Gaussian kernel \rightarrow

smooth functions (functions for which the Fourier coefficients decay rapidly)

$$\text{have small } \langle Y^*Y f, f \rangle = \|f\|_{\mathcal{H}}^2$$

Questions: ① what functions are "likely" to be found as minimizers?

② what are characteristics of the class of minimizers?

②: Sample complexity of minimization problems

①: Gaussian process interpretation of regression

: "Bayesian" perspective

function Complexity

PAC bound (probably approximately correct)
for finite hypothesis class:

$$\text{generalization error} \leq \frac{1}{m} \log \left(\frac{|H|}{\delta} \right)$$

$$\mathbb{E}_{(x,y) \sim \mathcal{D}} \ell(x,y, h(x)) \quad \text{w.h.p.}$$

$$\text{Training error} \\ \frac{1}{m} \sum_{(x,y) \in S} \ell(x,y, h(x))$$

Under the realizability
assumption,
training error = 0
for ERM.

For an ERM $h \in H$,

$$\mathbb{E}_{(x,y) \sim \mathcal{D}} \ell(x,y, h(x)) \leq \epsilon = \frac{1}{m} \log \frac{|H|}{\delta}$$

with Prob. $\geq 1 - \delta$

(over the randomness
in training data)

Sample complexity: function $(\epsilon, \delta) \rightarrow m(\epsilon, \delta)$
s.t. if trained with at least $m(\epsilon, \delta)$
 $m(\epsilon, \delta)$ samples, then, gen. error $R(h)$
 $< \epsilon$ with prob. $\geq 1 - \delta$.

$$\text{if } m = \frac{1}{\epsilon} \log \frac{|H|}{\delta}$$

then, $R(h) < \epsilon$ with prob $\geq 1 - \delta$.



Kernel regression: \mathcal{H} : RKHS

Vacuous bound

Rademacher complexity:

(one of many notions
to represent complexity
of function class)

Given $S \equiv \{(x_i, y_i)\}_{i \in [m]}$

$$\text{Rad}_S(\mathcal{H}) = \frac{1}{m} \mathbb{E}_{\sigma} \sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i h(x_i)$$

$$\sigma = [\sigma_1, \dots, \sigma_m]^T \in \mathbb{R}^m$$

σ_i iid Rademacher RV

$$\sigma_i = \begin{cases} +1 & \text{with } P_{\sigma} \frac{1}{2} \\ -1 & \text{with } P_{\sigma} \frac{1}{2} \end{cases}$$

$\text{Rad}(\mathcal{H}) \uparrow \Rightarrow$ can represent
noisy functions

Rademacher complexity of functions learned in kernel regression on RKHS

$$f(x) = \alpha \cdot \begin{bmatrix} \kappa(x_1, x) \\ \vdots \\ \kappa(x_m, x) \end{bmatrix} = \sum_{i=1}^m \alpha_i \kappa(x_i, x)$$

$$= \alpha \cdot \Phi(x)$$

$$\lambda \alpha^T G \alpha \leq \lambda \|\alpha\|_{G^{-1}}^2$$

$$\mathcal{H} = \left\{ x \mapsto \alpha \cdot \Phi(x) : \|\alpha\| < \Lambda, \text{Tr}(G) < \infty \right\}$$

Reg. term: $\lambda \alpha^T G \alpha$

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

$$\text{Rad}_S(\mathcal{H}) = \frac{1}{m} \mathbb{E}_{\sigma} \sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i h(x_i)$$

$$= \frac{1}{m} \mathbb{E}_{\sigma} \sup_{\|\alpha\| < \Lambda} \sum_{i=1}^m \sigma_i \alpha \cdot \phi(x_i)$$

$$= \frac{1}{m} \mathbb{E}_{\sigma} \sup_{\|\alpha\| < \Lambda} \alpha \cdot \sum_{i=1}^m \sigma_i \phi(x_i)$$

$$= \frac{1}{m} \mathbb{E}_{\sigma} \left\| \sum_{i=1}^m \sigma_i \phi(x_i) \right\| \quad (\text{Cauchy Schwarz})$$

$$\leq \frac{1}{m} \left(\mathbb{E}_{\sigma} \left\| \sum_{i=1}^m \sigma_i \phi(x_i) \right\|^2 \right)^{1/2}$$

$$(\sigma_1, \dots, \sigma_m) \quad (\text{Jensen's inequality})$$

$$\mathbb{E} X^2 \geq (\mathbb{E} X)^2$$

$$\left(\mathbb{E}_{\sigma} \left\| \sum_{i=1}^m \sigma_i \phi(x_i) \right\|^2 = \mathbb{E}_{\sigma} \left\| \sum_{i=1}^m \sigma_i^2 \phi(x_i) \right\|^2 + \mathbb{E}_{\sigma} \sum_{i \neq j} \sigma_i \sigma_j \langle \phi(x_i), \phi(x_j) \rangle \right)$$

$$\mathbb{E}_{\sigma} \sigma_i \sigma_j = 0$$

σ_i & σ_j are independent

$$\leq \frac{1}{m} \left(\mathbb{E}_{\sigma} \sum_{i=1}^m \sigma_i^2 \|\phi(x_i)\|^2 \right)^{1/2}$$

$$= \frac{1}{m} \left(\sum_{i=1}^m \|\phi(x_i)\|^2 \right)^{1/2} = \frac{1}{m} \|G\|_F$$

Recall

$$\phi(x_i) = [\kappa(x_i, x_1), \dots, \kappa(x_i, x_m)]$$

$$\|\phi(x_i)\|^2 = ?$$

$$\text{Rad}_S(\mathcal{H}) \leq \frac{1}{m} \|G\|_F$$

Thm:

Form of generalization bound:

$$R(h) \leq \hat{R}_S(h) + \text{Rad}_S(\mathcal{H}) +$$

gen error

$$3 \sqrt{\frac{\log^2 1/\delta}{2m}}$$

with Prob. $\geq 1 - \delta$.

• Proof: McDiarmid's inequality

$$R(h) - \hat{R}_S(h) \leq \text{Rad}_S(\mathcal{H}) + 3 \sqrt{\frac{\log^2 1/\delta}{2m}}$$

generalization gap

$$\leq \frac{1}{m} \|G\|_F + 3 \sqrt{\frac{\log^2 1/\delta}{2m}}$$

$$\frac{1}{m} \|G\|_F = \varepsilon$$

if Eigenvalues of G are \uparrow , $\text{Rad}(\mathcal{H}) \uparrow$.

(which converge to eigenvalues of T_{κ} as $m \rightarrow \infty$)