

10<sup>th</sup> & 12<sup>th</sup> presentations

- Regression, RKHS, statistical learning theory  
→ kernel methods, deep learning theory  
→ choice of kernels + scikit-learn  
→ generative modeling + kernel
- } Goal

Representer theorem:

Any minimizer of the "kernel regression" problem of the form

$$\min_{h \in \mathcal{H}} \hat{R}_S(h) + F(\|h\|_{\mathcal{H}})$$

(empirical risk)

where  $F: \mathbb{R}^+ \rightarrow \mathbb{R}$  is an increasing function has the form

$$h(x) = \sum_{i=1}^m \alpha_i x(x_i, x)$$

- In other words,  $h$  has a finite-dimensional representation in span  $\{x(x_i, \cdot)\}_{i=1}^m$  with  $\{x_i\}_{i=1}^m$  being training points
- $\hat{R}_S(h) \equiv$  empirical risk
- $S = \{(x_i, y_i)\}_{i=1}^m$
- Regularized ERM:  $\hat{R}_S(h) + F(\|h\|_{\mathcal{H}})$
- $\mathcal{H}$ : RKHS associated with PD kernel  $x$ .

Proof:

$$h = h_0 + h_{\perp}$$

$$\text{where } h_0 \in \text{span}\{x(x_i, \cdot)\}_{i=1}^m = H_0$$

$$h_{\perp} \in H_0^{\perp}$$

$$\begin{aligned} \hat{R}_S(h) &= \frac{1}{m} \sum_{i=1}^m \ell((x_i, y_i), h(x_i)) \\ &= \frac{1}{m} \sum_{i=1}^m \ell((x_i, y_i), \langle h, x(x_i, \cdot) \rangle) \quad (\text{reproducing property}) \\ &= \frac{1}{m} \sum_{i=1}^m \ell((x_i, y_i), \langle h_0 + h_{\perp}, x(x_i, \cdot) \rangle) \\ &= \frac{1}{m} \sum_{i=1}^m \ell((x_i, y_i), \langle h_0, x(x_i, \cdot) \rangle) \end{aligned}$$

- $F(\|h\|_{\mathcal{H}})$  is monotone increasing

$$F(\|h\|_{\mathcal{H}}^2) \text{ is also increasing}$$

(  $F$  quadratic fun. on  $[0, \infty)$  is increasing )

$$F(\|h\|_{\mathcal{H}}^2) = F(\|h_0\|_{\mathcal{H}}^2 + \|h_{\perp}\|_{\mathcal{H}}^2)$$

$$(\|h\|_{\mathcal{H}}^2 = \|h_0\|_{\mathcal{H}}^2 + \|h_{\perp}\|_{\mathcal{H}}^2 \text{ Pythagoras thm})$$

$$\geq F(\|h_0\|_{\mathcal{H}}^2)$$

if  $F$  is strictly increasing, then  $(h_{\perp} = 0 \text{ for minimizing})$

$$F(\|h\|_{\mathcal{H}}^2) > F(\|h_0\|_{\mathcal{H}}^2)$$

! solution when  $F$  is strictly increasing

Value of loss at  $h_0$  is smaller than value of loss at  $h$ , for every  $h \in \mathcal{H}$ .

How to use representation thm

WKT soln of

$$\min_{h \in \mathcal{H}} \hat{R}_S(h) + F(\|h\|_{\mathcal{H}}^2)$$

has the form

$$h(x) = \sum_{i=1}^m \alpha_i x(x_i, x)$$

$$\hat{R}_S(h) = \sum_{i=1}^m \ell((x_i, y_i), \sum_{j=1}^m \alpha_j x(x_i, x_j))$$

$$F(\|h\|_{\mathcal{H}}^2) = \sum_{i,j=1}^m \alpha_i \alpha_j x(x_i, x_j)$$

Can use any other regularizer, e.g.

$$\frac{\lambda}{2} \|\alpha\|^2$$

Typical kernel regression:

$$\min_{\alpha \in \mathbb{R}^m} \text{MSE}(\alpha) + \frac{\lambda}{2} \alpha^T G \alpha$$

$$\|h(x)\|_{\mathcal{H}}^2 = \langle h(x), h(x) \rangle_{\mathcal{H}}$$

$$= \left\langle \sum_{i=1}^m \alpha_i x(x_i, x), \sum_{j=1}^m \alpha_j x(x_j, x) \right\rangle_{\mathcal{H}}$$

$$= \sum_{i,j=1}^m \alpha_i \alpha_j \langle x(x_i, x), x(x_j, x) \rangle_{\mathcal{H}}$$

$$= \alpha^T G \alpha$$

$$G_{ij} = x(x_i, x_j)$$

$$m \times m$$

$$\text{MSE}(\alpha) = \frac{1}{m} \sum_{i=1}^m (y_i - h(x_i))^2$$

$$= \frac{1}{m} \sum_{i=1}^m \left( y_i - \sum_{j=1}^m \alpha_j x(x_i, x_j) \right)^2$$

$$= \frac{1}{m} \sum_{i=1}^m \left( y_i^2 + \left( \sum_{j=1}^m \alpha_j x(x_i, x_j) \right)^2 - 2 \sum_{j=1}^m y_i \alpha_j x(x_i, x_j) \right)$$

$$\min_{\alpha} \frac{-\lambda}{m} \sum_{i,j=1}^m y_i \alpha_j x(x_i, x_j)$$

$$+ \frac{1}{m} \sum_{i=1}^m \left( \sum_{j=1}^m \alpha_j x(x_i, x_j) \right)^2 + \lambda \alpha^T G \alpha$$

$$\min_{\alpha} \frac{1}{m} \|Y - G\alpha\|^2 + \lambda \alpha^T G \alpha$$

$$G[i, :] \alpha = h(x_i)$$

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \in \mathbb{R}^m$$

Linear regression for  $\alpha \in \mathbb{R}^m$ .

- Rmk: we are solving for function  $h \in \mathcal{H}$  that takes  $\mathcal{X}$  to scalars.

Inf-dim opt problem  $\rightarrow$  finite-dim

convex optimiz. in

dim = # training pts

and superficially, inde. of  $\dim(\mathcal{X})$ .

- Compare with linear reg.

$$\mathcal{H} = \{x \mapsto w^T x : w \in \mathbb{R}^d\}$$

$$\min_w \frac{1}{m} \|Y - Xw\|^2 + \lambda \|w\|^2$$

$$w_{\text{opt}} = \left( \frac{1}{m} X^T X + \lambda I \right)^{-1} X^T Y$$

$$X \in \mathbb{R}^{m \times d}$$

- Prob. for  $\alpha$ :

$$\min_{\alpha \in \mathbb{R}^m} \frac{1}{m} \|Y - G\alpha\|^2 + \lambda \alpha^T G \alpha$$

Solve:

$$\nabla_{\alpha} \left( \frac{1}{m} (Y - G\alpha)^T (Y - G\alpha) + \lambda \alpha^T G \alpha \right) = 0$$

$$\Rightarrow -\frac{1}{m} (Y - G\alpha)^T G + \lambda \alpha^T G = 0$$

$$\lambda \alpha^T G = \frac{1}{m} (Y - G\alpha)^T G$$

$$\lambda \alpha^T G = \frac{1}{m} Y^T G - \frac{1}{m} \alpha^T G^T G$$

$$\alpha^T \left( \lambda G + \frac{1}{m} G^T G \right) = \frac{1}{m} Y^T G$$

$$\left( \frac{1}{m} G^T G + \lambda G \right) \alpha = \frac{1}{m} G Y$$

$$\alpha = \left( \frac{1}{m} G^T G + \lambda G \right)^{-1} \frac{1}{m} G Y$$

- why is kernel regression effective?

$$\text{O/p: } h(x) = \sum_{i=1}^m \alpha_i x(x_i, x)$$

o.

o

o

o

o

o

o

o

## Feature maps

RKHS feature map:

$$x \rightarrow \kappa(x, \cdot)$$

$$\phi^{\text{RKHS}}(x) = \kappa(x, \cdot)$$

Mercer map

$$x \rightarrow \{\sqrt{\lambda_j} \psi_j(x)\}_{j \in \mathbb{N}}$$

where  $\lambda_j \in \mathbb{R}^+$  were eigenvalues of  $T_\kappa$  (H-S operator)

and  $\psi_j \in L^2(\mathcal{X}, \mu)$  are corresponding eigenfunctions

$$T_\kappa f(x) = \int f(y) \kappa(x, y) d\mu(y)$$

Last time:  $T_\kappa$  is compact on  $L^2(\mathcal{X}, \mu)$

and when  $\kappa$  is sym PD kernel,

$$\lambda_i \in \mathbb{R}^+$$

$$\kappa(x, x') = \sum_{j \in \mathbb{N}} \lambda_j \psi_j(x) \psi_j(x')$$

Want features s.t.

$$\begin{aligned} \kappa(x, x') &= \sum_{j=0}^{\infty} \sqrt{\lambda_j} \psi_j(x) \cdot \sqrt{\lambda_j} \psi_j(x') \\ &= \langle \phi^{\text{mer}}(x), \phi^{\text{mer}}(x') \rangle_{\ell^2} \end{aligned}$$

$$\phi^{\text{mer}}(x) = \{\sqrt{\lambda_j} \psi_j(x)\}_{j \in \mathbb{N}}$$

Recall: finite-dim. Mercer map (over fixed data pts)

( )  $\{\sqrt{\lambda_j} \psi_j\}_{j \in [m]}$   $\psi_j$  were rows of eigenvectors of Gram matrix  
 $\lambda_j$  are eigenvalues of Gram matrix

$$\begin{aligned} \kappa(x, x') &= \langle \phi^{\text{mer}}(x), \phi^{\text{mer}}(x') \rangle_{\ell^2(x)} \\ &= \langle \phi^{\text{RKHS}}(x), \phi^{\text{RKHS}}(x') \rangle_{\mathcal{H}} \end{aligned}$$

If eigenvalues of  $T_\kappa$  decay rapidly,

$$\phi^{\text{mer}}(x) = [\sqrt{\lambda_1} \psi_1(x), \sqrt{\lambda_2} \psi_2(x), \dots]$$

Even if  $\sup_{x \in \mathcal{X}} |\psi_j(x)| < \infty$  is false,

$$\sup_{x \in \mathcal{X}} |\sqrt{\lambda_j} \psi_j(x)| < \infty \text{ can be true}$$

Then, we may think of  $\mathcal{H}$  as a finite-dim. space

Even if  $\mathcal{H}$  is inf-dime., regularization with  $\|\cdot\|_{\mathcal{H}}$  makes opt prob finite dimensional

Empirical kernel map

$$x \rightarrow [k(x_1, x), \dots, k(x_m, x)] \\ \in \mathbb{R}^m$$

(dis. of RKHS map)

---

- Generalization bounds based on Rademacher <sup>complexity</sup>  
(Next time; prepare by reading bounds for  
finite hypothesis classes)
- Rademacher comp - one way of measuring  
size of hypothesis class.