

Nonparametric Learning

↳ STAT240: robust + non-parametric

→ Learning from some infinite-dimensional space.

→ parametric: linear regression w/ d -dimensional covariates

→ nonparametric: all smooth functions
- decision trees, neural networks

Still linear, but just w/ ∞ parameters

Today, Re-do basics of ordinary least squares + ridge regression, but in infinite dimensions.

Shift basis: parameter view \rightarrow function view.

Ordinary least squares:

$(x_1, y_1), \dots, (x_n, y_n)$

$$L(\beta) = \frac{1}{n} \sum_{i=1}^n (y_i - \langle \beta, x_i \rangle)^2$$

$$\beta^* = (X^T X)^{-1} X^T y = \underbrace{\left(\frac{1}{n} \sum_{i=1}^n x_i x_i^T \right)^{-1}}_S \left(\frac{1}{n} \sum_{i=1}^n x_i y_i \right)$$

$$L(\beta) - L(\beta^*) = (\beta - \beta^*)^T S (\beta - \beta^*)$$

Ridge regression:

$$L_\lambda(\beta) = \frac{\lambda}{n} \|\beta\|_2^2 + L(\beta)$$

$$\hat{\beta}_\lambda = (\lambda n \mathbf{I} + X^T X)^{-1} X^T y$$

$$= (\lambda \mathbf{I} + S)^{-1} \left(\frac{1}{n} \sum_{i=1}^n x_i y_i \right)$$

Parameter norm

$$\|\beta\|_2^2$$

Error norm

$$(\beta - \beta^*)^T S (\beta - \beta^*)$$

Function view: $f_\beta(x) = \beta^T x$ instead of β

$$f_\beta \in \mathbb{R}^n \quad (f_\beta(x_1), \dots, f_\beta(x_n))$$

$$L(f) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$L(f) - L(f^*) = \frac{1}{n} \|f - f^*\|_2^2$$

Question: what is $f^* = f_{\beta^*}$? What is $\|\beta\|_2^2$? $f_{\hat{\beta}_\lambda}$?

$$f_\beta = X\beta, \quad X \in \mathbb{R}^{n \times p}$$

$$f^* = f_{\beta^*} = X\beta^* = \underbrace{X(X^T X)^{-1} X^T}_{\text{projection onto linear span of } X} y$$

f^* : projection of y onto linear span of X

$$\| \beta \|_2^2 = \beta^T \beta$$

$$\beta = (X^T X)^{-1} X^T f \rightarrow f = X\beta$$

$$(X^T X)^{-1} X^T X \beta = \beta$$

$$\beta^T \beta = f^T \underbrace{X (X^T X)^{-1} X^T}_{\text{pseudoinverse}} f$$

$$= (X X^T)^{\dagger} f$$

$$K = X X^T$$

$$K_{ij} = \langle x_i, x_j \rangle$$

"kernel matrix"

$$\Rightarrow \boxed{\| \beta \|_2^2 = \begin{cases} f^T K^{\dagger} f & \text{if } f \in \text{span}(K) \\ \infty & \text{if } f \notin \text{span}(K) \end{cases}}$$

$$\operatorname{argmin}_f \frac{1}{n} \|y - f\|_2^2 \Rightarrow (f = y)$$

$$\operatorname{argmin}_{f: f \in \operatorname{span}(K)} \frac{\lambda}{n} R(f) + \frac{1}{n} \|y - f\|_2^2$$

	Complexity norm	Error norm
Parameter view	$\ \beta\ _2^2$	$(\beta - \beta^*)^T S (\beta - \beta^*)$
Function view	$f^T K^T f$	$\frac{1}{n} \ f - f^*\ _2^2$

Ridge regression

$$\operatorname{argmin}_f \frac{\lambda}{n} f^T K^T f + \frac{1}{n} \|f - y\|_2^2$$

$$K^T = K^{-1}$$

$$\Rightarrow \lambda K^T f + f - y = 0$$

$$\Rightarrow (\lambda K^T + I) f = y$$

$$= (\lambda K^{-1} + I)^{-1} = ((\lambda I + K) K^{-1})^{-1}$$

$$f = (\lambda K^T + I)^{-1} y$$

$$= K (\lambda I + K)^{-1} y$$

$$\rightarrow \hat{f}_\lambda = K (\lambda I + K)^{-1} y$$

$\hookrightarrow y$, kernel matrix K

Kernel ridge regression

$$f_{\beta}(x) = \beta^T X \quad \beta = (X^T X)^{-1} X^T f$$

$$= X^T \beta \quad K = X X^T$$

$$= X^T \underbrace{(X^T X)^{-1}} X^T f \quad K^T = X (X^T X)^{-2} X^T$$

$$= X^T \underbrace{X^T X (X^T X)^{-2} X^T}_{K^T} f$$

$$= \underbrace{X^T X^T}_{K_x^T} \underbrace{K^T f}_{K_x^T K^T f}$$

$$X \in \mathbb{R}^{n \times p}$$

$$x \in \mathbb{R}^{p \times 1}$$

$$(Xx)_i = \langle x_i, X \rangle$$

$$f = \hat{f} \lambda$$

$$K_x^T K_x^T K (\lambda I + K)^{-1} y$$

$$= \boxed{K_x^T (\lambda I + K)^{-1} y} \quad (\text{if } K \text{ invertible})$$

Reproducing

We only

Kernel

need

Hilbert

to

care

about

K

$\langle x_i, x_j \rangle$

Assume we have positive definite kernel function $k(x, x')$

$$T = \{x_1, \dots, x_n\}$$

$$K_{ij} = k(x_i, x_j)$$

then K is positive semi-definite,
and it is strictly positive definite if the x_i
are distinct. \hookrightarrow implies \mathcal{H} is infinite-dimensional
otherwise $\text{rank}(K) \leq p$

Def 3 Given K , the corresponding reproducing Kernel Hilbert space
 \mathcal{H} is all functions $f: X \rightarrow \mathbb{R}$ s.t.

$$\sup_T \|f\|_{K[T]} < \infty.$$

$$\hookrightarrow \|f\|_{K[T]} = \underbrace{\begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{bmatrix}^T K^{-1} \begin{bmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{bmatrix}}_{\text{"}\|y\|_Z^2\text{"}}$$

Also define the norm

$$\|f\|_{\mathcal{H}} = \sup_T \|f\|_{K[T]}.$$

Lemma 1 (Restrictions are contractions)

For any f and sets $T \subseteq T'$, $\|f\|_{K[T]} \leq \|f\|_{K[T']}$.

Lemma 2 (Isometric extension)

Let $T = \{x_1, \dots, x_n\}$ and $y_1, \dots, y_n \in \mathbb{R}$.
Then there exists an $f \in \mathcal{H}$ such that $f(x_i) = y_i \forall i$
and $\|f\|_{\mathcal{H}}^2 = \|y\|_{K[T]}^2 = y^T K^{-1} y$.

Pf of Lemma 2.

Define $f(x) = K_x^T K^{-1} y$

• $f(x_i) = y_i$; $f(x_i) = K_{x_i}^T K^{-1} y$
 $= K_i^T K^{-1} y = e_i^T y = y_i$

• Boundry norm of f

$$\|f\|_{\mathcal{H}} = \sup_{T'} \|f\|_{K[T']}$$

$$T' = \{x'_1, \dots, x'_n\}$$

$$\underbrace{f[T']}^T \underbrace{K[T']}^{-1} f[T']$$

$$\begin{bmatrix} f(x'_1) \\ \vdots \\ f(x'_n) \end{bmatrix}$$

$$f[T'] = \begin{bmatrix} K_{x'_1} \\ \vdots \\ K_{x'_n} \end{bmatrix}^T K^{-1} y$$

$$\underbrace{\hspace{10em}}_{K[T', T]}$$

$$K[T', T]_{ij} = k(x'_i, x_j)$$

$$\|f\|_{K[T']}^2 = y^T K^{-1} K[T, T'] K[T']^{-1} K[T', T] K^{-1} y$$

$$K_{11} = K[T]$$

$$K_{12} = K[T, T']$$

$$K_{21} = \dots$$

$$K_{22}$$

$$\begin{bmatrix} K_{11} & K_{12} \\ K_{21} & K_{22} \end{bmatrix} \succeq 0$$

Kernel matrix on $T \cup T'$

$$\hookrightarrow = y^T K_{11}^{-1} K_{12} K_{22}^{-1} K_{21} K_{11}^{-1} y$$

wants: $\leq y^T K_{11}^{-1} y$ (original $\|y\|_K^2$)

$$\Rightarrow \text{wants: } K_{11}^{-1} K_{12} K_{22}^{-1} K_{21} K_{11}^{-1} \leq K_{11}^{-1} \checkmark$$

$$K_{12} K_{22}^{-1} K_{21} \leq K_{11}$$

\hookrightarrow Holds via Schur complement lemma

$$\begin{bmatrix} K_{22} & K_{21} \\ K_{12} & K_{11} \end{bmatrix} \succeq 0 \iff \begin{bmatrix} A & B \\ B^T & C \end{bmatrix} \succeq 0$$

$$\iff B^T A^{-1} B \leq C$$

and $A \succeq 0$ \square

Implications of Lemma 2 for kernel ridge regression.

$$\hat{f}_\lambda = \underset{f}{\operatorname{argmin}} \frac{\lambda}{n} \|f\|_{K(T)}^2 + \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2$$

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

\uparrow depends on $T = \{x_1, \dots, x_n\}$

by Lemma 2

$$= \underset{f}{\operatorname{argmin}} \frac{\lambda}{n} \|f\|_{\mathcal{H}}^2 + \frac{1}{n} \sum_{i=1}^n (f(x_i) - y_i)^2$$

\uparrow \uparrow

can extend from T to X but only pay cost $\|f\|_{KCS}$

only cares about f on T

only cares about f on T

Bias-variance decomposition

$\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$

$$Y_i = f^*(x_i) + \varepsilon_i$$

\hat{f} : bias and variance of \hat{f}

$$\text{Bias}^2 = \frac{1}{n} \sum_{i=1}^n \left(\mathbb{E}_{\varepsilon} [f(x_i)] - f^*(x_i) \right)^2$$

$$\text{Var} = \frac{1}{n} \sum_{i=1}^n \text{Var}_{\varepsilon} \left[\hat{f}(x_i) \right]$$

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\varepsilon} \left(\hat{f}(x_i) - f^*(x_i) \right)^2 = \text{Bias}^2 + \text{Var}$$

Ordinary least squares; analyzed these in context of robust standard errors

\hat{f}_{λ} : ridge regression

$$\hat{f}_0 = K(\lambda I + K)^{-1} y \Big|_{\lambda=0} = K K^{-1} y = y$$

$$\text{Var}[\hat{f}_0(x_i)] = \text{Var}[y_i] = \sigma^2$$

Overall variance: $\sigma^2 \rightarrow$ Doesn't $\rightarrow 0$ as $n \rightarrow \infty$.
 Bias: 0