

Spring 2018: STA 6448
Advanced Probability and Inference II
Lecture 24

Yun Yang

- Non-parametric least squares

Estimation via constrained least-squares

- ▶ Standard non-parametric regression model:

$$y_i = f^*(x_i) + v_i, \quad v_i = \sigma w_i \sim \mathcal{N}(0, \sigma^2), \quad \text{for } i = 1, \dots, n.$$

- ▶ We consider estimate f^* is by constrained least-squares

$$\hat{f} \in \operatorname{argmin}_{f \in \mathcal{F}} \left\{ \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i))^2 \right\},$$

where \mathcal{F} is a suitably chosen subset of functions.

- ▶ Typically, we choose \mathcal{F} to be a compact subset of some ambient function class \mathcal{G} , for example, a ball of radius R in some norm $\|\cdot\|_{\mathcal{G}}$.
- ▶ For computational reasons, it can be convenient to use regularized estimators of the form

$$\hat{f} \in \operatorname{argmin}_{f \in \mathcal{G}} \left\{ \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda_n \|f\|_{\mathcal{G}}^2 \right\}.$$

Example: Linear regression

- ▶ For a given vector $\theta \in \mathbb{R}^d$, define the function

$$f_{\theta}(x) = \langle \theta, x \rangle.$$

- ▶ For a compact set $\mathcal{C} \subset \mathbb{R}^d$, define

$$\mathcal{F}_{\mathcal{C}} = \{f_{\theta} : \mathbb{R}^d \rightarrow \mathbb{R} \mid \theta \in \mathcal{C}\}.$$

- ▶ Constrained least-square:

$$\hat{\theta} \in \operatorname{argmin}_{\theta \in \mathcal{C}} \left\{ \frac{1}{2n} \|y - X\theta\|_2^2 \right\}.$$

- ▶ Examples:

- ▶ *Ridge regression*: $\mathcal{C} = \{\theta \in \mathbb{R}^d \mid \|\theta\|_2^2 \leq R_2\}.$
- ▶ *Lasso*: $\mathcal{C} = \{\theta \in \mathbb{R}^d \mid \|\theta\|_1 \leq R_1\}.$

Example: Cubic smoothing spline

- ▶ For a given radius $R > 0$, consider the class of twice continuously differentiable functions $f : [0, 1] \rightarrow \mathbb{R}$,

$$\mathcal{F}(R) := \left\{ f : [0, 1] \rightarrow \mathbb{R} \mid \int_0^1 (f''(x))^2 dx \leq R \right\}.$$

- ▶ This constraint can be understood as a Hilbert norm bound in a second-order Sobolev space.
- ▶ For this function class, the penalized non-parametric least squares estimate is given by

$$\hat{f} \in \operatorname{argmin}_f \left\{ \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda_n \int_0^1 (f''(x))^2 dx \right\}.$$

- ▶ It can be shown that any minimizer f is a cubic spline.
- ▶ In the limit as $R \rightarrow 0$, the cubic spline fit \hat{f} becomes a linear function.

Example: Kernel ridge regression

- ▶ Let \mathbb{H} be a Hilbert space, equipped with norm $\|\cdot\|_{\mathbb{H}}$.
- ▶ For some radius $R > 0$, consider the constrained least-square estimator

$$\hat{f} \in \operatorname{argmin}_{\|f\|_{\mathbb{H}} \leq R} \left\{ \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i))^2 \right\}.$$

- ▶ In practice, its dual form, the penalized least-square estimator is commonly used

$$\hat{f} \in \operatorname{argmin}_{f \in \mathbb{H}} \left\{ \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda_n \|f\|_{\mathbb{H}}^2 \right\}.$$

- ▶ In particular, we assume \mathbb{H} to be a *Reproducing kernel Hilbert space* (RKHS).

Reproducing kernel Hilbert space (RKHS)

Definition

A reproducing kernel Hilbert space \mathbb{H} is a Hilbert space of real-valued functions on \mathcal{X} such that for each $x \in \mathcal{X}$, the evaluation functional $L_x : \mathbb{H} \rightarrow \mathbb{R}, f \mapsto f(x)$ is bounded.

- ▶ When L_x is a bounded linear functional, the Riesz representation implies that there must exist some element R_x of the Hilbert space \mathbb{H} such that

$$f(x) = L_x(f) = \langle f, R_x \rangle_{\mathbb{H}} \quad \text{for all } f \in \mathbb{H}.$$

- ▶ This element R_x of \mathbb{H} is known as the representer of evaluation at $x \in \mathcal{X}$.
- ▶ The boundedness of R_x ensures that convergence of a sequence of functions in an RKHS implies pointwise convergence.

Examples

- ▶ The space of all linear functions $f_\beta(\cdot) = \langle \cdot, \beta \rangle$ over \mathbb{R}^m under the inner product

$$\langle f_\beta, f_{\beta'} \rangle_{\mathbb{H}} = \langle \beta, \beta' \rangle.$$

is an RKHS, whose representer of evaluation R_x is the function $R_x(z) = \langle x, z \rangle$.

- ▶ The space $\mathcal{L}^2[0, 1]$ is not an RKHS.

- ▶ The first order Sobolev space $\mathbb{H}^1[0, 1] =$

$$\{f : f(0) = 0, f \text{ is absolutely continuous with } f' \in L^2[0, 1]\}$$

is an RKHS under the inner product

$$\langle f_1, f_2 \rangle_{\mathbb{H}} = \int_0^1 f_1'(x) f_2'(x) dx,$$

whose representer of evaluation R_x is the function $R_x(z) = \min\{x, z\}$.

Examples

- ▶ More generally, consider the higher-order Sobolev space $\mathbb{H}^\alpha[0, 1]$ of real-valued functions on $[0, 1]$ that are α -times differentiable (almost everywhere), with $f^{(\alpha)}$ being Lebesgue integrable, and $f(0) = \dots = f^{(\alpha)} = 0$.
- ▶ Define the inner-product

$$\langle f_1, f_2 \rangle_{\mathbb{H}} = \int_0^1 f_1^{(\alpha)}(x) f_2^{(\alpha)}(x) dx.$$

- ▶ $\mathbb{H}^\alpha[0, 1]$ is an RKHS, and the representer of evaluation is

$$R_x(y) = \int_0^1 \frac{(x-z)_+^{\alpha-1}}{(\alpha-1)!} \frac{(y-z)_+^{\alpha-1}}{(\alpha-1)!} dz.$$

- ▶ This can be seen from the Taylor expansion formula

$$f(x) = \sum_{\ell=0}^{\alpha-1} f^{(\ell)}(0) \frac{x^\ell}{\ell!} + \int_0^1 f^{(\alpha)}(z) \frac{(x-z)_+^{\alpha-1}}{(\alpha-1)!} dz.$$

Kernel functions

Definition

Positive semidefinite kernel function A symmetric bivariate function $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty)$ is positive semidefinite (PSD) if for all integers $n \geq 1$ and elements $\{x_i\}_{i=1}^n \subset \mathcal{X}$, the $n \times n$ matrix K with elements $K_{ij} := K(x_i, x_j)$ is positive semidefinite.

PSD kernel from representer of evaluation

- ▶ Define \mathcal{K} via $\mathcal{K}(x, x') = \langle R_x, R_{x'} \rangle_{\mathbb{H}}$.
- ▶ \mathcal{K} is symmetric.
- ▶ \mathcal{K} is PSD: for any vector $\alpha \in \mathbb{R}^n$, we have

$$\begin{aligned}\alpha^T K \alpha &= \sum_{j,k} \alpha_j \alpha_k \mathcal{K}(x_j, x_k) = \left\langle \sum_j \alpha_j R_{x_j}, \sum_k \alpha_k R_{x_k} \right\rangle_{\mathbb{H}} \\ &= \left\| \sum_j \alpha_j R_{x_j} \right\|_{\mathbb{H}}^2 \geq 0.\end{aligned}$$

Reproducing kernel

- ▶ Indeed, for any $x \in \mathbb{H}$, the object $\mathcal{K}(\cdot, x)$ can be identified with R_x as an element of our Hilbert space, since

$$\mathcal{K}(x', x) = \langle R_{x'}, R_x \rangle_{\mathbb{H}} = R_x(x'), \quad \text{for all } x \in \mathcal{X}.$$

- ▶ Consequently, the function $\mathcal{K}(\cdot, x)$ satisfies the *reproducing kernel property*, namely,

$$\langle \mathcal{K}(\cdot, x), f \rangle_{\mathbb{H}} = f(x), \quad \text{for all } f \in \mathbb{H}.$$

- ▶ Conversely, given a positive semidefinite kernel function \mathcal{K} , we can define an associated function space

$$\widetilde{\mathbb{H}} := \left\{ f(\cdot) = \sum_{i=1}^n \alpha_i \mathcal{K}(\cdot, x_i) : n \in \mathbb{N} \text{ and } \{x_i\}_{i=1}^n \subset \mathcal{X} \right\}.$$

Reproducing kernel

- ▶ Given two functions $f, \bar{f} \in \widetilde{\mathbb{H}}$, where $\bar{f} = \sum_{i=1}^{\bar{n}} \bar{\alpha}_i \mathcal{K}(\cdot, \bar{x}_i)$, their inner product is defined by

$$\langle f, \bar{f} \rangle_{\widetilde{\mathbb{H}}} = \sum_{j=1}^n \sum_{k=1}^{\bar{n}} \alpha_j \bar{\alpha}_k \mathcal{K}(x_j, \bar{x}_k).$$

- ▶ The PSD property of \mathcal{K} implies $\|f\|_{\widetilde{\mathbb{H}}} \geq 0$, for all $f \in \widetilde{\mathbb{H}}$.
- ▶ We can define a Hilbert space \mathbb{H} as the completion of $\widetilde{\mathbb{H}}$ with respect to the inner product $\|\cdot\|_{\widetilde{\mathbb{H}}}$.

Theorem

Given any reproducing kernel Hilbert space \mathbb{H} of functions $f : \mathcal{X} \rightarrow \mathbb{R}$, there exists a unique positive semidefinite kernel function $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. Conversely, given any positive semidefinite kernel function \mathcal{K} , we can define an RKHS in which \mathcal{K} acts as the representer of evaluation.

Example: Kernel ridge regression, continued

- Recall the KRR estimator

$$\hat{f} \in \operatorname{argmin}_{f \in \mathbb{H}} \left\{ \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda_n \|f\|_{\mathbb{H}}^2 \right\}.$$

- Let $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ denote the reproducing kernel associated with the RKHS \mathbb{H} .

Theorem (Representer theorem)

Any solution \hat{f} of the KRR optimization problem takes the form

$$\hat{f}(\cdot) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\alpha}_i \mathcal{K}(\cdot, x_i).$$

Example: Kernel ridge regression, continued

- Define the empirical kernel matrix $K \in \mathbb{R}^{n \times n}$, with $K_{ij} = n^{-1} \mathcal{K}(x_i, x_j)$, and recall

$$\hat{f}(\cdot) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\alpha}_i \mathcal{K}(\cdot, x_i).$$

- Then, we can write

$$(\hat{f}(x_1), \dots, \hat{f}(x_n))^T = \sqrt{n} K \hat{\alpha},$$

where $\hat{\alpha} = (\hat{\alpha}_1, \dots, \hat{\alpha}_n)^T$.

- Solving the KRR optimization problem is inequivalent to solving the following quadratic programming

$$\hat{\alpha} \in \operatorname{argmin}_{\alpha \in \mathbb{R}^n} \left\{ \frac{1}{2n} \|y - \sqrt{n} K \alpha\|_2^2 + \lambda_n \underbrace{\alpha^T K \alpha}_{\|f\|_{\mathbb{H}}^2} \right\}.$$

Example: Convex regression

- ▶ Now suppose $f^* : \mathcal{C} \rightarrow \mathbb{R}$ is known to be a convex function over its domain \mathcal{C} , where \mathcal{C} is some convex and open subset of \mathbb{R}^d .
- ▶ It is natural to consider the least-squares estimator with a convexity constraint,

$$\hat{f} \in \operatorname{argmin}_{f \text{ is convex}} \left\{ \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i))^2 \right\}.$$

- ▶ Although this optimization problem is infinite-dimensional, we can convert it to an equivalent finite-dimensional problem.
- ▶ The convexity constraint implies there exist sub-gradient vectors $\{\tilde{z}_i\}_{i=1}^n$, such that for all $i = 1, \dots, n$,

$$f(x) \geq f(x_i) + \langle \tilde{z}_i, x - x_i \rangle \quad \text{for all } x \in \mathcal{C}.$$

Example: Convex regression

- ▶ Since the cost function depends only on the values $\tilde{y}_i = f(x_i)$, the optimum does not depend on the function behavior elsewhere.
- ▶ It suffices to solve the optimization problem

$$\min_{\{(\tilde{y}_i, \tilde{z}_i)\}_{i=1}^n} \frac{1}{2n} \sum_{i=1}^n (y_i - \tilde{y}_i)^2$$

such that $\tilde{y}_j \geq \tilde{y}_i + \langle \tilde{z}_i, x_j - x_i \rangle$ for all $i, j = 1, \dots, n$.

- ▶ An optimal solution $\{(\hat{y}_i, \hat{z}_i)\}_{i=1}^n$ can be used to define an estimate $\hat{f} : \mathcal{C} \rightarrow \mathbb{R}$ via

$$\hat{f}(x) = \max_{i=1, \dots, n} \{\hat{y}_i + \langle \hat{z}_i, x - x_i \rangle\}.$$

- ▶ \hat{f} is convex, and by the feasibility of the solution $\{(\hat{y}_i, \hat{z}_i)\}_{i=1}^n$, we are guaranteed that $\hat{f}(x_i) = \hat{y}_i$.