

Nonparametric Least Square Estimator

Luo Fangzhi

Sun Yat-sen University

Examples of Nonparametric Least Square

Kernel Ridge Regression : Let H be a reproducing kernel Hilbert space, equipped with the norm $\|\cdot\|_H$.

A constraint nonparametric least square estimator is by minimizing

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

subject to $\|f\|_{\mathbb{H}} < C$. Under some condition, this is equivalent to solve

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

for a corresponding regularization parameter $\lambda_n > 0$. As mentioned in last chapter, if the kernel of \mathbb{H} is \mathcal{K} , and the corresponding kernel function is $K(s, t)$, then the solution to f is $\hat{f}(\cdot) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\alpha}_i \mathcal{K}(\cdot, x_i)$, where $\hat{\alpha} := (\mathbf{K} + \lambda_n \mathbf{I}_n)^{-1} \frac{\mathbf{y}}{\sqrt{n}}$.

Problem Setup

- obs : (x_i, y_i)
- model : $y_i = f^*(x_i) + v_i, v_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$

We want to evaluate f^* by

$$\operatorname{argmin}_{f \in \mathcal{F}} \sum_i (y_i - f(x_i))^2. \quad (2)$$

The first intuition is the difficulty of evaluate f^* is from the complexity of \mathcal{F} .

The quality of evaluation is estimated by

$$\frac{1}{n} \sum_i (\hat{f}(x_i) - f^*(x_i))^2. \quad (3)$$

Remark that this is evaluation converge to $\mathbb{E}_{X,Y} [\hat{f}(x) - f^*(x)]^2$.

The Nonasymptotic Bound

Theorem

for any c , $\exists \delta_c$ such that, $\forall u > \delta_c$, $\frac{\mathbb{E}\mathcal{G}(u, \mathcal{F}^*)}{u} < c$. Then $\forall c_1$ subject to $2\sigma(c + c_1) > \delta_c$,

$$\mathbb{P}[\|\hat{f} - f^*\|_n > 2\sigma(c + c_1)] \leq \exp(-nc_1^2) \quad (4)$$

The Nonasymptotic Bound

Our goal is to establish a nonasymptotic bound of $\sum_i [\hat{f}(x_i) - f^*(x_i)]^2$. First tool is the basic inequality

$$\sum_i [y_i - f^*(x_i)]^2 \geq \sum_i [y_i - \hat{f}(x_i)]^2. \quad (5)$$

What follows can be obtained by simple transformation of the above inequality.

$$\frac{1}{2n} \sum_i [\hat{f}(x_i) - f^*(x_i)]^2 \leq \frac{\sigma}{n} \sum_i \omega_i [\hat{f}(x_i) - f^*(x_i)]. \quad (6)$$

Or equivalently,

$$\sum_i [\hat{f}(x_i) - f^*(x_i)]^2 \leq 2\sigma \sum_i \omega_i [\hat{f}(x_i) - f^*(x_i)]. \quad (7)$$

The right handside is very similar to a sub-Gaussian. But some observation is the right hand side is not L-Lipschitz continuous. So we cannot simple apply Thm 2.26.

The Nonasymptotic Bound

Moreover, if we let $\mathcal{G}(\delta, \mathcal{F}^*) = \sup_{g \in \mathcal{F}^*, \|g\|_n < \delta} \frac{1}{n} \sum_i \omega_i g(x_i)$, where $\mathcal{F}^* = \mathcal{F} - f^*$, $\|g\|_n = (\frac{1}{n} \sum_i g(x_i)^2)^{1/2}$. Then

$$\frac{1}{n} \sum_i [\hat{f}(x_i) - f^*(x_i)]^2 \leq \frac{2\sigma}{n} \sum_i \omega_i [\hat{f}(x_i) - f^*(x_i)] \leq 2\sigma \lim_{\delta \rightarrow \infty} \mathcal{G}(\delta, \mathcal{F}^*). \quad (8)$$

An important observation is $\mathcal{G}(t, \mathcal{F}^*)$ is increasing linearly or no more faster linearly with respect to a constant $C(\delta)$ for any $t > \delta$ if \mathcal{F}^* is star-shaped. Specifically,

$$\frac{\mathcal{G}(\delta, \mathcal{F}^*)}{\delta} \geq \frac{\mathcal{G}(u, \mathcal{F}^*)}{u}, \quad \frac{\mathbb{E}\mathcal{G}(\delta, \mathcal{F}^*)}{\delta} \geq \frac{\mathbb{E}\mathcal{G}(u, \mathcal{F}^*)}{u}, \quad \forall \delta < u. \quad (9)$$

The Nonasymptotic Bound

Let $Z_\delta(r) = \sup_{g \in \mathcal{F}^*, \|g\|_n < \delta} \frac{1}{n} \sum_i r_i g(x_i)$. For fixed δ , $Z_\delta(r)$ is L-Lipschitz continuous as a function of r . The textbook claims that the Lipschitz constant is at most $\frac{\delta}{\sqrt{n}}$. **But i can't prove this.** Consequently,

$$\mathbb{P}[\mathcal{G}(\delta, \mathcal{F}^*) - \mathbb{E}(\mathcal{G}(\delta, \mathcal{F}^*)) > t] \leq \exp(-\frac{nt^2}{\delta^2}), \forall t \quad (10)$$

for any $\|\hat{f} - f^*\|_n = \delta$,

$$\delta^2 \leq 2\sigma\mathcal{G}(\delta, \mathcal{F}^*) \leq 2\sigma(\mathbb{E}\mathcal{G}(\delta, \mathcal{F}^*) + t) \quad (11)$$

with prob $(1 - \exp(-\frac{nt^2}{\delta^2}))$.

So with prob $(1 - \exp(-\frac{nt^2}{\delta^2}))$,

$$\delta \leq 2\sigma\left(\frac{\mathbb{E}\mathcal{G}(\delta, \mathcal{F}^*)}{\delta} + \frac{t}{\delta}\right). \quad (12)$$

The Nonasymptotic Bound

Fortunately, for any c , $\exists \delta_c$, $\forall u > \delta_c$, $\frac{\mathbb{E}\mathcal{G}(u, \mathcal{F}^*)}{u} < c$, so $\forall \delta > \delta_c$,

$$\delta \leq 2\sigma\left(\frac{\mathbb{E}\mathcal{G}(\delta, \mathcal{F}^*)}{\delta} + \frac{t}{\delta}\right) \leq 2\sigma\left(c + \frac{t}{\delta}\right). \quad (13)$$

with prob $(1 - \exp(-\frac{nt^2}{\delta^2}))$. Let $t = c_1\delta$, then

$$\delta \leq 2\sigma(c + c_1). \quad (14)$$

with prob $(1 - \exp(-nc_1^2))$.

Determination of δ_c

The problem remained is the quantity of δ_c . It's related to \mathcal{F}^* . If δ_c is very big, we can only bound δ when δ_c is very big, then the bound is meaningless. Specifically we let $c = \frac{u}{2\sigma}$.

We want to know, for $c = \frac{u}{2\sigma}$, which δ_c is valid subject to $\forall u > \delta_c$, $\frac{\mathbb{E}\mathcal{G}(u, \mathcal{F}^*)}{u} < \frac{u}{2\sigma}$. Remark that here 2σ is not important.

Determination of δ_c

For any given $\delta \in (0, \sigma]$, let set $\mathbb{B}(\delta, \mathcal{F}^*) = \{g | g \in \mathcal{F}^*, \|g\|_n < \delta\}$. The minimal $\frac{\delta^2}{4\sigma}$ -covering of $\mathbb{B}(\delta, \mathcal{F}^*)$ is finite and is assumed to consist of $\{g^1, \dots, g^M\}$.

By the property of covering set, for any $g \in \mathbb{B}(\delta, \mathcal{F}^*)$ we have

$$\begin{aligned}
 \left| \frac{1}{n} \sum_i \omega_i g(x_i) \right| &\leq \left| \frac{1}{n} \sum_i \omega_i g^j(x_i) \right| + \left| \frac{1}{n} \sum_i \omega_i (g(x_i) - g^j(x_i)) \right| \\
 &\leq \max_j \left| \frac{1}{n} \sum_i \omega_i g^j(x_i) \right| + \sqrt{\frac{\sum_i \omega_i^2}{n}} \sqrt{\frac{\sum_i (g(x_i) - g^j(x_i))^2}{n}} \\
 &\leq \max_j \left| \frac{1}{n} \sum_i \omega_i g^j(x_i) \right| + \sqrt{\frac{\sum_i \omega_i^2}{n}} \frac{\delta^2}{4\sigma}
 \end{aligned} \tag{15}$$

Determination of δ_c

Then

$$\begin{aligned}
 \sup_{g \in \mathcal{F}^*} \{ \mathbb{E} | \frac{1}{n} \sum_i \omega_i g(x_i) | \} &\leq \mathbb{E} \mathcal{G}(\delta, \mathcal{F}^*) \\
 &\leq \mathbb{E} \max_j | \frac{1}{n} \sum_i \omega_i g^j(x_i) | + \mathbb{E} \sqrt{\frac{\sum_i \omega_i^2}{n}} \frac{\delta^2}{4\sigma} \\
 &\leq \mathbb{E} \max_j | \frac{1}{n} \sum_i \omega_i g^j(x_i) | + \frac{\delta^2}{4\sigma}
 \end{aligned} \tag{16}$$

So $\frac{\mathbb{E} \mathcal{G}(u, \mathcal{F}^*)}{u} < \frac{u}{2\sigma}$ exists if $\mathbb{E} \max_j | \frac{1}{n} \sum_i \omega_i g^j(x_i) | < \frac{\delta^2}{4\sigma}$

Determination of δ_c

By definition, for any j , $\frac{1}{n} \sum_i \omega_i g^j(x_i)$ is sub-Gaussian with $L = \frac{\delta^2}{2\sqrt{n}\sigma}$. The theorem 5.22 can be applied.

Let $Z = \max_j \frac{1}{n} \sum_i \omega_i g^j(x_i)$,

$$\mathbb{E}Z \leq \frac{16}{\sqrt{n}} \int_{\frac{\delta^2}{4\sigma}}^{\delta} \sqrt{\log N_n(t; \mathbb{B}_n(\delta; \mathcal{F}^*))} dt. \quad (17)$$

So $\frac{\mathbb{E}\mathcal{G}(u, \mathcal{F}^*)}{u} < \frac{u}{2\sigma}$ exists if

$$\frac{16}{\sqrt{n}} \int_{\frac{\delta^2}{4\sigma}}^{\delta} \sqrt{\log N_n(t; \mathbb{B}_n(\delta; \mathcal{F}^*))} dt \leq \frac{\delta^2}{4\sigma} \quad (18)$$

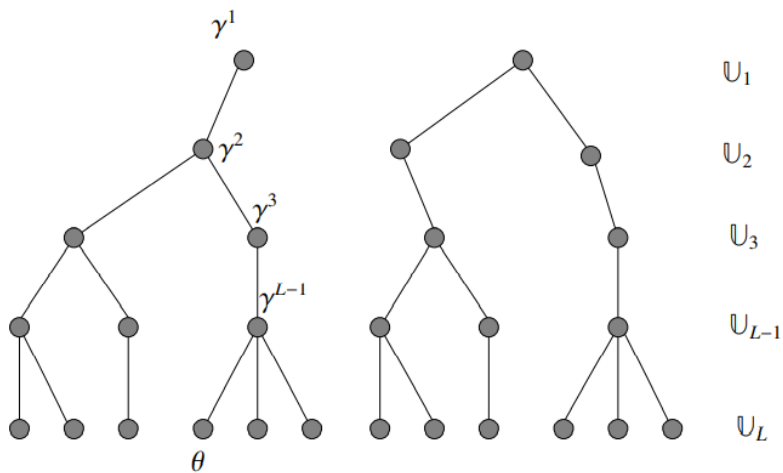
Let $\mathbb{U} = \{g^1, \dots, g^M\}$ let \mathbb{U}_m be a minimal $\epsilon_m = D2^{-m}$ covering set of U in the metric ρ_X , where we allow for any element of $\mathbb{B}(\delta, \mathcal{F}^*)$ to be used. Define the mapping $\pi_m : \mathbb{U} \rightarrow \mathbb{U}_m$ via

$$\pi_m(g) = \arg \min_{\beta \in \mathbb{U}_m} \rho_X(g, \beta),$$

Using this notation, we can decompose the random variable X_θ into a sum of increments in terms of an associated sequence $(\gamma^1, \dots, \gamma^L)$, where we define $\gamma^L = \theta$ and $\gamma^{m-1} := \pi_{m-1}(\gamma^m)$ recursively for $m = L, L-1, \dots, 2$. By construction, we then have the chaining relation

$$X_g - X_{\gamma^1} = \sum_{m=2}^L (X_{\gamma^m} - X_{\gamma^{m-1}})$$

and hence $|X_g - X_{\gamma^1}| \leq \sum_{m=2}^L \max_{\beta \in \mathbb{U}_m} |X_\beta - X_{\pi_{m-1}(\beta)}|$.



$|X_g - X_{\gamma^1}| \leq \sum_{m=2}^L \max_{\beta \in \mathbb{U}_m} |X_\beta - X_{\pi_{m-1}(\beta)}|$. From exercise 2.12 we know

$$\mathbb{E} \left[\max_{\gamma, \tilde{\gamma} \in \mathbb{U}_1} |X_\gamma - X_{\tilde{\gamma}}| \right] \leq 2D \sqrt{\log N(D/2)}.$$

Similarly, for each $m = 2, 3, \dots, L$, the set U_m has $N(D2^{-m})$ elements, and, moreover, $\max_{\beta \in \mathbb{U}_m} \rho_X(\beta, \pi_{m-1}(\beta)) \leq D2^{-(m-1)}$, whence

$$\mathbb{E} \left[\max_{\beta \in \mathbb{U}_m} |X_\beta - X_{\pi_{m-1}(\beta)}| \right] \leq 2D2^{-(m-1)} \sqrt{\log N(D2^{-m})}.$$

So

$$\mathbb{E} \left[\max_{g, \tilde{g} \in \mathbb{U}} |X_g - X_{\tilde{g}}| \right] \leq 4 \sum_{m=1}^L D2^{-(m-1)} \sqrt{\log N(D2^{-m})}.$$

Since the metric entropy $\log N(t)$ is non-increasing in t , we have

$$D2^{-(m-1)} \sqrt{\log N(D2^{-m})} \leq 4 \int_{D2^{-(m+1)}}^{D2^{-m}} \sqrt{\log N(u)} du,$$

Hence $2\mathbb{E} [\max_{g, \tilde{g} \in \mathbb{U}} |X_g - X_{\tilde{g}}|] \leq 32 \int_{\delta/4}^D \sqrt{\log N(u)} du$.