CS861: Theoretical Foundations of Machine Learning

Chapter 2: Nonparametric Methods

Kirthevasan Kandasamy UW-Madison

Introduction

Nonparametric methods: Learn in a function class which cannot be completely characterized by a finite number of scalar parameters.

Examples in regression:

- Nonparametric: the class of all Lipschitz-continuous functions, $\{f; |f(x) f(x')| \le L||x x'||_2\}.$
- ▶ Parametric: Linear class, $\{f; f(x) = \theta^{\top} x, \text{ for some } \theta \in \mathbb{R}^d\}$.

Contents

- 1. Nonparametric regression
- 2. Nonparametric density estimation

Ch 2.1: Nonparametric regression

Model: Let \mathcal{F} be the class of bounded L-Lipschitz functions in [0,1].

$$\mathcal{F} = \{f: [0,1] \to [0,1], |f(x_1) - f(x_2)| \le L|x_1 - x_2|\}$$

We observe an i.i.d dataset $\{(X_1, Y_1), \dots, (X_n, Y_n)\}$ drawn i.i.d from $P \in \mathcal{P}$ where,

$$\mathcal{P} = \bigg\{P; \quad 0 < \alpha_0 \leq p(x) \leq \alpha_1 < \infty, \text{ where } p(x) \text{ is the marginal of } x$$
 the regression function $f(\cdot) = \mathbb{E}[Y|X = \cdot] \in \mathcal{F},$
$$\operatorname{Var}(Y|X = x) \leq \sigma^2, \text{ for all } x \in [0,1] \bigg\}.$$

Nonparametric regression (cont'd)

We wish to estimate the regression function in the following loss. For any $g \in \mathbb{R}^{[0,1]}$,

$$L(g,P) = \int (f(x) - g(x))^2 p(x) dx$$
, where, $f(\cdot) = \mathbb{E}_P[Y|X = \cdot]$.

The minimax risk is,

$$R_n^{\star} = \inf_{\widehat{f}} \sup_{P \sim \mathcal{P}} \mathbb{E}_{S \sim P^n}[L(\widehat{f}(S), P)] = \inf_{\widehat{f}} \sup_{P \sim \mathcal{P}} \mathbb{E}_{S \sim P^n} \left[\int (\widehat{f}(S)(x) - f(x))^2 p(x) dx \right].$$

We will show that the minimax rate is $\Theta(n^{-2/3})$.

- Lower bound via Fano's method.
- Upper bound via the Nadaraya-Watson estimator.

Ch 2.1.1: Lower bounds for nonparametric regression

To apply our lower bound techniques, we cannot write the above loss in the form $\Phi \circ \rho$ where ρ is a metric Hence, let us consider $\mathcal{P}'' = \{P \in \mathcal{P}; p(x) = 1\} \subset \mathcal{P}$. Then, we can lower bound the minimax risk by

$$R_n^{\star} \geq \inf_{\widehat{f}} \sup_{P \subset \mathcal{P}''} \mathbb{E}_{S \sim P} \left[\underbrace{\int (f(x) - \widehat{f}(x))^2 dx}_{\Phi \circ \rho(f, \widehat{f})} \right]$$

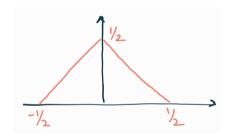
where, $\rho(f,g) = \|f - g\|_2$ and $\Phi(t) = t^2$.

Recall, Local Fano method. Let S be an i.i.d dataset from some distribution $P \in \mathcal{P}$. Let $\{P_1,\ldots,P_N\} \subset \mathcal{P}$ such that $N \geq 16$, $\rho(\theta(P_j),\theta(P_k)) \geq \delta$, and $\mathrm{KL}(P_j,P_k) \leq \log(N)/4n$ for all $j \neq k$. Then, $R_n^\star \geq \frac{1}{2}\Phi\left(\frac{\delta}{2}\right)$.

We will apply Fano's method in the following four steps: 1) Constructing alternatives. 2) Lower bound the minimum separation. 3) Upper bound the maximum KL. 4) Obtain final bound.

Step 1: Constructing Alternatives

Let us consider the following function,



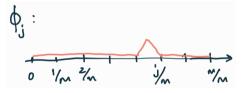
$$\psi(x) = \begin{cases} x + 1/2, & x \in [-1/2, 0], \\ -x + 1/2, & x \in [0, 1/2], \\ 0 & \text{otherwise.} \end{cases}$$

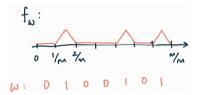
Note that, ψ is 1-Lipschitz, and $\int \psi^2 = 1/12$.

Step 1: Constructing Alternatives (cont'd)

Let h > 0 be a parameter to be chosen later. Let m = 1/h. Let $\omega \in \{0,1\}^m$. Let us define f_ω as follows:

$$\phi_j(x) = Lh\psi\left(\frac{x-(j-1/2)h}{h}\right), \qquad f_\omega(x) = \sum_{i=1}^m \omega_i\phi_i(x).$$





We will show that each f_{ω} is L-Lipschitz:

- f_{ω} is 0 at the boundaries (1/m, 2/m, ...).
- It is sufficient to check Lipschitzness within each bump ϕ_i .

- The gradient of
$$\phi_j$$
 is $|\phi_j'(t)| = \left| Lh\psi'\left(\frac{x-(j-1/2)h}{h}\right) \frac{1}{h} \right| = L$.

Step 1: Constructing Alternatives (cont'd)

Alternatives: Let Ω_m be the Gilbert-Varshamov-pruned m-hypercube. Define

$$\begin{split} \mathcal{F}' &= \{f_{\omega}; \omega \in \Omega_m\}, \\ \mathcal{P}' &= \{P; p(x) \text{ is uniform }, \ f(\cdot) = \mathbb{E}[Y|X = \cdot] \in \mathcal{F}', \\ Y|X &= x \sim \mathcal{N}(f(x), \sigma^2)\}. \end{split}$$

As $\mathcal{F}' \subset \mathcal{F}$, we have $\mathcal{P}' \subset \mathcal{P}'' \subset \mathcal{P}$.

Let P_{ω} be the distribution in \mathcal{P}' whose regression function is f_{ω} . Let p_{ω} be its density.

Step 2: Lower bounding the separation

For any $\omega, \omega' \in \Omega_m$, we have

$$\rho^{2}(f_{\omega}, f_{\omega'}) = \int_{0}^{1} (f_{\omega} - f_{\omega'})^{2}$$

$$= \sum_{j=1}^{m} \int_{\frac{j-1}{m}}^{\frac{j}{m}} (\omega_{j} \phi_{j} - \omega'_{j} \phi_{j})^{2}$$

$$= \sum_{j=1}^{m} \mathbb{1}(\omega_{j} \neq \omega'_{j}) \int_{\frac{j-1}{m}}^{\frac{j}{m}} \phi_{j}^{2}$$

Now, let us consider the integral,

$$\int_{\frac{j-1}{m}}^{\frac{j}{m}} \phi_j^2 = \int_{\frac{j-1}{m}}^{\frac{j}{m}} L^2 h^2 \psi^2 \left(\frac{x - (j-1/2)h}{h} \right) dx = \int_{-1/2}^{1/2} L^2 h^3 \psi^2(u) du = \frac{L^2 h^3}{12}.$$

Step 2: Lower bounding the separation (cont'd)

Recall, Gilbert-Varshamov bound. In the Gilbert-Varshamov pruned m-hypercube Ω_m of $\{0,1\}^m$, we have, (i) $|\Omega_m| \geq 2^{m/8}$, (ii) $H(\omega,\omega') \geq m/8$, $\forall \omega,\omega' \in \Omega_m$, (iii) $\mathbf{0}_m \in \Omega_m$.

Therefore,

$$\rho^2(f_{\omega},f_{\omega'})=H(\omega,\omega')\frac{L^2h^3}{12}.$$

By the Gilbert-Varshamov bound, $H(\omega, \omega') \ge \frac{m}{8} = \frac{1}{8h}$.

Therefore, the separation can be lower bounded by,

$$\min_{\omega,\omega'} \rho(f_{\omega}, f_{\omega'}) \geq \frac{Lh}{\sqrt{96}} \stackrel{\Delta}{=} \delta$$

٠

Step 3: Upper bounding the KL

Let $p_{\omega}, p_{\omega'}$ be the joint densities for any $\omega, \omega' \in \Omega_m$, respectively. Then,

$$\begin{split} \operatorname{KL}(P_{\omega},P_{\omega'}) &= \int_{\mathcal{X}\times\mathcal{Y}} p_{\omega} \log\left(\frac{p_{\omega}}{p_{\omega'}}\right) \\ &= \int_{0}^{1} \int_{-\infty}^{\infty} \overbrace{p_{\omega}(x)}^{=1} p_{\omega}(y|x) \log\left(\underbrace{\frac{p_{\omega}(x)}{p_{\omega'}(x)} p_{\omega}(y|x)}_{p_{\omega'}(x)} \right) \mathrm{d}y \mathrm{d}x \\ &= \int_{0}^{1} \operatorname{KL}\left(\mathcal{N}(f_{\omega}(x),\sigma^{2}), \mathcal{N}(f_{\omega'}(x),\sigma^{2})\right) \mathrm{d}x = \frac{1}{2\sigma^{2}} \int_{0}^{1} \left(f_{\omega}(x) - f_{\omega'}(x)\right)^{2} \mathrm{d}x \\ &= \frac{1}{2\sigma^{2}} \sum_{j=1}^{m} \mathbb{1}(\omega_{j} \neq \omega'_{j}) \int_{\frac{j-1}{m}}^{\frac{j}{m}} \left(\omega_{j}\phi_{j} - \omega'_{j}\phi_{j}\right)^{2} \\ &= \frac{H(\omega,\omega')}{2\sigma^{2}} \cdot \frac{L^{2}h^{3}}{12}. \quad \text{Similar calculations to separation.} \end{split}$$

Step 3: Upper bounding the KL (cont'd)

As
$$\max_{\omega,\omega'} H(\omega,\omega') \leq m = 1/h$$
, we have

$$\max_{\omega,\omega'} \mathrm{KL}(P_{\omega},P_{\omega'}) \leq \frac{L^2 h^2}{24\sigma^2}.$$

Step 4: Obtain final bound

Recall, Local Fano method. Let S be an i.i.d dataset from some distribution $P \in \mathcal{P}$. Let $\{P_1,\ldots,P_N\} \subset \mathcal{P}$ such that $N \geq 16$, $\delta \geq \rho(\theta(P_j),\theta(P_k))$, and $\mathrm{KL}(P_j,P_k) \leq \log(N)/4n$ for all $j \neq k$. Then, $R_n^\star \geq \frac{1}{2}\Phi\left(\frac{\delta}{2}\right)$.

Recall, Gilbert-Varshamov bound. In the Gilbert-Varshamov pruned *m*-hypercube Ω_m of $\{0,1\}^m$, we have, (i) $|\Omega_m| \geq 2^{m/8}$, (ii) $H(\omega,\omega') \geq m/8$, $\forall \, \omega,\omega' \in \Omega_m$, (iii) $\mathbf{0}_m \in \Omega_m$.

We need $\max_{p_{\omega},p_{\omega'}} \mathrm{KL}(P_{\omega},P_{\omega'}) \leq \frac{\log(|\mathcal{P}'|)}{4n}$. As $|\mathcal{P}'| \geq 2^{m/8}$, it is sufficient if,

$$\frac{L^2h^2}{24\sigma^2} \le \frac{\log(2^{m/8})}{4n} = \frac{\log(2)}{32nh}.$$

Therefore, choose h as follows,

$$h = \left(\frac{3\log(2)}{4}\right)^{1/3} \frac{\sigma^{2/3}}{n^{1/3}L^{2/3}}.$$

Step 4: Obtain final bound (cont'd)

This gives the following δ ,

$$\delta = \frac{Lh}{\sqrt{96}} = C_1 \frac{L^{1/3} \sigma^{2/3}}{n^{1/3}}.$$

By the local Fano method, we have obtain the following lower bound on the minimax risk,

$$R_n^{\star} \geq \frac{1}{2} \Phi\left(\frac{\delta}{2}\right) = C_2 \frac{L^{2/3} \sigma^{4/3}}{n^{2/3}}.$$

We require $N \ge 16$, which is satisfied if we have a sufficient number of samples,

$$N \ge 16 \iff 2^{m/8} \ge 16 \iff \frac{m}{8} = \frac{1}{8h} \ge 4 \iff h = C_3 \frac{\sigma^{2/3}}{n^{1/3} L^{2/3}} \le \frac{1}{32} \iff n \ge C_4 \frac{\sigma^2}{L^2}.$$

Ch 2.2: Upper bound

We will define the following estimator \hat{f} . Let h > 0 be a bandwidth parameter whose value we will specify shortly.

$$N(t) = \sum_{i=1}^{n} \mathbb{1}(X_i \in (t-h, t+h)),$$

$$\overline{f}(t) = \begin{cases} \frac{1}{2} & \text{if } N(t) = 0, \\ \frac{1}{N(t)} \sum_{i=1}^{n} Y_i \mathbb{1}(X_i \in (t-h, t+h)) & \text{if } N(t) > 0. \end{cases}$$

$$\widehat{f}(t) = \text{clip}(\overline{f}(t), 0, 1).$$

Good events

To bound err_t define $G_t = \{N(t) \geq \alpha_0 nh\}$. to be the "good event" in which a sufficient number of points fall within the [t-h,t+h] interval. We have, by Hoeffding's inequality,

$$\mathbb{P}(G_t^c) = \mathbb{P}\left(\sum_{i=1}^n \mathbb{1}(X_i \in (t-h, t+h) < \alpha_0 nh)\right)$$

$$\leq \mathbb{P}\left(\sum_{i=1}^n (\mathbb{1}(X_i \in (t-h, t+h) - P([t-h, t+h])) < -\alpha_0 nh)\right)$$
as $P([t-h, t+h]) = \int_{t-h}^{t+h} p(t) dt \geq 2\alpha_0 h$.
$$\leq \exp(-2\alpha_0^2 nh^2). \quad \text{Using Hoeffding's}$$
 (1)

Bounding the Risk

To analyze this estimator, we will first upper bound the risk as follows,

$$egin{aligned} R(P,\widehat{f}) &= \mathbb{E}\left[\int_0^1 (\widehat{f}(t) - f(t))^2 p(t) \mathrm{d}t
ight] \ &\leq lpha_1 \, \mathbb{E}\left[\int_0^1 (\widehat{f}(t) - f(t))^2 \mathrm{d}t
ight] \qquad ext{As } p(t) \leq lpha_1. \ &= lpha_1 \int_0^1 \underbrace{\mathbb{E}\left[(\widehat{f}(t) - f(t))^2
ight]}_{\mathrm{err}_t(\widehat{f})} \mathrm{d}t \end{aligned}$$

Let us now write,

$$\begin{split} \operatorname{err}_t(\widehat{f}) &= \mathbb{E}\left[(\widehat{f}(t) - f(t))^2\right] \\ &= \mathbb{E}\left[(\widehat{f}(t) - f(t))^2 \middle| G_t\right] \underbrace{\mathbb{P}(G_t)}_{\leq 1} + \underbrace{\mathbb{E}\left[(\widehat{f}(t) - f(t))^2 \middle| G_t^c\right]}_{\leq 1 \text{ as } f \text{ is bounded in } [0, 1]} \underbrace{\mathbb{P}(G_t^c)}_{< e^{-2\alpha_0^2 nh^2}} \end{split}$$

Bounding the error under the good event

To upper bound $\mathbb{E}[(\widehat{f}(t) - f(t))^2 | G_t]$, let us denote $A_i = \mathbb{I}(X_i \in (t - h, t + h))$ and expand $(\widehat{f}(t) - f(t))^2$ as follows:

$$\begin{split} \left(\widehat{f}(t) - f(t)\right)^2 &\leq \left(\overline{f}(t) - f(t)\right)^2 & \text{As } \widehat{f}(t) = \operatorname{clip}(\overline{f}(t), 0, 1) \\ &= \left(\frac{1}{N(t)} \sum_{i=1}^n \left(A_i Y_i - f(t)\right)\right)^2 \\ &= \left(\underbrace{\frac{1}{N(t)} \sum_{i=1}^n A_i \left(Y_i - f(X_i)\right)}_{v} + \underbrace{\frac{1}{N(t)} \sum_{i=1}^n A_i \left(f(X_i) - f(t)\right)}_{b}\right)^2 \\ &= v^2 + b^2 + 2bv. \end{split}$$

Bounding the variance

Let us decompose the variance v at t as follows:

$$\mathbb{E}[v^{2}|G_{t}] = \mathbb{E}\left[\mathbb{E}\left[\left(\frac{1}{N(t)}\sum_{i=1}^{n}A_{i}(Y_{i}-f(X_{i}))\right)^{2}\middle|G_{t},X_{1},\ldots,X_{n}\right]\right]$$

$$\stackrel{(a)}{=}\mathbb{E}\left[\mathbb{E}\left[\frac{1}{N(t)^{2}}\sum_{i=1}^{n}(A_{i}Y_{i}-f(X_{i}))^{2}\middle|G_{t},X_{1},\ldots,X_{n}\right]\right]$$

$$=\mathbb{E}\left[\frac{1}{N(t)^{2}}\sum_{i=1}^{n}A_{i}\operatorname{Var}(Y_{i}|X_{i})\middle|G_{t}\right] \qquad \text{Note } \mathbb{E}[Y_{i}|X_{i}] = f(X_{i})$$

$$\leq \mathbb{E}\left[\frac{\sigma^{2}N(t)}{N(t)^{2}}\middle|G_{t}\right] \leq \frac{\sigma^{2}}{\alpha_{0}nh}. \qquad \text{As } N(t) \geq \alpha_{0}nh \text{ under } G_{t}.$$

Above, conditioning on X_1, \ldots, X_n allows us to claim that the cross-terms are 0 in (a).

Bounding the bias

Next, let us consider the bias b. For any $X_i \in (t - h, t + h)$, we have $|f(X_i) - f(t)| \le L|X_i - t| \le Lh$.

Hence,

$$|b| = \left| \frac{1}{N(t)} \sum_{i=1}^n A_i \left(f(X_i) - f(t) \right) \right| \le \frac{1}{N(t)} \sum_{i=1}^n A_i \left| f(X_i) - f(t) \right|$$
 $\le \frac{1}{N(t)} \sum_{i=1}^n A_i Lh = Lh.$

Therefore, $\mathbb{E}[b^2|G_t] \leq L^2h^2$.

Cross-term

Finally, let us consider the cross-term. As $\mathbb{E}[Y_i|X_i] = f(X_i)$, we have,

$$\mathbb{E}[bv|G] = \mathbb{E}\left[b \cdot \mathbb{E}_{Y}\left[\frac{1}{N(t)}\sum_{i=1}^{n}A_{i}\left(Y_{i} - f(X_{i})\right)\middle|G_{t}, X_{1}, \ldots, X_{n}\right]\right] = 0.$$

Putting it together

We have shown,

$$R(P, \widehat{f}) \leq \alpha_{1} \int_{0}^{1} \underbrace{\mathbb{E}\left[(\widehat{f}(t) - f(t))^{2}\right]}_{\text{err}_{t}(\widehat{f})} dt$$

$$\operatorname{err}_{t}(\widehat{f}) = \mathbb{E}\left[\underbrace{(\widehat{f}(t) - f(t))^{2}}_{=(b+v)^{2}} \middle| G_{t}\right] \underbrace{\mathbb{P}(G_{t})}_{\leq 1} + \underbrace{\mathbb{E}\left[(\widehat{f}(t) - f(t))^{2}\middle| G_{t}^{c}\right]}_{\leq 1} \underbrace{\mathbb{P}(G_{t}^{c})}_{\leq 1}$$

$$|b| \leq \frac{1}{N(t)} \sum_{i=1}^{n} A_{i} L h = L h,$$

$$\mathbb{E}[v^{2}|G_{t}] \leq \frac{\sigma^{2}}{nh},$$

$$\mathbb{E}[bv|G_{t}] = 0.$$

Putting it together

This gives us,

$$\operatorname{err}_t(\widehat{f}) \leq L^2 h^2 + \frac{\sigma^2}{nh} + e^{-2\alpha_0 nh^2}.$$

And hence,

$$R(P,\widehat{f}) \leq \alpha_0 \left(L^2 h^2 + \frac{\sigma^2}{nh} + e^{-2\alpha_0 nh^2} \right).$$

Choosing $h = \frac{\sigma^{2/3}}{L^{2/3}n^{1/3}}$, we get

$$R(P,\widehat{f}) \leq 2\alpha_1 \frac{\sigma^{4/3} L^{2/3}}{n^{2/3}} + \alpha_1 \exp\left(-2\alpha_0 \frac{\sigma^{4/3} n^{1/3}}{L^{4/3}}\right).$$

This shows us that the minimax rate for this problem is $\Theta(n^{-2/3})$.

The Nadaraya-Watson (Kernel) Estimator

An estimator of the following form,

$$\widehat{f}(t) = \sum_{i=1}^n w_i(t) Y_i, \qquad \text{where,} \qquad w_i(t) = \begin{cases} \frac{K\left(\frac{t-X_i}{h}\right)}{\sum_{j=1}^n K\left(\frac{t-X_j}{h}\right)} & \text{if } \sum_{j=1}^n K\left(\frac{t-X_j}{h}\right) > 0, \\ 0 & \text{otherwise} \end{cases}$$

Here,

- $ightharpoonup K: \mathbb{R} \to \mathbb{R}$ is called a smoothing kernel.
- ▶ In our estimator, $K(t) = \mathbb{1}(|t| \le 1)$. This is sufficient for Lipschtiz smoothness.
- But other kernel choices lead to better rates under stronger guarantees. E.g. In the Hölder class $\mathcal{H}(\beta,L)$ in \mathbb{R}^d , where all $(\beta-1)$ partial derivatives are L-Lipschitz, the minimax rate is $\Theta(n^{\frac{-2\beta}{2\beta+d}})$.

Ch 2.2: Nonparametric density estimation

Model: Let \mathcal{F} be the class of bounded L-Lipschitz functions in [0,1],

$$\mathcal{F} = \{f: [0,1] \to [0,B], |f(x_1) = f(x_2)| \le L|x_1 - x_2|\}.$$

Let $\mathcal P$ be all continuous distributions whose pdf is in $\mathcal F$.

We observe n samples $S = \{X_1, \dots, X_n\}$ drawn i.i.d from $P \in \mathcal{P}$ and wish to estimate the density p in the L_2 loss,

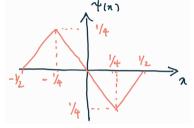
$$\Phi \circ
ho(p_1, p_2) = \|p_1 - p_2\|_2^2 = \int (p_1 - p_2)^2.$$

We will show $R_n^{\star} \in \Theta(n^{-2/3})$ where,

$$R_n^{\star} = \inf_{\widehat{\rho}} \sup_{p \in \mathcal{F}} \mathbb{E}_S \left[\|\widehat{p}(s) - p\|_2^2 \right].$$

Ch 2.2.1: Lower bounds for density estimation

Step 1: Constructing Alternatives. Let us consider the following function,



$$\psi(x) = \begin{cases} x + 1/2, & \text{if } x \in [-1/2, -1/4], \\ -x, & \text{if } x \in [-1/4, 1/4], \\ x - 1/2, & \text{if } x \in [1/4, 1/2], \\ 0 & \text{otherwise} \end{cases}$$

You can check that.

- $\blacktriangleright \psi$ is 1-Lipschitz.
- $\blacktriangleright \int \psi = 0$
- $\blacktriangleright \ \ \tfrac{-1}{4} \le \psi(t) \le \tfrac{1}{4}.$
- $\int \psi^2 = 1/48.$

Step 1: Constructing Alternatives (cont'd)

Let h > 0 be a parameter to be chosen later. Let m = 1/h.

Let $\omega \in \{0,1\}^m$. Let us define f_ω as follows:

You can verify that

- \triangleright Each p_{c} , is L-Lipschitz (similar calculation as before).
- ▶ Each p_{ω} is a valid density ($\int p_{\omega} = 1$ and $0 \ge p_{\omega} \ge B$) provided that $h < \frac{4}{7} \min(B - 1, 1)$.

Alternatives: Let $\mathcal{F}' = \{f_{\omega}; \omega \in \Omega_m\}$ where Ω_m is the Gilbert-Varshamov-pruned *m*-hypercube.

Step 2: Minimum Separation

For any $\omega, \omega' \in \Omega_d$, we have

$$\|p_{\omega} - p_{\omega'}\|_{2}^{2} = \sum_{j=1}^{m} \mathbb{1}(\omega_{j} \neq \omega'_{j}) L^{2} h^{3} \int_{-1/2}^{1/2} \psi^{2}(u) du$$
$$= \frac{L^{2} h^{3}}{48} H(\omega, \omega').$$

By the Gilbert-Varshamov Lemma, $\min_{\omega \neq \omega'} H(\omega, \omega') \geq m/8 = 1/(8h)$.

Therefore, we have

$$\min_{p_{\omega}, p'_{\omega} \in \mathcal{H}'} \|p_{\omega} - p_{\omega'}\|_2 \ge \sqrt{\frac{L^2 h^3}{48}} \sqrt{\min_{\omega, \omega'} H(\omega, \omega)} = \frac{Lh}{8\sqrt{6}} \stackrel{\Delta}{=} \delta \text{separation}$$
 (2)

Step 3: Upper bounding the maximum KL

To upper bound the KL divergence,

$$\mathrm{KL}(p_{\omega},p_{\omega'}) = \int p_{\omega} \log \left(\frac{p_{\omega}}{p_{\omega'}}\right) = \sum_{j=1}^{m} \underbrace{\int_{(j-1)/m}^{j/m} (1+\omega_{j}\phi_{j}) \log \left(\frac{1+\omega_{j}\phi_{j}}{1+\omega'_{j}\phi_{j}}\right)}_{A(\omega_{j},\omega'_{j})} \mathbb{1}(\omega_{j} \neq \omega'_{j})$$

Next, you can show (after some algebra), that when $\omega \neq \omega_i'$,

$$A(\omega_j, \omega_j') \le L^2 h^3 \int_{-1/2}^{1/2} \psi^2 = \frac{L^2 h^3}{48}.$$

Hint. Use inequalities $\log(1+x) \ge x - x^2$ for $x \ge -0.68$ and $\log(1+x) \le x$.

Therefore, $\mathrm{KL}(p_{\omega},p_{\omega'}) \leq \frac{L^2h^3}{48}H(\omega,\omega')$. As, $\max_{\omega,\neq\omega'}H(\omega,\omega')\leq m=1/h$, we have

$$\max_{p_{\omega},p_{\omega'}\in\mathcal{H}'}\mathrm{KL}(p_{\omega},p_{\omega'})\leq \frac{L^2h^4}{48}.$$

Step 4: Final Bound

We require the maximum KL to be smaller than $\frac{\log(|\mathcal{F}'|)}{4n}$. But we know $|\mathcal{F}'| \geq 2^{m/8}$. So it is sufficient if

$$\frac{L^2h^4}{48} < \frac{\log(2^{m/8})}{4n} = \frac{\log(2)}{32nh}.$$

So choose,

$$h = \left(\frac{3\log(2)}{2} \frac{1}{nL^2}\right)^{1/3}.$$

By the local Fano method, we obtain the following lower bound on the minimax risk,

$$R_n^{\star} \geq \frac{1}{2} \Phi\left(\frac{\delta}{2}\right) = C_2 L^2 h^2 = C_3 \frac{L^{2/3}}{n^{2/3}}.$$

Requirements

Requirements:

- ▶ $h \le \frac{4}{L} \min(B 1, 1)$, for it to be a valid density.
- ▶ $\frac{1}{h} = m \ge 8$ for Gilbert-Varshamov.
- ▶ $2^{m/8} \ge 16$, which is satisfied if $h \le 1/32$.
- ▶ $h \le 2.72/L$, an inequality we used.

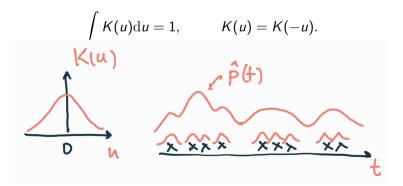
These conditions are satisfied if n is sufficiently large.

Ch 2.2.2: Upper bound via Kernel Density Estimation

A kernel density estimator \hat{p} with kernel K and bandwidth parameter h (to be chosen shortly) has the following form,

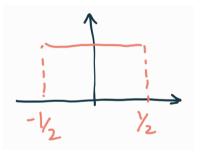
$$\widehat{p}(t) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{h} K\left(\frac{t - X_i}{h}\right),$$

Here, $K : \mathbb{R}_+ \to \mathbb{R}_+$ is a smoothing kernel with the following properties,



Boxcar Kernel

For our problem, we will use $K(u) = \mathbb{1}(|u| \le 1/2)$.



- ▶ This kernel is sufficient for Lipschitz functions.
- ▶ In HW2, you will analyse KDE for densities whose first derivative is Lipschitz. Then, we need to use a different kernel.

Bias-Variance Decomposition

Let us decompose the risk as follows,

$$\mathbb{E}\left[\|p-\widehat{\rho}\|_{2}^{2}\right]$$

$$= \mathbb{E}\left[\int (p-\widehat{\rho})^{2}\right]$$

$$= \mathbb{E}\left[\int (p-\mathbb{E}\widehat{\rho})^{2} + \int (\mathbb{E}\widehat{\rho}-\widehat{\rho})^{2} + 2\int (p-\mathbb{E}\widehat{\rho})(\mathbb{E}\widehat{\rho}-\widehat{\rho})\right]$$

$$= \int \underbrace{(p(t)-\mathbb{E}\widehat{\rho}(t))^{2}}_{\text{bias}^{2}(t)} dt + \int \underbrace{\mathbb{E}\left[(\widehat{\rho}(t)-\mathbb{E}\widehat{\rho}(t))^{2}\right]}_{\text{var}(t)} dt + 2\int (p(t)-\widehat{\rho}(t))\underbrace{\mathbb{E}\left[\widehat{\rho}(t)-\mathbb{E}\widehat{\rho}(t)\right]}_{=0} dt$$

Bounding the Bias

Let us first bound the bias,

$$bias(t) = \mathbb{E}[\widehat{p}(t)] - p(t) = \mathbb{E}_{X \sim P} \left[\frac{1}{h} K \left(\frac{X - t}{h} \right) \right] - p(t)$$
$$= \int \frac{1}{h} K \left(\frac{x - t}{h} \right) p(x) dx - p(t) \int K(u) du$$
$$= \int K(u) (p(t + uh) - p(t)) du.$$

Therefore, noting that $\int |u| K(u) du = 1/4$, we have

$$|\mathrm{bias}(t)| \leq \int K(u)|p(t+uh)-p(t)|\mathrm{d}u \leq \int K(u)Lh|u|\mathrm{d}u = \frac{1}{4}Lh.$$

Bounding the Variance

Let us next bound the variance.

$$\operatorname{var}(t) = \operatorname{\mathbb{V}ar}\left(\widehat{p}(t)\right) = \operatorname{\mathbb{V}ar}\left(\frac{1}{n}\sum_{i=1}^{n}\frac{1}{h}K\left(\frac{t-X_{i}}{h}\right)\right)$$

$$= \frac{1}{n}\operatorname{\mathbb{V}ar}_{X\sim P}\left(\frac{1}{h}K\left(\frac{t-X}{h}\right)\right)$$
as the terms are i.i.d
$$\leq \frac{1}{n}\mathbb{E}\left[\frac{1}{h^{2}}K^{2}\left(\frac{X-t}{h}\right)\right]$$

$$= \frac{1}{nh^{2}}\int K^{2}\left(\frac{x-t}{h}\right)p(x)\mathrm{d}x$$

$$= \frac{1}{nh}\int K^{2}(u)p(x+uh)\mathrm{d}u$$

$$\leq \frac{B}{nh}\int K^{2}(u)\mathrm{d}u = \frac{B}{nh}.$$
as $\int K^{2} = 1$

Putting It All Together

Putting these bounds together we get,

$$\mathbb{E}\left[\|p-\widehat{p}\|_{2}^{2}\right] = \int \operatorname{bias}^{2}(t) + \int \operatorname{var}(t)$$

$$\leq \int_{0}^{1} \frac{1}{4}L^{2}h^{2} + \int_{0}^{1} \frac{B}{nh}$$

$$\leq \frac{1}{4}L^{2}h^{2} + \frac{B}{nh}$$

Choosing $h = n^{-1/3}L^{-1/3}$, we get

$$\mathbb{E}\left[\|\boldsymbol{\rho}-\widehat{\boldsymbol{\rho}}\|_2^2\right]\in\mathcal{O}\left(\frac{L^{2/3}}{n^{2/3}}\right).$$

Boundary correction

- ▶ The bounds for bias(t) and var(t) are valid only within the interval [h, 1-h].
- ▶ In the analysis above, we may ignore this issue (i.e., assume the worst-case error) and still obtain the correct convergence rates.
- ► For higher-order smoothness assumptions, however, this boundary error may no longer be negligible.
- ▶ A common remedy is to mirror the data at the boundaries, which is effective when the underlying density has zero derivative at the boundary.