

HOMEWORK ASSIGNMENT 2

1 Data-driven bandwidth selection for histograms

We consider the task of unsupervised learning with $X_1, \dots, X_n \in [0, 1]$ iid distributed according to P . It is assumed that P has a density with respect to the Lebesgue measure and this density, denoted by f is unknown. The goal is to estimate f by some estimator \hat{f} so that the Mean Integrated Squared Error :

$$\text{MISE}_f(\hat{f}) = \mathbf{E} \left[\|\hat{f} - f\|_2^2 \right], \quad \left(\|g\|_2^2 \triangleq \int_0^1 g^2(x) dx \right)$$

be small.

We have seen in the lectures that the histogram estimator

$$\hat{f}_{n,h}(x) = \frac{1}{h} \sum_{k=1}^K \hat{p}_k \mathbb{1}_{C_k}(x), \quad \text{where} \quad \hat{p}_k = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{C_k}(X_i)$$

has a risk equal to

$$\text{MISE}_f(\hat{f}_{n,h}) = \|f\|_2^2 + \frac{1}{nh} - \frac{n+1}{nh} \sum_{k=1}^K p_k^2, \quad \text{where} \quad p_k = \int_{C_k} f(x) dx.$$

Furthermore, we have shown that

$$\hat{J}_n(h) = \frac{2}{(n-1)h} - \frac{n+1}{(n-1)h} \sum_{k=1}^K \hat{p}_k^2$$

is an unbiased estimator of the quantity $J(h) = \text{MISE}_f(\hat{f}_{n,h}) - \|f\|_2^2$.

Based on this, we have proposed the following rule for selecting among M candidates h_1, \dots, h_M of possible bandwidths :

$$\hat{h} \in \arg \min_{h_1, \dots, h_M} \hat{J}_n(h).$$

We have demonstrated that the resulting *adaptive* estimator $\hat{f}_{n,\hat{h}}$ satisfies the oracle inequality

$$\text{MISE}_f(\hat{f}_{n,\hat{h}}) \leq \min_{h \in \{h_1, \dots, h_M\}} \text{MISE}_f(\hat{f}_{n,h}) + 2\sqrt{M} \sup_{h>0} V_n(h)^{1/2},$$

where $V_n(h)$ is the variance of the zero-mean random variable

$$\xi_n(h) = \|\hat{f}_{n,h} - f\|_2^2 - \|f\|_2^2 - \hat{J}_n(h).$$

The aim of the present exercise is to find a relatively tight upper bound on $V_n(h)$. **In what follows, we assume that f is bounded and denote by $\|f\|_\infty$ the supremum of f .**

1. Show that $\xi_n(h) = \zeta_n^{(1)}(h) + \zeta_n^{(2)}(h)$, where

$$\zeta_n^{(1)}(h) = \frac{2n}{(n-1)h} \sum_{k=1}^K (\hat{p}_k^2 - \mathbf{E}[\hat{p}_k^2]),$$

$$\zeta_n^{(2)}(h) = \frac{2}{h} \sum_{k=1}^K p_k (p_k - \hat{p}_k).$$

2. Let us set

$$Z_i = \frac{2}{h} \sum_{k=1}^K p_k (p_k - \mathbb{1}_{C_k}(X_i)), \quad i = 1, \dots, n.$$

- (a) Check that $\tilde{\zeta}_n^{(2)}(h) = (1/n) \sum_{i=1}^n Z_i$. What can you say about the dependency between the random variables Z_1, \dots, Z_n ?
- (b) Show that the variance of Z_1 is upper-bounded by

$$\frac{4}{h^2} \sum_{k=1}^K p_k^3 (1 - p_k).$$

- (c) Deduce from previous questions an upper bound of the form

$$\text{Var}[\tilde{\zeta}_n^{(2)}(h)] \leq \frac{C \|f\|_\infty^2}{n}$$

with a constant C to be made precise.

3. We study now the first term $\tilde{\zeta}_n^{(1)}(h)$.

- (a) Let us denote $T = \{(i, j) \in \{1, \dots, n\}^2 \text{ s.t. } i < j\}$ and introduce the random variables

$$U_t = \sum_{k=1}^K \left(\mathbb{1}_{C_k}(X_i) \mathbb{1}_{C_k}(X_j) - p_k^2 \right), \quad \text{if } t = (i, j).$$

Show that $\mathbf{E}[U_t] = 0$ for every t and that

$$\tilde{\zeta}_n^{(1)}(h) = \frac{4}{nh(n-1)} \sum_{t \in T} U_t.$$

- (b) Let $t = (i, j)$ and $t' = (i', j')$ be elements of T such that all indices (i, i', j, j') are different. What is the value of the covariance between U_t and $U_{t'}$?
- (c) Check that for every $t = (i, j) \in T$:

$$\mathbf{E} \left[\left(\sum_{k=1}^K \mathbb{1}_{C_k}(X_i) \mathbb{1}_{C_k}(X_j) \right)^2 \right] = \sum_{k=1}^K \mathbf{E} \left[\left(\mathbb{1}_{C_k}(X_i) \mathbb{1}_{C_k}(X_j) \right)^2 \right] = \sum_{k=1}^K p_k^2.$$

Compute the variance of U_t .

- (d) Let us define the set

$$\mathcal{W} = \left\{ (i, j, i', j') \in T \times T \text{ s.t. } (i, j) \neq (i', j') \text{ and } \{i, j\} \cap \{i', j'\} \neq \emptyset \right\}.$$

Check that the cardinality of \mathcal{W} is less than $n(n-1)^2$.

- (e) Show that for every pair $(t, t') \in \mathcal{W}$, it holds that

$$\mathbf{E}[U_t U_{t'}] \leq \sum_{k=1}^K p_k^3 \leq h^2 \|f\|_\infty^2.$$

- (f) Show that

$$\text{Var}[\tilde{\zeta}_n^{(1)}(h)] = \frac{16}{h^2 n^2 (n-1)^2} \left\{ \sum_{(t, t') \in \mathcal{W}} \mathbf{E}[U_t U_{t'}] + \sum_{t \in T} \mathbf{E}[U_t^2] \right\}.$$

Deduce from preceding questions that

$$\text{Var}[\tilde{\zeta}_n^{(1)}(h)] \leq \frac{16 \|f\|_\infty^2}{n} + \frac{8 \|f\|_\infty}{hn(n-1)}$$

4. Show that $\|f\|_\infty \geq 1$. Explain why it is reasonable to assume that $h > 1/n$ and show that

$$V_n(h)^{1/2} = (\text{Var}[\tilde{\zeta}_n(h)])^{1/2} \leq \frac{9 \|f\|_\infty}{\sqrt{n}}.$$