

Mathematics of Machine Learning

Lecture 05/09/17

Outline

- 1 - Non-parametrics: an introduction
- 2 - Lower bounds
- 3 - Stochastic bandit models
- 4 - Classification: a theory
- 5 - Machine Learning and optimization
- 6 - Individual sequences

Prof: Aurélien Garivier

Sebastien Gerchinovitz

#Lectures: 6

1 Non-parametrics: an introduction

I Density Estimation

Literature: Tsybakov "Nonparametric Statistics" "Introduction to non-parametric estimation"

Model: $X_1, \dots, X_n \stackrel{iid}{\sim} P$, $P \in \mathcal{M}$ (family of proba. distributions)

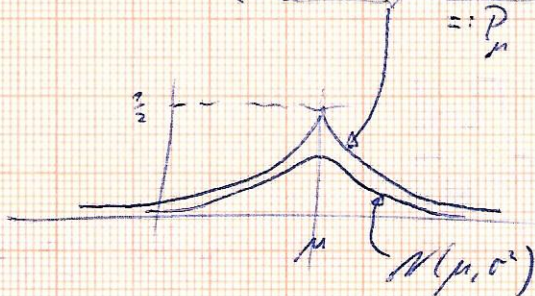
Problem: We don't know P , we want to estimate it.

Example (parametric): $\mathcal{M} := \{N(\mu, \sigma^2), \mu \in \mathbb{R}\}$, $\sigma^2 > 0$ is ^{given} known

→ estimate μ :

$$\hat{\mu}_n := \bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i$$

Example $\mathcal{M} = \left\{ \frac{1}{2} \exp(-|x-\mu|) d\lambda, \mu \in \mathbb{R} \right\}$ ← Lebesgue measure on \mathbb{R}
=: P_μ



- We could estimate μ with the Maximum-likelihood-method or the moment method:

• The expectation of P_μ is $E_\mu[X_i] = \mu$

• Law of large numbers yields $\bar{X}_n \rightarrow E_\mu[X_i] = \mu$
(LLN)

$$\hat{\mu}_n = \bar{X}_n$$

• for every estimator $\hat{\mu}_n$, we define the quadratic risk

$$R(\hat{\mu}_n, \mu) = \mathbb{E}_\mu[(\hat{\mu}_n - \mu)^2] \quad \text{and one can compute}$$

$$\begin{aligned} \mathbb{E}_\mu[(\hat{\mu}_n - \mu)^2] &= \mathbb{E}_\mu[(\hat{\mu}_n - \mathbb{E}_\mu[\hat{\mu}_n]) + (\mathbb{E}_\mu[\hat{\mu}_n] - \mu)]^2 \\ &= \mathbb{E}_\mu[(\hat{\mu}_n - \mathbb{E}_\mu[\hat{\mu}_n])^2] \\ &\quad + (\mathbb{E}_\mu[\hat{\mu}_n] - \mu)^2 \\ &\quad + \underbrace{\mathbb{E}_\mu[\hat{\mu}_n - \mathbb{E}_\mu[\hat{\mu}_n]]}_{=0} \cdot (\mathbb{E}_\mu[\hat{\mu}_n] - \mu) \end{aligned}$$

$$= \text{Var}(\hat{\mu}_n) + \underbrace{(\mathbb{E}_\mu[\hat{\mu}_n] - \mu)^2}_{=: (\text{bias}(\hat{\mu}_n))^2} \quad (\text{"variance-bias-decomposition"})$$

Exp: $R(\bar{X}_n, \mu) = \frac{\text{Var}(X_1)}{n}$

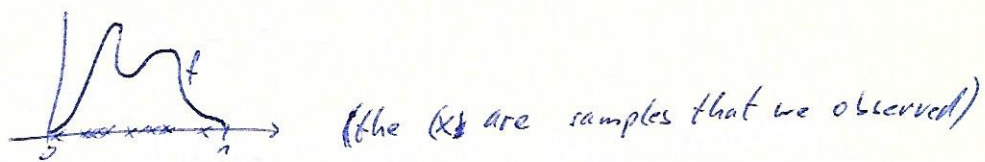
Remark:

One can prove: $\exists C > 0: \forall \hat{\mu}_n, \mathbb{E}_\mu[(\hat{\mu}_n - \mu)^2] \geq \frac{C}{n}$

Coming to non-parametrics:

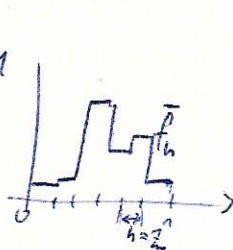
exp: $\mathcal{M} = \{dP(x) = f(x)dx, f \in C^2([0,1])\}$

$f \in \mathcal{M}$ may look like



→ How to estimate f ?

→ Idea: make an histogram:



that means we introduce a (regular) partition of $[0,1]$: C_1, \dots, C_m with $C_j = [\frac{j-1}{m}, \frac{j}{m}]$, $j=1, \dots, m$

$$\bar{f}_h(x) := \sum_{j=1}^m \frac{p_j}{h} \mathbb{1}_{C_j}(x) \quad \text{where } h = \frac{1}{m}, p_j := P(X_1 \in C_j)$$

$$\begin{aligned} \rightsquigarrow \frac{p_j}{h} &= \frac{1}{h} \int_{C_j} f(x) dx && \rightarrow \text{We estimate } \bar{f}_h \text{ instead of } f \text{ by} \\ \hat{f}_h(x) &= \frac{1}{nh} \sum_{i=1}^n \left(\sum_{j=1}^m \mathbb{1}_{C_j}(x_i) \mathbb{1}_{C_j}(x) \right) \end{aligned}$$

$$\mathcal{M}_h = \{f: [0,1] \rightarrow \mathbb{R}_+, f(x) = \sum_{j=1}^m \alpha_j \mathbb{1}_{C_j}(x), (\alpha_1, \dots, \alpha_m) \in \mathbb{R}^m\}$$

$$\hat{x}_j = \frac{1}{nh} \sum_{i=1}^n \mathbb{1}_{C_j}(X_i) \quad \mathcal{B}(p_j) \text{ (Bernoulli-distr. with parameter } p_j)$$

$$\mathbb{E}[\hat{x}_j] = \frac{1}{h} \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{C_j}(X_i)\right] = \frac{p_j}{h} = \alpha_j$$

$$\rightarrow \forall x \in [0,1] : \mathbb{E}[\hat{f}_h(x)] = \bar{f}_h(x)$$

Risk at $x \in [0,1]$

$$MSE_f(x, h) = \mathbb{E}_f \left[(\hat{f}_h(x) - f(x))^2 \right] = \underbrace{\left(\mathbb{E}_f[\hat{f}_h(x)] - f(x) \right)^2}_{\text{squared bias}} + \underbrace{\text{Var}_f[\hat{f}_h(x)]}_{\text{variance}}$$

$$\hat{f}_h(x) = \frac{\hat{p}_j(x)}{h} = \frac{1}{nh} \sum_{i=1}^n \mathbb{1}_{C_j(x)}(X_i) = \frac{Z_j}{nh} \quad \text{where } Z_j \sim \mathcal{B}(n, p_j(x))$$

the j for which holds that $x \in C_j$ binomial distr.

$$\rightarrow \mathbb{E}[\hat{f}_h(x)] = \frac{p_j(x)}{h}$$

$$\begin{aligned} \text{Var}(\hat{f}_h(x)) &= \frac{np_j(x)(1-p_j(x))}{(nh)^2} \\ &= \frac{p_j(x)(1-p_j(x))}{nh} \end{aligned}$$

$$\rightarrow \text{if we want to have } MSE_f(x, h) \xrightarrow{h \rightarrow 0} 0$$

$$\text{we need to choose } h = h_n \text{ and } h_n \xrightarrow{h \rightarrow 0} 0$$

Integrated risk

$$MISE_f(h) := \int_0^1 MSE_f(x, h) dx \stackrel{\text{Fubini}}{=} \mathbb{E}_f \left[\int_0^1 (\hat{f}_h(x) - f(x))^2 dx \right]$$

$$\text{Variance: } \int_0^1 \text{Var}_f(\hat{f}_h(x)) dx = \sum_{j=1}^m \int_{C_j} \text{Var}_f(\hat{f}_h(x)) dx = \sum_{j=1}^m \frac{p_j(1-p_j)}{nh} = \frac{1}{nh} \sum_{j=1}^m p_j - \frac{1}{n} \sum_{j=1}^m p_j^2$$

$\sum_{j=1}^m p_j = 1$

$$\text{Bias term: } \int_0^1 \left(\mathbb{E}_f[\hat{f}_h(x)] - f(x) \right)^2 dx = \sum_{j=1}^m \int_{C_j} \left(\frac{p_j}{h} - f(x) \right)^2 dx$$

$$= \sum_{j=1}^m \left\{ \frac{p_j^2}{h} - 2 \frac{p_j}{h} \int_{C_j} f(x) dx + \int_{C_j} f(x)^2 dx \right\}$$

$$= \int_0^1 f(x)^2 dx - \frac{1}{h} \sum_{j=1}^m p_j^2$$

Lemma: If $X_1, \dots, X_n \stackrel{iid}{\sim}$ with density f on $[0,1]$ and if \hat{f}_h is the histogram estimator with $m = \frac{1}{h}$ bins, then

$$MISE_f(h) = \mathbb{E}_f \left[\int_0^1 (\hat{f}_h(x) - f(x))^2 dx \right] = \int_0^1 \underbrace{f(x)^2}_{\int dx} + \frac{1}{nh} - \frac{2}{nh} \sum P_j^2$$

Question: How does it behave when $n \rightarrow \infty$ and $h_n \rightarrow 0$

$$\int_{c_j} f^2(x) dx - \frac{1}{h} P_j^2 = \int_{c_j} \left(f(x) - \frac{1}{h} \int_{c_j} f(u) du \right)^2 dx = \frac{1}{h^2} \int_{c_j} \left(\int_{c_j} (f(x) - f(u)) du \right)^2 dx$$

As f is assumed to be twice continuously differentiable,

$$f(u) - f(x) = (u-x) f'(a_j) + O(h^2) \quad \text{where } a_j = \frac{j-1}{m}$$

$$\int_{c_j} f^2(x) dx - \frac{1}{h} P_j^2 = \frac{f'(a_j)^2}{h^2} \int_{c_j} \left(\int_{c_j} (x-u) du \right)^2 dx + O(h^4)$$

Now: $\int_{c_j} \left(\int_{c_j} x-u du \right) dx = h^2 \int_0^1 \left(\int_0^1 (y-z) dz \right)^2 dy = \frac{h^5}{12}$

$$\Rightarrow \int_{c_j} f^2(x) dx - \frac{P_j^2}{h} = \frac{h^3}{12} f'(a_j)^2 + O(h^4) \\ = \frac{h^2}{12} \int_{c_j} f'(x)^2 dx + O(h^4) \quad \text{and}$$

$$MISE_f(h) = \sum_{j=1}^m \left(\int_{c_j} f^2(x) dx - \frac{P_j^2}{h} \right) + \frac{1}{nh} - \frac{1}{nh} \sum P_j^2 \\ = \frac{h^2}{12} \int_0^1 f'(x)^2 dx + O(h^3) + \frac{1}{nh} + O\left(\frac{1}{n}\right)$$

Theorem: if X_1, \dots, X_n is a sample of a probab. distr. with density in $C^2([0,1])$ and if $h_n \xrightarrow{n \rightarrow \infty} 0$, then the histogram estimator \hat{f}_{h_n} satisfies

$$MISE_f(\hat{h}_n) = \frac{h_n^2}{12} \int_0^1 f'(x)^2 dx + \frac{1}{nh_n} + O(h_n^3) + O\left(\frac{1}{n}\right) \quad \text{as } n \rightarrow \infty$$

\rightarrow in order to find the best number of bins, we need to minimize

$$g(h) = \frac{h^2}{12} \int_0^1 f'(x)^2 dx + \frac{1}{nh} \quad \xrightarrow{\text{obtain}} \quad h_{opt} = \left(\frac{12}{n} \int_0^1 f'(x)^2 dx \right)^{-\frac{1}{2}}$$

for that choice, $MISE_f(h_{opt}) = \left(\frac{12}{n} \int_0^1 f'(x)^2 dx \right)^{\frac{2}{3}} \cdot \frac{1}{n^{\frac{1}{3}}} + O\left(\frac{1}{n}\right)$.

Problems: - can we do better? (faster convergence) concerning the upper bound of MSE
 - we don't know $f' \Rightarrow$ we don't know h_{opt}

Memo: $MISE_f(h) = \int_0^1 f^2(x) dx + \frac{1}{nh} - \frac{n+1}{nh} \sum_{j=1}^m p_j^2$

We want to find h minimizing $MISE_f(h)$ just using X_1, \dots, X_n (and not anything unknown of f)

Remark: $\arg \min_h MSE_f(h) = \arg \min_h \left(MSE_f(h) - \int_0^1 f^2(x) dx \right)$
 $= \arg \min_h \left(\frac{1}{nh} - \frac{n+1}{nh} \sum_{j=1}^m p_j^2 \right)$

Idea: estimate $J_f(h) = MSE_f(h) - \int_0^1 f^2(x) dx$
 $= \frac{1}{nh} - \frac{n+1}{nh} \sum_{j=1}^m p_j^2$ by

$J_f(h) = \frac{1}{nh} - \frac{n+1}{nh} \sum_{j=1}^m \hat{p}_j^2$ } not good because biased:

$E[\hat{p}_j^2] = E[\hat{p}_j]^2 + \text{Var}[\hat{p}_j] = p_j^2 + \frac{p_j(1-p_j)}{n} = p_j^2 \left(1 - \frac{1}{n}\right) + \frac{p_j}{n}$

$\rightarrow E\left[\hat{p}_j^2 - \frac{\hat{p}_j}{n}\right] = p_j^2 \left(\frac{n-1}{n}\right) \Leftrightarrow E\left[\frac{n}{n-1} \left(\hat{p}_j^2 - \frac{\hat{p}_j}{n}\right)\right] = p_j^2$
 $= \frac{n}{n-1} \hat{p}_j^2 - \frac{\hat{p}_j}{n-1}$

$\rightarrow J_f(h) = \frac{1}{nh} - \frac{n+1}{nh} \sum_{j=1}^m \left(\frac{n}{n-1} \hat{p}_j^2 - \frac{\hat{p}_j}{n-1} \right)$

$= \frac{1}{nh} - \frac{n+1}{(n-1)h} \sum_{j=1}^m \hat{p}_j^2 + \frac{n+1}{n(n-1)h} \sum_{j=1}^m \hat{p}_j$

$= -\frac{n+1}{(n-1)h} \sum_{j=1}^m \hat{p}_j^2 + \frac{1}{nh} \left(\frac{n+1}{n-1} \right) = \frac{2}{(n-1)h} - \frac{n+1}{(n-1)h} \sum_{j=1}^m \hat{p}_j^2$

Prop: $E[J_f(h)] = J_f(h)$ "cv" = "cross-validation"

Alg: initialize: $m=1, m_{cv}=1, J_{cv} = +\infty$

while $m < n$ do

$J = \frac{2m}{n-1} - \frac{(n+1)m}{n-1} \sum_{j=1}^m \left(\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{C_j}(K_{ij}) \right)^2$

if $J < J_{cv}$ then

$m_{cv} = m$
 $J = J$

(Algo to find the best bin-width for our histogram)

endif
 $m = m + 1$
 return $\hat{h} = \frac{1}{mcr}$

Find estimator \hat{f}_h

Other ideas to estimate f :

coming from histograms:



if x is at the center of a bin: only

$$\hat{f}_h^{hist}(x) = \frac{1}{nh} \sum_{i=1}^n \mathbb{1}\{|X_i - x| < \frac{h}{2}\}$$

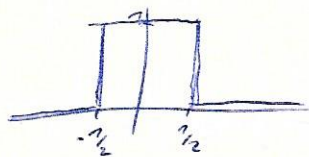
(Put the center of a bin the histogram estimator coincides with the rectangular kernel estimator!)

$$\hat{f}_h^R(x) = \frac{1}{nh} \sum_{i=1}^n \mathbb{1}\{|X_i - x| < \frac{h}{2}\} \quad \forall x \in [0, 1]$$

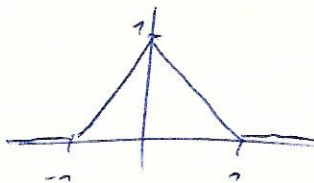
Def: $\hat{f}_h^K(z) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{X_i - z}{h}\right)$ "Kernel estimator", where K is a

Kernel fct. and h a width parameter

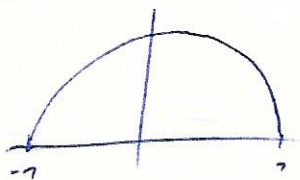
Kernel exp.: • rectangular kernel:



• triangular kernel

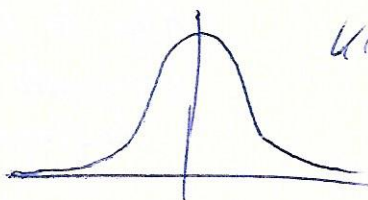


• Epanechnikov



$$K(u) = \frac{3}{4} (1 - u^2) \mathbb{1}_{[-1, 1]}(u)$$

• Gaussian kernel



$$K(u) = \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}}$$

Exercise if $K(u) \geq 0 \forall u \in \mathbb{R}$ and if $\int_{-\infty}^{\infty} K(u) du = 1$, then \hat{f}_h^K is a density function.

Remark: For now we only saw positive kernels, but there are also kernels that are partly negative (difference to densities)

What is the risk of the kernel density estimator?

Assumptions on the kernel

$$(1) \int_{-\infty}^{\infty} K(u) du = 1$$

$$(2) \int_{-\infty}^{\infty} u K(u) du = 0 \quad (\Leftarrow K \text{ is even})$$

$$(3) \int_{-\infty}^{\infty} u^2 |K(u)| du < \infty$$

$$(4) \int_{-\infty}^{\infty} K(u)^2 du < \infty$$

Prop. a) Under assumptions (1), (2), (3), if f is a bounded density with a bounded second derivative, then

$$\text{Bias}(\hat{f}_h^K) \leq C_1 h^2 \quad \text{where } C_1 = \frac{1}{2} \sup_z |f''(z)| \cdot \int_{-\infty}^{\infty} u^2 |K(u)| du$$

b) In addition, under conditions (1), (4)

$$\text{Var}(\hat{f}_h^K(x)) \leq \frac{C_2}{nh} \quad \text{where } C_2 = \sup_z f(z) \int_{-\infty}^{\infty} K(u)^2 du$$

Proof: a)
$$\begin{aligned} \mathbb{E}_f[\hat{f}_h^K(x)] &= \frac{1}{nh} \sum_{i=1}^n \mathbb{E}_f[K(\frac{X_i - x}{h})] = \frac{1}{h} \mathbb{E}_f[K(\frac{X_1 - x}{h})] = \frac{1}{h} \int_{-\infty}^{\infty} K(\frac{y-x}{h}) f(y) dy \\ &\stackrel{y=x+uh}{=} \int_{-\infty}^{\infty} K(u) f(x+uh) du \end{aligned}$$

Using a Taylor expansion:

$$\begin{aligned} E_f \left[\hat{f}_h^K(x) \right] &= \int_{-\infty}^{\infty} K(u) f(x+uh) du = \int_{-\infty}^{\infty} K(u) \left(f(x) + uh f'(x) + \frac{(uh)^2}{2} f''(\xi_u) \right) du \\ &= \int_{-\infty}^{\infty} K(u) du + h f'(x) \int_{-\infty}^{\infty} u K(u) du + \frac{h^2}{2} \int_{-\infty}^{\infty} u^2 K(u) f''(\xi_u) du \end{aligned}$$

where $\xi_u \in [x, x+uh]$

Thus, $| \text{Bias}(\hat{f}_h^K(x)) | = | E[\hat{f}_h^K(x)] - f(x) |$

$$\begin{aligned} &\leq \frac{h^2}{2} \left| \int_{-\infty}^{\infty} u^2 K(u) f''(\xi_u) du \right| \\ &\leq \frac{h^2}{2} \int_{-\infty}^{\infty} u^2 |K(u)| \cdot |f''(\xi_u)| du \\ &\leq h^2 \cdot \underbrace{\frac{1}{2} \sup_x |f''(z)| \int_{-\infty}^{\infty} u^2 |K(u)| du}_{= C_1} \end{aligned}$$

b) $| \text{Var}(\hat{f}_h^K(x)) = \frac{1}{(nh)^2} \text{Var}_f \left[\sum_{i=1}^n K\left(\frac{X_i - x}{h}\right) \right] = \frac{1}{nh^2} \text{Var} \left[K\left(\frac{X_1 - x}{h}\right) \right]$

$$\begin{aligned} &\leq \frac{1}{nh^2} E \left[K\left(\frac{X_1 - x}{h}\right)^2 \right] \\ &= \frac{1}{nh^2} \int_{-\infty}^{\infty} K\left(\frac{y-x}{h}\right)^2 f(y) dy \\ &= \frac{1}{nh} \int_{-\infty}^{\infty} K(u)^2 f(x+uh) du \leq \frac{1}{nh} \sup |f| \int_{-\infty}^{\infty} K(u)^2 du \end{aligned}$$

□

Consequence

$$MSE(\hat{f}_h^K(x)) \leq C_1 h^4 + \frac{C_2}{nh}$$

→ we can optimize in h to find the best bandwidth

$$h_{opt} = \left(\frac{C_2}{4C_1} \right)^{1/5} n^{-1/5} \quad \text{which leads to}$$

$$MSE(\hat{f}_{h_{opt}}^K(x)) \leq C n^{-4/5}$$

Remarks

• $n^{-4/5}$ is better than $n^{-2/3}$

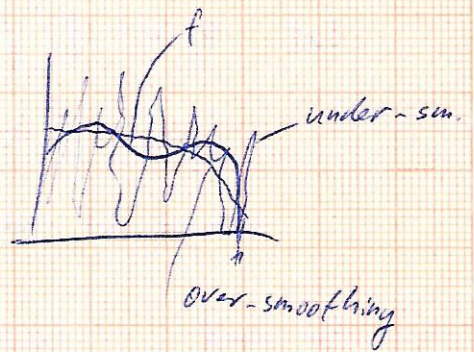
• $n^{-4/5}$ is optimal with those assumptions on f

• if we assume that f is ~~in $C^k(\mathbb{R}^d)$~~ $C^k(\mathbb{R}^d)$ then one can show that the optimal rate is

$$n^{-\frac{2k}{2k+d}} \quad (\text{we did } k=2, d=1)$$

• h too small \rightarrow "under smoothing"

• h too big \rightarrow "over-smoothing"



Proposition $\hat{J}(h) = \|\hat{f}_h^k\|_2^2 - \frac{2}{n(n-1)h} \sum_{i=0}^n \sum_{j \neq i}^n K\left(\frac{x_i - x_j}{h}\right)$ is an unbiased estimator of $J(h) = \text{MISE}(\hat{f}_h^k) - \|f\|_2^2$

Proof:
$$\begin{aligned} \mathbb{E}_f[\hat{J}(h)] &= \mathbb{E}_f[\|\hat{f}_h^k\|_2^2] - \frac{2}{n(n-1)h} \sum_{i=0}^n \sum_{j \neq i}^n \mathbb{E}_f\left[K\left(\frac{x_i - x_j}{h}\right)\right] \\ &= \mathbb{E}_f[\|\hat{f}_h^k\|_2^2] - \frac{2}{n(n-1)h} \sum_{i=0}^n \sum_{j \neq i}^n \int_{\mathbb{R}^2} K\left(\frac{x-y}{h}\right) f(x)f(y) dx dy \\ &= \mathbb{E}_f[\|\hat{f}_h^k\|_2^2] - \frac{2}{h} \int_{\mathbb{R}^2} K\left(\frac{x-y}{h}\right) f(x)f(y) dx dy \quad (*) \end{aligned}$$

But
$$\begin{aligned} J(h) &= \mathbb{E}_f[\|\hat{f}_h^k - f\|_2^2] - \|f\|_2^2 \\ &= \mathbb{E}_f[\|\hat{f}_h^k\|_2^2] - 2\mathbb{E}_f\left[\int_{\mathbb{R}^2} \hat{f}_h^k(x) f(x) dx\right] \\ &= \mathbb{E}_f[\|\hat{f}_h^k\|_2^2] - 2 \int_{\mathbb{R}^2} \mathbb{E}_f[\hat{f}_h^k(x)] f(x) dx \\ &= \mathbb{E}_f[\|\hat{f}_h^k\|_2^2] - 2 \int_{\mathbb{R}^2} \frac{1}{h} \int_{\mathbb{R}^2} K\left(\frac{y-x}{h}\right) f(y) f(x) dy dx \\ &\stackrel{(*)}{=} \mathbb{E}_f[\hat{J}(h)] \quad \square \end{aligned}$$

Also: - compute $\hat{h} = \arg\min_{h} \hat{J}(h)$ of possible values (for exp. on a grid of h)

- return $\hat{f}_{\hat{h}}^k$

Alternative method for choosing the bandwidth h : Lepski's method

Fix $x \in (0, 1)$, $\text{supp}(f) \subseteq [0, 1]$

$$\hat{f}_{h,h}^{\text{LC}}(x) = \frac{1}{2nh} \sum_{i=0}^n \mathbb{1}_{[x-h, x+h]}(X_i)$$

"locally constant"

Let us consider a finite sequence $h_1 < h_2 < \dots < h_M$ of possible bandwidths.

Let $e_j(x) = \mathbb{E}[\hat{f}_j(x)] - f(x)$ "approximation error" (= bias)

$\xi_j(x) = \hat{f}_j(x) - \mathbb{E}[\hat{f}_j(x)]$ "estimation error" (= variation)

where we use the notation $\hat{f}_j = \hat{f}_{h_j}^{\text{LC}}$. (h is fixed)

For small j , $|e_j(x)| < |\xi_j(x)|$.

For large j , $|e_j(x)| > |\xi_j(x)|$. Let $\delta > 0$.

Assume that you can choose a sequence $\epsilon_1 < \epsilon_2 < \dots < \epsilon_M$ such that

$\mathbb{P}(|\xi_j(x)| \leq \epsilon_j, \forall j=1, \dots, M) \geq 1 - \delta$. Then build the interval

$$I_j = [\hat{f}_j(x) - 2\epsilon_j, \hat{f}_j(x) + 2\epsilon_j].$$

We want to find the

bandwidth h_{j^*} for which $|e_{j^*}(x)| \approx |\xi_{j^*}(x)|$.

$$\text{Let } j^* = \min \{j : |e_j| \geq \epsilon_j\} - 1$$

$= \max \{j : |e_j| \leq \epsilon_j\}$. Then for all $j < j^*$, with high probab.

$$|\hat{f}_j(x) - f(x)| \leq |e_j(x)| + |\xi_j(x)| \leq |e_j(x)| + \epsilon_j \leq 2\epsilon_j \quad \text{which is equivalent to}$$

$$f(x) \in [\hat{f}_j(x) - 2\epsilon_j, \hat{f}_j(x) + 2\epsilon_j] = I_j$$

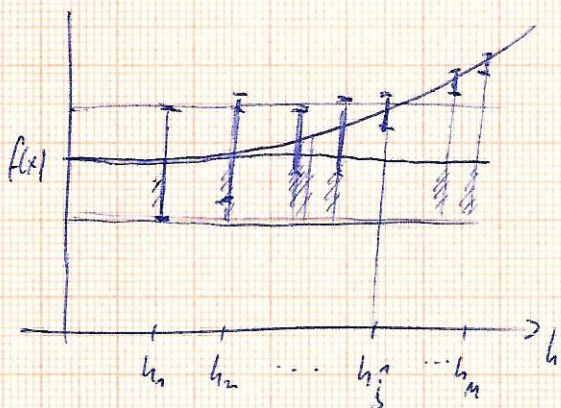
thus, $\mathbb{P}(f(x) \in \bigcap_{j=1}^{j^*} I_j) \geq 1 - \delta \Rightarrow \bigcap_{j=1}^{j^*} I_j \neq \emptyset$ and we define:

$$\hat{j} = \max \{J = 1, \dots, M : \bigcap_{j=1}^J I_j \neq \emptyset\}$$

$$\text{and } \hat{f}(x) = \hat{f}_{\hat{j}}(x)$$

Lemma: on $\Omega_0 = \left\{ \max_{1 \leq j \leq M} \left| \frac{f_j(x)}{e_j} \right| \leq \gamma \right\}$ it holds that

$$|f_j^*(x) - f(x)| \leq b e_{j^*}$$



Assumption: $\forall x \in [0,1], f(x) \leq L$ and L is known.

Parameter: ε (level of confidence)

Initialization: $j=0, D=[0,L]$

while $D \neq \emptyset$ and $j < M$:

$$j = j+1, h = h_j$$

$$\hat{f} = \frac{1}{2nh} \sum_{i=1}^n \mathbb{1}_{[x-h, x+h]}(X_i)$$

$$e = \sqrt{\frac{3L}{2nh} \ln\left(\frac{2M}{\varepsilon}\right)}$$

$$I = [\hat{f} - 2e, \hat{f} + 2e]$$

$$D = D \cap I$$

endwhile

return: $\hat{j} = j-1$

$$\hat{h} = h^{\hat{j}}$$

$$\hat{f} = \hat{f}_{h, h^{\hat{j}}}$$

Theorem: if $h_1 < h_2 < \dots < h_M$ and if f is a density on $[0,1]$ bounded by L , if $e_j = \sqrt{\frac{3L}{2nh_j} \ln\left(\frac{2M}{1-\varepsilon}\right)}$ then with probab. at least $1-\varepsilon$

$$j^* = \max \left\{ j \in \{1, \dots, M\} : \left| \mathbb{E} \left[\hat{f}_{h, h_j}^*(x) \right] - f(x) \right| \leq e_j \right\}$$

and if $h_j = h_n \cdot a^{j-1}$ ($a > 1$), if $|\mathbb{E}[\hat{f}_{n,h}^{\text{LC}}(x)] - f(x)| \leq Ch^\beta$

then with probab. $\geq 1 - \varepsilon$

$$|\hat{f}_n(x) - f(x)| \leq \text{Constant} \cdot n^{-\frac{\beta}{2\beta+1}} \cdot \left(\ln \frac{n}{\varepsilon}\right)^{\frac{\beta}{2\beta+1}}$$