

ON A NON-PARAMETRIC ESTIMATION OF THE REGRESSION FUNCTION

by

P. MAJOR

The regression function is usually estimated by a polynomial. In this paper an other method is investigated, similar to one of the methods of density estimation.

It is the following method: Let (ξ_i, η_i) $i=1, 2, \dots$ be independent identically distributed two-dimensional random variables, $P(0 \leq \eta_i \leq 1) = 1$. Let us divide the interval $[0, 1]$ into a_n equal parts. We estimate the function $R(x) = E(\xi | \eta = x)$ by the following $R_n(x)$

$$(1) \quad R_n(x) = \frac{\sum_{\substack{\eta_i \in \left[\frac{k}{a_n}, \frac{k+1}{a_n}\right] \\ (i \leq n)}} \xi_i}{\text{the number of } \eta_i - s \text{ } (i \leq n) \text{ falling into the interval } \left[\frac{k}{a_n}, \frac{k+1}{a_n}\right]}$$

$$x \in \left[\frac{k}{a_n}, \frac{k+1}{a_n}\right] \quad (k = 1, 2, \dots, a_n).$$

We investigate the order of magnitude of $R_n(x) - R(x)$ in supremum norm. We introduce some notations. $f(x)$ is the density function of η_i ,

$$\sigma^2(x) = E([\xi_i - E(\xi_i | \eta_i)]^2 | \eta_i = x).$$

Our first statement is the following

THEOREM 1. *Let us suppose that $f(x)$, $\sigma(x)$, $R(x)$ are differentiable, $|f'(x)| \leq K$, $|\sigma'(x)| \leq K$, $|R'(x)| \leq K$ and $f(x) > K'$, $\sigma(x) > K'$ for appropriate $K, K' > 0$ and for every x . Let us suppose further that $E(e^{t_0 |\xi|} | \eta = x) \leq c$ for appropriate $t_0 > 0$ and c for every $0 \leq x \leq 1$. If $n^\alpha < a_n < n^\beta$, $\frac{1}{3} < \alpha < \beta < 1$ then*

$$(2) \quad P \left(\sup_x \sqrt{\frac{nf(x)}{a_n}} \cdot \frac{|R_n(x) - R(x)|}{\sigma(x)} < \sqrt{2 \log a_n - \log \log a_n + y} \right) \rightarrow \exp \left(-\frac{e^{-y/2}}{\sqrt{\pi}} \right)$$

$$(3) \quad P \left(\sup_x \sqrt{\frac{nf(x)}{a_n}} \frac{R_n(x) - R(x)}{\sigma(x)} < \sqrt{2 \log a_n - \log \log a_n + y} \right) \rightarrow \exp \left(-\frac{e^{-y/2}}{2\sqrt{\pi}} \right)$$

as $n \rightarrow \infty$.

If the functions $f(x)$, $R(x)$, $\sigma(x)$ are smooth enough, it is worth substituting the estimation $R_n(x)$ by

$$(4) \quad \bar{R}_n(x) = R_n\left(\frac{k-\frac{1}{2}}{a_n}\right) + \left[R_n\left(\frac{k+\frac{1}{2}}{a_n}\right) - R_n\left(\frac{k-\frac{1}{2}}{a_n}\right)\right]\left(a_n x - k - \frac{1}{2}\right)$$

for

$$k - \frac{1}{2} < a_n x < k + \frac{1}{2}.$$

Then we have our.

THEOREM 1'. *If the conditions of Theorem 1 are fulfilled and, in addition $f(x)$, $R(x)$ and $\sigma(x)$ are differentiable twice with bounded second derivatives, then the relations (2) and (3) hold for $n^\alpha < a_n < n^\beta$, $\frac{1}{5} < \alpha < \beta < 1$, if we substitute $R_n(x)$ by $\bar{R}_n(x)$.*

Now if we want to construct a "confidence strip" i.e. we want to construct a region T in the plane such that $P(R(x) \in T) > 1 - \varepsilon$ with a prescribed ε , then we are interested in whether $f(x)$ and $\sigma(x)$ in the formulae (2) and (3) can be substituted by their estimations. The answer is in the affirmative.

THEOREM 2. *Let*

$$f_n(x) = \frac{a_n}{n} \left\{ \text{the number of } \eta_i - s \ (i \leq n) \text{ for which } \eta_i \in \left[\frac{k}{a_n}, \frac{k+1}{a_n} \right] \right\}$$

$$\sigma_n^2(x) = \frac{a_n}{nf_n(x)} \sum_{\substack{\eta_i \in \left[\frac{k}{a_n}, \frac{k+1}{a_n} \right] \\ i \leq n}} \xi_i^2 - \left[\frac{a_n}{nf_n(x)} \sum_{\substack{\eta_i \in \left[\frac{k}{a_n}, \frac{k+1}{a_n} \right] \\ i \leq n}} \xi_i \right]^2$$

if the condition of Theorem 1 holds

$$x \in \left[\frac{k}{a_n}, \frac{k+1}{a_n} \right] \quad k = 0, 1, \dots, a_n - 1.$$

Then

$$(2') \quad P \left(\sup_x \sqrt{\frac{nf_n(x)}{a_n}} \cdot \frac{|R_n(x) - R(x)|}{\sigma_n(x)} < \sqrt{2 \log a_n - \log \log a_n + y} \right) \rightarrow \exp \left(\frac{-e^{-y/2}}{\sqrt{\pi}} \right)$$

and

$$(3') \quad P \left(\sup_x \sqrt{\frac{nf_n(x)}{a_n}} \cdot \frac{|R_n(x) - R(x)|}{\sigma_n(x)} < \sqrt{2 \log a_n - \log \log a_n + y} \right) \rightarrow \exp \left(\frac{-e^{-y/2}}{2\sqrt{\pi}} \right)$$

if.

As an immediate consequence of Theorem 1 we get that

$$\sup_x \sqrt{\frac{nf(x)}{2a_n \log a_n}} \cdot \frac{|R_n(x) - R(x)|}{\sigma(x)} \rightarrow 1$$

where the symbol \Rightarrow means convergence in probability. The question arises whether this expression is also convergent with probability 1.

THEOREM 3. *Under the conditions of theorem 1 we have*

$$\mathbb{P} \left(\limsup_{n \rightarrow \infty} \sqrt{\frac{nf(x)}{2a_n \log a_n}} \frac{|R_n(x) - R(x)|}{\sigma(x)} = 1 \right) = 1$$

$$\mathbb{P} \left(\limsup_{n \rightarrow \infty} \sqrt{\frac{nf(x)}{2a_n \log a_n}} \frac{R_n(x) - R(x)}{\sigma(x)} = 1 \right) = 1$$

if

$$n^\alpha < a_n < n^\beta, \quad \frac{1}{3} < \alpha < \beta < 1.$$

For proving our theorems we need some lemmas:

LEMMA 1. *Let $\xi_{n1}, \xi_{n2} \dots \xi_{nn}$, $n=1, 2 \dots$ be a double array of random variables independent and identically distributed within a row. Let $E\xi_{n1}=0$, $D^2\xi_{n1}=1$, $Ee^{t|\xi_{n1}|} < c$ for appropriate $t>0$ and c . Then*

$$\mathbb{P}(\xi_{n1} + \dots + \xi_{nn} > \sqrt{n} \cdot x_n) \sim 1 - \Phi(x_n) \quad \text{for } x_n = o(\sqrt[n]{n})$$

where $\Phi(x)$ is the standard normal distribution function.

The proof is essentially the same as in [2] page 517, and we omit it.

LEMMA 2. *Let $(\xi_i^{(n)}, \eta_i^{(n)})$ be pairs of random variables independent and identically distributed for fixed n ,*

$$\mathbb{P}(\eta_i^{(n)} = k) = p_k^{(n)}, \quad k = 1, 2, \dots, a_n,$$

$$\sum_{k=1}^{a_n} p_k^{(n)} = 1, \quad n^\varepsilon < a_n < n^{1-\varepsilon}, \quad \varepsilon > 0,$$

$$\mathbb{E}(\xi_i^{(n)} | \eta_i^{(n)} = k) = 0, \quad \mathbb{E}(\xi_i^{(n)2} | \eta_i^{(n)} = k) = \sigma_k^{(n)2}, \quad \frac{c_1}{a_n} < p_k^{(n)} < \frac{c_2}{a_n}$$

$$\mathbb{E}(e^{t|\xi_i^{(n)}|} | \eta_i^{(n)}) \leq c, \quad \sigma_k^{(n)} > c'$$

for appropriate $c_1, c_2, c, t > 0$, c' for every k and n . Then

$$\mathbb{P} \left(\max_k \frac{\sqrt{np_k^{(n)}} \cdot \sum_{i=1}^n \xi_i^{(n)} \cdot I_{(\eta_i^{(n)}=k)}}{\sum_{i=1}^n I_{(\eta_i^{(n)}=k)}} < \sqrt{2 \log a_n - \log \log a_n + y} \right) \rightarrow \exp \left(\frac{-e^{-y/2}}{2\sqrt{\pi}} \right)$$

and

$$\mathbb{P} \left(\max_k \frac{\sqrt{np_k^{(n)}}}{\sigma_k^{(n)}} \cdot \left| \frac{\sum_{i=1}^n \xi_i^{(n)} \cdot I_{(\eta_i^{(n)}=k)}}{\sum_{i=1}^n I_{(\eta_i^{(n)}=k)}} \right| < \sqrt{2 \log a_n - \log \log a_n + y} \right) \rightarrow \exp \left(\frac{-e^{-y/2}}{\sqrt{\pi}} \right)$$

where $I_{(\eta_i^{(n)}=k)} = 1$ if $\eta_i^{(n)} = k$ and 0 otherwise.

PROOF OF LEMMA 2. In the following we omit the superscript n . Let us introduce the notations $s_j = \sum_{i=1}^n \xi_i I_{(\eta_i=j)}$ and $v_j = \sum_{i=1}^n I_{(\eta_i=j)}$. Then

$$\begin{aligned} \mathbb{P}(s_1 < x_1, \dots, s_{a_n} < x_{a_n} | v_1 = t_1, \dots, v_{a_n} = t_{a_n}) &= \\ &= \mathbb{P}(s_1 < x_1 | v_1 = t_1) \dots \mathbb{P}(s_{a_n} < x_{a_n} | v_{a_n} = t_{a_n}) \end{aligned}$$

Let $\tilde{\xi}_{ij}$, $i=1, 2 \dots j=1, 2 \dots a_n$ be — for fixed j — independent identically distributed random variables with distribution function $\mathbb{P}(\tilde{\xi}_{ij} < x) = \mathbb{P}(\xi_1 = x_1 | \eta_1 = j)$. Then

$$\mathbb{P}(s_j < x_j | v_j = t_j) = \mathbb{P}(\tilde{\xi}_{1j} + \dots + \tilde{\xi}_{tj,j} < x_j)$$

Now we prove that

$$(5) \quad \mathbb{P}(|v - np_j| > \sqrt{np_j} n^{\frac{1-\varepsilon}{3}}) \leq 2 \exp \left(-\frac{n^{\frac{2}{3}(1-\varepsilon)}}{3} \right)$$

Really

$$\begin{aligned} \mathbb{P}(v_j - np_j > \sqrt{np_j} n^{\frac{1-\varepsilon}{3}}) &= \mathbb{P} \left(\sum_{i=1}^n [I_{(\eta_i=j)} - p_j] > \sqrt{np_j} n^{\frac{1-\varepsilon}{3}} \right) = \\ &= \mathbb{P} \left(\exp \left(n^{\frac{1-\varepsilon}{3}} \sum_{i=1}^n \frac{I_{(\eta_i=j)} - p_j}{\sqrt{np_j}} \right) > \exp(n^{\frac{2}{3}(1-\varepsilon)}) \right) \leq \\ &\leq \frac{\left[\mathbb{E} \left[\exp \left(n^{\frac{1-\varepsilon}{3}} \frac{I_{(\eta_i=j)} - p_j}{\sqrt{np_j}} \right) \right] \right]^n}{\exp(n^{\frac{2}{3}(1-\varepsilon)})} \leq \frac{\left[1 + \frac{2}{3} n^{\frac{2}{3}(1-\varepsilon)} \mathbb{E} \left(\frac{I_{(\eta_1=j)} - p_j}{\sqrt{np_j}} \right)^2 \right]^n}{\exp(n^{\frac{2}{3}(1-\varepsilon)})} \leq \\ &\leq \frac{\exp \left[\frac{2}{3} (1-p_j) n^{\frac{2}{3}(1-\varepsilon)} \right]}{\exp(n^{\frac{2}{3}(1-\varepsilon)})} \leq \exp \left(-\frac{n^{\frac{2}{3}(1-\varepsilon)}}{3} \right) \end{aligned}$$

and similarly

$$\mathbb{P}(v_j - np_j < -\sqrt{np_j} n^{\frac{1-\varepsilon}{3}}) = \exp \left(-\frac{n^{\frac{2}{3}(1-\varepsilon)}}{3} \right).$$

Thus (5) is valid.

Now, by lemma 1 and by the relation $\frac{1}{\sqrt{2\pi}} e^{-x^2/2} \sim 1 - \Phi(x)$ for $x \rightarrow \infty$ we have

$$\begin{aligned} & \mathbb{P}\left(\frac{\sqrt{np_j}}{\sigma_j} \cdot \frac{s_j}{v_j} > \sqrt{2 \log a_n - \log \log a_n + y} \mid v_j = t\right) = \\ &= \mathbb{P}\left(\frac{\tilde{\xi}_{1j} + \dots + \tilde{\xi}_{tj}}{\sigma_j \sqrt{t}} > \sqrt{\frac{t}{np_j}} \sqrt{2 \log a_n - \log \log a_n + y}\right) \sim \\ &\sim \sqrt{\frac{np_j}{2\pi t}} \cdot \frac{1}{\sqrt{2 \log a_n - \log \log a_n + y}} \cdot e^{-\frac{np_j}{2t}(2 \log a_n - \log \log a_n + y)} \sim \frac{e^{-y/2}}{2\sqrt{\pi} a_n} \end{aligned}$$

if

$$|t - np_j| < \sqrt{np_j} n^{\frac{1-\varepsilon}{3}}$$

On the other hand from (5) we have

$$\mathbb{P}(|v_k - np_k| > \sqrt{np_k} n^{\frac{1-\varepsilon}{3}} \text{ for some } k) \leq 2a_n \exp\left(-\frac{n^{\frac{2}{3}(1-\varepsilon)}}{3}\right) = o(1)$$

Thus

$$\begin{aligned} & \mathbb{P}\left(\max_k \frac{\sqrt{np_k}}{\sigma_k} \cdot \frac{s_k}{v_k} < \sqrt{2 \log a_n - \log \log a_n + y}\right) = \\ &= \left[1 - \frac{e^{-y/2}}{2\sqrt{\pi} a_n} (1 + o(1))\right]^{a_n} \mathbb{P}(|v_k - np_k| < \sqrt{np_k} n^{\frac{1-\varepsilon}{3}} \text{ for every } k) + o(1) = \\ &= \exp\left(-\frac{e^{-y/2}}{2\sqrt{\pi}}\right) + o(1) \end{aligned}$$

as we have stated.

The second relation of lemma 2 can be proved in the same way.

LEMMA 3. *Under the conditions of lemma 2*

$$\mathbb{P}\left(\max_k \frac{\sqrt{np_k^{(n)}}}{\sigma_k^{(n)}} \cdot \frac{\sum_{i=1}^n \zeta_i^{(n)} I_{(\eta_i^{(n)}=k)}}{\sum_{i=1}^n I_{(\eta_i^{(n)}=k)}} < \sqrt{2 \log a_n - 3 \log \log a_n - 2 \log 2\sqrt{\pi} c}\right) = \frac{1}{a_n^{c+o(1)}}$$

$$\mathbb{P}\left(\max_k \frac{\sqrt{np_k^{(n)}}}{\sigma_k^{(n)}} \cdot \left| \frac{\sum_{i=1}^n \zeta_i^{(n)} I_{(\eta_i^{(n)}=k)}}{\sum_{i=1}^n I_{(\eta_i^{(n)}=k)}} \right| < \sqrt{2 \log a_n - 3 \log \log a_n - 2 \log 2\sqrt{\pi} c}\right) = \frac{1}{a_n^{c+o(1)}}$$

The proof is similar to that of Lemma 2.

LEMMA 4. Under the conditions of Lemma 2 we have

$$\begin{aligned} \mathbb{P} \left(\max_k \frac{\sqrt{np_k^{(n)}}}{\sigma_k^{(n)}} \left(1 + o\left(\frac{1}{\log a_n}\right) \right) \cdot \frac{\sum_{i=k}^n I_{(\eta_i^{(n)}=k)} \left[\zeta_i^{(n)} + o\left(\sqrt{\frac{a_n}{n \log a_n}}\right) \right]}{\sum_{i=1}^n I_{(\eta_i^{(n)}=k)}} < \sqrt{2 \log a_n - \log \log a_n + y} \right) \rightarrow \exp\left(-\frac{e^{-y/2}}{2\sqrt{\pi}}\right) \\ \mathbb{P} \left(\max_k \frac{\sqrt{np_k^{(n)}}}{\sigma_k^{(n)}} \left(1 + o\left(\frac{1}{\log a_n}\right) \right) \cdot \left| \frac{\sum_{i=1}^n I_{(\eta_i^{(n)}=k)} \left[\zeta_i^{(n)} + o\left(\sqrt{\frac{a_n}{n \log a_n}}\right) \right]}{\sum_{i=1}^n I_{(\eta_i^{(n)}=k)}} \right| < \sqrt{2 \log a_n - \log \log a_n + y} \right) \rightarrow \exp\left(-\frac{e^{-y/2}}{\sqrt{\pi}}\right). \end{aligned}$$

Lemma 4 is an easy consequence of Lemma 2.

PROOF OF THEOREM 1. Let us divide the interval $[0,1]$ into a_n equal parts. Obviously

$$m_k^{(n)} = \mathbb{E} \left(\zeta_i \middle| \frac{k-1}{a_n} \leq \eta_i < \frac{k}{a_n} \right) = \frac{\int_{\frac{k-1}{a_n}}^{\frac{k}{a_n}} R(y)f(y) dy}{\int_{\frac{k-1}{a_n}}^{\frac{k}{a_n}} f(y) dy}$$

and

$$D_k^{2(n)} = D^2 \left(\zeta_i \middle| \frac{k-1}{a_n} \leq \eta_i < \frac{k}{a_n} \right) = \frac{\int_{\frac{k-1}{a_n}}^{\frac{k}{a_n}} [\sigma^2(y) + R^2(y)]f(y) dy}{\int_{\frac{k-1}{a_n}}^{\frac{k}{a_n}} f(y) dy} - [m_k^{(n)}]^2.$$

Let us define the random variables $(\zeta_i^{(n)}, \eta_i^{(n)})$ $i=1, 2 \dots n$ in the following way:

$$\eta_i^{(n)} = k \quad \text{and} \quad \zeta_i^{(n)} = \frac{\zeta_i - m_k^{(n)}}{D_k^{(n)}} \cdot \sigma \left(\frac{k - \frac{1}{2}}{a_n} \right)$$

if

$$\frac{k-1}{a_n} \leq \eta_i < \frac{k}{a_n} \quad i = 1, 2, \dots, n \quad k = 1, 2, \dots, a_n.$$

Then the conditions of Lemma 2 are satisfied for the sequence $(\xi_i^{(n)}, \eta_i^{(n)})$. Indeed,

$p_k^{(n)} = \int_{\frac{k-1}{a_n}}^{\frac{k}{a_n}} f(y) dy$, thus $\frac{c_1}{a_n} < p_k^{(n)} < \frac{c_2}{a_n}$ since $f(x)$ is bounded both from above and below. $E(\xi_i^{(n)} | \eta_i^{(n)} = k) = 0$, $(\sigma_k^{(n)})^2 = D^2(\xi_i^{(n)} | \eta_i^{(n)} = k) = \sigma^2 \left(\frac{k - \frac{1}{2}}{a_n} \right)$ is bounded from below. $\frac{D_k^{(n)}}{\sigma \left(\frac{k - \frac{1}{2}}{a_n} \right)}$ is bounded from below and

$$E \left(\exp \left(t_0 \frac{D_k^{(n)}}{\sigma \left(\frac{k - \frac{1}{2}}{a_n} \right)} | \xi_i^{(n)} \right) \middle| \eta_i^{(n)} = k \right) \equiv c \exp(t_0 m_k^{(n)})$$

is also bounded as required. By the Taylor formula

$$f(y) = f(x) + (y - x) \cdot f'(x + \theta(y - x))$$

$$R(y) = R(x) + (y - x) \cdot R'(x + \theta'(y - x)) \quad 0 < \theta, \theta', \theta'' < 1$$

$$\sigma(y) = \sigma(x) + (y - x) \cdot \sigma'(x + \theta''(y - x))$$

The first derivatives of $f(x)$, $R(x)$, $\sigma(x)$ are bounded and so it can be seen easily that $m_k^{(n)} = R(x) + O\left(\frac{1}{a_n}\right)$, $p_k^{(n)} = \frac{f(x)}{a_n} + O\left(\frac{1}{a_n^2}\right)$, $D^2(\xi_i^{(n)} | \eta_i^{(n)} = k) = \sigma^2(x) + O\left(\frac{1}{a_n}\right)$ if $x \in \left[\frac{k-1}{a_n}, \frac{k}{a_n}\right]$. Since $n^\alpha < a_n$ and $\alpha > \frac{1}{3}$, $\frac{1}{a_n} = o\left(\sqrt{\frac{a_n}{n \log a_n}}\right)$ and

$$\sqrt{\frac{nf(x)}{a_n}} \frac{1}{\sigma(x)} = \frac{\sqrt{np_k^{(n)}}}{\sigma_k^{(n)}} \left(1 + o\left(\frac{1}{\log a_n}\right) \right) \quad \text{for } x \in \left[\frac{k-1}{a_n}, \frac{k}{a_n}\right].$$

Hence Lemma 4 can be applied, and we get

$$\begin{aligned} P \left(\sup_k \sup_{k \in \left[\frac{k-1}{a_n}, \frac{k}{a_n}\right]} \sqrt{\frac{nf(x)}{a_n}} \frac{\sum_{i=1}^n I_{(\eta_i^{(n)} = k)} [\xi_i - R(x)]}{\sigma(x) \sum_{i=1}^n I_{(\eta_i^{(n)} = k)}} < \sqrt{2 \log a_n - \log \log a_n + y} \right) \rightarrow \\ \rightarrow \exp \left(-\frac{e^{-y/2}}{2\sqrt{\pi}} \right) \end{aligned}$$

which is identical to formula (3). The second statement of lemma 4 gives formula (4).

The PROOF of Theorem 1' is almost the same. By the Taylor expansion up to 2 terms

$$\begin{aligned} R(x) &= R\left(\frac{k-\frac{1}{2}}{a_n}\right) + R'\left(\frac{k-\frac{1}{2}}{a_n}\right)\left(x - \frac{k-\frac{1}{2}}{a_n}\right) + \\ &\quad + \frac{1}{2} R''\left(\frac{k-\frac{1}{2}}{a_n} + \theta\left(x - \frac{k-\frac{1}{2}}{a_n}\right)\right)\left(x - \frac{k-\frac{1}{2}}{a_n}\right)^2 \\ f(x) &= f\left(\frac{k-\frac{1}{2}}{a_n}\right) + f'\left(\frac{k-\frac{1}{2}}{a_n}\right)\left(x - \frac{k-\frac{1}{2}}{a_n}\right) + \\ &\quad + \frac{1}{2} f''\left(\frac{k-\frac{1}{2}}{a_n} + \theta'\left(x - \frac{k-\frac{1}{2}}{a_n}\right)\right)\left(x - \frac{k-\frac{1}{2}}{a_n}\right)^2 \\ \sigma(x) &= \sigma\left(\frac{k-\frac{1}{2}}{a_n}\right) + \sigma'\left(\frac{k-\frac{1}{2}}{a_n}\right)\left(x - \frac{k-\frac{1}{2}}{a_n}\right) + \\ &\quad + \frac{1}{2} \sigma''\left(\frac{k-\frac{1}{2}}{a_n} + \theta''\left(x - \frac{k-\frac{1}{2}}{a_n}\right)\right)\left(x - \frac{k-\frac{1}{2}}{a_n}\right)^2 \end{aligned}$$

where $0 < \theta, \theta', \theta'' < 1$. It implies that

$$m_k^{(n)} = R\left(\frac{k-\frac{1}{2}}{a_n}\right) + O\left(\frac{1}{a_n^2}\right), \quad p_k^{(n)} = \frac{f\left(\frac{k-\frac{1}{2}}{a_n}\right)}{a_n} + O\left(\frac{1}{a_n^3}\right)$$

and

$$D^2(\zeta_i^{(n)} | \eta_i^{(n)} = k) = \sigma^2\left(\frac{k-\frac{1}{2}}{a_n}\right) + O\left(\frac{1}{a_n^2}\right),$$

and these imply, just as in Theorem 1, the limit relation

$$\begin{aligned} \mathbb{P} \left(\sup_k \sqrt{\frac{nf\left(\frac{k-\frac{1}{2}}{a_n}\right)}{a_n}} \frac{R_n\left(\frac{k-\frac{1}{2}}{a_n}\right) - R\left(\frac{k-\frac{1}{2}}{a_n}\right)}{\sigma\left(\frac{k-\frac{1}{2}}{a_n}\right)} < \sqrt{2 \log a_n - \log \log a_n + y} \right) &\rightarrow \\ &\rightarrow \exp\left(\frac{-e^{-y/2}}{2\sqrt{\pi}}\right). \end{aligned}$$

Now by

$$R(x) = R\left(\frac{k-\frac{1}{2}}{a_n}\right) + \left[R\left(\frac{k+\frac{1}{2}}{a_n}\right) - R\left(\frac{k-\frac{1}{2}}{a_n}\right) \right] \left(a_n x - k - \frac{1}{2}\right) + O\left(\frac{1}{a_n^2}\right)$$

for $x \in \left[\frac{k-\frac{1}{2}}{a_n}, \frac{k+\frac{1}{2}}{a_n}\right]$ and (4) the statement of Theorem 1' is valid.

PROOF OF THEOREM 2. Since we have already proved Theorem 1, it is sufficient to prove that

$$\mathbb{P} \left(\sup_x \frac{\sqrt{f(x)}}{\sigma(x)} \frac{\sigma_n(x)}{\sqrt{f_n(x)}} > 1 + \frac{\delta_n}{\sqrt{\log a_n}} \right) \rightarrow 0$$

and

$$\mathbb{P} \left(\inf_x \frac{\sqrt{f(x)}}{\sigma(x)} \frac{\sigma_n(x)}{\sqrt{f_n(x)}} > 1 - \frac{\delta_n}{\sqrt{\log a_n}} \right) \rightarrow 0$$

for an appropriate sequence $\delta_n \rightarrow 0$, i.e. $\frac{\sqrt{f(x)}}{\sigma(x)}$ and $\frac{\sqrt{f_n(x)}}{\sigma_n(x)}$ are near to each other. We shall estimate $f_n(x) - f(x)$ and $\sigma_n(x) - \sigma(x)$. This will be similar to the proof of Lemma 2.

$$f(x) - a_n p_k^{(n)} = O \left(\frac{1}{a_n} \right) \quad x \in \left[\frac{k+1}{a_n}, \frac{k}{a_n} \right]$$

and

$$f_n(x) - a_n p_k^{(n)} = \frac{a_n}{n} \sum_{i=1}^n [I_{(\eta_i^{(n)}=k)} - \mathbb{E} I_{(\eta_i^{(n)}=k)}].$$

Thus, exactly the same way as in proving estimation (5) we get that

$$\mathbb{P} \left(|f_n(x) - f(x)| > n^{\frac{1-\alpha}{3}} \frac{\sqrt{n}}{a_n \sqrt{p_k^{(n)}}} \right) = O(\exp(-n^{\frac{1-\alpha}{3}})).$$

$n^{\frac{1-\alpha}{3}} \frac{\sqrt{n}}{a_n \sqrt{p_k^{(n)}}} < K \cdot n^{-\frac{1-\alpha}{6}}$ for an appropriate constant K . Thus

$$\mathbb{P} \left(\sup_x |f_n(x) - f(x)| > K \cdot n^{-\frac{1-\alpha}{6}} \right) \leq a_n O(\exp(-n^{\frac{1-\alpha}{3}})) \rightarrow 0$$

We have

$$\sigma^2(x) - D_k^{(n)} = O \left(\frac{1}{a_n} \right) \quad \text{and} \quad \sigma_n^2(x) = \frac{\sum_{i=1}^n \xi_i^2 I_{(\eta_i^{(n)}=k)}}{\sum_{i=1}^n I_{(\eta_i^{(n)}=k)}} - \left[\frac{\sum_{i=1}^n \xi_i I_{(\eta_i^{(n)}=k)}}{\sum_{i=1}^n I_{(\eta_i^{(n)}=k)}} \right]^2$$

for

$$x \in \left[\frac{k-1}{a_n}, \frac{k}{a_n} \right].$$

Now we estimate

$$\mathbb{P} \left(\left| \frac{\sum_{i=1}^n \xi_i I_{(\eta_i^{(n)}=k)}}{\sum_{i=1}^n I_{(\eta_i^{(n)}=k)}} - m_k^{(n)} \right| > \frac{n^{\frac{1-\alpha}{3}}}{\sqrt{np_k^{(n)}}} \right).$$

Let $\tilde{\xi}_1, \dots, \tilde{\xi}_n \dots n=1, 2, \dots$ be independent identically distributed random variables having the distribution

$$\mathbb{P}(\tilde{\xi}_i < x) = \mathbb{P}(\xi_1 - m_k^{(n)} < x | \eta_1^{(n)} = k)$$

Then

$$\begin{aligned} \mathbb{P}\left(\left|\frac{\sum_{i=1}^n \xi_i I_{(\eta_i^{(n)}=k)}}{\sum_{i=1}^n I_{(\eta_i^{(n)}=k)}} - m_k^{(n)}\right| > \frac{n^{\frac{1-\alpha}{3}}}{\sqrt{np_k^{(n)}}} \left|\sum_{i=1}^n I_{(\eta_i^{(n)}=k)} = t\right.\right) = \\ = \mathbb{P}\left(\left|\frac{\bar{\xi}_1 + \dots + \bar{\xi}_t}{\sqrt{t}}\right| < n^{\frac{1-\alpha}{3}} \sqrt{\frac{t}{np_k^{(n)}}}\right) = O\left(\exp(-n^{\frac{1-\alpha}{3}})\right) \end{aligned}$$

for $|t - np_k^{(n)}| < n^{\frac{1-\alpha}{3}} \sqrt{np_k^{(n)}}$. On the other hand

$$\mathbb{P}\left(\left|\sum_{i=1}^n I_{(\eta_i^{(n)}=k)} - np_k^{(n)}\right| > n^{\frac{1-\alpha}{3}} \sqrt{np_k^{(n)}}\right) = O\left(\exp(-n^{\frac{1-\alpha}{3}})\right)$$

Thus

$$\mathbb{P}\left(\left|\frac{\sum_{i=1}^n \xi_i I_{(\eta_i^{(n)}=k)}}{\sum_{i=1}^n I_{(\eta_i^{(n)}=k)}} - m_k^{(n)}\right| > \frac{n^{\frac{1-\alpha}{3}}}{\sqrt{np_k^{(n)}}} \text{ for some } k\right) \equiv a_n O\left(\exp(-n^{\frac{1-\alpha}{3}})\right) \rightarrow 0$$

The relation

$$\mathbb{P}\left(\left|\frac{\sum_{i=1}^n \xi_i^2 I_{(\eta_i^{(n)}=k)}}{\sum_{i=1}^n I_{(\eta_i^{(n)}=k)}} - \frac{\int_{\frac{k-1}{a_n}}^{\frac{k}{a_n}} [\sigma^2(y) + R^2(y)] f(y) dy}{\int_{\frac{k-1}{a_n}}^{\frac{k}{a_n}} f(y) dy}\right| > \frac{n^{\frac{1-\alpha}{3}}}{\sqrt{np_k^{(n)}}} \text{ for some } k\right) \rightarrow 0$$

can be proved in the same way. From these estimations it follows that $|f_n(x) - f(x)|$ and $|\sigma_n(x) - \sigma(x)|$ are less than n^{-c} for appropriate $c > 0$ with a probability near 1. Thus Theorem 2 is valid.

PROOF OF THEOREM 3. First we prove that

$$\mathbb{P}\left(\limsup_x \sqrt{\frac{n f(x)}{2a_n \log a_n}} \frac{|R_n(x) - R(x)|}{\sigma(x)} > 1 - \varepsilon\right) = 1$$

for arbitrary $\varepsilon > 0$.

It is easy to check the validity of the conditions of Lemma 2 for the sequence $(\xi_i^{(n)}, \eta_i^{(n)})$, therefore by Lemma 3 we can state that

$$\begin{aligned} \mathbb{P} \left(\max_k \sqrt{np_k^{(n)}} \left| \frac{\sum_{i=1}^n \xi_i^{(n)} I_{(\eta_i^{(n)}=k)}}{\sigma \left(\frac{k-\frac{1}{2}}{a_n} \right) \sum_{i=1}^n I_{(\eta_i^{(n)}=k)}} \right| > \sqrt{2 \log a_n - 3 \log \log a_n - 2 \log \sqrt{\pi} c} \right) = \\ = \frac{1}{a_n^{c+o(1)}} \end{aligned}$$

Thus, choosing n_0 large enough, by the Borel—Cantelli lemma we obtain

$$\begin{aligned} \max_k \sqrt{np_k^{(n)}} \left| \frac{\sum_{i=1}^n \xi_i^{(n)} I_{(\eta_i^{(n)}=k)}}{\sigma \left(\frac{k-\frac{1}{2}}{a_n} \right) \sum_{i=1}^n I_{(\eta_i^{(n)}=k)}} \right| < \\ < \sqrt{2 \log a_n - 3 \log \log a_n - 2 \log \sqrt{\pi} c} \quad \text{for } n > n_0(\omega) \end{aligned}$$

with probability 1. Since

$$\sqrt{np_k^{(n)}} = \sqrt{\frac{nf(x)}{a_n}} \left(1 + O \left(\frac{1}{a_n} \right) \right), \quad \sigma \left(\frac{k-\frac{1}{2}}{a_n} \right) = \sigma(x) + O \left(\frac{1}{a_n} \right)$$

and

$$\xi_i^{(n)} I_{(\eta_i^{(n)}=k)} = \left[\xi_i - R(x) + O \left(\frac{1}{a_n} \right) \right] I_{(\eta_i^{(n)}=k)} \left[1 + O \left(\frac{1}{a_n} \right) \right]$$

if $x \in \left[\frac{k-1}{a_n}, \frac{k}{a_n} \right]$, our statement follows.

The proof of the relation

$$\mathbb{P} \left(\limsup_x \sqrt{\frac{nf(x)}{2a_n \log a_n}} \frac{|R_n(x) - R(x)|}{\sigma(x)} < 1 - \varepsilon \right) = 1$$

is more involved.

Set

$$\bar{\xi}_i = \frac{\xi_i - \mathbb{E}(\xi_i | \eta_i)}{\mathbb{D}(\xi_i | \eta_i)}.$$

Let us define the random variable

$$T_n(x) = \sum_{\substack{i \leq n \\ \eta_i \leq x}} \bar{\xi}_i \quad 0 \leq x \leq 1.$$

We want to prove that

(6)

$$\mathbb{P}(A_n) = \mathbb{P} \left(\sup_{\substack{y-x \leq \frac{1}{a_n} \\ 0 \leq x < y \leq 1}} \frac{1}{\sqrt{f(x)}} |T_n(x) - T_n(y)| > \left(1 + \varepsilon + \frac{\sqrt{\varepsilon}}{2}\right) \sqrt{\frac{2n \log a_n}{a_n}} \right) \leq \frac{12}{\varepsilon \cdot a_n^\varepsilon}.$$

First we estimate, for fixed x , the probability

$$\mathbb{P} \left(\sup_{x \leq y \leq x + \frac{3+\varepsilon}{3a_n}} \frac{1}{\sqrt{f(x)}} |T_n(x) - T_n(y)| > (1 + \varepsilon) \sqrt{\frac{2n \log a_n}{a_n}} \right).$$

Let us order the η_i -s falling into the interval $[x, x + \frac{3+\varepsilon}{3a_n}]$ in increasing order $x \leq \eta_1^* < \eta_2^* \dots < \eta_t^* < x + \frac{3+\varepsilon}{3a_n}$ and denote by $\bar{\xi}_i^*$ the $\bar{\xi}_i$ which is the pair of η_j^* . Then

$$\sup_{x < y < x + \frac{3+\varepsilon}{3a_n}} |T_n(y) - T_n(x)| = \sup_{k \leq t} \left| \sum_{j=1}^k \bar{\xi}_j^* \right|$$

We can state that

$$(7) \quad \mathbb{P} \left(\sup_{k \leq t} \sum_{j=1}^k \bar{\xi}_j^* > (1 + \varepsilon) \sqrt{\frac{2nf(x) \log a_n}{a_n}} \middle| \begin{array}{l} \text{the number of } \eta_i \text{-s } i \leq n \text{ falling} \\ \text{into the interval } [x, x + \frac{3+\varepsilon}{3a_n}] \\ \text{is } t, \eta_1^* = x_1, \dots \eta_t^* = x_t \end{array} \right) =$$

$$= \mathbb{P} \left(\sup_{k \leq t} \sum_{j=1}^k \zeta_j > (1 + \varepsilon) \sqrt{\frac{2nf(x) \log a_n}{a_n}} \right)$$

where $\zeta_1, \zeta_2, \dots, \zeta_t$ are independent and $\mathbb{P}(\zeta_j < y) = \mathbb{P}(\zeta_1 < y | \eta_1 = x_j)$. Now the sequence $\sum_{j=1}^k \zeta_j$, $k = 1, 2, \dots, t$ forms a martingale, so $\exp(t \sum_{j=1}^k \zeta_j)$ is a semimartingale for arbitrary t , and by a well-known martingale inequality see [1] p. 314.,

$$(8) \quad \mathbb{P} \left(\sup_{k \leq t} \sum_{j=1}^k \zeta_j > (1 + \varepsilon) \sqrt{\frac{2nf(x) \log a_n}{a_n}} \right) =$$

$$= \mathbb{P} \left(\sup_{k \leq t} \exp \left(\sqrt{\frac{2a_n \log a_n}{nf(x)}} \sum_{j=1}^k \zeta_j \right) > \exp(2(1 + \varepsilon) \log a_n) \right) \leq$$

$$\leq \frac{\prod_{i=1}^t \mathbb{E} \exp \left(\sqrt{\frac{2a_n \log a_n}{nf(x)}} \zeta_i \right)}{a_n^{2(1+\varepsilon)}}$$

Now

$$(9) \quad \begin{aligned} \mathbb{E} \exp \left(\sqrt{\frac{2a_n \log a_n}{nf(x)}} \zeta_i \right) &\equiv 1 + \mathbb{E} \sqrt{\frac{2a_n \log a_n}{nf(x)}} \zeta_i + \\ &+ \mathbb{E} \left[\frac{a_n \log a_n}{nf(x)} \zeta_i^2 \exp \left| \sqrt{\frac{2a_n \log a_n}{nf(x)}} \zeta_i \right| \right] \equiv 1 + \left(1 + \frac{\varepsilon}{2} \right) \frac{a_n \log a_n}{nf(x)} \equiv \\ &\equiv \exp \left[\left(1 + \frac{\varepsilon}{2} \right) \frac{a_n \log a_n}{nf(x)} \right] \end{aligned}$$

It can be proved in the same way as relation (5) that

$$\mathbb{P} \left(\left| t - \frac{n \left(1 + \frac{\varepsilon}{3} \right)}{a_n} f(x) \right| > n^{\frac{1-\alpha}{3}} \sqrt{\frac{n}{a_n} f(x)} \right) = O(\exp(-n^{\frac{1}{3}(1-\alpha)}))$$

$\left(t \text{ is the number of } \eta_i\text{'s } (i \leq n) \text{ falling into the interval } \left[x, x + \frac{3+\varepsilon}{3a_n} \right] \right)$, and so by (7), (8) and (9)

$$\mathbb{P} \left(\sup_{x < y < x + \frac{3+\varepsilon}{3a_n}} (T_n(y) - T_n(x)) > (1 + \varepsilon) \sqrt{\frac{2nf(x) \log a_n}{a_n}} \right) < \frac{1}{a_n^{1+\varepsilon}}$$

Using this relation for the sequence $(-\xi_1, \eta_1), \dots, (-\xi_n, \eta_n)$ we get that

$$\mathbb{P} \left(\sup_{x < y < x + \frac{3+\varepsilon}{3a_n}} |T_n(x) - T_n(y)| > (1 + \varepsilon) \sqrt{\frac{2nf(x) \log a_n}{a_n}} \right) < \frac{2}{a_n^{1+\varepsilon}}$$

In the same way (x is fixed) we find

$$\mathbb{P} \left(\sup_{x < y < x + \frac{\varepsilon}{3a_n}} |T_n(y) - T_n(x)| > \frac{\sqrt{\varepsilon}}{2} \sqrt{\frac{2nf(x) \log a_n}{a_n}} \right) < \frac{2}{a_n^{1+\varepsilon}}$$

Thus

$$\mathbb{P} \left(\sup_{x < y < x + \frac{3+\varepsilon}{3a_n}} \frac{1}{\sqrt{f(x)}} |T_n(y) - T_n(x)| > (1 + \varepsilon) \sqrt{\frac{2n \log a_n}{a_n}} \right) < \frac{6}{\varepsilon a_n^\varepsilon}$$

$$x = \frac{k\varepsilon}{3a_n} \left(k = 0, 1, \dots, \frac{3a_n}{\varepsilon} \right)$$

and

$$\mathbb{P} \left(\sup_{x < y < x + \frac{\varepsilon}{3a_n}} \frac{1}{\sqrt{f(x)}} |T_n(y) - T_n(x)| > \frac{\sqrt{\varepsilon}}{2} \sqrt{\frac{2n \log a_n}{a_n}} \right) < \frac{6}{\varepsilon a_n^\varepsilon}$$

$$x = \frac{k\varepsilon}{3a_n} \left(k = 0, 1, \dots, \frac{3a_n}{\varepsilon} \right).$$

Let us now consider an arbitrary interval $[u, v]$ such that $0 < v - u \leq \frac{1}{a_n}$. Choose the integer k in such a way that for the number $x = \frac{k\varepsilon}{3a_n}$ the inequality $0 < u - x < \frac{\varepsilon}{3a_n}$ hold. Then

$$|T_n(v) - T_n(u)| \leq \sup_{x < y < x + \frac{3+\varepsilon}{3a_n}} |T_n(y) - T_n(x)| + \sup_{x < y < x + \frac{\varepsilon}{3a_n}} |T_n(y) - T_n(x)|,$$

and therefore (6) follows from the last two inequalities. Choosing a $c > \frac{6}{\varepsilon}$ we get

$\mathbb{P}(A_{l^c}) < \frac{1}{l^2}$, and by the Borel—Cantelli lemma

$$(10) \quad \lim_{l \rightarrow \infty} A_{l^c} = \emptyset$$

Given a number n it determines a unique l for which $l^c \leq n < (l+1)^c$. Then $n - l^c < c(l+1)^{c-1} < 2cn^{c-1}$. By similar calculations as before we get

$$\mathbb{P}\left(\sum_{i=l^c}^n \bar{\xi}_i I_{(\eta_i^{(n)}=k)} > \varepsilon \sqrt{\frac{2nf(x) \log a_n}{a_n}}\right) < \exp(-Dn^{1/c})$$

with appropriate $D > 0$ for $k = 1, 2, \dots, a_n$, $x \in \left[\frac{k-1}{a_n}, \frac{k}{a_n}\right]$ and

$$\mathbb{P}(B_n) = \mathbb{P}\left(\sup_{\substack{k \\ x \in \left[\frac{k-1}{a_n}, \frac{k}{a_n}\right]}} \left| \frac{1}{\sqrt{f(x)}} \sum_{i=l^c}^n \bar{\xi}_i I_{(\eta_i^{(n)}=k)} \right| > \varepsilon \sqrt{\frac{2n \log a_n}{a_n}}\right) < a_n \exp(-Dn^{1/c})$$

This yields the relation

$$(11) \quad \lim B_n = \emptyset.$$

By the definition of A_n and B_n and relations (10) and (11) we have with probability 1

$$(12) \quad \left| \sup \frac{1}{\sqrt{f(x_k)}} \sum_{i=1}^n \bar{\xi}_i I_{(\eta_i^{(n)}=k)} \right| < \left(1 + 2\varepsilon + \frac{\sqrt{\varepsilon}}{2}\right) \sqrt{\frac{2n \log a_n}{a_n}}$$

for

$$x_k \in \left[\frac{k-1}{a_n}, \frac{k}{a_n}\right]$$

for every large n .

An easy computation shows that with probability one

$$(13) \quad \sum_{i=1}^n I_{(\eta_i^{(n)}=k)} > (1-\varepsilon) \frac{nf(x_k)}{a_n}, \quad x_k \in \left[\frac{k-1}{a_n}, \frac{k}{a_n}\right]$$

for every large n and $k=1, 2, \dots, a_n$. Since $R(x)$, $\sigma(x)$ and $f(x)$ are smooth enough

$$\mathbb{P} \left(\overline{\lim}_{n \rightarrow \infty} \sup_x \sqrt{\frac{nf(x)}{2a_n \log a_n}} \frac{|R_n(x) - R(x)|}{\sigma(x)} < \frac{1+2\varepsilon + \frac{\sqrt{\varepsilon}}{2}}{1-\varepsilon} \right) = 1$$

from the relations (12) and (13). Since it is true for arbitrary small $\varepsilon > 0$, the first statement of theorem 3 is valid. The second statement can be proved in the same way.

LITERATURE

- [1] DOOB, J. L.: *Stochastic Processes*, Wiley, New York (1953).
- [2] FELLER, W.: *An Introduction to Probability Theory and Its Applications* 2, Wiley, New York.
- [3] RÉVÉSZ, P.: On empirical density function, *Periodica Math. Hungar.* 2 (1972) 85—110.
- [4] Надарая, Э. А.: Замечания о непараметрических оценках плотности вероятности и кривой регрессии. *Теория вероят. и ее примен.* XV. 1 (1970) 139—142.

Mathematical Institute of the Hungarian Academy of Sciences, Budapest

(Received April 15, 1973)