

ON THE ESTIMATION OF DISTRIBUTION FUNCTION ON INDIRECT SAMPLE

Elizbar Nadaraya¹, Petre Babilua², Grigol Sokhadze³

Iv. Javakhishvili Tbilisi State University, Tbilisi, Georgia

¹elizbar.nadaraya@tsu.ge, ²petre.babilua@tsu.ge, ³grigol.sokhadze@tsu.ge

1. Let X_1, X_2, \dots, X_n be a sample of independent observations of a non-negative random value X with a distribution function $F(x)$. In problems of the theory of censored observations the sample values are pairs of random values $Y_i = (X_i \wedge t_i)$ and $Z_i = I(Y_i = X_i)$, $i = \overline{1, n}$, where t_i are given numbers ($t_i \neq t_j$ for $i \neq j$) or random values independent of X_i , $i = \overline{1, n}$. Throughout the paper $I(\cdot)$ denotes the indicator of the set A .

We will consider here several different cases: the observer has an access only to random values $\xi_i = I(X_i < t_i)$, $t_i = c_F \frac{2i-1}{2n}$, $i = \overline{1, n}$, $c_F = \inf \{x : F(x) = 1\} < \infty$.

The problem consists in estimating distribution functions $F(x)$ by the sample $\xi_1, \xi_2, \dots, \xi_n$. Such a problem arises, for example, in corrosion investigations (see [1] where an experiment connected with corrosion is described).

We will consider estimates for $F(x)$ that are analogous to regression curve estimates of Nadaraya-Watson type and have the form

$$\hat{F}_n(x) = F_{n1}(x)F_{n2}(x), \quad F_{n1}(x) = \sum_{i=1}^n K\left(\frac{x-t_i}{h}\right)\xi_i, \quad F_{n2}(x) = \left(\sum_{i=1}^n K\left(\frac{x-t_i}{h}\right)\right)^{-1},$$

where $K(x)$ is some weight function (kernel), $\{h = h(n)\}$ is a sequence of positive numbers converging to zero.

Lemma. Assume that

1⁰. $K(x)$ is some distribution density with bounded variation and $K(x) = K(-x)$, $x \in R = (-\infty, \infty)$. If $nh \rightarrow \infty$, then

$$\frac{1}{nh} \sum_{j=1}^n K^{m_1-1}\left(\frac{x-t_j}{h}\right) F^{m_2-1}(t_j) = \frac{1}{c_F h} \int_0^{c_F} K^{m_1-1}\left(\frac{x-u}{h}\right) F^{m_2-1}(u) du + O\left(\frac{1}{nh}\right) \quad (2)$$

uniformly with respect to $x \in [0, c_F]$, m_1, m_2 are natural.

Without loss of generality we assume below that the interval $[0, c_F] = [0, 1]$.

Theorem 1. Let $F(x)$ be continuous and the conditions of the lemma be fulfilled. Then the estimate $\hat{F}_n(x)$ is asymptotically unbiased and consistent at all points $x \in [0, 1]$. Moreover, $\hat{F}_n(x)$ is distributed asymptotically normally, i.e.

$$\sqrt{nh}(\hat{F}_n(x) - EF_n(x))\sigma^{-1}(x) \xrightarrow{d} N(0,1),$$

$$\sigma^2(x) = F(x)(1-F(x)) \int K^2(u) du,$$

where d denotes convergence in distribution, and $N(0,1)$ a random value having a normal distribution with mean 0 and variance 1.

2. Uniform consistency. We define the conditions under which the estimate $\hat{F}_n(x)$ converges uniformly in probability (almost surely) to a true $F(x)$.

Following E. Parzen [2], we introduce the Fourier transform of $K(x)$

$$\varphi(t) = \int_{-\infty}^{\infty} e^{itx} K(x) dx$$

and assume that

2⁰. $\varphi(t)$ is absolutely integrable. Then $F_{n1}(x)$ can be written in the form

$$F_{n1}(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-iu\frac{x}{h}} \varphi(u) \frac{1}{nh} \sum_{j=1}^n \xi_j e^{iu\frac{t_j}{h}} du .$$

Thus

$$F_{n1}(x) - EF_{n1}(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-iu\frac{x}{h}} \varphi(u) \frac{1}{nh} \sum_{j=1}^n (\xi_j - F(t_j)) e^{iu\frac{t_j}{h}} du .$$

Denote

$$d_n = \sup_{x \in \Omega_n} |\hat{F}_n(x) - EF_n(x)|, \quad \Omega_n = [h^\alpha, 1 - h^\alpha], \quad 0 < \alpha < 1 .$$

Theorem 2. Let $K(x)$ satisfy conditions 1⁰ and 2⁰.

(a) Let $F(x)$ be continuous and $n^{\frac{1}{4}}h(n) \rightarrow \infty$, then

$$D_n = \sup_{x \in \Omega_n} |\hat{F}_n(x) - F(x)| \xrightarrow{p} 0 ;$$

(b) If $\sum_{n=1}^{\infty} n^{-\frac{p}{4}}h(n) < \infty$, $p > 4$, then $D_n \rightarrow 0$ almost surely.

Proof. We have

$$\sup_{x \in \Omega_n} \left(1 - \frac{1}{h} \int_0^1 K\left(\frac{x-u}{h}\right) du \right) \leq \int_{-\infty}^{-h^{\alpha-1}} K(u) du + \int_{h^{\alpha-1}}^{\infty} K(u) du \rightarrow 0 \quad (1)$$

This and this lemma imply

$$\sup_{x \in \Omega_n} |F_{n2}(x) - 1| \rightarrow 0 \quad (2)$$

i.e., due to uniform convergence, for any $\varepsilon_0 > 0$, $0 < \varepsilon_0 < 1$, and sufficiently large $n \geq n_0$ we have $F_{n2}(x) \geq 1 - \varepsilon_0$ uniformly with respect to $x \in \Omega_n$.

Denote

$$A_n = \sup_{u \in R} \left| \frac{1}{nh} \sum_{j=1}^n \eta_j e^{iu\frac{t_j}{h}} \right|, \quad \eta_j = \xi_j - F(t_j).$$

Then

$$d_n^p \leq \frac{(1 - \varepsilon_0)^{-p}}{(2\pi)^p} A_n^p \left| \int_{-\infty}^{\infty} |\varphi(u)| du \right|^p, \quad p > 4. \quad (3)$$

Note that

$$\begin{aligned}
 EA_n^p &= \frac{1}{(nh)^p} E \sup_{u \in R} \left| \sum_{j=1}^n \eta_j^2 + \sum_{k \neq j} \eta_k \eta_j \cos \left(\left(\frac{t_k - t_j}{h} \right) u \right) \right|^{\frac{p}{2}} = \\
 &= \frac{1}{(nh)^p} E \sup_{u \in R} \left| \sum_{j=1}^n \eta_j^2 + \sum_{k \neq j} \eta_k \eta_j \cos \left(\frac{k-j}{nh} u \right) \right|^{\frac{p}{2}} = \\
 &= \frac{1}{(nh)^p} E \sup_{u \in R} \left| \sum_{j=1}^n \eta_j^2 + \sum_{\substack{m=-n+1 \\ m \neq 0}}^{n-1} \sum_{\substack{j=1 \\ m \neq 0}}^{n-|m|} \eta_j \eta_{j+m} \cos \left(\frac{m}{nh} u \right) \right|^{\frac{p}{2}}.
 \end{aligned}$$

From this, by the inequality

$$\left| \sum_{j=1}^n a_j \right|^q \leq n^{q-1} \sum_{j=1}^n |a_j|^q, \quad q \geq 1,$$

we have

$$EA_n^p \leq \frac{2^{\frac{p-1}{2}}}{(nh)^p} E \left[\sum_{i=1}^n \eta_i^2 \right]^{\frac{p}{2}} + \frac{2^{\frac{p-1}{2}}}{(nh)^p} E \sup_{u \in R} \left| \sum_{\substack{m=-n+1 \\ m \neq 0}}^{n-1} \cos \left(\frac{m}{nh} u \right) \sum_{j=1}^{n-|m|} \eta_j \eta_{j+m} \right|^{\frac{p}{2}} = C_{n1} + C_{n2}. \quad (4)$$

Let us estimate C_{n1} and C_{n2} :

$$C_{n1} \leq \frac{2^{\frac{p-1}{2}}}{n^{\frac{p+1}{2}} h^p} \sum_{j=1}^n E |\eta_j|^p = \frac{2^{\frac{p-1}{2}}}{n^{\frac{p+1}{2}} h^p} \sum_{j=1}^n \left[(1-F(t_j))^p F(t_j) + F^p(t_j) (1-F(t_j)) \right] \leq c_2 \frac{1}{n^2 h^p}. \quad (5)$$

Further, using Whittle's inequality [3] for moments of quadratic form, we obtain

$$C_{n2} \leq \frac{2^{\frac{p-1}{2}} (2n-3)^{\frac{p-1}{2}}}{(nh)^p} \sum_{\substack{m=-n+1 \\ m \neq 0}}^{n-1} E \left| \sum_{j=1}^{n-|m|} \eta_j \eta_{j+m} \right|^{\frac{p}{2}},$$

thus

$$E \left| \sum_{j=1}^{n-|m|} \eta_j \eta_{j+m} \right|^{\frac{p}{2}} \leq c(p) (n-|m|)^{\frac{p}{4}},$$

where $c(p)$ depends only on p and $E |\eta_j|^p \leq 1$.

Thus

$$\sum_{\substack{m=-n+1 \\ m \neq 0}}^{n-1} E \left| \sum_{j=1}^{n-|m|} \eta_j \eta_{j+m} \right|^{\frac{p}{2}} \leq 2c(p) \sum_{m=1}^{n-1} m^{\frac{p}{4}} = O \left(n^{\frac{p+1}{4}} \right)$$

and

$$C_{n2} = O \left(\frac{1}{n^4 h^p} \right). \quad (6)$$

After combining the relations (3), (4), (5) and (6), we obtain

$$Ed_n^p = O\left(\frac{1}{n^4 h^p}\right), \quad p > 4.$$

Therefore

$$P\left\{\sup_{x \in \Omega_n} |\hat{F}_n(x) - EF_n(x)| \geq \varepsilon\right\} \leq \frac{C_3}{\varepsilon^p n^4 h^p}. \quad (7)$$

Further we obtain

$$\sup_{x \in \Omega_n} |EF_n(x) - F(x)| \leq \frac{1}{1 - \varepsilon_0} \left(\sup_{x \in \Omega_n} |EF_{n1}(x) - F(x)| + \sup_{x \in \Omega_n} |1 - F_{n2}(x)| \right). \quad (8)$$

The second summand in the right-hand part of (8) tends, by virtue of (2), to zero, while the first summand is estimated as follows:

$$\sup_{x \in \Omega_n} |EF_{n1}(x) - F(x)| \leq S_{n1} + S_{n2} + O\left(\frac{1}{nh}\right), \quad (9)$$

$$S_{n1} = \sup_{0 \leq x \leq 1} \left| \frac{1}{h} \int_0^1 (F(y) - F(x)) K\left(\frac{x-y}{h}\right) dy \right|,$$

$$S_{n2} = \sup_{x \in \Omega_n} \left(1 - \frac{1}{h} \int_0^1 K\left(\frac{x-y}{h}\right) dy \right),$$

and, by virtue of (1),

$$S_{n2} \rightarrow 0. \quad (10)$$

Now let us consider S_{n1} . Note that

$$\begin{aligned} S_{n1} &\leq \sup_{0 \leq x \leq 1} \int_0^1 |F(y) - F(x)| \frac{1}{h} K\left(\frac{x-y}{h}\right) dy = \sup_{0 \leq x \leq 1} \int_{x-1}^x |F(x-u) - F(x)| \frac{1}{h} K\left(\frac{u}{h}\right) du \leq \\ &\leq \sup_{0 \leq x \leq 1} \int_{-\infty}^{\infty} |F(x-u) - F(x)| \frac{1}{h} K\left(\frac{u}{h}\right) du. \end{aligned} \quad (11)$$

Assume that $\delta > 0$ and divide the integration domain in (11) into two domains $|u| \leq \delta$ and $|u| > \delta$. Then

$$\begin{aligned} S_{n1} &\leq \sup_{0 \leq x \leq 1} \int_{|u| \leq \delta} |F(x-u) - F(x)| \frac{1}{h} K\left(\frac{u}{h}\right) du + \sup_{0 \leq x \leq 1} \int_{|u| > \delta} |F(x-u) - F(x)| \frac{1}{h} K\left(\frac{u}{h}\right) du \leq \\ &\leq \sup_{x \in R} \sup_{|u| \leq \delta} |F(x-u) - F(x)| + 2 \int_{|u| \geq \frac{\delta}{h}} K(u) du. \end{aligned} \quad (12)$$

By a choice of $\delta > 0$ the first summand in the right-hand part of (12) can be made arbitrarily small. After choosing $\delta > 0$ and making n tend to infinity, we obtain that the second summand tends to zero.

Thus

$$\lim_{n \rightarrow \infty} S_{n1} = 0. \quad (13)$$

Finally, the proof of the theorem follows from the relations (7)-(10) and (13).

Remarks.

- 1) If $K(x) = 0$, $|x| \geq 1$ and $\alpha = 1$, i.e., $\Omega_n = [h, 1-h]$, then $S_{n2} = 0$.
- 2) Under the conditions of Theorem 2,

$$\sup_{x \in [a,b]} |\hat{F}_n(x) - F(x)| \rightarrow 0$$

in probability (almost surely) for any fixed interval $[a,b] \subset [0,1]$ since there exists n_0 such that $[a,b] \subset \Omega_n$, $n \geq n_0$.

Let us assume that $h = n^{-\gamma}$, $\gamma > 0$. The conditions of Theorem 2 are fulfilled:

$$n^{\frac{1}{4}} h_n \rightarrow \infty \text{ if } \gamma < \frac{1}{4}$$

and

$$\sum_{n=1}^{\infty} n^{-\frac{p}{4}} h_n^{-p} < \infty \text{ if } \gamma < \frac{1}{4} - \frac{1}{p}, p > 4.$$

3. Estimation of moments. In the considered problem there naturally arises the question of estimation of integral functional of $F(x)$, for example, of moments μ_m , $m \geq 1$:

$$\mu_m = m \int_0^1 t^{m-1} (1 - F(t)) dt.$$

As estimates for μ_m we will consider the statistics

$$\hat{\mu}_{nm} = 1 - \frac{m}{n} \sum_{j=1}^n \xi_j \frac{1}{h} \int_h^{1-h} t^{m-1} K\left(\frac{t-t_j}{h}\right) F_{n2}^{-1}(t) dt.$$

Theorem 3. Let $K(x)$ satisfy condition I^0 and, in addition to this, $K(x) = 0$ outside the interval $[-1,1]$. If $nh \rightarrow \infty$ as $n \rightarrow \infty$, then $\hat{\mu}_{nm}$ is an asymptotically unbiased, consistent estimate for μ_m and, moreover,

$$\frac{\sqrt{n}(\hat{\mu}_{nm} - E\hat{\mu}_{nm})}{\sigma} \xrightarrow{d} N(0,1), \quad \sigma^2 = m^2 \int_0^1 t^{2m-2} F(t)(1 - F(t)) dt.$$

References

1. Manjgaladze, K.V. On one estimate of a distribution function and its moments. (Russian) *Bull. Acad. Sci. Georgian SSR* **124** (1980), No. 2, pp. 261-268.
2. Parzen, E. On estimation of a probability density function and mode. *Ann. Math. Statist.* **33** (1962), pp. 1065-1076.
3. Whittle, P. Bounds for the moments of linear and quadratic forms in independent variables. *Teor. Veroyatnost. i Primenen.* **5** (1960), pp. 331-335.