

МАТЕМАТИКА

Н. В. СМИРНОВ

О ПОСТРОЕНИИ ДОВЕРИТЕЛЬНОЙ ОБЛАСТИ ДЛЯ ПЛОТНОСТИ РАСПРЕДЕЛЕНИЯ СЛУЧАЙНОЙ ВЕЛИЧИНЫ

(Представлено академиком А. Н. Колмогоровым 14 VII 1950)

Пусть произведено  $n$  независимых наблюдений случайной величины  $\xi$ , распределенной с непрерывной плотностью  $f(x)$ . Обычно практикуемый прием приближенной оценки неизвестной функции  $f(x)$  по данным выборки на сегменте  $[a, b]$  заключается в построении „гистограммы“ частот  $f_n^*(x)$ . Подразделяя сегмент  $[a, b]$  на  $s$  отрезков  $\Delta_1, \Delta_2, \dots, \Delta_s$  равной длины  $h = \frac{b-a}{s}$  и подсчитав численности  $m_1, m_2, \dots, m_s$  наблюдений, попавших на последовательные отрезки  $\Delta_k$ , полагают:

$$f_n^*(x) = \frac{m_k}{nh} \text{ при } x \in \Delta_k. \quad (1)$$

Эту эмпирическую кривую сравнивают с теоретической гистограммой  $\bar{f}(x)$ :

$$\bar{f}(x) = \frac{p_k}{h} = \frac{1}{h} \int_{\Delta_k} f(x) dx. \quad (2)$$

Легко обнаруживается из элементарных соображений, что  $f_n^*(x)$  сходится по вероятности к  $\bar{f}(x)$  в каждой точке  $x \in [a, b]$ . Если при возрастании  $n$  надлежащим образом увеличивать и число  $s$  отрезков подразделения, то можно притти к более точным выводам, впервые полученным В. И. Гливенко (1). В настоящей заметке мы формулируем теоремы, которые при известных условиях могут быть использованы для более точной и надежной оценки плотности.

Теорема 1. Если непрерывная плотность распределения  $f(x)$  удовлетворяет на  $[a, b]$  условиям:

$$\min \{f(x)\} = \mu > 0; \quad a \leq x \leq b, \quad (\text{A})$$

$$\int_a^b f(x) dx = 1 - \alpha < 1, \quad (\text{B})$$

и при возрастании  $n$  и  $s$  имеем:

$$\overline{\lim}_{n \rightarrow \infty} \frac{s^3 (\ln s)^3}{n} < \infty, \quad (\text{C})$$

то для любого  $\lambda, -\infty < \lambda < \infty$ ,

$$P_n = P \left\{ \text{Max} \frac{|f_n^*(x) - f(x)|}{V\bar{f}(x)} \leqslant \frac{l_s + \frac{\lambda}{l_s}}{Vnh} \right\} = \\ = \left[ 2 \Phi \left( l_s + \frac{\lambda}{l_s} \right) \right]^s + O \left( \frac{V \ln s}{s} \right) + O \left( \frac{s^{3/2}}{Vn} \right) = e^{-2e^{-\lambda}} + O \left( \frac{1}{V \ln s} \right) \quad (3)$$

и

$$P_n \rightarrow e^{-2e^{-\lambda}} \quad (n \rightarrow \infty); \quad (3^1)$$

при этом

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_0^x e^{-x^2/2} dx$$

и  $l_s$  есть корень уравнения

$$\frac{1}{2} - \Phi(x) = \frac{1}{s}. \quad (4)$$

В формулировке теоремы 1 теоретическую гистограмму  $\bar{f}(x)$  нельзя заменить плотностью  $f(x)$ , так как (даже при наличии ограниченной производной  $f'(x)$ ) систематическое расхождение между  $\bar{f}(x)$  и  $f(x)$  значительно превосходит случайную погрешность. Однако в качестве приближающей эмпирической функции берут также „полигон частот“ и соответственно заменяют  $\bar{f}(x)$  теоретическим полигоном.

Эти функции определяются следующим образом.

Пусть  $x_k (k = 1, 2, \dots, s)$  — середины промежутков  $\Delta_k$  и  $a_k$  — правые концы этих промежутков; в каждом сегменте  $[x_k, x_{k+1}]$  полагаем:

$$\varphi_n^*(x) = \frac{m_k + m_{k+1}}{2nh} + (x - a_k) \frac{(m_{k+1} - m_k)}{nh^2}. \quad (5^1)$$

и

$$\varphi_s(x) = \frac{p_k + p_{k+1}}{2h} + (x - a_k) \frac{(p_{k+1} - p_k)}{h^2}, \quad (5^2)$$

так что

$$\varphi_n^*(x_k) = \frac{m_k}{nh}; \quad \varphi_n^*(x_{k+1}) = \frac{m_{k+1}}{nh}; \quad \varphi_s(x_k) = \frac{p_k}{h}; \quad \varphi_s(x_{k+1}) = \frac{p_{k+1}}{h}.$$

Пусть еще

$$\psi_s(x) = \frac{\sqrt{\frac{p_k}{h}} + \sqrt{\frac{p_{k+1}}{h}}}{2} + (x - a_k) \left( \sqrt{\frac{p_{k+1}}{h}} - \sqrt{\frac{p_k}{h}} \right), \quad (6)$$

$$\psi_s(x_k) = \sqrt{\frac{p_k}{h}}; \quad \psi_s(x_{k+1}) = \sqrt{\frac{p_{k+1}}{h}}. \quad (6^1)$$

Тогда имеет место

Теорема 2. При условиях (A), (B) и (C) теоремы 1 для вероятности

$$P'_n = P \left\{ \text{Max} \frac{|\varphi_n^*(x) - \varphi_s(x)|}{\psi_s(x)}; \quad a + \frac{h}{2} \leqslant x \leqslant b - \frac{h}{2} \right\}$$

справедливы утверждения (3) и (3<sup>1</sup>) теоремы 1.

Из теоремы 2 легко получается следующая

**Теорема 3.** Если  $f(x)$  имеет на  $[a, b]$  ограниченную вторую производную и выполняются условия (A), (B) и (C) теоремы 1 и, кроме того,

$$\frac{n \ln s}{s^5} \rightarrow 0 \quad (n \rightarrow \infty), \quad (\text{D})$$

то

$$P \left\{ \max \frac{|\varphi_n^*(x) - f(x)|}{Vf(x)} \leqslant \frac{l_s + \frac{\lambda}{l_s}}{Vnh} \right\} \xrightarrow[n \rightarrow \infty]{} e^{-2e^{-\lambda}}. \quad (7)$$

Условиям (C) и (D) можно удовлетворить, полагая, например,  $s = n^\alpha$ ,  $\frac{1}{5} < \alpha < \frac{1}{3}$ .

Теорема 3 дает возможность строить доверительные области с заданным коэффициентом довери  $\beta$  для оценки неизвестной плотности  $f(x)$  по данным выборки.

Определив при заданном  $\beta$  ( $0 < \beta < 1$ )  $\lambda = \lambda_\beta$ , удовлетворяющее уравнению

$$e^{-2e^{-\lambda}} = \beta$$

(или, более точно, из уравнения  $\left[2\Phi\left(l_s + \frac{\lambda}{l_s}\right)\right]^s = \beta$ ), находим

$$\frac{l_s + \lambda_\beta / l_s}{Vnh} = t_\beta.$$

Неравенство

$$\max \frac{|\varphi_n^*(x) - f(x)|}{Vf(x)} \leqslant t_\beta$$

означает, что кривая  $f(x)$  на всем сегменте  $[a, b]$  не выходит за границы полосы  $\mathfrak{U}_\beta$ , ограниченной кривыми:

$$y_1(x) = \varphi_n^*(x) + \frac{t_\beta^2}{2} - t_\beta \sqrt{\varphi_n^*(x) + \frac{t_\beta^2}{4}}$$

и

$$y_2(x) = \varphi_n^*(x) + \frac{t_\beta^2}{2} + t_\beta \sqrt{\varphi_n^*(x) + \frac{t_\beta^2}{4}}.$$

Полоса  $\mathfrak{U}_\beta$  покрывает кривую плотности  $y = f(x)$  при больших  $n$  и  $s$  при условиях теоремы 3 с вероятностью, сколь угодно близкой к  $\beta$ .

Поступило  
4 IV 1950

#### ЦИТИРОВАННАЯ ЛИТЕРАТУРА

<sup>1</sup> В. И. Глиденко, Курс теории вероятностей, 1939, стр. 180.