

## НЕПАРАМЕТРИЧЕСКАЯ ОЦЕНКА СМЕСИ ПЛОТНОСТЕЙ ВЕРОЯТНОСТИ, ОСНОВАННАЯ НА ТЕХНОЛОГИИ РАЗМНОЖЕНИЯ СТАТИСТИЧЕСКИХ ДАННЫХ

*Исследована непараметрическая оценка смеси плотностей вероятности, синтез которой основан на технологии размножения статистических данных. Установлены условия ее асимптотической несмещенности и состоятельности. На этой основе проведено сравнение свойств предлагаемой оценки смеси плотностей вероятности с непараметрической оценкой Розенблатта–Парзена.*

*Ключевые слова:* плотность вероятности, непараметрическая оценка, размножение данных, асимптотические свойства.

Большинство статистических методов обработки информации ориентировано на представительные обучающие выборки. Однако при решении прикладных задач часто имеется ограниченный объем наблюдений – короткая либо малая выборка, что обуславливается нестационарностью объекта исследования, высокой стоимостью и сложностью получения дополнительной информации. Получаемые на их основе решающие правила не всегда обеспечивают приемлемые результаты, так как информации малых выборок недостаточно для оценивания вероятностных характеристик изучаемых закономерностей.

Проблемы малых выборок можно разрешить с помощью технологий обработки информации, основанных на бутстреп-методах. Ниже на основе результатов аналитических исследований будет обоснована эффективность его применения при непараметрическом оценивании плотностей вероятности.

**Синтез непараметрической оценки смеси плотностей вероятности, основанной на технологии бутстреп-метода.** Пусть  $V = (x^i, i = 1, n)$  – выборка из  $n$  независимых наблюдений случайной величины  $x = (x_v, v = 1, k)$  с плотностью вероятности  $p(x)$ , вид которой априори неизвестен.

Сформируем на основе исходной выборки  $N$  групп наблюдений выборку  $V_j = (x^i, i \in I_j)$ , где  $I_j$  – множество номеров элементов из  $V$ , составляющих  $j$ -ю группу. Количество элементов в группах одинаково и равно  $\bar{n} = n - n'$ . Каждая пара групп  $V_j, V_t, j, t = 1, N, j \neq t$  отличается  $n'$  элементами. Количество групп элементов  $N = n/n'$ .

По каждой выборке  $V_j$  построим непараметрические оценки плотностей вероятности [1; 2]:

$$\bar{p}_j(x) = \frac{1}{\bar{n} \prod_{v=1}^k c_v} \sum_{i \in I_j} \prod_{v=1}^k \Phi \left( \frac{x_v - x_v^i}{c_v} \right), \quad (1)$$

$$j = \overline{1, N},$$

где  $\Phi(u)$  – ядерные функции, удовлетворяющие условиям  $H$ :

$$\begin{aligned} \Phi(u) &= \Phi(-u), \quad 0 \leq \Phi(u) < \infty, \\ \int \Phi(u) du &= 1, \quad \int u^2 \Phi(u) du = 1, \\ \int u^m \Phi(u) du &< \infty, \quad 0 \leq m < \infty; \end{aligned}$$

$c_v = c_v(\bar{n})$  – коэффициенты размытости ядерных функций, значения которых убывают с ростом  $\bar{n}$ . Здесь и далее бесконечные пределы интегрирования опускаются.

В качестве приближения  $p(x)$  по статистической выборке  $V$  примем смесь непараметрических оценок  $\bar{p}_j(x)$  плотности вероятности

$$\bar{\bar{p}}(x) = \frac{1}{N} \sum_{j=1}^N \bar{p}_j(x). \quad (2)$$

Статистика (2) построена в соответствии с бутстреп-методом и допускает использование технологии параллельных вычислений.

Исследуем асимптотические свойства оценки плотности вероятности (2) в условиях, когда  $k = 1$ .

*Теорема.* Пусть  $p(x)$  и первые две ее производные ограничены и непрерывны; ядерные функции  $\Phi(u)$  удовлетворяют условиям нормированности, положительности и симметричности  $H$ ; последовательность  $c(\bar{n}) = c$  коэффициентов размытости ядерных функций такова, что при  $\bar{n} \rightarrow \infty$  значения  $c \rightarrow 0$ , а при  $\bar{n}c \rightarrow \infty$  и  $\frac{1}{\bar{n}} \rightarrow 0$ ,  $\frac{n'}{\bar{n}^2} \rightarrow 0$ . Тогда при конечных значениях  $N$  непараметрическая оценка  $\bar{\bar{p}}(x)$  смеси плотности вероятности  $p(x)$  обладает свойством асимптотической несмещенности и состоятельности.

**Доказательство.** По определению

$$\begin{aligned} M(\bar{\bar{p}}(x)) &= \frac{1}{N} \times \\ &\times \sum_{j=1}^N M(\bar{p}_j(x)) = \frac{1}{c} \int \Phi \left( \frac{x-t}{c} \right) p(t) dt = \\ &= \int \Phi(u) p(x-cu) du, \end{aligned}$$

где  $M$  – знак математического ожидания.

Разлагая  $p(x-cu)$  в ряд Тейлора и ограничиваясь первыми двумя членами ряда при  $\bar{n} \rightarrow \infty$ , имеем

$$W_1 = M(\bar{\bar{p}}(x) - p(x)) \sim \frac{p^{(2)}(x)}{2} c^2, \quad (3)$$

где  $p^{(2)}(x)$  – вторая производная плотности вероятности  $p(x)$  по  $x$ . Отсюда из условия  $c = c(\bar{n})$  и  $c \rightarrow 0$  при  $\bar{n} \rightarrow \infty$  следует свойство асимптотической несмещенности непараметрической оценки смеси плотностей вероятности (2).

Для доказательства сходимости  $\bar{\bar{p}}(x)$  в среднеквадратическом отклонении рассмотрим выражение

$$\begin{aligned}
 & M \int (p(x) - \bar{p}(x))^2 dx = \\
 & = M \int \left[ \frac{1}{N} \sum_{j=1}^N (p(x) - \bar{p}_j(x)) \right]^2 dx = \\
 & = \frac{1}{N^2} M \left[ \sum_{j=1}^N \int (p(x) - \bar{p}_j(x))^2 dx + \right. \\
 & \left. + \sum_{\substack{j=1 \\ i=1 \\ i \neq j}}^N \sum_{i=1}^N \int (p(x) - \bar{p}_j(x)) \times \right. \\
 & \left. \times (p(x) - \bar{p}_i(x)) dx \right]. \quad (4)
 \end{aligned}$$

Найдем асимптотическое выражение функционала

$$\begin{aligned}
 & M \int (p(x) - \bar{p}_j(x))(p(x) - \bar{p}_i(x)) dx = \\
 & = \int p^2(x) dx - M \int \bar{p}_i(x) p(x) dx - \\
 & - M \int \bar{p}_j(x) p(x) dx + M \int \bar{p}_j(x) \bar{p}_i(x) dx. \quad (5)
 \end{aligned}$$

Преобразуем его последнее слагаемое:

$$\begin{aligned}
 & M \int \bar{p}_j(x) \bar{p}_i(x) dx = \\
 & = \frac{1}{\bar{n}^2 c^2} \int \left[ \sum_{i \in I_j} M \Phi^2 \left( \frac{x - x^i}{c} \right) + \right. \\
 & \left. + \sum_{i \in I_j} \sum_{v \in I_i \setminus I_j} M \Phi \left( \frac{x - x^i}{c} \right) M \Phi \left( \frac{x - x^v}{c} \right) \right] dx,
 \end{aligned}$$

которое при достаточно большом объеме  $\bar{n}$  элементов в группах  $V_j$ ,  $j = \bar{1}, \bar{N}$ , может быть представлено в виде

$$\begin{aligned}
 & \frac{\bar{n} - n'}{\bar{n}^2 c} \|\Phi(u)\|^2 + \frac{\bar{n}^2 - (\bar{n} - n')}{\bar{n}^2} \times \\
 & \times \int (p(x) + c^2 p^{(2)}(x)/2)^2 dx, \quad (6)
 \end{aligned}$$

где  $\|\Phi(u)\|^2 = \int \Phi^2(u) du$ .

Заметим, что при  $\bar{n} \rightarrow \infty$

$$M \int \bar{p}_j(x) p(x) dx \sim \|p(x)\|^2 + \frac{c^2}{2} \int p^{(2)}(x) p(x) dx,$$

где  $\|p(x)\|^2 = \int p^2(x) dx$ . Тогда асимптотическое выражение для функционала (5) соответствует выражению

$$\begin{aligned}
 & \frac{\bar{n} - n'}{\bar{n}^2 c} \|\Phi(u)\|^2 + \frac{c^4}{4} \|p^{(2)}(x)\|^2 + \\
 & + \frac{\bar{n} - n'}{\bar{n}^2} \int (p(x) + c^2 p^{(2)}(x)/2)^2 dx. \quad (7)
 \end{aligned}$$

С учетом (7) и справедливости при  $\bar{n} \rightarrow \infty$  утверждения [2]

$$M \int (p(x) - \bar{p}_j(x))^2 dx \sim \frac{\|\Phi(u)\|^2}{\bar{n}c} + \frac{c^4 \|p^{(2)}(x)\|^2}{4},$$

запишем асимптотическое выражение для (4):

$$\begin{aligned}
 & M \int (p(x) - \bar{p}(x))^2 dx \sim \frac{1}{N} \times \\
 & \times \left( \frac{\|\Phi(u)\|^2}{\bar{n}c} + \frac{c^4 \|p^{(2)}(x)\|^2}{4} \right) + \\
 & + \frac{N-1}{N} \left[ \frac{\bar{n} - n'}{\bar{n}^2 c} \|\Phi(u)\|^2 + \frac{c^4 \|p^{(2)}(x)\|^2}{4} + \right. \\
 & \left. + \frac{\bar{n} - n'}{\bar{n}^2} \int (p(x) + c^2 p^{(2)}(x)/2)^2 dx \right].
 \end{aligned}$$

Отсюда, пренебрегая величинами малости  $0\left(\frac{1}{\bar{n}}\right)$ ,  $0\left(\frac{n'}{\bar{n}^2}\right)$ , получим

$$\begin{aligned}
 & M \int (p(x) - \bar{p}(x))^2 dx \sim \frac{\|\Phi(u)\|^2}{N \bar{n}c} \times \\
 & \times \left( 1 + \frac{(\bar{n} - n')(N-1)}{\bar{n}} \right) + \frac{c^4 \|p^{(2)}(x)\|^2}{4}. \quad (8)
 \end{aligned}$$

Нетрудно заметить, что при выполнении условий  $c \rightarrow 0$ ,  $\bar{n}c \rightarrow \infty$  при  $\bar{n} \rightarrow \infty$  оценка плотности вероятности (2) сходится в среднеквадратическом отклонении к  $p(x)$ , а с учетом свойства ее асимптотической несмещенности является состоятельной.

**Сравнение асимптотических свойств статистики (2) и непараметрической оценки Розенблатта–Парзена.** Определим минимальное значение  $W_2$  выражения (8) при оптимальных значениях  $\bar{c}$  коэффициентов размытости ядерных функций непараметрических оценок  $\bar{p}_j(x)$ ,  $j = \bar{1}, \bar{N}$ , составляющих их смесь  $\bar{p}(x)$  (2).

В принятых допущениях значение

$$\bar{c} = \left( \frac{\|\Phi(u)\|^2}{\bar{n} \|p^{(2)}(x)\|^2} \right)^{\frac{1}{5}}.$$

Тогда

$$\begin{aligned}
 & W_2 = \left[ \left( \frac{\|\Phi(u)\|^2}{n - n'} \right)^4 \|p^{(2)}(x)\|^2 \right]^{\frac{1}{5}} \times \\
 & \times \left( \frac{1}{N} \left( 1 + \frac{n - n'}{n} (N-1) \right) + \frac{1}{4} \right) = \\
 & = \left[ \left( \frac{\|\Phi(u)\|^2}{n - n'} \right)^4 \|p^{(2)}(x)\|^2 \right]^{\frac{1}{5}} \left( \frac{n - n'}{n} + \frac{1}{4} \right). \quad (9)
 \end{aligned}$$

Если  $n' = 0$ , то  $W_2$  совпадает с минимальным значением асимптотического выражения среднеквадратического отклонения

$$W'_2 = \frac{5}{4} \left[ \left( \frac{\|\Phi(u)\|^2}{n} \right)^4 \|p^{(2)}(x)\|^2 \right]^{\frac{1}{5}}$$

для оценки плотности вероятности типа Розенблатта–Парзена

$$\tilde{p}(x) = \frac{1}{nc} \sum_{i=1}^n \Phi \left( \frac{x - x^i}{c} \right) \quad (10)$$

при оптимальных значениях  $c = c^*$ . При этом

$$\frac{W_2}{W'_2} = \frac{4}{5} \left( \frac{n}{n - n'} \right)^{\frac{4}{5}} \left( \frac{n - n'}{n} + \frac{1}{4} \right). \quad (11)$$

По аналогии сравним главные дисперсионные составляющие

$$W_3 = \frac{n - n'}{n} \left[ \left( \frac{\|\Phi(u)\|^2}{n - n'} \right)^4 \|p^{(2)}(x)\|^2 \right]^{\frac{1}{5}},$$

$$W'_3 = \left[ \left( \frac{\|\Phi(u)\|^2}{n} \right)^4 \|p^{(2)}(x)\|^2 \right]^{\frac{1}{5}}$$

непараметрических оценок плотностей вероятностей (2), (10). Их отношение будет следующим

$$\frac{W_3}{W'_3} = \left( \frac{n-n'}{n} \right)^{\frac{1}{5}}.$$

Так как объем статистических данных при синтезе составляющих (1) статистики (2) меньше, чем при формировании непараметрической оценки плотности вероятности (10), то очевидно, что ее смещение меньше по сравнению с оценкой плотности вероятности (2). Действительно, отношение их минимальных асимптотических выражений смещений при оптимальных значениях коэффициентов размытости имеет вид

$$\frac{W_1}{W'_1} = \left( \frac{n}{n-n'} \right)^{\frac{2}{5}} > 1. \quad (13)$$

Значения отношений (11), (12), (13) при  $n' = \alpha n$  приведены в таблице.

Использование статистики (2) позволяет несколько улучшить эффективность оценивания плотности вероятности по сравнению с оценкой Розенблатта–Парзена (10).

С уменьшением  $n' = \alpha n$  растет количество групп наблюдений  $V_j, j = 1, N$ , но снижается уровень их разнообразия. При этом состав групп наблюдений незначительно отличается. Поэтому в данных условиях аппроксимационные свойства статистик (2) и (10) практически одинаковы.

При увеличении  $n'$  уменьшается количество составляющих смеси непараметрических оценок плотностей

вероятности (2), что приводит к снижению ее аппроксимационных свойств, несмотря на рост разнообразия групп наблюдений  $V_j, j = 1, N$ . При этом уменьшение дисперсии смеси (2) объясняется различными темпами изменения значений ее смещения и среднеквадратического отклонения от искомой плотности вероятности.

Таким образом, использование технологии бутстреп-метода позволяет повысить эффективность оценивания плотностей вероятности. Получаемые при этом непараметрические оценки смеси плотностей вероятности обладают повышенными аппроксимационными свойствами, что особенно наблюдается в снижении их дисперсии и среднеквадратического отклонения. Определены условия их асимптотической сходимости и преимущества перед традиционной непараметрической оценкой Розенблатта–Парзена. Следует ожидать более значительного преимущества предлагаемой методики оценивания плотностей вероятности в условиях малых выборок.

### Библиографический список

1. Parzen, E. On estimation of a probability density function and mode / E. Parzen // Ann. Math. Statistic. 1962. Vol. 33. P. 1065–1076.
2. Епанечников, В. А. Непараметрическая оценка многомерной плотности вероятности // Теория вероятности и ее применения. 1969. Т. 14. Вып. 1. С. 156–161.

Зависимость отношений  $W_i/W'_i, i = \overline{1, 3}$  от значений  $n'$

$\frac{W_i}{W'_i}$	$\alpha$						
	0	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{20}$	$\frac{1}{40}$
(11)	1	1,044	1,002	0,998	0,999	0,999	0,999
(12)	1	0,87	0,94	0,97	0,987	0,989	0,995
(13)	1	1,32	1,12	1,054	1,024	1,020	1,008

A. V. Lapko, V. A. Lapko

## NONPARAMETRIC ESTIMATION OF A MIX PROBABILITY DENSITY, BASED ON TECHNOLOGY DUPLICATION OF STATISTICAL DATA

*The nonparametric estimation of the mix probability density which synthesis is based on technology duplication of statistical data is investigated. The asymptotic unbiasedness conditions and solvencies are determined. On this basis the properties comparison of the offered estimation of the mix probability density with the nonparametric estimation such as Rozenblatt–Parzen is carried out.*

*Keywords: probability density, nonparametric estimation, data duplication, asymptotic properties.*

© Ланко А. В., Ланко В. А., 2009