

УДК 519.7

А. В. Лапко, В. А. Лапко

(Красноярск)

**НЕПАРАМЕТРИЧЕСКИЕ МЕТОДИКИ АНАЛИЗА
МНОЖЕСТВ СЛУЧАЙНЫХ ВЕЛИЧИН**

Предлагается методика анализа множеств случайных величин в задачах восстановления стохастических зависимостей и распознавания образов, основанная на оценивании вероятностных законов распределения элементов множеств и их преобразовании с помощью непараметрических процедур. Исследованы асимптотические свойства моделей. Полученные результаты имеют актуальное значение при обработке больших массивов статистических данных.

Введение. Существующий парадокс формирования исходных данных при обнаружении статистических закономерностей состоит в сопоставлении значений случайных величин, вероятность которого невелика. Более обоснованными являются постановки задач, когда обучающие выборки состояются из пар множеств случайных величин.

Пусть состояние исследуемого объекта характеризуется множествами $X \subset \mathbb{R}^k$ и $Y \subset \mathbb{R}^l$ независимых наблюдений случайных величин x и y , взаимосвязь между которыми определяется неизвестным преобразованием

$$R: X \rightarrow Y.$$

Априорную информацию составляют N пар множеств $(X^i, Y^i, i = \overline{1, N})$, где множеству X^i соответствует вполне определенное множество Y^i . Назовем X входным, а Y выходным множеством. Подобные условия встречаются при исследовании объектов, параметры которых многократно измеряются в течение короткого интервала времени; оценивании качества партии изделий по случайно выбранному их ограниченному набору; обработке больших массивов статистических данных, которые с помощью декомпозиции преобразуются к обучающей выборке изучаемой задачи.

Предлагаются методы анализа множеств случайных величин в задачах обнаружения статистических закономерностей. Идея предлагаемого подхода заключается в замене операций над множествами менее трудоемкими и хорошо разработанными операциями над функциями либо над их параметрами с использованием методов локальной аппроксимации.

Регрессионная оценка плотности вероятности. Вычислительная эффективность непараметрических алгоритмов обработки информации во

многим определяется объемом статистических данных и снижается по мере его увеличения, что затрудняет построение систем принятия решений в условиях больших выборок.

Для обхода возникающих проблем рассмотрим методику сжатия априорной информации на основе ее декомпозиции и последующего анализа получаемых множеств случайных величин.

Пусть $x^i, i = \overline{1, n}$, – выборка из n независимых наблюдений случайной величины $x \in R^1$ с неизвестной плотностью вероятности $p(x)$, которая ограничена и непрерывна со всеми своими производными до второго порядка включительно. Причем

$$\int (p^{(2)}(x))^2 dx = \|p^{(2)}(x)\|^2 < \infty,$$

где $p^{(2)}(x)$ – вторая производная $p(x)$. Эти условия, налагаемые на $p(x)$, обозначим через G_2 . Для упрощения записи бесконечные пределы интегрирования здесь и в дальнейшем опускаются.

Разобьем область определения $p(x)$ на N непересекающихся интервалов длиной 2β и сформируем множества случайных величин $X^j, j = \overline{1, N}$. В качестве характеристики X^j примем частоту P^j попадания случайной величины x в j -й интервал и его центр z^j . На основе полученной информации составим статистическую выборку $(z^j, p^j = P^j / (2\beta), j = \overline{1, N})$, где центры z введенных интервалов имеют равномерный закон распределения $p(z) = (2\beta N)^{-1}$, а объем полученной выборки N может быть значительно меньше исходной n .

В качестве приближения по эмпирическим данным искомой плотности $p(x)$ примем непараметрическую оценку условного математического ожидания

$$\bar{p}(x) = c^{-1} \sum_{i=1}^N P^i \Phi \left(\frac{x - x^i}{c} \right), \quad (1)$$

где $\Phi(u) \in H$ – ядерная функция, удовлетворяющая условиям регулярности $H[1]$:

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

$c = c(n)$ – убывающая с ростом n последовательность положительных чисел (коэффициентов размытости).

Статистику (1) нетрудно получить, подставляя в выражение условного математического ожидания непараметрическую оценку совместной плотности вероятности случайных величин (P, z) и известную плотность $p(z)$. Регрессионная оценка плотности $p(x)$ является нормированной функцией, т. е. удовлетворяет основному свойству плотности вероятности.

Асимптотические свойства $\bar{p}(x) \forall x \in R^1$ определяет

Теорема 1. Пусть $p(x)$ удовлетворяет условиям G_2 , $\Phi(u) \in H$, кроме того, последовательности $c > 0$, $\beta > 0$ таковы, что при $n, N \rightarrow \infty$, значения $c, \beta \rightarrow 0$, $\beta/c \rightarrow 0$, $nc \rightarrow \infty$.

Тогда смещение

$$W_1(\bar{p}(x)) = M\{\bar{p}(x) - p(x)\} \sim p^{(2)}(x)(\beta^2/3 + c^2)/2,$$

среднеквадратическое отклонение

$$\begin{aligned} W_2(\bar{p}(x)) = M\left\{\|\bar{p}(x) - p(x)\|^2\right\} &\sim \|\Phi(u)\|^2/nc + \|p^{(2)}(x)\|^2(\beta^2/3 + c^2)^2/4 + \\ &+ \left(\|\Phi(u)\|^2\beta/c\right) \left(2\|p(x)\|^2 + \beta^4\|p^{(2)}(x)\|^2/18\right) + \beta c\|u\Phi(u)\|^2 \times \\ &\times \left(2\|p^{(1)}(x)\|^2 + \beta^2\|p^{(2)}(x)\|^2/3\right) + \beta c^3\|p^{(2)}(x)\|^2\|u^2\Phi(u)\|^2/2. \end{aligned}$$

Сравнение асимптотических свойств $\bar{p}(x)$ и непараметрической оценки $\bar{p}_1(x)$ типа Парзена [2], восстанавливаемой по выборке объема N , позволяет получить оценку сверху разности среднеквадратических отклонений

$$W_2(\bar{p}(x)) - W_2(\bar{p}_1(x)) < \frac{\|\Phi(u)\|^2}{c} \left(\frac{1}{n} - \frac{1}{N} + 2\beta\right).$$

Смещение $\bar{p}(x)$ несколько выше, чем смещение традиционной оценки Парзена.

Преимущество регрессионной плотности вероятности $\bar{p}(x)$ заключается не только в повышении вычислительной эффективности непараметрических алгоритмов, синтезируемых на ее основе, но и в упрощении задачи оптимизации $\bar{p}(x)$ по параметрам c в режиме «скользящего экзамена» по выборке $(z'_i, P'_i/2\beta, i = \overline{1, N})$.

В многомерном случае при $x = (x_1, \dots, x_k)$ регрессионная оценка плотности вероятности запишется в виде

$$\bar{p}(x) = \left(\prod_{v=1}^k c_v\right)^{-1} \sum_{i=1}^{N^k} P'_i \prod_{v=1}^k \Phi\left(\frac{x_v - x'_v}{c_v}\right). \quad (3)$$

На основе предлагаемых оценок плотности вероятности (1), (3) возможна реализация последовательной вычислительной процедуры

$$\bar{A}_s(\bar{p}_s(x), r_{s+1}(x)), \quad s = \overline{1, m}, \quad (4)$$

каждый последующий $(s+1)$ -й этап которой реализуется непараметрическим алгоритмом $\bar{A}_s(\cdot)$ в некоторой окрестности предыдущего решения $r_s(x)$ и уточняет его $r_{s+1}(x)$. Причем синтез алгоритмов $\bar{A}_s(\cdot)$ осуществляется на основе непараметрических оценок плотности вероятности $\bar{p}_s(x)$, восстанавливаемых по выборкам ограниченного объема N_s , значительно меньшего исходного n .

Установлено, что при $N_s = N$, $s = \overline{1, m}$, с увеличением числа этапов s процедуры последовательной обработки информации среднеквадратическая ошибка оценивания $p(x)$ снижается пропорционально значению $(N/2)^{4(s-1)/5}$. Особенность структуры регрессионных оценок плотности вероятности позволяет решить проблему их доверительного оценивания по следующей методике:

– рассчитываются доверительные границы P_1^i, P_2^i для частоты P^i попадания случайной величины x в i -й интервал при некотором значении коэффициента доверия α ;

– формируются статистические выборки $(z^i, P_1^i, i = \overline{1, N}), (z^i, P_2^i, i = \overline{1, N})$;

– строятся верхняя $\tilde{p}_1(x)$ и нижняя $\tilde{p}_2(x)$ толерантные границы регрессионной оценки плотности вероятности

$$\tilde{p}_v(x) = c^{-1} \sum_{j=1}^N P_v^j \Phi \left(\frac{x - z^j}{c} \right), \quad v=1, 2.$$

По аналогии осуществляется доверительное оценивание других непараметрических решающих функций.

Применение метода декомпозиции обучающей выборки в задаче распознавания образов. Синтез непараметрических алгоритмов распознавания образов основан на оценивании байесовского уравнения разделяющей поверхности [3]. В случае двух классов $\Omega_1(x)$ и $\Omega_2(x)$ уравнение разделяющей поверхности имеет вид

$$f_{12}(x) = p_2(x) - p_1(x) \quad (5)$$

и может быть оценено с помощью непараметрических методов статистики по обучающей выборке $V = (x^i, \sigma(x^i), i = \overline{1, n})$.

Примем в качестве оценок плотностей вероятности $x = (x_1, \dots, x_k)$ в классах $\Omega_1(x)$ и $\Omega_2(x)$ процедуры типа (3). Тогда непараметрическая оценка $f_{12}(x)$ (5) представляется выражением

$$\bar{f}_{12}(x) = \left(\prod_{v=1}^k c_v \right)^{-1} \sum_{i=1}^{N^k} \sigma_{12}(x^i) P^i \prod_{v=1}^k \Phi \left(\frac{x_v - x_v^i}{c_v} \right), \quad (6)$$

где

$$\sigma_{12}(x^i) = \begin{cases} -1, & \text{если } x^i \in \Omega_1(x), \\ 1, & \text{если } x^i \in \Omega_2(x); \end{cases}$$

$N^k = N_j^k + N_{\bar{j}}^k$ ($N_j^k, N_{\bar{j}}^k$ – количество элементов выборки, на основании которых формируются оценки плотностей $\bar{p}_1(x), \bar{p}_2(x)$).

Асимптотическую сходимость $\bar{f}_{12}(x)$ нетрудно доказать, используя утверждения теоремы 1.

Оптимизация непараметрической оценки решающей функции (6) по параметрам $\beta_v, c_v, v=1, k$, осуществляется из условия минимума статистической оценки вероятности ошибки распознавания образов, формируемой по контрольной выборке $V_1 \subset V$. Проводя предварительную нормировку компонент $x_v, v=1, k$, относительно их среднего значения, размерность задачи оптимизации можно сократить до двух параметров: $\beta_v = \beta, c_v = c, v=1, k$.

По результатам статистического моделирования применение предлагаемых алгоритмов при снижении временных затрат на порядок обеспечивает увеличение ошибки распознавания образов не более чем на 5 % по сравнению с прямыми методами обработки. При случайной стратегии формирования выборки меньшего объема $N^k < n$ отмечается значительное увеличение ошибки распознавания (до 15 %).

Синтез и анализ непараметрической регрессии на основе метода декомпозиции выборки. Пусть переменные регрессии $y = \varphi(x)$ скаляры $x, y \in R^1$. В соответствии с методикой синтеза $\bar{p}(x)$ (1) нетрудно получить выражения для оценки плотности вероятности двумерной случайной величины (x, y) :

$$\bar{p}(x, y) = c^{-2} \sum_{i=1}^{N^2} P^i \Phi\left(\frac{x-x'}{c}\right) \Phi\left(\frac{y-y'}{c}\right),$$

где объем N^2 рабочей выборки $(x', y', P^i, i=1, N^2)$, сформированной на основе обучающей $V = (x^j, y^j, j=1, n)$, значительно меньше n .

Тогда с учетом $p(y/x) = 2N\beta p(x, y)$ непараметрическая оценка регрессии представляется в виде

$$\bar{y} = \varphi(x) = \frac{2N\beta}{c} \sum_{i=1}^{N^2} y^i P^i \Phi\left(\frac{x-x'}{c}\right). \quad (7)$$

Справедлива следующая

Теорема 2. Пусть

1) функции $\varphi(x), p(x, y), p(x)$ ограничены и непрерывны со всеми своими производными до второго порядка включительно;

2) ядерные функции $\Phi(u) \in H$;

3) последовательности $c = c(n) \rightarrow 0, \beta = \beta(n) \rightarrow 0$ при $n \rightarrow \infty$, а $nc \rightarrow \infty$.

Тогда смещение

$$M(\bar{\varphi}(x) - \varphi(x)) \sim \varphi(x) \left(\frac{p(x)}{p_1(x)} - 1 \right) + \left(\frac{\beta^2}{3} + c^2 \right) \times$$

$$\begin{aligned} & \times \left(\frac{\varphi(x)p^{(2)}(x)}{2p_1(x)} + \frac{\varphi^{(1)}(x)p^{(1)}(x)}{p_1(x)} + \varphi^{(2)}(x)p(x) \right) + \\ & + \frac{\varphi^{(2)}(x)p^{(2)}(x)}{6} \beta^2 + O(c^4, \beta^4, c^2\beta^2), \end{aligned}$$

квадратическое отклонение

$$\begin{aligned} M(\bar{\varphi}(x) - \varphi(x))^2 & \sim \varphi^2(x)(p(x)/p_1(x) - 1)^2 + (\beta^2/3 + c^2) \times \\ & \times [(p(x)/p_1(x) - 1)(\varphi^2(x)p^{(2)}(x) + 2\varphi(x)\varphi^{(1)}(x)p^{(1)}(x))/p_1(x) + \\ & + 2\varphi(x)p(x)\varphi^{(2)}(x)(p(x)/2p_1^2(x) - 1)] + \beta^2(\varphi(x)p^{(1)}(x) + p(x)\varphi^{(1)}(x) + \\ & + p^2(x))/(3p_1^2(x)) + c^2\varphi^2(x)p^2(x)/4 + \\ & + (\varphi^2(x)p(x))/(ncp_1^2(x)) \int \Phi^2(u)du + O(c^4, \beta^4, c^2\beta^2, \beta^2/nc, c/n). \end{aligned}$$

Таким образом, статистика (7) обладает свойством асимптотической несмещенности и сходимостью в квадратическом, если закон распределения исходной выборки наблюдений x равномерный, т. е. $p(x) = p_1(x)$, что возможно при проведении активного эксперимента с исследуемой зависимостью $y = \varphi(x)$.

Для устранения смещения умножим $\bar{y} = \bar{\varphi}(x)$ на $p_1(x)/\bar{p}(x)$, получим

$$\bar{\bar{y}} = \sum_{i=1}^{N^2} y' \alpha'(x) / \sum_{i=1}^{N^2} \alpha'(x), \quad (8)$$

где $\alpha'(x) = P' \Phi \left(\frac{x - x'}{c} \right)$.

В многомерном случае при $x = (x_1, \dots, x_k)$ структура оценки (8) не меняется, а весовая функция принимает вид

$$\alpha'(x) = P' \prod_{v=1}^k \Phi \left(\frac{x_v - x'_v}{c_v} \right), \quad i = \overline{1, N^k}.$$

Анализ множеств случайных величин при восстановлении стохастических зависимостей. Определим на элементах множеств X', Y' исходной обучающей выборки $(x', y', i = \overline{1, N})$ непараметрические оценки плотностей вероятности [4]:

$$\bar{p}_i(x) \forall x \in X', \quad \bar{p}_i(y) \forall y \in Y', \quad i = \overline{1, N}.$$

Тогда два множества X^i, X^j близки, если соответствующие им статистические оценки функций распределения $\bar{F}_i(x), \bar{F}_j(x)$ тождественны, т. е. справедлива гипотеза

$$H_0: \bar{F}_i(x) \equiv \bar{F}_j(x) \quad (9)$$

с некоторым уровнем доверия β .

Пусть вероятностные законы распределения элементов случайных множеств X^i, X^j оцениваются по выборкам объема N_i, N_j . В качестве критерия близости между оценками функций распределения $\bar{F}_i(x)$ и $\bar{F}_j(x)$ может быть взята ядерная мера

$$h(X^i, X^j) = \prod_{v=1}^k \Phi(\bar{F}_i(x_v), \bar{F}_j(x_v));$$

$$\Phi(\bar{F}_i(x_v), \bar{F}_j(x_v)) = \begin{cases} (D_\beta + c)^{-1}, & \text{если } \Delta \leq D_\beta, \\ \frac{D_\beta + c - \Delta}{D_\beta + c}, & \text{если } D_\beta \leq \Delta \leq D_\beta + c, \\ 0, & \text{если } \Delta > D_\beta; \end{cases} \quad (10)$$

$$\Delta = \max_{x_v} |\bar{F}_i(x) - \bar{F}_j(x)|.$$

Предложенная мера близости (10) основана на использовании критерия Смирнова для проверки гипотезы H_0 (9) с уровнем доверия β :

$$D_\beta = \sqrt{-\ln \frac{\beta}{2} \left(\frac{1}{N_i} + \frac{1}{N_j} \right)} / 2.$$

Превышение порогового значения D_β критерия Смирнова означает нарушение гипотезы H_0 .

На основе введенных понятий запишем непараметрическую модель преобразования случайных множеств

$$\bar{p}(y) = \frac{\sum_{i=1}^N \bar{p}_i(y) h(X, X^i)}{\sum_{i=1}^N h(X, X^i)}, \quad (11)$$

$$G: \bar{p}(y) \rightarrow \bar{Y}.$$

Оператор G является датчиком случайных величин, с помощью которого формируется оценка \bar{Y} множества Y , соответствующего X .

Выбор порогового значения c осуществляется из условия минимума ошибки оценивания $\bar{p}(y)$ решающим правилом (11) в режиме «скользящего экзамена».

Если плотность $p(y)$ представима в виде линейного функционала от $p_j(y)$, $j=1, \overline{N_1}$, $N_1 < N$:

$$p(y) = \frac{1}{N_1} \sum_{j=1}^{N_1} p_j(y), \quad (12)$$

то при ограниченных $p_j^{(2)}(y)$, $j=1, \overline{N_1}$, и $c_j^2(N_j) \rightarrow 0$ при $N_j \rightarrow \infty$, $j=1, \overline{N_1}$, непараметрическая статистика $\bar{p}(y)$ является несмещенной оценкой $p(y)$:

$$M(\bar{p}(y) - p(y)) \sim \frac{1}{2N_1} \sum_{j=1}^{N_1} c_j^2 p_j^{(2)}(y) + O(c^4).$$

При дополнительных условиях $N_j c_j \rightarrow \infty \forall N_j \rightarrow \infty$, $c_j \rightarrow 0$, $j=1, \overline{N_1}$, имеет место асимптотическая сходимость $\bar{p}(y)$ к $p(y)$ в среднеквадратическом

$$M \int (\bar{p}(y) - p(y))^2 dy \sim \frac{1}{N_1^2} \sum_{j=1}^{N_1} \left[\frac{\|\Phi(u)\|^2}{N_j c_j} + \frac{c_j^4 \|p_j^{(2)}(y)\|^2}{4} (1 + N_1) \right].$$

Заключение. Рассмотрены непараметрические модели анализа множеств случайных величин в задачах восстановления стохастических зависимостей и распознавания образов. Идея предлагаемого подхода заключается в замене операций над множествами операциями с вероятностными законами распределения их элементов либо с параметрами законов с использованием методов локальной аппроксимации. Полученные результаты позволяют с новых теоретических позиций решить проблемы доверительного оценивания непараметрических статистик и их оптимизации. Наиболее перспективным направлением применения разработанных непараметрических моделей является обработка больших массивов статистических данных, часто встречающихся при исследовании медико-биологических и экологических систем.

СПИСОК ЛИТЕРАТУРЫ

1. Епаненчиков В. А. Непараметрическая оценка многомерной плотности вероятности // Теория вероятности и ее применения. 1969. **14**, вып. 1. С. 156.
2. Parzen E. On estimation of a probability density function and mode // Ann. Math. Stat. 1962. **33**. P. 1065.
3. Цыпкин Я. З. Основы теории обучающихся систем. М.: Наука, 1970.
4. Лапко А. В., Лапко В. А., Соколов М. И., Ченцов С. В. Непараметрические системы классификации. Новосибирск: Наука, 2000.

Институт вычислительного моделирования СО РАН,
E-mail: lapko@ksc.krasn.ru

Поступила в редакцию
24 апреля 2002 г.