

УДК 519.7

А. В. Лапко, В. А. Лапко

## НЕПАРАМЕТРИЧЕСКАЯ РЕГРЕССИЯ В УСЛОВИЯХ ЧАСТИЧНЫХ СВЕДЕНИЙ О ВИДЕ ВОССТАНАВЛИВАЕМЫХ ЗАКОНОМЕРНОСТЕЙ

Предлагается непараметрическая модель стохастической зависимости, которая обеспечивает учёт априорных сведений о виде восстанавливаемых закономерностей. Исследуются её асимптотические свойства, проводится анализ результатов вычислительных экспериментов.

Ключевые слова: непараметрическая регрессия, стохастические зависимости, априорные сведения, асимптотические свойства.

Для наиболее полного учёта априорной информации о виде восстанавливаемых зависимостей и экспериментальных данных о её локальном поведении широко используются гибридные модели [1]. Традиционные гибридные модели сочетают в одном решающем правиле преимущество параметрических и непараметрических аппроксимаций. При этом единое решающее правило образуют параметрическая модель восстанавливаемой зависимости и корректирующая её функция непараметрического типа, которые строятся в одном и том же пространстве переменных. Полученные результаты были развиты для условий наличия частичных априорных сведений о виде восстанавливаемых зависимостей в ограниченном пространстве признаков [2]. Основная проблема применения гибридных моделей состоит в выборе вида корректирующей функции, которая является трудно формализуемой. Для её обхода предлагается использовать непараметрическую регрессию, синтез которой основан на обобщении априорной информации о виде восстанавливаемых зависимостей и экспериментальных данных об их локальном поведении.

Цель работы состоит в обосновании возможности учёта априорных сведений о виде восстанавливаемых закономерностей при синтезе непараметрических моделей стохастических зависимостей, основанных на оценках плотности вероятности типа Розенблатта–Парзена [3].

**Синтез модифицированной непараметрической регрессии.** Пусть об искомой однозначной зависимости

$$y = \psi(x) \quad \forall \quad x \in R^k \quad (1)$$

известно её описание

$$\bar{y}_1 = F(\bar{x}_1, \alpha) \quad \forall \quad \bar{x}_1 \in R^{k^2}, \quad k^2 < k$$

относительно некоторого ограниченного набора признаков из  $x = (x_1, \bar{x}_1)$ ,  $x_1 = (x_{1v}, v = \overline{1, k^1})$ ,  $k = k^1 + k^2$  и выборка  $V = (x_v^i, v = \overline{1, k}, y^i, i = \overline{1, n})$  экспериментальных данных, составленная из статистически независимых значений переменных  $x, y$  исследуемой зависимости (1). Параметры  $\alpha$  полинома  $F(\bar{x}_1, \alpha)$  заданы.

Задача состоит в построении модифицированной непараметрической модели  $\bar{y}(x)$  зависимости (1), совмещающей в одном решающем правиле всю имеющуюся априорную информацию.

На основании исходных экспериментальных данных сформируем промежуточную обучающую выборку

$$V^1 = (x_1^i, \bar{y}_1^i = F(\bar{x}_1^i, \alpha), y^i, i = \overline{1, n}).$$

Принимая в качестве оптимального решающего правила, в смысле минимума среднеквадратического отклонения, условное математическое ожидание  $\Phi(x)$  [4], построим непараметрическую регрессию

$$\bar{y}(x) = \sum_{i=1}^n y^i \beta_i(x), \quad (2)$$

где

$$\beta_i(x) = \frac{\prod_{v=1}^{k^1} \Phi\left(\frac{x_{1v} - x_{1v}^i}{c_v}\right) \Phi\left(\frac{\bar{y}_1 - \bar{y}_1^i}{c}\right)}{\sum_{i=1}^n \prod_{v=1}^{k^1} \Phi\left(\frac{x_{1v} - x_{1v}^i}{c_v}\right) \Phi\left(\frac{\bar{y}_1 - \bar{y}_1^i}{c}\right)}.$$

В статистике (2) ядерные функции  $\Phi(u)$  удовлетворяют условиям  $H$ :

$$\Phi(u) = \Phi(-u), \quad 0 \leq \Phi(u) < \infty,$$

$$\int \Phi(u) du = 1, \quad \int u^2 \Phi(u) du = 1,$$

$$\int u^m \Phi(u) du < \infty, \quad 0 \leq m < \infty;$$

$c = c(n)$ ,  $c_v = c_v(n)$ ,  $v = \overline{1, k^1}$  – коэффициенты размытости ядерных функций, значения которых убывают с ростом объёма  $n$  обучающей выборки. Здесь и далее бесконечные пределы интегрирования опускаются.

При оценивании зависимости в ситуациях  $x = (x_1, \bar{x}_1)$  сначала вычисляется  $\bar{y}_1 = F(\bar{x}_1, \alpha)$ , а затем по данным  $(x_1, \bar{y}_1)$  в соответствии со статистикой (2) определяется значение  $\bar{y}(x)$ .

Оптимизация модифицированной непараметрической регрессии (2) по коэффициентам размытости

ядерных функций  $c, c_v, v = \overline{1, k1}$  осуществляется в режиме «скользящего экзамена» из условия минимума статистической оценки среднеквадратической ошибки аппроксимации искомой зависимости.

**Асимптотические свойства непараметрической регрессии.** Без существенной потери общности будем считать, что в частичном наборе признаков  $x_{1v}, v = \overline{1, k1}$  их количество  $k1 = 1$ . В качестве ядерной функции примем функцию вида

$$\Phi(u) = \begin{cases} \frac{1}{2} & \forall |u| < 1 \\ 0 & \forall |u| \geq 1. \end{cases}$$

В этом случае непараметрическая регрессия (2) запишется как

$$y(x) = \frac{\sum_{i=1}^n y^i \Phi\left(\frac{x_1 - x_1^i}{c_1}\right) \Phi\left(\frac{\bar{y}_1 - \bar{y}_1^i}{c}\right)}{\sum_{i=1}^n \Phi\left(\frac{x_1 - x_1^i}{c_1}\right) \Phi\left(\frac{\bar{y}_1 - \bar{y}_1^i}{c}\right)}. \quad (3)$$

Тогда справедливо следующее утверждение.

*Теорема.* Пусть 1) частичные сведения  $\bar{y}_1 = F(\bar{x}_1, \alpha)$  восстанавливаемой зависимости (1) принадлежат к классу линейных полиномов; 2) функция  $y = \varphi(x)$  и плотности вероятностей  $p(x), p(x, y_1)$  ограничены вместе со своими производными до второго порядка включительно; 3) ядерные функции  $\Phi(u)$  являются положительными, симметричными и нормированными; 4) последовательности коэффициентов размытости  $c_1(n), c(n)$  ядерных функций таковы, что при  $n \rightarrow \infty$  их значения стремятся к нулю. Тогда непараметрическая регрессия (3) обладает свойством асимптотической несмещённости.

*Доказательство.* Представим модель (3) в виде

$$y(x) = \frac{(nc_1c)^{-1} \sum_{i=1}^n y^i \Phi\left(\frac{x_1 - x_1^i}{c_1}\right) \Phi\left(\frac{\bar{y}_1 - \bar{y}_1^i}{c}\right)}{(nc_1c)^{-1} \sum_{i=1}^n \Phi\left(\frac{x_1 - x_1^i}{c_1}\right) \Phi\left(\frac{\bar{y}_1 - \bar{y}_1^i}{c}\right)} = \frac{\bar{z}_1(x)}{\bar{z}_2(x)}. \quad (4)$$

Проведем преобразования

$$M \frac{\bar{z}_1(x)}{\bar{z}_2(x)} = M \left[ \frac{\bar{z}_1(x)}{M \bar{z}_2(x)} + \frac{\bar{z}_1(x)}{\bar{z}_2(x)} \frac{M \bar{z}_2(x) - \bar{z}_2(x)}{M \bar{z}_2(x)} \right], \quad (5)$$

где  $M$  – знак математического ожидания.

Ввиду ограниченности значений  $\bar{y}_1(x) = \frac{\bar{z}_1(x)}{\bar{z}_2(x)}$  свойства статистики (3) зависят от асимптотического поведения  $M(\bar{z}_1(x)), M(\bar{z}_2(x))$ . Вычислим

$$M(\bar{z}_2(x)) = (nc_1c)^{-1} \sum_{i=1}^n \int \dots \int \Phi\left(\frac{x_1 - x_1^i}{c_1}\right) \times \\ \times \Phi\left(\sum_{v=1}^{k2} \alpha_v \frac{\bar{x}_{1v} - \bar{x}_{1v}^i}{c}\right) p(x_1^i, \bar{x}_{1v}^i, v = \overline{1, k2}) dx_1^i d\bar{x}_{1v}^i \dots d\bar{x}_{1k2}^i.$$

Так как  $x_1^i, \bar{x}_{1v}^i, v = \overline{1, k2}$  являются значениями одних и тех же случайных величин  $t, t_v, v = \overline{1, k2}$  с плотностью вероятности  $p(t, t_v, v = \overline{1, k2})$ , то

$$M(\bar{z}_2(x)) = (c_1c)^{-1} \int \dots \int \Phi\left(\frac{x_1 - t}{c_1}\right) \times \\ \times \Phi\left(\sum_{v=1}^{k2} \frac{\alpha_v}{c} (\bar{x}_{1v} - t_v)\right) p(t, t_v, v = \overline{1, k2}) dt dt_1 \dots dt_{k2}.$$

Проведём замену переменных  $u = \frac{x_1 - t}{c_1}$ ,

$u_v = \frac{\alpha_v (\bar{x}_{1v} - t_v)}{c}$ . После несложных преобразований получим

$$M(\bar{z}_2(x)) = \frac{c^{k2-1}}{\prod_{v=1}^{k2} \alpha_v} \int \dots \int \Phi(u) \Phi\left(\sum_{v=1}^{k2} u_v\right) \times \\ \times p\left(x_1 - c_1u, \bar{x}_{1v} - \frac{c}{\alpha_v} u_v, v = \overline{1, k2}\right) du du_1 \dots du_{k2}. \quad (6)$$

Разложим функцию

$$p\left(x_1 - c_1u, \bar{x}_{1v} - \frac{c}{\alpha_v} u_v, v = \overline{1, k2}\right)$$

в ряд Тейлора в точке  $x = (x_1, x_{1v}, v = \overline{1, k2})$  и преобразуем (6) с учётом свойств:

$$\frac{1}{2^{k2-1}} \int_{-1}^1 \dots \int_{-1}^1 \Phi\left(\sum_{v=1}^{k2} u_v\right) du_1 \dots du_{k2} = 1, \\ \frac{1}{2^{k2-1}} \int_{-1}^1 \dots \int_{-1}^1 u_t \Phi\left(\sum_{v=1}^{k2} u_v\right) du_1 \dots du_{k2} = 0, \quad t = \overline{1, k2}, \\ \frac{1}{2^{k2-1}} \int_{-1}^1 \dots \int_{-1}^1 u_t^2 \Phi\left(\sum_{v=1}^{k2} u_v\right) du_1 \dots du_{k2} = \beta^2.$$

В результате при  $n \rightarrow \infty$  имеем

$$M(\bar{z}_2(x)) \sim \frac{2^{k2-1} c^{k2-1}}{\prod_{v=1}^{k2} \alpha_v} \times \\ \times \left[ p(x) + \frac{c_1^2}{2} p_{x_1}^{(2)}(x) + \frac{c^2 \beta^2}{2^{k2}} \sum_{v=1}^{k2} \frac{1}{\alpha_v^2} p_v^{(2)}(x) + o(c^4) \right], \quad (7)$$

где  $p_{x_1}^{(2)}(x), p_v^{(2)}(x)$  – вторые производные плотности вероятности  $p(x_1, \bar{x}_{1v}, v = \overline{1, k2})$  по переменным  $x_1, \bar{x}_{1v}, v = \overline{1, k2}$  соответственно.

Следуя приведённой технологии вычислений, найдём асимптотическое выражение для

$$\begin{aligned}
 M(\bar{z}_1(x)) &= (nc_1c)^{-1} \sum_{i=1}^n \int \dots \int y^i \Phi\left(\frac{x_1 - x_1^i}{c_1}\right) \Phi\left(\sum_{v=1}^{k_2} \alpha_v \frac{\bar{x}_{1v} - \bar{x}_{1v}^i}{c}\right) \times \\
 &\times p\left(y^i, x_1^i, \bar{x}_{1v}^i, v = \overline{1, k_2}\right) dy^i dx_1^i d\bar{x}_{11}^i \dots d\bar{x}_{1k_2}^i = \\
 &= \frac{c^{k_2-1}}{\prod_{v=1}^{k_2} \alpha_v} \int \dots \int \varphi\left(x_1 - c_1 u, \bar{x}_{1v} - \frac{c}{\alpha_v} u_v, v = \overline{1, k_2}\right) \times \\
 &\times \Phi(u) \Phi\left(\sum_{v=1}^{k_2} u_v\right) p\left(x_1 - c_1 u, \bar{x}_{1v} - \frac{c}{\alpha_v} u_v, v = \overline{1, k_2}\right) du du_1 \dots du_{k_2} \sim \\
 &\sim \frac{2^{k_2-1} c^{k_2-1}}{\prod_{v=1}^{k_2} \alpha_v} \left[ \varphi(x) p(x) + \frac{c_1^2}{2} (\varphi(x) p(x))_{x_1}^{(2)} + \right. \\
 &\left. + \frac{c^2 \beta^2}{2^{k_2}} \sum_{v=1}^{k_2} \frac{1}{\alpha_v^2} (\varphi(x) p(x))_{x_{1v}}^{(2)} + 0(c_1^2 c^2, c_1^4 c^4, v = \overline{1, k_2}) \right]. \quad (8)
 \end{aligned}$$

В выражении (8)  $(\varphi(x) p(x))_{x_1}^{(2)}$ ,  $(\varphi(x) p(x))_{x_{1v}}^{(2)}$  – вторые производные произведения двух функций по переменным  $x_1$ ,  $\bar{x}_{1v}, v = \overline{1, k_2}$  соответственно.

Подставим выражения (7) и (8) в (5), получим

$$\begin{aligned}
 M(\bar{y}(x)) &\sim M\left(\frac{\bar{z}_1(x)}{\bar{z}_2(x)}\right) \sim \\
 &\frac{\varphi(x) p(x) + \frac{c_1^2}{2} (\varphi(x) p(x))_{x_1}^{(2)} + \frac{c^2 \beta^2}{2^{k_2}} \sum_{v=1}^{k_2} \frac{1}{\alpha_v^2} (\varphi(x) p(x))_{x_{1v}}^{(2)}}{p(x) + \frac{c_1^2}{2} p_{x_1}^{(2)}(x) + \frac{c^2 \beta^2}{2^{k_2}} \sum_{v=1}^{k_2} \frac{1}{\alpha_v^2} p_v^{(2)}(x)}. \quad (9)
 \end{aligned}$$

Из анализа выражения (9) следует, что при  $c_1 = c_1(n) \rightarrow 0$ ,  $c = c(n) \rightarrow 0$  с ростом  $n \rightarrow \infty$  изучаемая статистика (3) обладает свойством асимптотической несмещённости.

*Замечания.* При  $k_2 = 2$  полученные результаты могут быть использованы при исследовании свойств статистических моделей, основанных на методе группового учёта аргументов [5]. Идея метода заключается в построении последовательности моделей

$$\bar{y}_j = \bar{\varphi}_j(x_j, \bar{y}_{j-1}), i = \overline{1, m}. \quad (10)$$

Ранее не используемая в моделях  $\bar{y}_t, t = \overline{1, j-1}$  компонента  $x_j$  вектора аргументов  $x$  обеспечивает в наборе с  $\bar{y}_{j-1}$  минимальное расхождение значений  $\bar{y}_j$  с экспериментальными данными. На каждом этапе процедуры (10) искомая зависимость оценивается в пространстве двух переменных  $(x_j, \bar{y}_{j-1})$ .

**Анализ результатов вычислительных экспериментов.** На основании данных вычислительных экс-

периментов сравнивалась эффективность статистики (2) и традиционной непараметрической регрессии

$$\tilde{y}(x) = \frac{\sum_{i=1}^n y^i \prod_{v=1}^k \Phi\left(\frac{x_v - x_v^i}{c_v}\right)}{\sum_{i=1}^n \prod_{v=1}^k \Phi\left(\frac{x_v - x_v^i}{c_v}\right)}. \quad (11)$$

В качестве искомой зависимости (1) использовался полином второй степени

$$\psi(x) = x_1^2 + 2x_2^2 + x_1 x_2 + x_3 + 0,5x_4 + 2x_5, \quad (12)$$

каждый аргумент которого принимает значения из интервала  $x_v \in [0, 1]$ ,  $v = \overline{1, 5}$  с равномерным законом распределения. Частичные сведения о восстанавливаемой зависимости в соответствии с условиями теоремы определяются линейным полиномом

$$\bar{y}_1 = x_3 + 0,5x_4 + 2x_5.$$

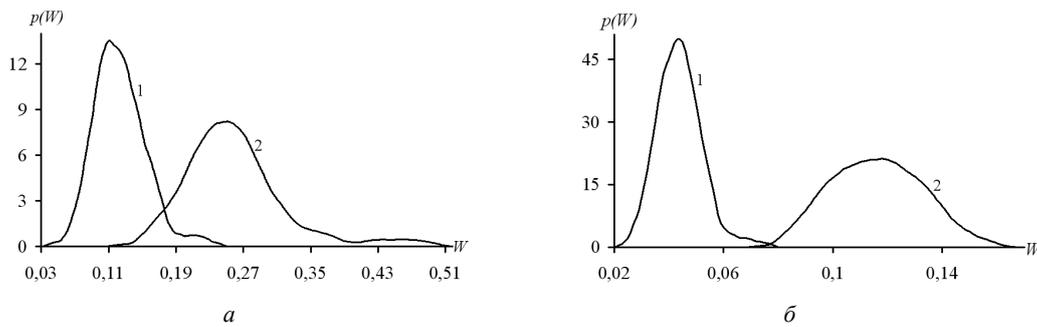
При формировании обучающей выборки  $V = (x_v^i, v = \overline{1, 5}, y^i, i = \overline{1, n})$  на значения функции (12) накладывалась аддитивная помеха

$$y^i = \psi(x^i) \left(1 + 2(\varepsilon^i - 0,5)r\right), \quad (13)$$

где  $\varepsilon$  – случайная величина с равномерным законом в диапазоне  $[0; 1]$ ;  $r$  – параметр, определяющий уровень шума.

При синтезе непараметрических моделей (2), (11) использовалась ядерная функция Епанечникова, а их оптимизация по коэффициентам размытости осуществлялась в режиме «скользящего экзамена» из условия минимума среднеквадратического критерия. При этом полагалось, что значения коэффициентов размытости  $c_v = c, v = \overline{1, 5}$  для непараметрической регрессии и  $\bar{c}_v = \bar{c}, v = \overline{1, 2}$  – для модели (2), так как интервалы изменения аргументов восстанавливаемой зависимости априори одинаковые. В качестве критерия эффективности моделей (2), (11) принимались среднеквадратические отклонения  $W_2, W_{11}$  их значений от функции (12), которые оценивались по контрольной выборке  $V_k$  объёма  $n_k = 10\,000$ . При этом ситуации из выборки  $V_k$ , в которых исследуемые непараметрические модели не идентифицируют значения функции (12), не участвуют в формировании критериев их эффективности. Доля таких ситуаций не превышает значений 0,06 от объёма контрольной выборки.

Вычислительные эксперименты при фиксированных условиях исследования осуществлялись 60 раз. По полученным результатам восстанавливались плотности вероятности  $p(W_2), p(W_{11})$  оценок среднеквадратических отклонений  $W_2, W_{11}$  соответственно моделей (2), (11) (см. рисунок).



Оценки плотностей вероятностей  $p(W_2)$ ,  $p(W_{11})$  среднеквадратических отклонений модифицированной непараметрической модели  $\bar{y}$  (кривая 1) и непараметрической регрессии  $\hat{y}$  (кривая 2). Условия эксперимента: объём обучающей выборки  $n = 50$  (а),  $n = 200$  (б); уровень шума в процедуре (13)  $r = 0,1$

Анализ данных вычислительных экспериментов показывает, что статистические оценки законов распределения значений критериев эффективности  $W_2$ ,  $W_{11}$  достоверно отличаются при различных объёмах обучающих выборок, причём интервалы изменения среднеквадратического критерия  $W_{11}$  непараметрической регрессии (11) характеризуются большими их значениями по сравнению с модифицированной регрессией (2). С ростом объёма обучающей выборки преимущество моделей (2) возрастает. Например, отношение  $R = \bar{W}_{11} / \bar{W}_2$  средних значений  $W_2$ ,  $W_{11}$  соответствует 2, 3 при  $n = 50$  и 3 при  $n = 200$ .

Эффективность модифицированной непараметрической модели (2) объясняется возможностью снижения её размерности за счёт использования априорных сведений о наличии линейной взаимосвязи между переменными исследуемой зависимости. Данное заключение согласуется с результатами исследования гибридных моделей стохастических зависимостей [1].

Традиционная непараметрическая регрессия, основанная на оценке плотности вероятности типа Розенблатта–Парзена, обобщена при построении статистических моделей в условиях наличия частичных сведений о виде восстанавливаемых зависимостей. Предлагаемая модифицированная непараметрическая

регрессия обладает свойством асимптотической несмещённости. Это позволяет аналитически обосновать возможность частичного сжатия пространства признаков на основе линейных преобразований без существенной потери полезной информации.

Перспективное направление дальнейших исследований состоит в развитии предлагаемого подхода на анализ свойств статистических моделей, основанных на методе группового учёта аргументов.

#### Библиографические ссылки

1. Лапко А. В., Лапко В. А. Гибридные модели стохастических зависимостей // Автометрия. 2002. № 5. С. 38–48.
2. Лапко В. А. Синтез и анализ гибридных моделей стохастических зависимостей в условиях наличия их частного описания // Автометрия. 2004. № 1. С. 51–59.
3. Parzen E. On estimation of a probability density function and mode // Ann. Math. Statistic. 1962. Vol. 33. P. 1065–1076.
4. Надарая Э. А. Непараметрические оценки кривой регрессии // Тр. ВЦ АН СССР. 1965. Вып. 5. С. 56–68.
5. Ивахненко А. Г. Непараметрический комбинированный алгоритм МГУА на операторах поиска аналогов // Автоматика. 1990. № 5. С. 14–27.

A. V. Lapko, V. A. Lapko

#### NONPARAMETRIC REGRESSION IN THE CONDITIONS OF PARTIAL DATA ON A MODE OF RESTORED LEGITIMACIES

*The nonparametric model of stochastic dependence which provides for registration of a priori data on a mode of restored legitimacies is offered. Its asymptotic properties are researched.*

*Keywords: nonparametric regression, stochastic dependences, a priori data, asymptotic properties.*

© Лапко А. В., Лапко В. А., 2011