

И. В. Г е т м а н с к а я

**СОСТОЯТЕЛЬНАЯ ОЦЕНКА ПАРАМЕТРА
ОДНОПАРАМЕТРИЧЕСКОЙ ПАРНОЙ
РЕГРЕССИИ**

Рассматривается реализация параметра однопараметрической парной регрессии в точке с координатами, соответствующими значениям фактора и отклика. Найдены ее основные числовые характеристики, с помощью которых определена состоятельная оценка регрессионного коэффициента.

Предположим, требуется построить математическую модель $Y = r(X)$, связывающую независимую переменную X и зависимую переменную Y , по результатам n наблюдений значений фактора X_i и отклика $y_i = Y_i + \varepsilon_i$, $i = \overline{1, n}$, содержащего в себе погрешность ε_i , возникающую либо из-за неучтенных факторов, либо из-за ошибок измерений.

Класс функций, в котором целесообразно искать наилучшую (в определенном смысле) аппроксимацию зависимости $Y = r(X)$, можно подобрать по внешнему виду экспериментальной зависимости y от X или исходя из физических соображений, связанных с существом решаемой задачи.

Предположим, \mathfrak{S} — класс допустимых моделей зависимостей. Аппроксимирующая функция $r(X)$ из класса \mathfrak{S} , называется линией регрессии. Если класс \mathfrak{S} задается некоторым параметрическим семейством функций $\{r(X, \theta)\}$, где θ — свободный параметр, то θ называется регрессионным коэффициентом, а аппроксимирующие функции $Y = r(X, \theta) \in \mathfrak{S}$ называют парными однопараметрическими регрессиями. В настоящий момент разработаны различные методы оценивания параметра θ . Наибольшее распространение среди них получили алгоритмы итерационного типа [1, 2], для которых "... первостепенное значение имеет удачный выбор начального приближения" [1]. Принято в качестве начального приближения использовать одно из значений параметра θ , реализующего зависимость $Y = r(X, \theta)$ в точках с координатами (X_i, y_i) , $i = \overline{1, n}$ [1].

В работе приводится доказательство того, что среднее значение реализаций $\hat{\theta}_i$ параметра θ регрессии $Y = r(X, \theta)$ в экспериментальных данных (X_i, y_i) с учетом поправки является состоятельной оценкой регрессионного коэффициента.

Регрессионный коэффициент, реализующий значение фактора и отклика. Пусть

$$Y = r(X, \theta) \in \Upsilon \subseteq \mathbb{R}^1 \quad (1)$$

— функция регрессии известного вида, определенная с точностью до подлежащего оценке регрессионного коэффициента $\theta \in \Theta \subseteq \mathbb{R}^1$. Здесь $X \in \aleph \subseteq \mathbb{R}^1$ — независимая переменная. Пусть (X, y) — пара наблюдений значения функции отклика y , полученных при значении объясняющей переменной X . Причем

$$y = Y + \varepsilon, \quad (2)$$

где погрешность ε — случайная величина с начальными моментами

$$E\varepsilon = 0, \quad E\varepsilon^2 =: \sigma^2. \quad (3)$$

Предположим, что в точках области $A \in \aleph \times \Upsilon$, $A = \{(X, y) \in \aleph \times \Upsilon : |y - Y| \leq \varepsilon_m\}$, где $|\varepsilon| < \varepsilon_m$, существует и единственная заданная неявно уравнением (1) функция $\theta = \theta(X, y)$, отображающая $A \xrightarrow{\theta(X, y)} D$, где $D \in \Theta$.

Теорема 1. Пусть выполняются условия (1–3), а также а) $|\varepsilon| < 1$ ($\varepsilon_m = 1$), кроме того, в области $\aleph \times D$: б) функция регрессии $Y = r(X, \theta)$ строго монотонная по переменной θ ; в) существуют r'_θ , $r''_{\theta\theta}$, $r'''_{\theta\theta\theta}$ и постоянная $q > 0$: $|r'_\theta| > q$. Тогда реализация $\overset{\circ}{\theta}$ параметра θ функции регрессии $Y = r(X, \theta)$ в точке $(X, y) \in A$ — смещенная оценка θ с приближением смещения

$$W \cong -\sigma^2 \frac{r''_{\theta\theta}(X, \theta)}{2(r'_\theta(X, \theta))^3}.$$

Доказательство. По определению [3] оценка $\overset{\circ}{\theta}$ параметра θ — несмещенная оценка, если выполняется равенство $E[\overset{\circ}{\theta}] = \theta$. По постановке задачи $\theta = \theta(X, Y)$ решенное относительно θ уравнение $R(X, Y, \theta) = 0$, где $R(X, Y, \theta) = Y - r(X, \theta)$. Так как $y = r(X, \overset{\circ}{\theta})$, то $\overset{\circ}{\theta} = \theta(X, y)$. Функцию $\theta(X, y)$ представим как функцию, зависящую от случайного аргумента ε : $\theta(X, y) = \theta(X, Y + \varepsilon) = \theta(\varepsilon)$.

По определению математического ожидания неслучайной функции θ от случайного аргумента ε имеем

$$E[\overset{\circ}{\theta}] = \int_{-1}^1 \theta(\varepsilon) dF_\varepsilon.$$

Выполнение условий б) и в) теоремы означает дифференцируемость функции $\theta(X, Y)$ по Y в точке $(X, Y) \in A$ до третьего порядка

включительно, что позволяет представить функцию $\theta(X, Y)$ по формуле Тейлора в ε -окрестности точки Y с точностью до бесконечно малой $o(y - Y)^2 = o(\varepsilon^2)$:

$$\theta(X, Y + \varepsilon) = \theta(X, Y) + \varepsilon\theta'_Y(X, Y) + \frac{(\varepsilon)^2}{2}\theta''_{YY}(X, Y) + o(\varepsilon^2).$$

В результате получим

$$E[\overset{\circ}{\theta}] = \int_{-1}^1 \left(\theta(X, Y) + \varepsilon\theta'_Y(X, Y) + \frac{(\varepsilon)^2}{2}\theta''_{YY}(X, Y) + o(\varepsilon^2) \right) dF_\varepsilon =$$

$$\theta(X, Y) \int_{-1}^1 dF_\varepsilon + \theta'_Y(X, Y) \int_{-1}^1 \varepsilon dF_\varepsilon + \frac{\theta''_{YY}(X, Y)}{2} \int_{-1}^1 \varepsilon^2 dF_\varepsilon + \int_{-1}^1 o(\varepsilon^2) dF_\varepsilon.$$

Поскольку $\int_{-1}^1 dF_\varepsilon = 1$, $\int_{-1}^1 \varepsilon dF_\varepsilon = E\varepsilon$, $\int_{-1}^1 \varepsilon^2 dF_\varepsilon = E\varepsilon^2$, а $\int_{-1}^1 o(\varepsilon^2) dF_\varepsilon = E[o(\varepsilon^2)]$, кроме того, $\theta = \theta(X, Y)$ — истинное значение θ , а по условию (3) $E\varepsilon = 0$, $E\varepsilon^2 =: \sigma^2$, то

$$E[\overset{\circ}{\theta}] = \theta + \sigma^2 \frac{\theta''_{YY}(X, Y)}{2} + E[o(\varepsilon^2)].$$

Обозначим $W = \sigma^2 \frac{\theta''_{YY}(X, Y)}{2}$. Пренебрегая третьим слагаемым, по-

лучим $\tilde{E}[\overset{\circ}{\theta}] \approx \theta + W$ — приближение $E[\overset{\circ}{\theta}]$.

Выразим смещение W через функцию регрессии $r(X, \theta)$. Для этого найдем значение $\theta''_{YY}(X, Y)$ по правилу производной неявно заданной $R(X, Y, \theta) = Y - r(X, \theta) = 0$ функции $\theta(X, Y)$ [4]:

$$\frac{\partial^2 \theta(X, Y)}{\partial Y^2} = \frac{-r''_{\theta\theta}}{(r'_\theta)^3}.$$

В результате получим $W = -\sigma^2 \frac{r''_{\theta\theta}(X, \theta)}{2(r'_\theta(X, \theta))^3}$. Что и требовалось доказать.

Теорема 2. При выполнении условий **теоремы 1** дисперсия $Var[\overset{\circ}{\theta}]$ регрессионного коэффициента $\overset{\circ}{\theta}$, реализующего функцию регрессии $Y = r(X, \theta)$ в точке (X, y) , приближенно равна

$$Var[\overset{\circ}{\theta}] \cong \left(\frac{\sigma}{r'_\theta(X, \theta)} \right)^2.$$

Доказательство. Дисперсию $Var \left[\overset{\circ}{\theta} \right]$ найдем по формуле

$$Var \left[\overset{\circ}{\theta} \right] = E \left[\left(\overset{\circ}{\theta} \right)^2 \right] - \left(E \left[\overset{\circ}{\theta} \right] \right)^2.$$

Начальный момент второго порядка $E \left[\left(\overset{\circ}{\theta} \right)^2 \right]$ случайной величины $\overset{\circ}{\theta}$ можно найти как начальный момент первого порядка неслучайной функции

$$\left(\overset{\circ}{\theta} \right)^2 = (\theta(X, y))^2 = (\theta(X, Y + \varepsilon))^2 = (\theta(\varepsilon))^2 = \varphi(\varepsilon)$$

от случайного аргумента ε [5], а именно:

$$E [\varphi(\varepsilon)] = \int_{-\infty}^{\infty} \varphi(\varepsilon) dF_{\varepsilon}.$$

Если

$$\theta(X, Y + \varepsilon) = \theta(X, Y) + \varepsilon \theta'_Y(X, Y) + \frac{(\varepsilon)^2}{2} \theta''_{YY}(X, Y) + o(\varepsilon^2)$$

— приближение функции $\theta(X, Y + \varepsilon)$ по формуле Тейлора в ε -окрестности точки Y , то

$$\begin{aligned} \theta^2(X, Y + \varepsilon) &= \theta^2(X, Y) + 2\varepsilon \theta(X, Y) \theta'_Y(X, Y) + \\ &+ \varepsilon^2 \left((\theta'_Y(X, Y))^2 + \theta(X, Y) \theta''_{YY}(X, Y) \right) + \psi(o(\varepsilon^2)), \end{aligned}$$

где $\psi(o(\varepsilon^2))$ включает в себя слагаемые, зависящие от степеней ε выше второй. Пренебрегая $\psi(o(\varepsilon^2))$ и интегрируя последнее выражение на интервале $(-1, 1)$, получим

$$E \left[\left(\overset{\circ}{\theta} \right)^2 \right] \cong \theta^2 + \sigma^2 \left((\theta'_Y(X, Y))^2 + \theta \theta''_{YY}(X, Y) \right).$$

В результате

$$\begin{aligned} Var \left[\overset{\circ}{\theta} \right] &\cong E \left[\left(\overset{\circ}{\theta} \right)^2 \right] - \left(E \left[\overset{\circ}{\theta} \right] \right)^2 = \\ &= \theta^2 + \sigma^2 \left((\theta'_Y(X, Y))^2 + \theta \theta''_{YY}(X, Y) \right) - \\ &- \left(\theta^2 + \sigma^2 \theta \theta''_{YY}(X, Y) \right) = \sigma^2 (\theta'_Y(X, Y))^2. \end{aligned}$$

Так как

$$\frac{\partial \theta(X, Y)}{\partial y} = \frac{1}{r'_\theta},$$

то приближенно

$$\text{Var} [\overset{\circ}{\theta}] \cong \frac{\sigma^2}{\left(r'_\theta(X, \theta)\right)^2}.$$

Состоятельная оценка регрессионного коэффициента. Пусть выполняются условия (1)–(3) и имеются результаты n независимых наблюдений значений фактора X_i и отклика

$$y_i = Y_i + \varepsilon_i, \quad i = \overline{1, n}, \quad (4)$$

где ε_i — реализации случайной величины ε .

Обозначим

$$\overset{\circ}{\theta} = \left(\overset{\circ}{\theta}_1, \overset{\circ}{\theta}_2, \dots, \overset{\circ}{\theta}_n\right),$$

где $\overset{\circ}{\theta}_i$ — значения регрессионного коэффициента θ , реализующие функцию регрессии $Y = r(X, \theta)$ в точках (X_i, y_i) , $i = \overline{1, n}$. Получить $\overset{\circ}{\theta}_i$ можно, решив систему уравнений $y_i = r\left(X_i, \overset{\circ}{\theta}_i\right)$, $i = \overline{1, n}$, относительно $\overset{\circ}{\theta}_i$. Необходимо обратить внимание на то, что в общем случае решение каждого из уравнений этой системы является не единственным. На практике получить единственное решение можно, сузив область возможных значений $\overset{\circ}{\theta}_i$, “... исходя из условий конкретной задачи” [6].

Пусть

$$\overset{\circ}{\bar{\theta}} := n^{-1} \sum_{i=1}^n \overset{\circ}{\theta}_i \quad (5)$$

— среднее значение ряда $\overset{\circ}{\theta}$.

Теорема 3. При выполнении условий (1–3), а в каждой точке (X_i, θ) , $i = \overline{1, n}$, условий (4) и условий **теоремы 2**

$$\tilde{\theta} = \overset{\circ}{\bar{\theta}} - n^{-1} \sum_{i=1}^n W_i \quad (6)$$

является состоятельной оценкой регрессионного коэффициента θ .

Здесь $W_i = -\sigma^2 \frac{r''_{\theta\theta}(X_i, \theta)}{2(r'_\theta(X_i, \theta))^3}$.

Доказательство. Для доказательства состоятельности оценки $\tilde{\theta}$ регрессионного коэффициента θ воспользуемся следующей теоремой: если оценка $\tilde{\theta}$ параметра θ — несмещенная и $\text{Var} [\tilde{\theta}] \rightarrow 0$ при $n \rightarrow \infty$, то $\tilde{\theta}$ — состоятельная оценка параметра θ .

Таким образом, необходимо доказать, что

- 1) $E [\tilde{\theta}] = \theta$;
- 2) $Var [\tilde{\theta}] \rightarrow 0$ при $n \rightarrow \infty$.

Доказательство п. 1.

$$E [\tilde{\theta}] = E \left[\overset{\circ}{\theta} - n^{-1} \sum_{i=1}^n W_i \right] = E \left[n^{-1} \sum_{i=1}^n \left(\overset{\circ}{\theta}_i - W_i \right) \right].$$

Исходя из свойств математического ожидания, с учетом того, что W_i — не случайная величина, имеем:

$$E \left[n^{-1} \sum_{i=1}^n \left(\overset{\circ}{\theta}_i - W_i \right) \right] = n^{-1} \sum_{i=1}^n \left(E [\overset{\circ}{\theta}_i] - W_i \right).$$

Так как для любого i выполняются все условия теоремы 1, то $E [\overset{\circ}{\theta}_i] = \theta + W_i$, откуда следует, что $E [\tilde{\theta}] = n^{-1}n\theta = \theta$.

Доказательство п. 2.

$$Var [\tilde{\theta}] = Var \left[n^{-1} \sum_{i=1}^n \left(\overset{\circ}{\theta}_i - W_i \right) \right].$$

Величины $\overset{\circ}{\theta}_i$, $i = \overline{1, n}$, — независимые, как значения неслучайной функции $\overset{\circ}{\theta}_i = \theta (X_i, Y_i + \varepsilon_i)$ независимых случайных аргументов ε_i . В результате $Var [\tilde{\theta}]$ найдем как дисперсию линейной комбинации независимых случайных величин:

$$Var \left[n^{-1} \sum_{i=1}^n \left(\overset{\circ}{\theta}_i - W_i \right) \right] = Var \left[n^{-1} \sum_{i=1}^n \overset{\circ}{\theta}_i \right] = n^{-2} \sum_{i=1}^n Var [\overset{\circ}{\theta}_i].$$

Обозначим $M = \max_i Var [\overset{\circ}{\theta}_i]$. По **теореме 2** $M = \max_i \sigma^2 (r'_\theta (X_i, \theta))^{-2}$.

Согласно условию **с) теоремы 1** $\forall i \exists$ постоянная $q > 0: |r'_\theta (X_i, \theta)| \geq q$.

Тогда $M \leq \left(\frac{\sigma}{q} \right)^2 < \infty$. Так как $0 \leq M < \infty$, то $0 \leq Var [\tilde{\theta}] \leq n^{-2}n \times M = n^{-1}M \xrightarrow{n \rightarrow \infty} 0$, откуда следует $Var [\tilde{\theta}] \xrightarrow{n \rightarrow \infty} 0$, что и требовалось доказать.

Вычислительный эксперимент оценивания регрессионного коэффициента однопараметрической парной регрессии. Вычислительный эксперимент проводился в компьютерной системе

математических символьных вычислений Maple. Модель значений отклика строилась по формуле (2). Регрессия задавалась функцией вида

$$r(X, \theta) = \text{ch}(X/\theta)$$

при значении параметра $\theta = 1$, заданной на интервале $\aleph = (0, 2, 2)$.

Моделирование погрешности отклика ε осуществлялось программой random статистического пакета Stats системы Maple. Значения ε строились как элементы выборки объема n равномерного распределения на интервале $[-\sqrt{3}\sigma, \sqrt{3}\sigma]$, $\sigma = 0,202$. На рис. 1 изображен один из результатов моделирования.

В таблице приведены результаты эксперимента оценивания по формулам (5), (6). Значения первого столбца соответствуют номеру эксперимента. Значения N соответствуют количеству найденных значений $\overset{\circ}{\theta}_i$ как решений уравнений $y_i = r(x_i, \overset{\circ}{\theta}_i)$, $i = \overline{1, n}$ в области $\Theta = [\theta/10, 2\theta]$ предполагаемых значений θ . Уравнения решались численным методом половинного деления, программно реализованным в среде Maple. Значения N' соответствуют количеству значений $\overset{\circ}{\theta}_i$, для которых выполняются условия *теоремы 1*. Значения $\varepsilon_{\overset{\circ}{\theta}} = \left| \overset{\circ}{\theta} - \theta \right|$ соответствуют абсолютной погрешности оценки по фор-

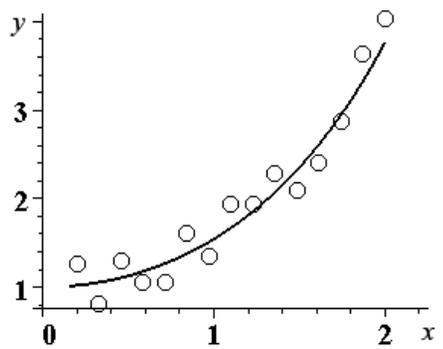


Рис. 1. Модель значений фактора и отклика:

— — линия регрессии $Y = r(X, \theta)$;
 \circ — точки с координатами (X_i, y_i) , где $y_i = Y_i + \varepsilon_i$

№	n	N	$\varepsilon_{\overset{\circ}{\theta}}$	N'	$\varepsilon_{\bar{\theta}}$
1	4	3	0,2309	2	0,0199
2	6	5	0,0603	3	0,0408
3	9	7	0,0090	5	0,0029
4	11	8	0,1841	5	0,0462
5	15	13	0,0277	9	0,0135
6	18	14	0,0513	10	0,0246
7	20	17	0,056	10	0,0007
8	26	21	0,0631	13	0,0005
9	31	22	0,0075	18	0,0058
10	41	33	0,0161	23	0,0146
11	51	49	0,0524	30	0,0111

муле (5) — $\overset{\circ}{\theta} = N^{-1} \sum_{i=1}^N \overset{\circ}{\theta}_i$, а $\varepsilon_{\tilde{\theta}} = |\tilde{\theta} - \theta|$ — оценкам по формуле (6) —

$\tilde{\theta} = (N')^{-1} \sum_{i=1}^{N'} \left(\overset{\circ}{\theta}_i - W_i \right)$. Величины $W_i = -\sigma^2 \frac{r''_{\theta\theta}(X_i, \theta)}{2(r'_{\theta}(X_i, \theta))^3}$, завися-

щие от оцениваемого параметра θ , были заменены их приближенными

значениями $\tilde{W}_i = -\sigma^2 \frac{r''_{\theta\theta}\left(X_i, \overset{\circ}{\theta}\right)}{2\left(r'_{\theta}\left(X_i, \overset{\circ}{\theta}\right)\right)^3}$.

Подробнее остановимся на результатах пятого эксперимента, модель которого приведена на рис. 1. На рис. 2 изображена поверхность $\theta(X, Y)$ неявно заданная уравнением $Y = r(X, \theta)$. Как видно из рис. 2, в точках с координатами (0,3285, 0,8149), (0,7142, 1,0465) уравнение $\theta = \theta(X, Y)$ в области предполагаемых значений $\theta = (0,1, 2)$ решений не имеет. Поэтому $N = n - 2 = 13$. Значение в ячейке $N' = N - 4 = 9$, так как четыре точки с координатами (0,2, 1.2646), (0,4571, 1,2908), (0,5857, 1,0466), (0,8428, 1,6071) были удалены по признаку невыполнения условий **теоремы 1**, а именно: W_i для них были равны 32,9426; 1,1597; 1,4169; 0,0891, что значительно превышает не только $\varepsilon_i^2 \approx \sigma^2 = 0,0408$, но и $\varepsilon_i \approx \sigma = 0,202$.

Из результатов оценивания, приведенных в таблице, видно, что при увеличении $N' \leq N$ прослеживается неустойчивая сходимости к истинному значению θ состоятельной оценки $\tilde{\theta}$. При этом она является более точной, чем смещенная оценка $\overset{\circ}{\theta}$.

Выводы. 1. Предложенная оценка регрессионного коэффициента однофакторной однопараметрической регрессии, ввиду доказанной ее состоятельности, может служить не только начальным приближением существующих итерационных методов оценивания, но и быть при некоторых условиях самостоятельной оценкой.

2. В процессе доказательства состоятельности оценки получены ее основные характеристики, а именно: математическое ожидание и дисперсия, с помощью которых возможно построение интервальных оценок различной степени надежности.

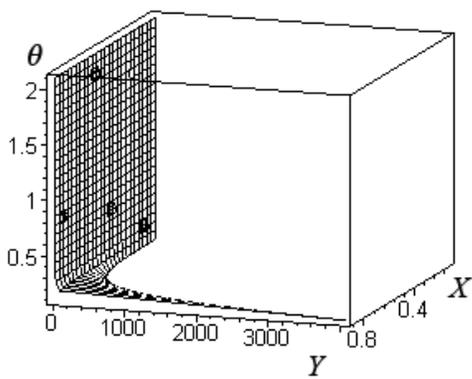


Рис. 2. Поверхность $\theta = \theta(X, Y)$, заданная уравнением $Y = r(X, \theta)$

3. Найденное значение дисперсии зависит не только от погрешностей экспериментальных данных, но и от частной производной регрессии. Это делает возможным осуществлять математически обоснованный поиск оптимальной регрессии в рамках не только одного, но и нескольких регрессионных классов.

4. Приведенные результаты получены без принятия допущений о каком-либо законе распределения погрешностей, поэтому предложенная оценка может быть применима для достаточно широкого класса задач обработки данных реальных экспериментов.

СПИСОК ЛИТЕРАТУРЫ

1. Айвазян С. А., Енюков И. С., Мешалкин Л. Д. Прикладная статистика: исследование зависимостей. – М.: Финансы и статистика, 1985. – 487 с.
2. Бард Й. Нелинейное оценивание параметров. – М.: Статистика, 1979. – 349 с.
3. М а т е м а т и ч е с к а я статистика: Учеб. для вузов / В.Б. Горяинов, И.В. Павлов, Г.М. Цветкова и др. Под ред. В.С. Зарубина, А.П. Крищенко. – М.: Изд-во МГТУ им. Н.Э. Баумана, 2001. – 424 с. (Сер. Математика в техническом университете. Вып. XVII).
4. К у д р я в ц е в Л. Д. Курс математического анализа: Учеб. для студентов ун-тов и вузов. В 3 т. Т. 1. – М.: Высш. школа., 1988. – 712 с.
5. В е н т ц е л ь Е. С. Теория вероятностей. – М.: Гос. изд-во физ.-мат. лит., 1962. – 564 с.
6. Г р е ш и л о в А. А. Математические методы построения прогнозов. – М.: Радио и связь, 1997. – 112 с.

Статья поступила в редакцию 16.01.2006

Ирина Васильевна Гетманская родилась в 1956 г., окончила в 1978 г. Казахский государственный университет им. С.М. Кирова. Старший преподаватель кафедры “Высшая математика” Калужского филиала МГТУ им. Н.Э. Баумана. Автор 6 научных работ в области математической физики и прикладной статистики.

I.V. Getmanskaya (b. 1956) graduated from the Kazakh State University n.a.S.M.Kirov in 1978. Senior teacher of “Higher Mathematics” department of the Bauman Moscow State Technical University. Author of 6 publications in the field of mathematical physics and applied statistics.

