

МНК-ОЦЕНКИ ПАРАМЕТРОВ ЛИНЕЙНОЙ ПАРНОЙ КОНФЛЮЭНТНОЙ МОДЕЛИ

Введение. В различных областях естествознания часто возникает задача исследования априори постулируемых функциональных зависимостей между количественными неслучайными переменными, значения которых определяются из эксперимента. Об экспериментальных ошибках при измерении всех переменных предполагается лишь, что они случайны, некоррелированы и распределены по некоторому вероятностному закону с нулевым математическим ожиданием и конечной постоянной дисперсией. В работе рассматривается часто используемый в такой ситуации метод наименьших квадратов (МНК) для нахождения оценок параметров в случае двухпараметрической линейной функциональной зависимости между изучаемыми переменными с ошибками в обеих переменных (простая линейная конфлюэнтная модель).

Постановка задачи. Предположим, что изучаются n объектов с характерными для каждого из них двумя признаками X, Y , (например, рост и вес для каждого из n человек или доход и расход для каждой из n семей). Обозначим через X_i, Y_i гипотетические неизвестные значения соответствующих признаков X, Y для каждого i -ого объекта, $i = \overline{1, n}$.

Пусть в отношении X и Y сделано предположение, что они связаны линейной зависимостью

$$Y = aX + b, \quad (1)$$

где коэффициенты a, b неизвестны и подлежат определению. Тогда соотношению (1) соответствует связь между значениями признаков для каждого из n объектов, имеющая вид

$$Y_i = aX_i + b, \quad i = \overline{1, n}. \quad (2)$$

Будем считать также, что для каждого из n изучаемых объектов независимо от других объектов одновременно измеряются два интересующих нас признака X, Y . Результаты измерений представляют собой набор $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, где x_i, y_i – значения признаков X, Y , соответственно, зафиксированных при измерении у i -ого объекта.

Однако при проведении измерений каждого из признаков в общем случае могут допускаться неточности. Тогда справедливы равенства:

$$x_i = X_i + e_{xi}, \quad y_i = Y_i + e_{yi}, \quad i = \overline{1, n}, \quad (3)$$

где e_{xi}, e_{yi} - ошибки измерений, рассматриваемые как значения некоррелированных случайных величин (с.в.) e_x, e_y с нулевыми математическими ожиданиями и конечными дисперсиями:

$$Ee_x = 0, Ee_y = 0, De_x = \sigma_x^2, De_y = \sigma_y^2, \sigma_x^2 + \sigma_y^2 \neq 0, \text{cov}(e_x, e_y) = 0.$$

Из соотношений (2), (3) получаем следующую «структурную» [1, Гл.29] связь между измерениями

$$y_i = ax_i + b + (e_{yi} - ae_{xi}), \quad i = \overline{1, n}. \quad (4)$$

При этом, измерения $(x_i, y_i), i = \overline{1, n}$, можно рассматривать как результат наблюдений с.в. (x, y) . В силу (3) имеем

$$Ex = X, Ey = Y, \quad (5)$$

а в силу (1) справедливо равенство

$$Ey = aEx + b. \quad (6)$$

Возникает следующая задача: в условиях пассивного эксперимента по результатам измерений $(x_i, y_i), i = \overline{1, n}$, построить статистические точечные оценки неизвестных признаков (X, Y) и коэффициентов a, b (2), для которых достигается минимальное значение суммы $\left\{ \sum_{i=1}^n \frac{e_{xi}^2}{\sigma_x^2} + \sum_{i=1}^n \frac{e_{yi}^2}{\sigma_y^2} \right\}$. Выбор последнего функционала связан с тем, что его минимизация равносильна максимизации функции правдоподобия в предположении двумерного нормального закона распределения с.в. (e_x, e_y) [1, Гл.29].

Ниже будет получен явный вид МНК-оценок для простой линейной конъюнктной модели, показана их состоятельность и эквивалентность

оценкам, получаемым по методу максимального правдоподобия при предположении нормальности ошибок [1, Гл.29].

Оценка параметров. В силу условия (2) справедливо равенство

$$\min_{\substack{a,b,X_i,Y_i,i=\overline{1,n} \\ Y_i=aX_i+b}} \left\{ \sum_{i=1}^n \frac{e_{x_i}^2}{\sigma_x^2} + \sum_{i=1}^n \frac{e_{y_i}^2}{\sigma_y^2} \right\} = \min_{a,b,X_i,i=\overline{1,n}} \left\{ \sum_{i=1}^n \frac{(x_i - X_i)^2}{\sigma_x^2} + \sum_{i=1}^n \frac{(y_i - aX_i - b)^2}{\sigma_y^2} \right\}. \quad (7)$$

Обозначим

$$S = S(a, b, X_1, X_2, \dots, X_n) = \sum_{i=1}^n \frac{(x_i - X_i)^2}{\sigma_x^2} + \sum_{i=1}^n \frac{(y_i - aX_i - b)^2}{\sigma_y^2}, \quad (8)$$

тогда задача минимизации (7) сводится к выбору в качестве оценок для a , b и $X_i, i = \overline{1, n}$ таких функций \hat{a} , \hat{b} и $\hat{X}_i, i = \overline{1, n}$ от измерений $x_i, y_i, i = \overline{1, n}$, для которых выполнено соотношение

$$S(\hat{a}, \hat{b}, \hat{X}_1, \hat{X}_2, \dots, \hat{X}_n) \leq S(a, b, X_1, X_2, \dots, X_n)$$

при любых допустимых значениях a , b и $X_i, i = \overline{1, n}$.

Для решения задачи сначала найдем множество всех возможных оценок a , b и $X_i, i = \overline{1, n}$, для которых функция $S(a, b, X_1, X_2, \dots, X_n)$ может иметь экстремальное значение, а потом из этого множества выберем те оценки, при которых S принимает наименьшее значение.

Утверждение 1. Решения \hat{a} , \hat{b} и $\hat{X}_i, i = \overline{1, n}$ задачи (7) удовлетворяют равенствам

$$\begin{cases} \hat{X}_i = \frac{\sigma_y^2 x_i + \sigma_x^2 \hat{a} (y_i - \hat{b})}{\sigma_y^2 + \sigma_x^2 \hat{a}^2}, i = \overline{1, n}. \\ \hat{a}^2 \sigma_x^2 s_{xy} - \hat{a} (\sigma_x^2 s_y^2 - \sigma_y^2 s_x^2) - \sigma_y^2 s_{xy} = 0, \\ \hat{b} = \bar{y} - \hat{a} \bar{x}. \end{cases} \quad (9)$$

где $\bar{x}, \bar{y}, s_x^2, s_y^2, s_{xy}$ - выборочные средние, дисперсии и ковариация измерений $x_i, y_i, i = \overline{1, n}$, определяемые как

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i, s_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2, s_y^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2,$$

$$s_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}).$$

Доказательство. Дифференцируемая функция S (8) может иметь экстремумы только при тех значениях a , b и $X_i, i = \overline{1, n}$, при которых частные производные $S(a, b, X_1, X_2, \dots, X_n)$ по всем $n+2$ неизвестным равны нулю. Для их определения имеем следующие системы алгебраических уравнений

$$\begin{cases} \frac{\partial S}{\partial X_i} = 0, \\ i = \overline{1, n}, \\ \frac{\partial S}{\partial a} = 0, \\ \frac{\partial S}{\partial b} = 0, \end{cases} \Leftrightarrow \begin{cases} X_i (\sigma_y^2 + \sigma_x^2 a^2) = \sigma_y^2 x_i + \sigma_x^2 (a y_i - ab), \\ i = \overline{1, n}, \\ \sum_{i=1}^n (y_i - a X_i - b) X_i = 0, \\ \sum_{i=1}^n (y_i - a X_i - b) = 0. \end{cases}$$

Выразим из первого уравнения X_i и подставим в оставшиеся два. После преобразований получим систему уравнений

$$\begin{cases} X_i = \frac{\sigma_y^2 x_i + \sigma_x^2 (a y_i - ab)}{\sigma_y^2 + \sigma_x^2 a^2}, \\ i = \overline{1, n}, \\ a^2 \sigma_x^2 s_{xy} - a (\sigma_x^2 s_y^2 - \sigma_y^2 s_x^2) - \sigma_y^2 s_{xy} = 0, \\ b = \bar{y} - a \bar{x}. \end{cases}$$

Что и требовалось доказать.

Таким образом, задача (7) может быть переформулирована в следующем виде: найти, среди всего множества оценок \hat{a} , \hat{b} и $\hat{X}_i, i = \overline{1, n}$, удовлетворяющих (9), те, которые дают наименьшее значение $S(a, b, X_1, X_2, \dots, X_n)$. Более того, после подстановки найденных

$\hat{X}_i = \frac{\sigma_y^2 x_i + \sigma_x^2 (\hat{a} y_i - \hat{a} \hat{b})}{\sigma_y^2 + \sigma_x^2 \hat{a}^2}$ в $S(a, b, X_1, X_2, \dots, X_n)$ и проведения простейших

преобразований вместо (7) приходим к следующей задаче:

найти, среди всего множества оценок \hat{a} , \hat{b} , удовлетворяющих

$$\begin{cases} \hat{a}^2 \sigma_x^2 s_{xy} - \hat{a} (\sigma_x^2 s_y^2 - \sigma_y^2 s_x^2) - \sigma_y^2 s_{xy} = 0, \\ \hat{b} = \bar{y} - \hat{a} \bar{x}, \end{cases} \quad (9')$$

те, которые дают наименьшее значение функционала

$$\tilde{S}(\hat{a}, \hat{b}) = \frac{\sum_{i=1}^n (y_i - \hat{a}x_i - \hat{b})^2}{\sigma_y^2 + \hat{a}^2 \sigma_x^2}. \quad (8')$$

Нетрудно показать, что система (9') является системой нормальных уравнений задачи

$$\min_{a,b} \tilde{S}(a,b) = \min_{a,b} \frac{\sum_{i=1}^n (y_i - ax_i - b)^2}{\sigma_y^2 + a^2 \sigma_x^2} \quad (7')$$

минимизации функционала $\tilde{S}(a,b)$ по множеству всех возможных значений параметров a , b .

Действительно, система частных производных задачи (7') имеет вид

$$\begin{cases} \frac{\partial \tilde{S}}{\partial a} = \frac{-2}{(\sigma_y^2 + a^2 \sigma_x^2)} \left\{ \sum_{i=1}^n x_i (y_i - ax_i - b) + \frac{a \sigma_x^2 \sum_{i=1}^n (y_i - ax_i - b)^2}{(\sigma_y^2 + a^2 \sigma_x^2)} \right\} = 0, \\ \frac{\partial \tilde{S}}{\partial b} = \frac{-2 \sum_{i=1}^n (y_i - ax_i - b)}{(\sigma_y^2 + a^2 \sigma_x^2)} = 0. \end{cases} \quad (9'')$$

После несложных, но трудоемких преобразований получаем (9').

Следствие 1. При $\sigma_x^2 \neq \sigma_y^2$, $\sigma_x^2 \neq 0$, $\sigma_y^2 \neq 0$, $s_{xy} \neq 0$ наименьшее значение функции \tilde{S} (8'), а следовательно S (8), задачи (7') ((7)), достигается на оценках, определяемых равенствами

$$\left\{ \begin{aligned} \hat{a} &= \frac{(\sigma_x^2 s_y^2 - \sigma_y^2 s_x^2) + \sqrt{(\sigma_x^2 s_y^2 - \sigma_y^2 s_x^2)^2 + 4\sigma_x^2 \sigma_y^2 s_{xy}^2}}{2\sigma_x^2 s_{xy}} = \\ &= \frac{\left(\frac{s_y}{s_x} - \frac{\sigma_y^2 s_x}{\sigma_x^2 s_y}\right) + \sqrt{\left(\frac{s_y}{s_x} - \frac{\sigma_y^2 s_x}{\sigma_x^2 s_y}\right)^2 + 4\frac{\sigma_y^2}{\sigma_x^2} r_{xy}^2}}{2r_{xy}}, \\ \hat{b} &= \bar{y} - \hat{a}\bar{x}, \end{aligned} \right. \quad (10)$$

где $r_{xy} = \frac{s_{xy}}{s_x s_y}$ - выборочный коэффициент корреляции.

Доказательство. Введем обозначения для корней квадратного трехчлена в (9) ((9'))

$$\hat{a}^{(+)} = \frac{(\sigma_x^2 s_y^2 - \sigma_y^2 s_x^2) + \sqrt{(\sigma_x^2 s_y^2 - \sigma_y^2 s_x^2)^2 + 4\sigma_x^2 \sigma_y^2 s_{xy}^2}}{2\sigma_x^2 s_{xy}},$$

$$\hat{a}^{(-)} = \frac{(\sigma_x^2 s_y^2 - \sigma_y^2 s_x^2) - \sqrt{(\sigma_x^2 s_y^2 - \sigma_y^2 s_x^2)^2 + 4\sigma_x^2 \sigma_y^2 s_{xy}^2}}{2\sigma_x^2 s_{xy}}.$$

Подставив $\hat{b} = \bar{y} - \hat{a}\bar{x}$ в $\hat{S}(\hat{a}, \hat{b})$ получим

$$\hat{S}(\hat{a}) = \frac{\sum_{i=1}^n ((y_i - \bar{y}) - \hat{a}(x_i - \bar{x}))^2}{\sigma_y^2 + \hat{a}^2 \sigma_x^2} = \frac{ns_y^2 + \hat{a}^2 ns_x^2 - 2\hat{a}ns_{xy}}{\sigma_y^2 + \hat{a}^2 \sigma_x^2}.$$

Сравним $\hat{S}(\hat{a}^{(+)})$ с $\hat{S}(\hat{a}^{(-)})$. Рассмотрим 2 случая. Пусть $s_{xy} > 0$. Тогда очевидно $\hat{a}^{(+)} > 0, \hat{a}^{(-)} < 0$ и поэтому $\hat{S}(\hat{a}^{(+)}) < \hat{S}(\hat{a}^{(-)})$. При $s_{xy} < 0$ имеем $\hat{a}^{(+)} < 0, \hat{a}^{(-)} > 0$, а значит, как и в предыдущем случае, $\hat{S}(\hat{a}^{(+)}) < \hat{S}(\hat{a}^{(-)})$.

Для полноты доказательства покажем теперь, что значения (10)

являются точкой экстремума функции $\tilde{S}(a, b) = \frac{\sum_{i=1}^n (y_i - ax_i - b)^2}{\sigma_y^2 + a^2 \sigma_x^2}$, в которой она принимает наименьшее значение. Для этого, вводя обозначения

$$A = \frac{\partial^2 \tilde{S}}{\partial a^2} \Big|_{\substack{a=\hat{a}, \\ b=\hat{b}}}, B = \frac{\partial^2 \tilde{S}}{\partial a \partial b} \Big|_{\substack{a=\hat{a}, \\ b=\hat{b}}}, C = \frac{\partial^2 \tilde{S}}{\partial b^2} \Big|_{\substack{a=\hat{a}, \\ b=\hat{b}}},$$

достаточно показать справедливость неравенств

$$A > 0, C > 0, D = B^2 - AC < 0.$$

Выражения для производных суть

$$\frac{\partial^2 \tilde{S}}{\partial a^2} = \frac{2}{(\sigma_y^2 + a^2 \sigma_x^2)} \left\{ \sum_{i=1}^n x_i^2 + \frac{4a\sigma_x^2 \sum_{i=1}^n x_i (y_i - ax_i - b)}{(\sigma_y^2 + a^2 \sigma_x^2)} - \frac{\sigma_x^2 (\sigma_y^2 - 3a^2 \sigma_x^2) \sum_{i=1}^n (y_i - ax_i - b)^2}{(\sigma_y^2 + a^2 \sigma_x^2)^2} \right\},$$

$$\frac{\partial^2 \tilde{S}}{\partial a \partial b} = \frac{2}{(\sigma_y^2 + a^2 \sigma_x^2)} \left\{ \sum_{i=1}^n x_i + \frac{2a\sigma_x^2 \sum_{i=1}^n (y_i - ax_i - b)}{(\sigma_y^2 + a^2 \sigma_x^2)} \right\},$$

$$\frac{\partial^2 \tilde{S}}{\partial b^2} = \frac{2n}{(\sigma_y^2 + a^2 \sigma_x^2)}.$$

Тогда, учитывая, что \hat{a} , \hat{b} (10) удовлетворяют (9''), получим

$$A = \frac{2n}{(\sigma_y^2 + \hat{a}^2 \sigma_x^2)} \left\{ \frac{2\sigma_x^2 s_{xy}^2}{(\sigma_x^2 s_y^2 - \sigma_y^2 s_x^2) + \sqrt{(\sigma_x^2 s_y^2 - \sigma_y^2 s_x^2)^2 + 4\sigma_x^2 \sigma_y^2 s_{xy}^2}} + (\bar{x})^2 \right\} > 0,$$

$$B = \frac{2n\bar{x}}{(\sigma_y^2 + \hat{a}^2 \sigma_x^2)},$$

$$C = \frac{2n}{(\sigma_y^2 + \hat{a}^2 \sigma_x^2)} > 0,$$

$$D = \frac{-8n^2 \sigma_x^2 s_{xy}^2}{(\sigma_y^2 + \hat{a}^2 \sigma_x^2)^2 \left\{ (\sigma_x^2 s_y^2 - \sigma_y^2 s_x^2) + \sqrt{(\sigma_x^2 s_y^2 - \sigma_y^2 s_x^2)^2 + 4\sigma_x^2 \sigma_y^2 s_{xy}^2} \right\}} < 0.$$

Что и требовалось доказать. Отметим, что вид оценок (10) эквивалентен оценкам, получаемым по методу максимального правдоподобия при предположении нормальности ошибок [1, Гл.29].

Следствие 2. В случае $\sigma_x^2 = \sigma_y^2$ задача (7) ((7')) имеет единственное решение, удовлетворяющее уравнениям

$$\begin{aligned} \hat{b}_{opt} &= \bar{y} - \hat{a}_{opt} \bar{x}, \\ \hat{a}_{opt}^2 s_{xy} - \hat{a}_{opt} (s_y^2 - s_x^2) - s_{xy} &= 0, \end{aligned}$$

и определяется равенствами

$$\begin{aligned} \hat{b}_{opt} &= \bar{y} - \hat{a}_{opt} \bar{x}, \\ \hat{a}_{opt} &= \frac{(s_y^2 - s_x^2) + \sqrt{(s_y^2 - s_x^2)^2 + 4s_{xy}^2}}{2s_{xy}} = \\ &= \frac{\left(\frac{s_y}{s_x} - \frac{s_x}{s_y} \right) + \sqrt{\left(\frac{s_y}{s_x} - \frac{s_x}{s_y} \right)^2 + 4r_{xy}^2}}{2r_{xy}}. \end{aligned} \quad (11)$$

Величины \hat{a}_{opt} и \hat{b}_{opt} являются оценками коэффициентов a , b «ортогональной» регрессии [2].

Следствие 3. В случае $\sigma_x^2 = 0$ решением задачи (7) ((7')) являются оценки «прямой» регрессии

$$\hat{b}_{np} = \bar{y} - \hat{a}_{np} \bar{x},$$

$$\hat{a}_{np} = \frac{s_{xy}}{s_x^2} = r_{xy} \frac{s_y}{s_x}. \quad (12)$$

Следствие 4. В случае $\sigma_y^2 = 0$ решением задачи (7) ((7')) являются оценки «обратной» регрессии

$$\hat{b}_{обр} = \bar{y} - \hat{a}_{обр} \bar{x},$$

$$\hat{a}_{обр} = \frac{s_y^2}{s_{xy}} = \frac{1}{r_{xy}} \frac{s_y}{s_x}. \quad (13)$$

Следствие 5. Оценки \hat{a} и \hat{b} (10) могут быть выражены через оценки «прямой» и «обратной» регрессий в виде

$$\hat{b} = \bar{y} - \hat{a} \bar{x},$$

$$\hat{a} = \frac{\left(s_y^2 - \frac{\sigma_y^2}{\sigma_x^2} s_x^2 \right) + \sqrt{\left(s_y^2 - \frac{\sigma_y^2}{\sigma_x^2} s_x^2 \right)^2 + 4 \frac{\sigma_y^2}{\sigma_x^2} s_{xy}^2}}{2s_{xy}} = \quad (14)$$

$$= \frac{\left(\hat{a}_{обр} - \frac{\sigma_y^2}{\sigma_x^2 \hat{a}_{np}} \right) + \sqrt{\left(\hat{a}_{обр} - \frac{\sigma_y^2}{\sigma_x^2 \hat{a}_{np}} \right)^2 + 4 \frac{\sigma_y^2}{\sigma_x^2}}}{2}.$$

Следствие 6. Если в качестве приближений дисперсий ошибок σ_x^2, σ_y^2 взять их выборочные аналоги s_x^2, s_y^2 , то оценки \hat{a} и \hat{b} (10) являются оценками «диагональной» регрессии [2]

$$\hat{b}_{диаг} = \bar{y} - \hat{a}_{диаг} \bar{x},$$

$$\hat{a}_{диаг} = \sqrt{\frac{s_y^2}{s_x^2}} = \sqrt{\hat{a}_{np} \hat{a}_{обр}}. \quad (15)$$

Таким образом, доказанное утверждение обобщает ранее известные оценки параметров линейной парной регрессии и дает новый явный вид МНК-оценок \hat{a} и \hat{b} (10) в случае наличия ошибок в обеих переменных X, Y .

Свойства оценок параметров. Остановимся лишь на некоторых свойствах оценок (10).

Утверждение 2. Справедливы соотношения

$$|\hat{a}_{np}| \leq |\hat{a}| \leq |\hat{a}_{обр}|,$$

причем равенства имеют место лишь в случае точной линейной зависимости переменных X, Y .

Доказательство аналогично доказательству соотношения $|\hat{a}_{np}| \leq |\hat{a}_{опт}| \leq |\hat{a}_{обр}|$ в [2].

Докажем теперь, аналогично тому как это сделано в [1, Гл.29], что оценки (10), как и (11) [2], являются состоятельными, в отличие от оценок (12), (13) в случае ошибок в обеих переменных X, Y ($\sigma_x^2 \neq 0, \sigma_y^2 \neq 0$).

Утверждение 3. При условии $\text{cov}(e_x, e_y) = 0$ и

$$S_X^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \xrightarrow{n \rightarrow \infty} \Sigma_X^2 \quad \text{оценки } \hat{a} \quad \text{и} \quad \hat{b} \quad (10) \quad \text{являются}$$

состоятельными.

Доказательство. Используя (4') и свойства дисперсий и ковариаций, а также учитывая, что X, Y не являются случайными величинами, имеем следующие разложения

$$s_x^2 = s_{X+e_x}^2 = S_X^2 + s_{e_x}^2 + 2s_{X, e_x},$$

$$s_y^2 = s_{aX+b+e_y}^2 = a^2 S_X^2 + s_{e_y}^2 + 2as_{X, e_y},$$

$$s_{xy} = s_{(X+e_x)(aX+b+e_y)} = aS_X^2 + s_{e_x, e_y} + s_{X, e_y} + as_{X, e_x}.$$

Воспользуемся тем, что выборочные дисперсии и ковариация с ростом объема выборки сходятся по вероятности к своим математическим ожиданиям, то есть стремятся к своим теоретическим аналогам, если последние существуют [3, Гл.10]. В нашем случае, поскольку $s_{X, e_x} \rightarrow 0, s_{X, e_y} \rightarrow 0, n \rightarrow \infty$, имеем:

$$s_x^2 \rightarrow \Sigma_X^2 + \sigma_x^2, s_y^2 \rightarrow a^2 \Sigma_X^2 + \sigma_y^2, s_{xy} \rightarrow a \Sigma_X^2, n \rightarrow \infty.$$

Теперь, подставляя в (10), имеем

$$\hat{a} \rightarrow \frac{\Sigma_X^2 (a^2 \sigma_x^2 - \sigma_y^2) + \sqrt{\Sigma_X^4 (a^2 \sigma_x^2 - \sigma_y^2)^2 + 4a^2 \sigma_x^2 \sigma_y^2 \Sigma_X^4}}{2a^2 \sigma_x^2 \Sigma_X^2} =$$

$$= \frac{(a^2 \sigma_x^2 - \sigma_y^2) + (a^2 \sigma_x^2 + \sigma_y^2)}{2a^2 \sigma_x^2} = a, n \rightarrow \infty.$$

Отсюда следует, что оценка \hat{a} состоятельная.

Докажем теперь состоятельность оценки \hat{b} . Воспользуемся тем фактом, что $\bar{y} \rightarrow Y = aX + b$, $\bar{x} \rightarrow X$, $n \rightarrow \infty$ [3, Гл.10]. Тогда $\hat{b} = \bar{y} - \hat{a}\bar{x} \rightarrow aX + b - aX = b$, $n \rightarrow \infty$. Оценка \hat{b} (10) - состоятельная.

Что и требовалось доказать.

Численные иллюстрации

Пример 1 (дисперсии ошибок заданы). Пусть истинные значения параметров модели (1) $a = 0.5$, $b = 4$ и $\{X_i, i = \overline{1,6}\} = \{6,7,8,9,10,11\}$. С помощью датчика случайных чисел нормального закона распределения получим выборку n значений ошибок $e_{x_i} \sim N(0,1)$ и n значений $e_{y_i} \sim N(0,3)$, а затем рассчитаем x_i, y_i по правилу (2), (3). Предполагая известными дисперсии ошибок $\sigma_x^2 = 1, \sigma_y^2 = 3$, для параметров a, b конфлюэнтной модели (4) по выборке $(x_i, y_i), i = \overline{1,n}$, были вычислены оценки (10), (12), (13), (15). При различных n полученные результаты приведены в табл.1.

Табл.1. Оценки параметров линейной модели при нормально распределенных ошибках.

Метод оценивания	n=300		n=30000	
	a=0.5	b=4	a=0.5	b=4
(10)	0.4985	4.0124	0.5	4.004
(12)	0.3694	5.1099	0.372	5.085
(13)	2.6019	-13.866	2.552	-13.439
(15)	0.9804	-0.0835	0.975	-0.035

Восстановленные по оценкам графики линейных приближений в случае $n = 300$ изображены на рис.1. «Крестиками» отмечены точки с истинными значениями $\{(X_i, Y_i), i = \overline{1,6}\}$, а «ромбиками» моделируемые точки с ошибками в обеих координатах $\{(x_i, y_i), i = \overline{1,6}\}$.

Как и ожидалось, наилучшее приближение к истинным значениям $a = 0.5, b = 4$ дают оценки (10).

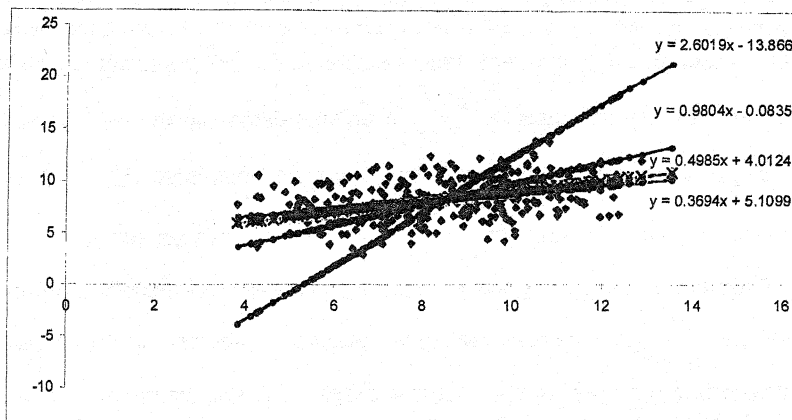


Рис.1

Пример 2 (дисперсии ошибок заданы). Аналогично первому примеру с помощью датчика случайных чисел получим выборку n значений равномерно распределенных ошибок e_{xi} , с $Ee_{xi} = 0, De_{xi} = 1$, и n значений ошибок e_{yi} , с $Ee_{yi} = 0, De_{yi} = 3$, а затем рассчитаем x_i, y_i по правилу (2), (3). Предполагая известными дисперсии ошибок $\sigma_x^2 = 1, \sigma_y^2 = 3$, для параметров a, b конъюэнтной модели (4) по выборке $(x_i, y_i), i = \overline{1, n}$, были вычислены оценки (10), (12), (13), (15). При различных n результаты расчетов приведены в табл.2.

Табл. 2. Оценки параметров линейной модели при равномерно распределенных ошибках.

Метод оценивания	n=300		n=30000	
	a=0.5	b=4	a=0.5	b=4
(10)	0.502	3.980	0.498	4.014
(12)	0.373	5.083	0.371	5.093
(13)	2.582	-13.695	2.557	-13.487
(15)	0.981	-0.087	0.975	-0.034

Как и в предыдущем примере, наилучшее приближение к истинным значениям $a = 0.5$, $b = 4$ дают оценки (10).

Пример 3 (дисперсии ошибок неизвестны, но известно их отношение). На практике редко встречается ситуация когда дисперсии ошибок известны. Однако для применения \hat{a} , \hat{b} (10) достаточно знать отношение дисперсий ошибок $\frac{\sigma_y^2}{\sigma_x^2}$ и использовать равенство (10) или (14). В рассмотренных выше примерах вместо значений $\sigma_x^2 = 1, \sigma_y^2 = 3$ достаточно задать их отношение $\frac{\sigma_y^2}{\sigma_x^2} = 3$ и результат будет тот же.

Примером, где дисперсии ошибок неизвестны, а известно лишь их отношение $\frac{\sigma_y^2}{\sigma_x^2}$, может служить задача оценивания линейной зависимости $Y = aX + b$ между ростом детей (Y) и средним ростом их родителей (X) [4]. При объеме выборки $n = 928$ данные группировались по 11 группам среднего роста у родителей и 14 группам роста у детей и приведены в табл.3. В верхней строке указан средний рост родителей в дюймах соответствующей группы, а по строкам – общее количество взрослых детей с указанным в левом столбце их росте в дюймах.

Табл.3. Данные о среднем росте родителей и их взрослых детей.

		Средний рост родителей (в дюймах)										
		63.5	64.5	65.5	66.5	67.5	68.5	69.5	70.5	71.5	72.5	73.5
Рост взрослых детей (в дюймах)	74.2							5	3	2	4	
	73.2						3	4	3	2	2	3
	72.2			1		4	4	11	4	9	7	1
	71.2			2		11	18	20	7	4	2	
	70.2			5	4	19	21	25	14	10	1	
	69.2	1	2	7	13	38	48	33	18	5	2	
	68.2	1		7	14	28	34	20	12	3	1	
	67.2	2	5	11	17	38	31	27	3	4		
	66.2	2	5	11	17	36	25	17	1	3		
	65.2	1	1	7	2	15	16	4	1	1		
	64.2	4	4	5	5	14	11	16				
	63.2	2	4	9	3	5	7	1	1			
	62.2		1		3	3						
	61.2	1	1	1			1		1			

Предполагая постоянной дисперсию роста взрослых людей и равной σ_y^2 , будем иметь дисперсию среднего роста родителей, как дисперсию половины суммы роста людей, равной $\frac{\sigma_y^2}{2} = \sigma_x^2$, а дисперсию роста детей равной σ_y^2 . Данному предположению $\frac{\sigma_y^2}{\sigma_x^2} = 2$ близко отношение выборочных дисперсий роста детей s_y^2 и среднего роста родителей s_x^2 равное $\frac{s_y^2}{s_x^2} = 1.97$. Для оценки параметров линейной модели применим формулы (12), (13), (15) и (10) с $\frac{\sigma_y^2}{\sigma_x^2} = 2$. Полученные результаты вычислений представлены в табл.4.

Табл.4. Оценки параметров линейной модели для данных в табл.3.

Метод оценивания	$n = 928$	
	(10)	1.3548
(12)	0.6455	24.004
(13)	3.0488	-140.151
(15)	1.4029	-27.726

На рис.2 в виде точек изображены исходные данные из табл.3 и их предполагаемые погрешности. По оценкам параметром, указанным в табл.4, построены прямые, соответствующие приближенным линейным моделям. Оценки (15) «диагональной» регрессии близки по значению с оценками (10) и поэтому соответствующий им график приближенной линейной модели не изображен на рис.2.

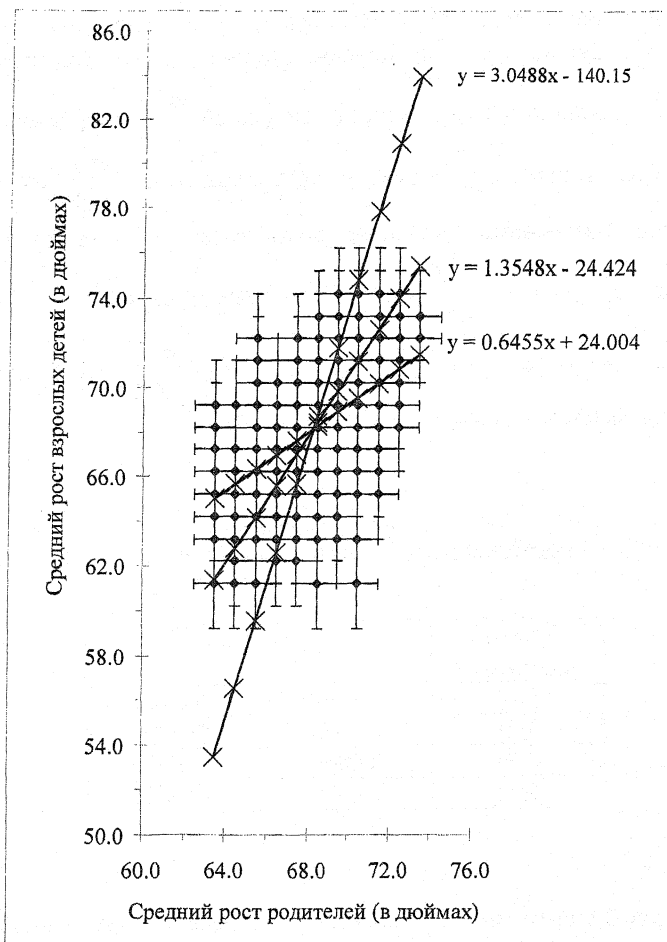


Рис.2.

Galton [4] по результатам своих исследований сделал вывод о «регрессивной» линейной зависимости роста детей от роста их родителей, что соответствует оценкам (12) «прямой» регрессионной модели. В то время как оценки (10) конфлюэнтной модели показывают «прогрессивную» линейную зависимость ростов.

Таким образом, полученные точечные МНК-оценки (10) для параметров простой линейной конфлюэнтной модели являются

наилучшими по сравнению с известными оценками «прямой», «обратной», «ортогональной», «диагональной» регрессии, а также эквивалентны оценкам, получаемым по методу максимального правдоподобия при предположении нормальности ошибок.

Литература

1. Кендалл М., Стьюарт А. Статистические выводы и связи. М.: Наука, 1973.
2. Маленко Э. Статистические методы эконометрии. М.: Статистика, 1975.
3. Кендалл М., Стьюарт А. Теория распределений. М.: Наука, 1966.
4. Galton F. Regression towards mediocrity in hereditary stature // J. of the Royal Anthropolog. Inst., 1885. V.15. P.228-262.