

**Метод наименьших расстояний для конглоэнтной модели с гомоскедастичными измерениями столбцов матрицы и правой части.**

Автором в работах [1,2,3] в самых общих предположениях проведено доказательство метода наименьших расстояний оценивания параметров конглоэнтных моделей пассивного эксперимента. В данной работе подробно исследуется, возможно, наиболее распространенный частный случай такой модели с различной по векторам измерений, но гомоскедастичной внутри каждого вектора, дисперсией. Результат иллюстрируется примером из эконометрии.

**1. Основная модель.**

Рассмотрим основное предположение [2,3] для наиболее типичного случая, когда все ошибки в столбцах конглоэнтной матрицы  $X$  и в векторе отклика  $y$  гомоскедастичны.

**Предположение 1.** Будем исходить из следующей линейной стохастической модели пассивного эксперимента:

$$\begin{cases} y = \Xi\beta + H\delta + e, & y \in R^n, \beta \in R^m, H \in R^{n \times k}, \delta \in R^k, Ee = 0, Eee^T = \sigma^2 I, \\ X = \Xi + C, & X \in R^{n \times m}, EC = 0, EeC^T = 0, Ec_i c_j^T = \delta_{ij} \mu_i^2 I, i, j = \overline{1, m}. \end{cases} \quad (1)$$

с матрицей  $[\Xi, H]$  линейного функционального уравнения

$$\vartheta = \Xi\beta + H\delta, \vartheta \in R^n, \Xi \in R^{n \times m}, H \in R^{n \times k}, \beta \in R^m, \delta \in R^k, \quad (2)$$

полного ранга ( $\delta_{ij}$  - символ Кронекера).

То есть в модели (1) пассивные наблюдения, которые будем считать нормально распределенными, имеют одинаковую дисперсию  $\sigma^2$  в  $y$  и одинаковые дисперсии  $\mu_i^2$  в каждом  $i$ -том столбце матрицы  $X$ . Тогда задача оценивания параметров будет иметь вид:

**Задача 1.** Зная одну реализацию  $\tilde{y}$  и  $\tilde{X}$  ( $\text{rank } \tilde{X} = m$ ):

$$\begin{cases} \tilde{y} = \Xi\beta + H\delta + \tilde{e}, \\ \tilde{X} = \Xi + \tilde{C}, \end{cases}$$

случайных величин  $u$  и  $X$ , точно теоретическую матрицу  $H$ , дисперсии  $\sigma^2, \mu_i^2, i = \overline{1, m}$ , оценить неизвестные элементы  $\vartheta, \Xi$  и неизвестные параметры  $\beta, \delta$  линейной стохастической модели (1) методом максимального правдоподобия.

Т.е., учитывая в данном случае тот факт, что ошибки измерений не зависят от параметров [3] получим, что параметры надо оценить таким образом, чтобы квадрат расстояния измерений  $\tilde{y}, \tilde{x}_i, i = \overline{1, m}$ , от искомым значений  $\vartheta, \xi_i, i = \overline{1, m}$ ,

$$S^2 = \sigma^{-2} (\tilde{y} - \vartheta)^T (\tilde{y} - \vartheta) + \sum_{i=1}^m \mu_i^{-2} (\tilde{x}_i - \xi_i)^T (\tilde{x}_i - \xi_i) \quad (3)$$

в линейных ограничениях  $\vartheta = \Xi\beta + H\delta$  (2) был бы минимален (иначе говоря, методом наименьших расстояний - м.н.р.).

**Теорема 1.** Оценки  $\hat{b}, \hat{d}$  параметров  $\beta, \delta$  линейной модели (1) из Задачи 1 доставляют минимум квадратичной форме

$$S^2 = \frac{\|\tilde{X}\beta + H\delta - \tilde{y}\|^2}{\sigma^2 + \sum_{i=1}^m \mu_i^2 \beta_i^2}, \quad (4)$$

или (как можно показать с помощью дифференцирования) вычисляются из следующей нормальной системы линейных алгебраических уравнений (с.л.а.у.):

$$\begin{cases} (\bar{X}^T \bar{X} - \text{diag}(\mu_1^2, \dots, \mu_m^2) \hat{S}^2) \hat{b} + \bar{X}^T \hat{N} \hat{d} = \bar{X}^T \bar{y}, \\ \hat{N}^T \bar{X} \hat{b} + \hat{N}^T \hat{N} \hat{d} = \hat{N}^T \bar{y}, \end{cases} \quad (5)$$

где  $\hat{S}^2$  - минимальное значение остаточной суммы квадратов  $S^2$  (4).

Теорема является частным случаем доказанного в [2,3] утверждения с более общим видом ковариационной матрицы.

## 2. Частный случай.

В качестве следствия рассмотрим построение оценки для линейной двупараметрической стохастической компонентной модели с различной для каждой переменной, но гомоскедастичной внутри измерений дисперсией:

$$\begin{cases} y = \beta_1 \xi_1 + \beta_2 \xi_2 + \delta + e, \quad y \in \mathbb{R}^n, \quad \beta \in \mathbb{R}^1, \quad \delta \in \mathbb{R}^1, \quad Ee = 0, \quad Eee^T = \sigma^2 I, \\ x_1 = \xi_1 + c_1, \quad x_1 \in \mathbb{R}^n, \quad Ec_1 = 0, \quad Eec_1^T = 0, \quad Ec_1 c_1^T = \mu_1^2 I, \\ x_2 = \xi_2 + c_2, \quad x_2 \in \mathbb{R}^n, \quad Ec_2 = 0, \quad Eec_2^T = 0, \quad Ec_2 c_2^T = \mu_2^2 I, \end{cases} \quad (6)$$

и одним свободным (регрессионным) параметром  $\delta$ .

Как следует из теоремы 1, оценки  $\hat{b}_1, \hat{b}_2, \hat{d}$  параметров  $\beta_1, \beta_2, \delta$  линейной модели (6) из Задачи 1 доставляют минимум квадратичной форме

$$S^2 = \frac{\|\beta_1 \bar{x}_1 + \beta_2 \bar{x}_2 + \delta - \bar{y}\|^2}{\sigma^2 + \mu_1^2 \beta_1^2 + \mu_2^2 \beta_2^2}, \quad (7)$$

или является решением следующей с.л.а.у.:

$$\begin{bmatrix} \sum_{i=1}^n \tilde{x}_{1i}^2 - \mu_1^2 \hat{S}^2 & \sum_{i=1}^n \tilde{x}_{1i} \tilde{x}_{2i} & \sum_{i=1}^n \tilde{x}_{1i} \\ \sum_{i=1}^n \tilde{x}_{1i} \tilde{x}_{2i} & \sum_{i=1}^n \tilde{x}_{2i}^2 - \mu_2^2 \hat{S}^2 & \sum_{i=1}^n \tilde{x}_{2i} \\ \sum_{i=1}^n \tilde{x}_{1i} & \sum_{i=1}^n \tilde{x}_{2i} & n \end{bmatrix} * \begin{bmatrix} b_1 \\ b_2 \\ d \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n \tilde{x}_{1i} \tilde{y}_i \\ \sum_{i=1}^n \tilde{x}_{2i} \tilde{y}_i \\ \sum_{i=1}^n \tilde{y}_i \end{bmatrix} \quad (8)$$

где  $\hat{S}^2$  - минимальное по параметрам  $\beta_1, \beta_2, \delta$  значение остаточной суммы квадратов  $S^2$  (7). Очевидно, что для решения с.л.а.у. (8) необходимо использовать итерационный процесс.

## 3. Пример построения оценки.

Рассмотрим пример из раздела регрессионного анализа книги [4], где цена поддержанного автомобиля Ниссан оценивается в функции года производства ("возраста")  $\xi_1$  и пробега  $\xi_2$ .

Рассмотрим сначала чистую регрессию в функции возраста. Для простой линейной модели имеем оценку метода наименьших квадратов [5] (м.н.к.) вида:

$$y = \delta_1 \eta_1 + \delta + e \cong -2,026 \eta_1 + 19,547, \quad y \in \mathbb{R}^{11}, \quad \hat{S}^2 = 1,58. \quad (0,28)$$

Далее, для вычисления м.н.р.-оценки, будем считать возраст автомобиля равномерно распределенным в течение года, откуда имеем [6] для его дисперсии величину

$$c_1 \in U_{-0,5,0,5}; \quad Ec_1 = 0; \quad \mu_1^2 = Dc_1 = Ec_1^2 - (Ec_1)^2 = 1/12.$$

и будем полагать, что равномерное распределение для 11 измерений достаточно хорошо аппроксимируется нормальным.

Ниссан	$\tilde{x}_1(\eta_1)$ : годы	$\tilde{x}_2(\eta_2)$ : мили/1000	Y: цены\$/1000
1	5	57	8,5
2	4	40	10,3
3	6	77	7
4	5	60	8,2
5	5	49	8,9
6	5	47	9,8
7	6	58	6,6
8	6	39	9,5
9	2	8	16,9
10	7	69	7
11	7	89	4,8

Таблица 1. Возраст автомобиля, его пробег и объявленная стоимость. В скобках приведены обозначения для регрессионного анализа.

Воспользуемся оценкой дисперсии стоимости из м.н.к. как начальным значением для итерационного процесса вычисления дисперсии цены  $y$ . Итерируя до ожидаемой величины степеней свободы  $\hat{S}^2=11-2=9$ , достаточно просто подобрать отношение дисперсий исходных данных. Откуда имеем м.н.р.-оценку

$$y = \beta_1 \xi_1 + \delta + e \cong -2,10 \hat{x}_1 + 20,0, \hat{S}^2 = 9,0; y \in \mathbb{R}^{11}, \sigma^2 = 1,2; \mu_1^2 = 1/12.$$

Расчет м.н.р. точнее оценивает стоимость двухлетнего автомобиля, чем то недостаточно точное описание (как это отмечено в[4]), что получено для этого автомобиля из м.н.к. В итоге м.н.р.-оценки цен полностью согласуются с исходными данными. Из нее следует, что цена нового автомобиля \$20.000 и срок его службы ~10 лет. На рис. 1 содержатся исходные данные и обе оценки.

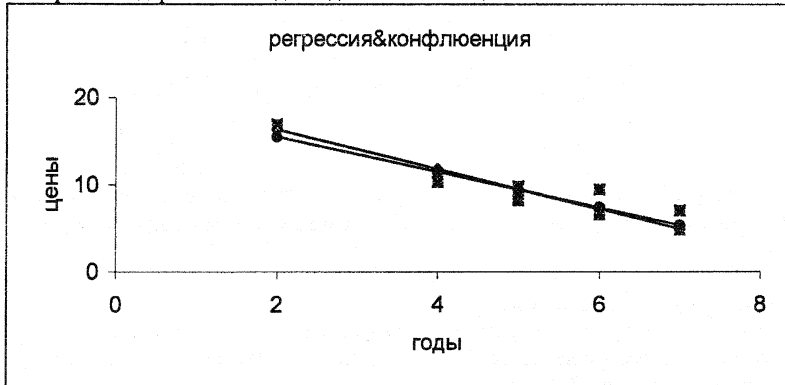


Рис. 1. Однопараметрическая модель. Квадратами обозначены исходные данные, точками - м.н.к.-оценка, ромбами - м.н.р.-оценка.

Рассмотрим теперь регрессию в функции пробега. Для простой линейной модели имеем м.н.к.-оценку

$$y = \delta_2 \eta_2 + \delta + e \cong -0,136 \eta_2 + 16,204, y \in \mathbb{R}^{11}, \hat{S}^2 = 1,2; \quad (9)$$

(0,016)

и м.н.р.-оценку (считаем ошибку пробега также равномерно распределенной)

$$y = \beta_2 \xi_2 + \delta + e \cong -0,136 \tilde{x}_2 + 16,205, \hat{S}^2 = 9,0; y \in \mathbb{R}^{11}, \sigma^2 = 1,2; \mu_2^2 = 1/12.$$

Здесь расчет м.н.р. практически совпадает с описанием, полученным из м.н.к.-расчетов, так как теоретическая дисперсия ошибки пробега очень мала, что не совсем согласуется с исходными данными.

**Замечание.** Действительно, показания спидометра автомобиля напрямую зависят от размера колес, а изменение радиуса колеса на один сантиметр вследствие износа протектора влечет за собой в среднем 4-5% ошибку в измерении пробега, откуда в среднем  $\mu_2^2 \approx 2,5$ . Поэтому просчитаем результат с такой дисперсией. Здесь также по величине степеней свободы достаточно просто рассчитывается отношение дисперсий исходных данных:

$$y = \beta_2 \xi_2 + \delta + e \cong -0,137 \bar{x}_2 + 16,24, \hat{S}^2 = 9,0; y \in \mathbf{R}^{11}, \sigma^2 = 1,1; \mu_2^2 = 2,5.$$

Такая смена дисперсии не сильно влияет на оценку, поскольку на самом деле дисперсия должна быть много больше для большого пробега и такой подход требует расчетов с гетероскедастичной дисперсией [2,3], выходящих за рамки данной работы, однако проведем их, имеем оценку

$$y = \beta_2 \xi_2 + \delta + e \cong -0,14 \bar{x}_2 + 16,4, \hat{S}^2 = 9; y \in \mathbf{R}^{11}, \sigma^2 = 1,06; \mu_2^2 = (0,05 \bar{y})^2; (10)$$

которая, хотя и лучше согласуется с исходными данными, показывает тем не менее, что только пробег, также как и только возраст, не дает полного описания стоимости автомобиля. На рис. 2 также содержатся исходные данные и оценки (9) и (10).

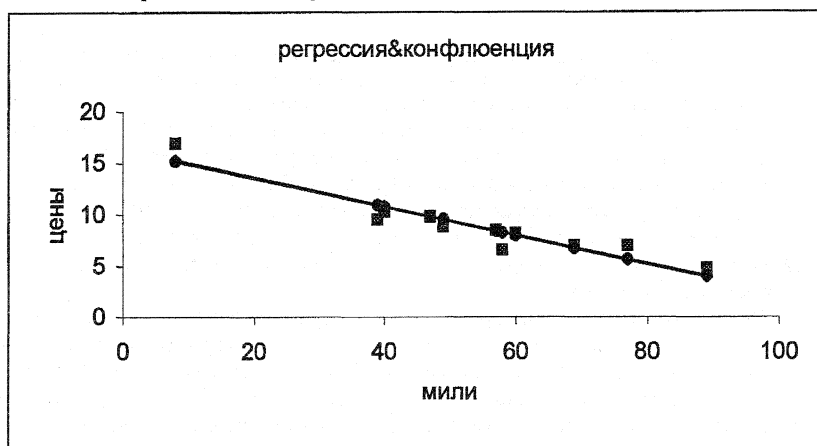


Рис. 2. Однопараметрическая модель. Квадраты - данные, ромбы - м.н.р.-оценка, точки - м.н.к.-оценка.

Рассмотрим, наконец, регрессию по обоим переменным. Для двухпараметрической линейной модели имеем м.н.к.-оценку

$$y = \delta_1 \eta_1 + \delta_2 \eta_2 + \delta + e \cong -0,95 \underset{(0,39)}{\eta_1} - 0,082 \underset{(0,026)}{\eta_2} + 18,3; y \in \mathbf{R}^{11}, \hat{S}^2 = 0,78;$$

и м.н.р.-оценку

$$y = \beta_1 \xi_1 + \beta_2 \xi_2 + \delta + e \cong -0,99 \bar{x}_1 - 0,088 \bar{x}_2 + 18,4; \hat{S}^2 = 8,0; y \in \mathbf{R}^{11}, \sigma^2 = 1; \mu_1^2 = 0,034; \mu_2^2 = 2,5$$

И м.н.к. и м.н.р.-оценки согласуются с исходными данными. На рис. 3 также содержатся данные и обе оценки и видно, что шесть срединных точек (наиболее правильно представляющих эксперимент) лучше приближаются м.н.р.

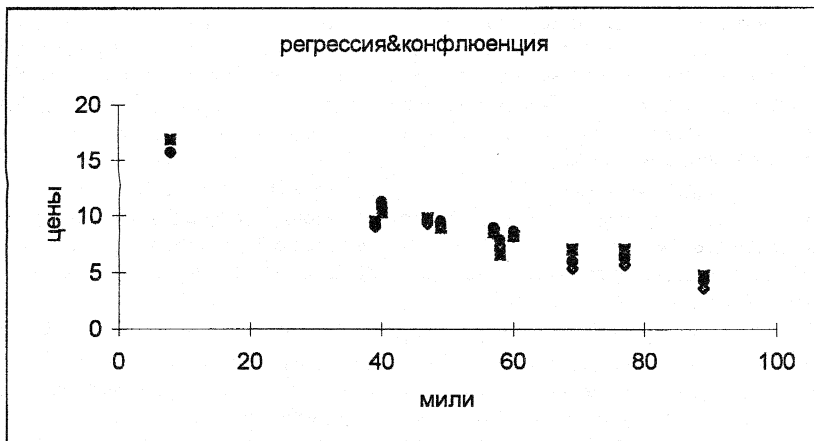


Рис. 3. Двупараметрическая линейная модель (на диаграмме в проекции плоскости регрессии на плоскость цены\*пробег). Квадраты - данные, ромбы - м.н.р.-оценка, точки - м.н.к.-оценка.

Поскольку возраст и километраж сильно коррелированы (в принципе они не могут эффективно работать вместе в этой модели), то ранг матрицы (8) чувствителен к заданию величины дисперсии и м.н.р.-оценка быстро теряет устойчивость, когда дисперсии независимых переменных выравниваются. По величине степеней свободы достаточно просто рассчитывается отношение дисперсий исходных данных.

Достаточно очевидно, что модель в самом деле нелинейна по этим параметрам (поскольку цена автомобиля не равна нулю через 9-10 лет) и требует более детального рассмотрения в этом плане.

В данной работе автором подробно изучен метод построения оценок параметров для стохастических линейных конфлюентно-регрессионных моделей с различными гомоскедастичными дисперсиями - класса моделей, с одной стороны, достаточно полно описывающих события, встречающиеся в созерцательных исследованиях, с другой стороны, оценки метода наименьших расстояний для которых достаточно легко рассчитать в обычном табличном процессоре типа Excel. На приведенных примерах видно, что использование различной дисперсии может значительно улучшить простейшее описание данных задачи и позволяет полнее провести апостериорное исследование исходных данных и их причинно-следственных связей.

### Литература.

1. Меченов А.С. О частично приближенных системах линейных алгебраических уравнений // Ж. вычисл. матем. и матем. физ., 1991, т. 31, №6, с. 790-799.
2. Меченов А.С. О подходе максимального правдоподобия к оценке параметров линейных функциональных соотношений. // Численные методы в математической физике. М.: Московский университет. 1996. С. 153-159.
3. Меченов А.С. О конфлюентном подходе в регрессионном анализе. // Методы математического моделирования. М.: Московский университет. 1998. С. 42-53.
4. Weiss N.A. Introductory statistics, New York. 1995.
5. Gauss K.F. Theoria motus corporum coelestium in sectionibus conicis solum ambientium, 1809, Hamburgi. (Русский перевод Н.Ф. Булаевского: Гаусс К.Ф. Избранные геодезические произведения Т.1-2. М.: Изд-во геодезической литературы. 1957. С. 152).
6. Боровков А.А. Теория вероятностей. М.: Наука, 1986, 432 с.