

УДК 519.233.5+519.654

О СВЯЗИ МЕЖДУ КОЭФФИЦИЕНТАМИ ПРОСТОЙ И МНОЖЕСТВЕННОЙ РЕГРЕССИОННЫХ МОДЕЛЕЙ

В. Г. Панов, А. Н. Вараксин

Аннотация. Рассмотрена задача нахождения зависимости между коэффициентами множественной регрессионной модели и коэффициентами простых корреляционных моделей, описывающих зависимость между предикторными переменными. Показано, что искомая зависимость описывается матричным равенством.

Ключевые слова: линейная регрессия, простая регрессия, коэффициенты регрессии, метод наименьших квадратов.

Исследование зависимости переменной-отклика в условиях многофакторного воздействия является актуальной и трудной задачей математической статистики. При этом, как правило, получение точной зависимости переменной-отклика от переменных-факторов (называемых также *предикторами*) оказывается слишком сложным и практически бесполезным даже в случае функциональной зависимости. Кроме того, отсутствие функциональной зависимости между переменными является типичным, если эти переменные имеют вероятностный характер. В этом случае описание ожидаемой зависимости носит приближенный характер и выражается некоторым функциональным уравнением между средними величинами.

Точная постановка соответствующей задачи рассматривается в регрессионном анализе. Наиболее простая и разработанная модель регрессионного анализа предполагает линейную статистическую зависимость между стохастическими переменной-откликом Y и переменными-предикторами X_1, X_2, \dots, X_n . Тогда зависимость Y от X_1, X_2, \dots, X_n можно моделировать линейным регрессионным уравнением, т. е. уравнением множественной линейной регрессии:

$$y = b_0 + \sum_{k=1}^n b_k x_k. \quad (1)$$

Постоянные b_0, b_1, \dots, b_n являются решениями вариационной задачи минимизации среднеквадратичного отклонения (\mathbf{E} — оператор математического ожидания):

$$\min_{b_0, b_1, \dots, b_n} \mathbf{E} \left(Y - b_0 - \sum_{k=1}^n b_k X_k \right)^2. \quad (2)$$

Таким образом, задача нахождения уравнения множественной регрессии переменной-отклика на заданный набор переменных-предикторов меняется, как только меняется набор предикторов или сама переменная-отклик. В регрессионном анализе зависимость между коэффициентами регрессионных моделей,

построенных по разным системам предикторов, обычно не рассматривается, так как ожидать какой-либо зависимости в этих моделях в общем случае не приходится.

Однако учет статистической зависимости как между предикторами, так и между откликом и частью предикторов делает проблему поиска связи между коэффициентами линейных регрессионных моделей более содержательной.

А именно, рассмотрим следующую задачу. Пусть имеется набор из $n+1$ случайных величин X_1, X_2, \dots, X_n, Y , у которых, вообще говоря, произвольные распределения (предполагается, однако, что среди случайных величин X_1, X_2, \dots, X_n нет постоянных). В рассматриваемых дальше линейных регрессионных моделях переменные X_1, X_2, \dots, X_n являются объясняющими (предикторными), а Y — откликом. Для оценки зависимости отклика Y от X_1, X_2, \dots, X_n рассмотрим линейное регрессионное уравнение, т. е. уравнение множественной линейной регрессии (1).

Будем считать, что для любой пары объясняющих переменных X_i, X_j найдено уравнение линейной регрессии

$$x_i = c_{0ij} + c_{ij}x_j, \quad i, j = 1, 2, \dots, n, \quad (3)$$

где коэффициенты c_{0ij}, c_{ij} являются решениями вариационной задачи

$$\min_{c_{0ij}, c_{ij}} \mathbf{E}(X_i - c_{0ij} - c_{ij}X_j)^2, \quad i, j = 1, 2, \dots, n. \quad (4)$$

Допустим также, что известно уравнение простой регрессии Y по каждой из переменных X_k :

$$y = a_{0k} + a_k x_k, \quad k = 1, 2, \dots, n. \quad (5)$$

Как и раньше, коэффициенты a_{0k}, a_k являются решениями задачи минимизации функционала среднеквадратичного уклонения

$$\min_{a_{0k}, a_k} \mathbf{E}(Y - a_{0k} - a_k X_k)^2, \quad k = 1, 2, \dots, n. \quad (6)$$

При обычных условиях [1, 2] каждая из этих экстремальных задач имеет единственное решение, причем для задачи (4) выполняются равенства $c_{0ii} = 0$ и $c_{ii} = 1$.

Рассматриваемый в предлагаемой статье вопрос состоит в исследовании связи между коэффициентами введенных выше регрессионных моделей: существует ли какая-либо связь между коэффициентами $\{a_k\}, \{b_k\}, \{c_{ij}\}$?

Предполагаемая связь регрессионных коэффициентов возникла как гипотеза при обработке экспериментальных данных в проблемах медико-биологического мониторинга [3, 4]. В качестве исследуемой переменной Y рассматривалась некоторая интегральная характеристика, которую можно назвать заболеваемостью населения, а в качестве факторов X_1, X_2, \dots, X_n — загрязнение атмосферного воздуха различными токсикантами. В этих работах отмечалось, что в задачах мониторинга модели множественной регрессии должны рассматриваться в условиях согласованного изменения предикторов $\{X_i\}$. В работе [4] предложены два варианта согласованного изменения $\{X_i\}$: пропорциональное и корреляционное. В данной работе изучаются соотношения между коэффициентами при корреляционной зависимости предикторов.

С более формальной точки зрения задача состоит в исследовании зависимости между моделями множественной регрессии отклика на предикторы,

парных корреляций предикторов и простых корреляций отклика на каждый из предикторов.

Перед формулировкой основного утверждения докажем вспомогательные леммы. Для случайной величины X обозначим через \bar{X} среднее значение (математическое ожидание) этой случайной величины. Так как будем рассматривать только дискретные случайные величины, ниже во всех случаях среднее значение соответствующей случайной величины существует. Для заданного набора предикторов X_1, X_2, \dots, X_n и отклика Y введем следующую вспомогательную матрицу порядка $n + 3$:

$$A = \begin{pmatrix} 1 & \bar{X}_1 & \bar{X}_2 & \dots & \bar{X}_n & \bar{Y} & 0 \\ \bar{X}_1 & \bar{X}_1^2 & \bar{X}_1\bar{X}_2 & \dots & \bar{X}_1\bar{X}_n & \bar{X}_1\bar{Y} & 0 \\ \bar{X}_2 & \bar{X}_2\bar{X}_1 & \bar{X}_2^2 & \dots & \bar{X}_2\bar{X}_n & \bar{X}_2\bar{Y} & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \bar{X}_n & \bar{X}_n\bar{X}_1 & \bar{X}_n\bar{X}_2 & \dots & \bar{X}_n^2 & \bar{X}_n\bar{Y} & 0 \\ 0 & \bar{X}_1 & \bar{X}_2 & \dots & \bar{X}_n & \bar{Y} & 1 \\ 0 & \bar{X}_k\bar{X}_1 & \bar{X}_k\bar{X}_2 & \dots & \bar{X}_k\bar{X}_n & \bar{X}_k\bar{Y} & \bar{X}_k \end{pmatrix}. \quad (7)$$

Предпоследняя $(n + 2)$ -я строка матрицы A получается из ее первой строки перестановкой первого и последнего элементов. Таким же образом из $(k + 1)$ -й ее строки получается последняя $(n + 3)$ -я строка.

Лемма 1. Матрица A вырождена.

Доказательство. Раскроем определитель матрицы A по последнему столбцу:

$$\det(A) = (-1)^{2n+5} \begin{vmatrix} 1 & \bar{X}_1 & \bar{X}_2 & \dots & \bar{X}_n & \bar{Y} \\ \bar{X}_1 & \bar{X}_1^2 & \bar{X}_1\bar{X}_2 & \dots & \bar{X}_1\bar{X}_n & \bar{X}_1\bar{Y} \\ \bar{X}_2 & \bar{X}_2\bar{X}_1 & \bar{X}_2^2 & \dots & \bar{X}_2\bar{X}_n & \bar{X}_2\bar{Y} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \bar{X}_n & \bar{X}_n\bar{X}_1 & \bar{X}_n\bar{X}_2 & \dots & \bar{X}_n^2 & \bar{X}_n\bar{Y} \\ 0 & \bar{X}_k\bar{X}_1 & \bar{X}_k\bar{X}_2 & \dots & \bar{X}_k\bar{X}_n & \bar{X}_k\bar{Y} \end{vmatrix} + \bar{X}_k \begin{vmatrix} 1 & \bar{X}_1 & \bar{X}_2 & \dots & \bar{X}_n & \bar{Y} \\ \bar{X}_1 & \bar{X}_1^2 & \bar{X}_1\bar{X}_2 & \dots & \bar{X}_1\bar{X}_n & \bar{X}_1\bar{Y} \\ \bar{X}_2 & \bar{X}_2\bar{X}_1 & \bar{X}_2^2 & \dots & \bar{X}_2\bar{X}_n & \bar{X}_2\bar{Y} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \bar{X}_n & \bar{X}_n\bar{X}_1 & \bar{X}_n\bar{X}_2 & \dots & \bar{X}_n^2 & \bar{X}_n\bar{Y} \\ 0 & \bar{X}_1 & \bar{X}_2 & \dots & \bar{X}_n & \bar{Y} \end{vmatrix}.$$

В первом определителе вычитаем последнюю строку из $(k + 1)$ -й строки, а во втором вычитаем последнюю строку из первой строки:

$$\det(A) = - \begin{vmatrix} 1 & \bar{X}_1 & \bar{X}_2 & \dots & \bar{X}_n & \bar{Y} \\ \bar{X}_1 & \bar{X}_1^2 & \bar{X}_1\bar{X}_2 & \dots & \bar{X}_1\bar{X}_n & \bar{X}_1\bar{Y} \\ \bar{X}_k & 0 & 0 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & \bar{X}_k\bar{X}_1 & \bar{X}_k\bar{X}_2 & \dots & \bar{X}_k\bar{X}_n & \bar{X}_k\bar{Y} \end{vmatrix}$$

$$\begin{aligned}
 & + \overline{X_k} \begin{vmatrix} 1 & 0 & 0 & \dots & 0 & 0 \\ \overline{X_1} & \overline{X_1^2} & \overline{X_1 X_2} & \dots & \overline{X_1 X_n} & \overline{X_1 Y} \\ \overline{X_2} & \overline{X_2 X_1} & \overline{X_2^2} & \dots & \overline{X_2 X_n} & \overline{X_2 Y} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \overline{X_n} & \overline{X_n X_1} & \overline{X_n X_2} & \dots & \overline{X_n^2} & \overline{X_n Y} \\ 0 & \overline{X_1} & \overline{X_2} & \dots & \overline{X_n} & \overline{Y} \end{vmatrix} \\
 = & (-1)^{k+1} \overline{X_k} \begin{vmatrix} \overline{X_1} & \overline{X_2} & \dots & \overline{X_n} & \overline{Y} \\ \overline{X_1^2} & \overline{X_1 X_2} & \dots & \overline{X_1 X_n} & \overline{X_1 Y} \\ \dots & \dots & \dots & \dots & \dots \\ \overline{X_{k-1} X_1} & \overline{X_{k-1} X_2} & \dots & \overline{X_{k-1} X_n} & \overline{X_{k-1} Y} \\ \overline{X_{k+1} X_1} & \overline{X_{k+1} X_2} & \dots & \overline{X_{k+1} X_n} & \overline{X_{k+1} Y} \\ \dots & \dots & \dots & \dots & \dots \\ \overline{X_n X_1} & \overline{X_n X_2} & \dots & \overline{X_n^2} & \overline{X_n Y} \\ \overline{X_k X_1} & \overline{X_k X_2} & \dots & \overline{X_k X_n} & \overline{X_k Y} \end{vmatrix} \\
 & + \overline{X_k} \begin{vmatrix} \overline{X_1^2} & \overline{X_1 X_2} & \dots & \overline{X_1 X_n} & \overline{X_1 Y} \\ \overline{X_2 X_1} & \overline{X_2^2} & \dots & \overline{X_2 X_n} & \overline{X_2 Y} \\ \dots & \dots & \dots & \dots & \dots \\ \overline{X_n X_1} & \overline{X_n X_2} & \dots & \overline{X_n^2} & \overline{X_n Y} \\ \overline{X_1} & \overline{X_2} & \dots & \overline{X_n} & \overline{Y} \end{vmatrix}.
 \end{aligned}$$

В первом определителе последовательными перестановками поднимаем последнюю строку до тех пор, пока она не станет $(k + 1)$ -й. Во втором определителе последнюю строку поднимаем на место первой. Получаются одинаковые определители с множителями $(-1)^{k+1+n-k} = (-1)^{n+1}$ и $(-1)^n$ соответственно. Лемма доказана.

Для доказательства основной теоремы нам понадобится теорема Лапласа о разложении определителя, которую удобно сформулировать, как в [5] (см. также [6]).

Пусть $Q_{r,n}$ — совокупность всех строго возрастающих последовательностей длины r из $(1, 2, \dots, n)$; $A(\alpha) = A(\alpha_1, \alpha_2, \dots, \alpha_r)$ — определитель подматрицы матрицы A , образованной пересечением строк с номерами из $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_r)$ и столбцов с номерами из $\beta = (\beta_1, \beta_2, \dots, \beta_r)$, где α, β — наборы из $Q_{r,n}$; $A(\alpha^*)$ — определитель подматрицы матрицы A , стоящей на пересечении строк с номерами из α^* и столбцов с номерами из β^* , где α^*, β^* — наборы из $Q_{n-r,n}$, дополнительные для наборов α и β соответственно: $\alpha \cap \alpha^* = \emptyset$, $\beta \cap \beta^* = \emptyset$, $\alpha \cup \alpha^* = \beta \cup \beta^* = \{1, 2, \dots, n\}$; для любого набора $\beta = (\beta_1, \beta_2, \dots, \beta_r)$ обозначим $|\beta| = \sum_{i=1}^r \beta_i$.

Теорема Лапласа. Пусть A — квадратная матрица порядка n , $1 \leq r \leq n$, α — фиксированная последовательность из $Q_{r,n}$. Тогда

$$\det(A) = (-1)^{|\alpha|} \sum_{\beta \in Q_{r,n}} (-1)^{|\beta|} A(\alpha, \beta) \cdot A(\alpha^*, \beta^*).$$

Возьмем $r = n + 1$, $\alpha = (1, 2, \dots, n + 1)$, $\alpha \in Q_{n+1, n+3}$, и применим к матрице (7) порядка $n + 3$ теорему Лапласа. Учитывая лемму 1, получим равенство

$$\sum_{\beta \in Q_{n+1, n+3}} (-1)^{|\beta|} A\left(\begin{matrix} 1, 2, \dots, n+1 \\ \beta_1, \beta_2, \dots, \beta_{n+1} \end{matrix}\right) \cdot A\left(\begin{matrix} n+2, n+3 \\ * * \end{matrix}\right) = 0. \quad (8)$$

Звездочками в этом равенстве обозначены такие номера (расположенные в порядке возрастания) от 1 до $n + 3$, которые вместе с набором β образуют все множество $\{1, 2, \dots, n + 3\}$.

Лемма 2. *Количество наборов β из множества $Q_{n+1, n+3}$, для которых произведение определителей в равенстве (8) отлично от нуля, не более $n + 1$.*

ДОКАЗАТЕЛЬСТВО. Если в набор β не входит 1, то набор, дополнительный для β , начинается с 1 и, следовательно,

$$A \begin{pmatrix} n+2 & n+3 \\ 1 & * \end{pmatrix} = 0.$$

С другой стороны, в набор β не может входить число $n + 3$, так как в силу строения матрицы A

$$A \begin{pmatrix} 1, 2, \dots, n+1 \\ \beta_1, \beta_2, \dots, n+3 \end{pmatrix} = 0.$$

Следовательно, интересующие нас наборы β должны иметь вид

$$(1, \beta_2, \dots, \beta_n, n+1) \quad \text{или} \quad (1, \beta_2, \dots, \beta_n, n+2).$$

Так как первый набор определен однозначно и равен $(1, 2, \dots, n, n+1)$, остается доказать, что наборов второго типа будет ровно n . Из этих наборов выделим случай

$$\beta = (1, 2, \dots, n, n+2).$$

Для остальных наборов рассматриваемого типа существует единственное число i от 2 до n , (включительно) такое, что

$$\begin{aligned} \beta_j &= j & \text{при } j = 1, 2, \dots, i-1, \\ \beta_i &= i+1 & \text{при } j = i, \\ \beta_j &= j+1 & \text{при } j = i+1, i+2, \dots, n+1. \end{aligned} \tag{9}$$

Обозначим через \mathbf{b}_i последовательность, определяемую условиями (9). Таким образом, наборов второго типа будет ровно столько, сколько значений может принимать индекс i в условиях (9). Следовательно, всего наборов второго типа будет $n - 1 + 1 = n$, что и доказывает лемму.

Для формулировки основного результата введем следующие матрицы:

$$\mathcal{A} = (a_1, a_2, \dots, a_n), \quad \mathcal{B} = (b_1, b_2, \dots, b_n), \quad \mathcal{C} = (c_{ij})_{n \times n},$$

где $\{a_k\}$, $\{b_k\}$, $\{c_{ij}\}$ — коэффициенты уравнений (1), (3), (5), полученные как решение вариационных задач (2), (4), (6).

Теорема. *Коэффициенты a_k, b_i, c_{ik} связаны равенством*

$$a_k = \sum_{i=1}^n b_i \cdot c_{ik} \tag{10}$$

или, в матричных обозначениях,

$$\mathcal{A} = \mathcal{B} \cdot \mathcal{C}.$$

где \mathbf{b}_i — набор (9), $i = 2, 3, \dots, n$. Подставляя эти выражения в равенство (8), получаем

$$\begin{aligned}
0 &= \sum_{\beta \in Q_{n+1, n+3}} (-1)^{|\beta|} A \begin{pmatrix} 1, 2, \dots, n+1 \\ \beta_1, \beta_2, \dots, \beta_{n+1} \end{pmatrix} \cdot A \begin{pmatrix} n+2, n+3 \\ * * \end{pmatrix} \\
&= (-1)^{\frac{(n+1)(n+2)}{2}} \Delta \cdot (-\delta_{1n}) + (-1)^{\frac{(n+1)(n+2)}{2}+1} \Delta_n \cdot (-d_n) \\
&\quad + \sum_{i=2}^n (-1)^{\frac{(n+2)(n+3)}{2}-i} (-1)^{n-i+1} \Delta_{i-1} \cdot (-d_{i-1k}) \\
&= (-1)^{\frac{(n+1)(n+2)}{2}+1} \left[\Delta \cdot \delta_{1k} - \Delta_n \cdot d_{nk} - \sum_{i=2}^n (-1)^{n+1-i} (-1)^{n-i+1} \Delta_{i-1} \cdot d_{i-1k} \right] \\
&= (-1)^{\frac{(n+1)(n+2)}{2}+1} \left[\Delta \cdot \delta_{1k} - \Delta_n \cdot d_{nk} - \sum_{i=2}^n \Delta_{i-1} \cdot d_{i-1k} \right].
\end{aligned}$$

Следовательно,

$$\Delta \cdot \delta_{1k} = \Delta_n \cdot d_{nk} + \sum_{i=2}^n \Delta_{i-1} \cdot d_{i-1k}.$$

Деля это равенство на $\delta_k \cdot \Delta$ и учитывая равенства (12), получим

$$a_k = \sum_{i=1}^n b_i \cdot c_{ik},$$

что и требовалось доказать.

Следствие 1. Если предикторы X_1, X_2, \dots, X_n попарно не коррелированы (т. е. $c_{ik} = 0$ для всех $i \neq k$), то $a_k = b_k$. Таким образом, коэффициенты множественной регрессии в этом случае совпадают с коэффициентами парных корреляций.

Доказательство следует из того, что в случае нулевых попарных корреляций предикторов матрица \mathcal{C} становится единичной.

Следствие 2. Если матрица \mathcal{C} обратима, то строку коэффициентов $\mathcal{B} = (b_1, b_2, \dots, b_n)$ уравнения (1) можно выразить через коэффициенты простых корреляций Y по предикторам X равенством

$$\mathcal{B} = A \cdot \mathcal{C}^{-1}.$$

Свободный член b_0 можно найти, например, из первого уравнения (11).

ЗАМЕЧАНИЯ. 1. Как показывает следствие 2, с помощью полученного соотношения (10) можно вычислить коэффициенты множественной регрессии отклика Y на множество предикторов X_1, X_2, \dots, X_n через коэффициенты парных корреляций \mathcal{A} и \mathcal{C} . В регрессионном анализе существует другой метод сведения уравнения множественной регрессии к серии уравнений парных корреляций, изложенный, например, в [8, разд. 4.1]. Однако в этом методе существенную роль играет вычисление регрессии остатков на остатки, которое в данном случае совсем не требуется. Кроме того, сама процедура вычисления коэффициентов множественной регрессии выполняется последовательно в несколько шагов (тем больше, чем больше предикторов), что приводит к накоплению вычислительной ошибки. В доказанном выражении это достигается с помощью одного матричного произведения.

2. Сам метод из [8] является методом нахождения коэффициентов множественной регрессии и не может быть распространен на вычисление каких-либо других регрессионных коэффициентов (например, парных) по известным коэффициентам множественной регрессии. В то же время равенство $\mathcal{A} = \mathcal{B} \cdot \mathcal{C}$ дает возможность вычисления коэффициентов регрессии отклика на каждый предиктор при известной матрице коэффициентов регрессий между предикторами и коэффициентов множественной регрессии отклика на все предикторы.

ЛИТЕРАТУРА

1. Магнус Я. Р., Нейдеккер Х. Матричное дифференциальное исчисление с приложениями к статистике и эконометрике. М.: Физматлит, 2002.
2. Козлов М. В., Прохоров А. В. Введение в математическую статистику. М.: Изд-во Моск. ун-та, 1987.
3. Вараксин А. Н., Маслакова Т. А., Чуканов В. Н., Антонов К. Л. Регрессионная модель зависимости заболеваемости населения от степени загрязнения атмосферного воздуха // Экологические системы и приборы. 2004. № 4. С. 52–55.
4. Маслакова Т. А., Вараксин А. Н., Чуканов В. Н. Интерпретация прогностических регрессионных моделей в области медико-экологического мониторинга // Экологические системы и приборы. 2008. № 2. С. 6–9.
5. Маркус М., Минк Х. Обзор по теории матриц и матричных неравенств. М.: Едиториал УРСС, 2004.
6. Гантмахер Ф. Р. Теория матриц. М.: Наука, 1988.
7. Линник Ю. В. Метод наименьших квадратов и основы математико-статистической теории обработки наблюдений. М.: Наука, 1962.
8. Дрейпер Н., Смит Г. Прикладной регрессионный анализ. М.: Финансы и статистика, 1986. Кн. 1.

Статья поступила 7 мая 2008 г., окончательный вариант — 5 декабря 2008 г.

Панов Владимир Григорьевич
Уральский гос. университет, кафедра математического анализа и теории функций,
пр. Ленина, 51, Екатеринбург 620083
vladimir.panov@usu.ru

Вараксин Анатолий Николаевич
лаборатория математического моделирования
Института промышленной экологии УрО РАН,
ул. Софьи Ковалевской, 20-а, Екатеринбург 620219, ГСП-594
varaksin@ecko.uran.ru