

УДК 519.24

Л.Г. Раскин, О.В. Серая

Национальный технический университет «ХПИ», Харьков

ИНФОРМАЦИОННЫЕ ПРОБЛЕМЫ КАНОНИЧЕСКОГО РЕГРЕССИОННОГО АНАЛИЗА

Рассмотрены технологии решения задач регрессионного анализа для случаев, когда нарушаются традиционные исходные предпосылки. При этом разработаны методы оценивания параметров уравнения регрессии, если независимые переменные – случайные величины или заданы нечетко. Предложено решение задачи компараторной идентификации для нечетких исходных данных.

Ключевые слова: регрессионный анализ, нарушение предпосылок, нечеткие исходные данные, компараторная идентификация.

Введение

Регрессионный анализ – мощный и эффективный статистический метод построения математических моделей, описывающих зависимость между показателем функционирования анализируемой системы y и обуславливающими, объясняющими независимыми переменными (факторами) F_1, F_2, \dots, F_m . С целью выявления этой связи проводится серия экспериментов, в которой каждому опыту $(F_{j1}, F_{j2}, \dots, F_{jm})$ ставится в соответствие его результат – значение зависимой переменной y_j , $j = 1, 2, \dots, n$. Искомая связь обычно описывается полиномом Колмогорова-Габор, который в простейшем случае имеет вид:

$$y_j = x_0 + F_{j1}x_1 + F_{j2}x_2 + \dots + F_{jm}x_m + \varepsilon_j \quad (1)$$

Здесь F_{ji} – значение i -й независимой переменной в j -м опыте, $i = 0, 1, 2, \dots, m$, $j = 1, 2, \dots, n$.

В матричной форме соотношение (1) имеет вид $F\mathbf{X} = \mathbf{Y}$, где

$$F = \begin{pmatrix} 1 & F_{11} & F_{12} & \dots & F_{1m} \\ 1 & F_{21} & F_{22} & \dots & F_{2m} \\ \dots & \dots & \dots & \dots & \dots \\ 1 & F_{n1} & F_{n2} & \dots & F_{nm} \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} x_0 \\ x_1 \\ \dots \\ x_m \end{pmatrix}, \quad \mathbf{Y} = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix}.$$

В каноническом регрессионном анализе делаются следующие основные предположения.

1. Значения независимых переменных F_i измеряются без ошибок, $i = 1, 2, \dots, m$.

2. Зависимая переменная y в каждом опыте оценивается со случайной ошибкой ε_j , которая нормально распределена с нулевым математическим ожиданием и известной дисперсией σ^2 .

3. Случайные ошибки ε_j в разных опытах не коррелированы.

В этих предложениях оценки неизвестных коэффициентов x_0, x_1, \dots, x_m регрессионного полинома (1) получают методом наименьших квадратов

(МНК), минимизируя сумму квадратов отклонений значений результирующей переменной y_j от соответствующих значений $\sum_{i=0}^m x_i F_{ji}$, предсказываемых моделью (1).

Постановка задачи. При решении многих практических задач возникают ситуации, когда исходные предпосылки классического регрессионного анализа не верны. При этом, естественно, шаблонное применения МНК может привести к грубым ошибкам. Поставим задачу разработки методик оценивания коэффициентов уравнения регрессии для разных вариантов нарушений канонических предпосылок.

Основные результаты

А. Оценивание параметров регрессионной модели по данным экспериментов со случайными ошибками в значениях независимых переменных (нарушена предпосылка 1).

Будем считать, что наблюдаемое значение F_{ji} есть нормально распределенная случайная величина с математическим ожиданием, равным истинному значению w_{ji} , и дисперсией σ_{ji}^2 , а наблюдаемое значение Y_j есть нормально распределенная случайная величина с математическим ожиданием, равным истинному значению v_j , и дисперсией σ_j^2 .

При этом, как легко показать [1], оценки регрессионных коэффициентов, получаемых с использованием МНК, оказываются смещенными. В [2] доказано, что величина смещения зависит от исходных данных и может быть очень большой. Известные методы снижения уровня смещения [3,4], сопряжены с необходимостью проведения сложных вычислений и приводят к приемлемым результатам, если ошибки оценивания независимых переменных малы.

В связи с этим рассмотрим другой подход к решению задачи оценивания регрессионных коэф-

коэффициентов в условиях, когда значения независимых переменных оцениваются с ошибкой.

Принятые предположения с учетом принципа инверсии позволяют записать законы распределения неизвестных истинных значений независимых и зависимой переменных в виде

$$\mu(w_{ji}) = \frac{1}{\sqrt{2\pi\sigma_{ji}}} \exp\left\{-\frac{(w_{ji} - F_{ji})^2}{2\sigma_{ji}^2}\right\}, \quad (2)$$

$$\mu(v_j) = \frac{1}{\sqrt{2\pi\sigma_j}} \exp\left\{-\frac{(v_j - y_j)^2}{2\sigma_j^2}\right\},$$

$$i = 1, 2, \dots, m, \quad j = 1, 2, \dots, n. \quad (3)$$

Если модель (1) адекватна, то для истинных значений наблюдаемых величин должны выполняться равенства:

$$x_0 + w_{11}x_1 + w_{12}x_2 + \dots + w_{1m}x_m = v_1,$$

$$\dots \dots \dots \quad (4)$$

$$x_0 + w_{n1}x_1 + w_{n2}x_2 + \dots + w_{nm}x_m = v_n.$$

Поскольку параметры системы (4) – случайные величины, то соотношения (2) - (4) задают стохастическую систему линейных алгебраических уравнений. Теперь поступаем следующим образом. Решим порождаемую (4) детерминированную систему линейных алгебраических уравнений, заменив случайные параметры (w_{ji}) и (v_i) их математическими ожиданиями. При этом получим

$$\sum_{i=0}^m F_{ji}x_i = y_j, \quad j = 1, 2, \dots, n. \quad (5)$$

Эта переопределенная система в матричной форме имеет вид

$$FX = Y, \quad F = (F_{ji}),$$

$$F_{j0} = 1, \quad i = 0, 1, 2, \dots, m, \quad j = 1, 2, \dots, n,$$

решается путем минимизации функционала

$$J = (FX - Y)^T (FX - Y),$$

а искомый вектор – решение имеет вид

$$X^{(0)} = (F^T F)^{-1} F^T Y. \quad (6)$$

Назовем это решение модальным.

Введем далее набор переменных $Z_j(X)$:

$$Z_1(X) = w_{11}x_1 + w_{12}x_2 + \dots + w_{1m}x_m,$$

$$\dots \dots \dots \quad (7)$$

$$Z_n(X) = w_{n1}x_1 + w_{n2}x_2 + \dots + w_{nm}x_m,$$

задающих предсказываемые моделью значения зависимой переменной в каждом опыте. Введенные переменные – это случайные величины, плотности распределения которых совместно с $(\mu(w_{ij}))$ определяются набором X:

$$\phi(Z_j) = \frac{1}{\sqrt{2\pi D_j}} \exp\left\{-\frac{(Z_j - m_j)^2}{2D_j}\right\},$$

$$m_j = \sum_{i=0}^m F_{ji}x_i, \quad D_j = \sum_{i=0}^m \sigma_{ji}^2 x_i^2, \quad (8)$$

$$\sigma_{j0}^2 = 0, \quad j = 1, 2, \dots, n.$$

Теперь детерминированным решением стохастической задачи (4) будем называть набор $X = (x_i), i = 1, 2, \dots, m$, минимизирующий сумму дисперсий случайных величин $Z_j(x)$ и наименее уклоняющийся от $X^{(0)}$. Смысл этого критерия понятен. Его использование обеспечивает получение набора детерминированных оценок регрессионных коэффициентов x_1, x_2, \dots, x_m , для которых плотности распределения случайных величин Z_1, Z_2, \dots, Z_n , определяющих предсказываемые моделью значения зависимой переменной в каждом опыте, наименее размыты, и имеющих математические ожидания, максимально близкие к нулю. Формальное описание критерия имеет вид:

$$L(x) = \sum_{j=1}^n D_j(x) + \sum_{i=0}^m (x_i - x_i^{(0)})^2 =$$

$$= \sum_{i=0}^m \left[x_i^2 \left(\sum_{j=1}^n \sigma_{ji}^2 \right) + (x_i - x_i^{(0)})^2 \right] = \quad (9)$$

$$= \sum_{i=0}^m \left[x_i^2 D_{\Sigma i} + (x_i - x_i^{(0)})^2 \right].$$

Найдем набор X, минимизирующий (9), методом множителей Лагранжа, дополнив модель задачи условием нормировки

$$\sum_{i=0}^m x_i = 1.$$

Запишем функцию Лагранжа

$$\Phi(x) = L(x) - \lambda \left(\sum_{i=0}^m x_i - 1 \right).$$

Далее имеем

$$\frac{d\Phi(x)}{dx_i} = 2x_i(D_{\Sigma i} + 1) - 2x_i^{(0)} - \lambda = 0,$$

откуда

$$x_i = \frac{2x_i^{(0)} + \lambda}{2(D_{\Sigma i} + 1)} = \frac{\lambda}{2} \cdot \frac{1}{D_{\Sigma i} + 1} + \frac{x_i^{(0)}}{D_{\Sigma i} + 1}.$$

Значение $\lambda/2$ найдем из условия нормировки

$$\sum_{i=0}^m x_i = \frac{\lambda}{2} \sum_{i=0}^m \frac{1}{D_{\Sigma i} + 1} + \sum_{i=0}^m \frac{x_i^{(0)}}{D_{\Sigma i} + 1} = 1,$$

$$\frac{\lambda}{2} = \left(1 - \sum_{i=0}^m \frac{x_i^{(0)}}{D_{\Sigma i} + 1} \right) / \sum_{i=0}^m \frac{1}{D_{\Sigma i} + 1}.$$

$$\sum_{i=1}^m x_i V_{ji} + x_{m+1} = 0, j = 1, 2, \dots, n - 1. \quad (15)$$

При этом, с целью исключения тривиального решения $x_1 = x_2 = \dots = x_m = x_{m+1} = 0$ к системе (15) следует добавить ещё одно уравнение

$$\sum_{i=1}^m x_i = 1, \quad (16)$$

являющееся условием нормировки для коэффициентов уравнения регрессии (11), а также неравенство $x_{m+1} > 0$. Получаемое решение используется для расчета «полезности» объектов по формуле (11).

Г. Компараторная идентификация для случая, когда независимые переменные заданы нечетко.

Во многих практических ситуациях значения характеристик объектов не могут быть определены точно. Для описания этих значений используем технологию нечеткой математики [5,8]. Будем считать, что значение i -й характеристики j -го объекта есть нечеткое число r_{ji} с функцией принадлежности

$$\mu_{ji}(r_{ji}), i = 1, 2, \dots, m, j = 1, 2, \dots, n.$$

Поскольку характеристики объектов – нечеткие числа, то нечеткими является и результаты расчетов значений уровня «полезности» объектов

$$Q_j(x) = \sum_{i=1}^m x_i r_{ji}, i = 1, 2, \dots, m. \quad (17)$$

Пусть, для определенности, r_{ji} – нечеткие числа с треугольной функцией принадлежности, то есть

$$\mu_{ij}(r_{ij}) = \begin{cases} 0, & r_{ji} < a_{ji}, \\ \frac{r_{ji} - a_{ji}}{c_{ji} - a_{ji}}, & a_{ji} \leq r_{ji} < c_{ji}, \\ \frac{b_{ji} - r_{ji}}{b_{ji} - c_{ji}}, & c_{ji} \leq r_{ji} < b_{ji}, \\ 1, & r_{ji} \geq b_{ji}. \end{cases} \quad (18)$$

Поскольку треугольные нечеткие числа – частный случай чисел с функцией принадлежности (L – R)- типа, то для приближенных расчетов можно использовать правила выполнения операций, принятых для нечетких чисел этого типа.

С учетом этого, определим функции принадлежности нечетких чисел $Q_i(X)$:

$$\mu_i(Q_i(X)) = \begin{cases} 0, & Q_j(X) < A_i, \\ \frac{Q_i(X) - A_i}{C_i - A_i}, & A_i \leq Q_i(X) < C_i, \\ \frac{B_i - Q_i(X)}{B_i - C_i}, & C_i \leq Q_i(X) < B_i, \\ 0, & Q_i(X) \geq B_i \end{cases} \quad (19)$$

$$A_i = \sum_{j=1}^n x_j a_{ji}, C_i = \sum_{j=1}^n x_j c_{ji}, B_i = \sum_{j=1}^n x_j b_{ji}.$$

Если в рассматриваемом случае ранжировка (11) объектов по уровню их «полезности» сохранилась, то естественным аналогом (13) будут нечеткие неравенства

$$\begin{aligned} \zeta_1(X) &= Q_2(X) - Q_1(X) = \\ &= x_1(r_{21} - r_{11}) + x_2(r_{22} - r_{12}) + \\ &\quad + \dots + x_m(r_{2m} - r_{1m}) = \\ &= \sum_{i=1}^m x_i w_{ji} < 0, \end{aligned} \quad (20)$$

$$\begin{aligned} \zeta_{n-1}(X) &= Q_n(X) - Q_{n-1}(X) = \\ &= \sum_{i=1}^m x_i w_{n-1,i} < 0, \\ w_{ji} &= r_{j+1,i} - r_{ji}, \\ j &= 1, 2, \dots, n - 1, i = 1, 2, \dots, m. \end{aligned} \quad (21)$$

Поставим задачу отыскания неотрицательного набора $X = (x_1, x_2, \dots, x_n)$, обеспечивающего выполнение неравенств (20).

Система неравенств (20) с добавлением положительной переменной x_{m+1} преобразуется в нечеткую систему линейных алгебраических уравнений

$$\sum_{i=1}^m x_i w_{ij} + x_{m+1} = 0, j = 1, 2, \dots, n - 1, \quad (22)$$

где w_{ji} - нечеткие числа, функция принадлежности которых определяется с учетом (22).

Используя стандартное описание функции принадлежности треугольных чисел L – R -типа в виде

$$r_{ji} = \langle c_{ji}, \alpha_{ji}, \beta_{ji} \rangle_{L-R},$$

где c_{ji} - мода числа r_{ji} , $\alpha_{ji} = c_{ji} - a_{ji}$, $\beta_{ji} = b_{ji} - c_{ji}$, и правила выполнения операций, принятых для нечетких чисел (L – R)-типа, запишем функцию принадлежности числа w_{ji} .

Имеем

$$\begin{aligned} w_{ji} &= r_{j,i+1} - r_{ij} = \\ &= \langle c_{j,i+1} - c_{ij}, c_{ij} - a_{ij} + b_{j,i+1} - c_{j,i+1}, \\ & b_{ij} - c_{ij} + c_{j,i+1} - a_{j,i+1} \rangle_{LR} = \langle m_{ji}, \alpha_{ji}, \beta_{ji} \rangle_{LR} \end{aligned}$$

Здесь m_{ji} - мода числа w_{ji} , α_{ji} и β_{ji} - левый и правый коэффициент нечеткости.

Введем нечеткие числа

$$Z_j(X) = \sum_{i=1}^m x_i w_{ij} + x_{m+1}, j = 1, 2, \dots, n - 1,$$

и запишем функции их принадлежности:

$$\mu(Z_j(X)) = \begin{cases} 0, & Z_j(X) < a_1 \\ \frac{Z_j(x) - a_1}{\sum_{i=1}^m x_i \alpha_{ji}}, & a_1 \leq Z_j(X) \leq a_2, \\ \frac{a_3 - Z_j(X)}{\sum_{i=1}^m x_i \beta_{ji}}, & a_2 \leq Z_j(x) \leq a_3, \\ 0, & Z_j(x) \geq a_3; \end{cases} \quad (23)$$

$$a_1 = \sum_{i=1}^m x_i (m_{ji} - \alpha_{ji}) + x_{m+1}, \quad a_2 =$$

$$= \sum_{i=1}^m x_i m_{ji} + x_{m+1},$$

$$a_3 = \sum_{j=1}^n x_j (m_{ji} + \beta_{ji}) + x_{n+1}.$$

Преобразуем вытекающую из (22) систему нечетких уравнений

$$Z_j(X) = 0, \quad j = 1, 2, \dots, n-1,$$

в обычную систему линейных алгебраических уравнений задав, нечеткие числа W_{ji} равными их модальным значениям. При этом получим

$$\sum_{i=1}^m x_i m_{ji} + x_{m+1} = 0, \quad j = 1, 2, \dots, n-1. \quad (24)$$

С целью исключения тривиального решения $x_i = 0, i = 1, 2, \dots, m+1$, системы (24) добавим к ней ещё одно нормирующее уравнение (6).

Пусть набор $X^{(0)} = (x_0^{(0)}, x_1^{(0)}, \dots, x_m^{(0)}, x_{m+1}^{(0)})$ есть решение системы (24), (16).

Используем введенное в [6] определение «четкое решение нечеткой системы линейных алгебраических уравнений». В соответствии с этим определением четким решением системы уравнений (24), (16) будем называть набор $X = (x_1, x_2, \dots, x_{n+1})$, минимизирующий сумму площадей фигур, ограниченных функциями принадлежности $\mu(Z_j)$ нечетких чисел

Z_1, Z_2, \dots, Z_{n-1} , и наименее уклоняющийся от $X^{(0)}$.

Смысл этого определения понятен. Его использование обеспечивает получение набора четких чисел $(x_1, x_2, \dots, x_{n+1})$, максимально близкого к модальному $X^{(0)}$, и для которого функции принадлежности нечетких чисел Z_1, Z_2, \dots, Z_{m-1} наименее размыты. В качестве критерия компактности функций принадлежности нечетких чисел $Z_j, j = 1, 2, \dots, n-1$, могут быть использованы квадраты длины интервалов – носителей соответствующих нечетких чисел. Тогда мера качества решения системы (24), (16) будет иметь вид:

$$J = \sum_{i=1}^n x_i^2 \left(\sum_{i=1}^{m-1} (\beta_{ij} + \alpha_{ij}) \right)^2 + \sum_{i=1}^m (x_i - x_i^{(0)})^2 \quad (25)$$

Минимизация (25) совместно с условием нормировки (16) дает искомый набор

$$X = (x_1, \dots, x_{m+1}).$$

Список литературы

1. Вучков И. Прикладной линейный регрессионный анализ / Пер. с болг. И. Вучков, Л. Бояджиева, Б. Солаков. – М.: Финансы и статистика, 1987. – 239 с.
2. Rao J. Subrahmaniam K. Combining independent estimators and estimation in linear regression with unequal variances // *Biometrics*. – 1997. – V. 27. – P. 971
3. Durbin J. Errors in variables / J. Durbin // *Rev. Int. Stat. Inst.* – 1954. – V. 22. – P. 23.
4. Madansky A. The biting of straight lines when both variables are subject to error / A. Madansky // *J. Amer. Stat. Assoc.* – 1959. – V. 54. – P. 173
5. Раскин Л.Г. Нечеткая математика. Основы теории. Приложение / Л.Г. Раскин, О.В. Серая. – Х.: Парус, 2008. – 352 с.
6. Ведение в нормативную теорию принятия решений / В.В. Крючковский, Э.Г. Петров, Н.А. Соколова, В.Е. Ходаков. – Херсон : Гринь Д.С., 2013. – 284 с.
7. Зуховицкий С.И. Линейное и выпуклое программирование / С.И. Зуховицкий, Л.И. Авдеева. – М.: Наука, 1967. – 460 с.
8. Zadeh L.A. Fuzzy sets / Zadeh L.A. // *Inf. Conf.* – 1965. – 8. – P. 338-353.

Поступила в редколлегию 25.08.2015

Рецензент: д-р техн. наук, проф. С.М. Порошин, Национальный технический университет «ХПИ», Харьков.

ІНФОРМАЦІЙНІ ПРОБЛЕМИ КАНОНІЧНОГО РЕГРЕСІЙНОГО АНАЛІЗУ

Л.Г. Раскін, О.В. Сіра

Розглянуто технології вирішення задач регресійного аналізу для випадків, коли порушуються традиційні вихідні передумови. При цьому розроблені методи оцінювання параметрів рівняння регресії, якщо незалежні змінні - випадкові величини або задані нечітко. Запропоновано вирішення завдання компараторної ідентифікації для нечітких вихідних даних.

Ключові слова: регресійний аналіз, порушення передумов, нечіткі вихідні дані, компараторного ідентифікація

INFORMATION PROBLEMS OF CANONICAL REGRESSION ANALYSIS

L.G. Raskin, O.W. Seraja

The technology solutions for the problems of regression analysis where the traditional presuppositions are violated. At the same time methods for estimating the parameters of the regression equation when the independent variables - the random fuzzy variables are developed. A solution of the problem of fuzzy comparator for identifying the source data is offered.

Keywords: regression analysis, violation of assumptions, fuzzy source data, comparator identification.