

ИНФОРМАЦИОННЫЕ ТЕХНОЛОГИИ И ТЕЛЕКОММУНИКАЦИИ
INFORMATION TECHNOLOGY AND TELECOMMUNICATIONS

УДК 519.862.6

DOI: 10.21822/2073-6185-2022-49-3-32-38

Оригинальная статья/Original Paper

**Программа оценивания с помощью метода наименьших квадратов неэлементарных
линейных регрессий с двумя переменными**

М.П. Базилевский, Д.В. Карбушева

Иркутский государственный университет путей сообщения,
664074, г. Иркутск, ул. Чернышевского, 15, Россия

Резюме. Цель. Целью исследования является разработка программы приближенного оценивания специфицированных на основе производственной функции Леонтьева регрессионных моделей (неэлементарных регрессий с двумя переменными) и её применение для моделирования уровня безработицы в Иркутской области. **Метод.** Оценивание неэлементарных регрессий осуществляется с помощью метода наименьших квадратов. Для нахождения приближенных оценок использован ранее разработанный алгоритм, предполагающий решение весьма трудоемкой вычислительной задачи. **Результат.** На основе этого алгоритма в среде программирования Delphi была разработана специальная программа. Программа предусматривает работу в ручном и автоматическом режимах. В ручном режиме по заданным критериям определяются оценки параметров модели, сумма квадратов остатков, коэффициент детерминации, критерий Стьюдента, Дарбина-Уотсона и для каждой переменной номера срабатываний компонент бинарной операции по выборке. В автоматическом режиме определяются наилучшие оценки неэлементарной регрессии по критериям: суммы квадратов остатков, коэффициента детерминации, Стьюдента и Дарбина-Уотсона. При этом строятся графики всех основных характеристик в зависимости от ключевого параметра модели. С помощью разработанной программы построена модель уровня безработицы в Иркутской области. **Вывод.** Построенная с помощью разработанной программы модель оказалась лучше, чем традиционная модель множественной линейной регрессии. Программа является универсальной и может применяться для решения конкретных прикладных задач анализа данных.

Ключевые слова: регрессионная модель, метод наименьших квадратов, производственная функция Леонтьева, неэлементарная линейная регрессия, коэффициент детерминации, безработица

Для цитирования: М.П. Базилевский, Д.В. Карбушева. Программа оценивания с помощью метода наименьших квадратов неэлементарных линейных регрессий с двумя переменными. Вестник Дагестанского государственного технического университета. Технические науки. 2022; 49(3):32-38. DOI:10.21822/2073-6185-2022-49-3-32-38

**The program for estimation non-elementary linear regressions
with two variables using ordinary least squares**

M.P. Bazilevskiy, D.V. Karbusheva

Irkutsk State Transport University,
15 Chernyshevskogo Str., Irkutsk 664074, Russia

Abstract. Objective. The aim of this article is to develop a program for approximate estimation of regression models specified on the basis of the Leontief production function (non-elementary regressions with two variables) and use it for modeling the unemployment rate in the Irkutsk region. **Method.** Estimation of non-elementary regressions is carried out using ordinary least squares method. To find approximate estimates, we used a previously developed algorithm that involves solving a very laborious computational problem. **Result.** Based on this algorithm, a special program was developed in the Delphi programming environment. The program provides for work in manual and automatic

modes. In manual mode, according to the specified criteria, the estimates of the model parameters, the residual sum of squares, the coefficient of determination, the Student's criterion, Durbin-Watson's criterion and, for each variable, the number of the binary operation components triggerings on the sample, are determined. In automatic mode, the best estimates of non-elementary regression are determined according to the criteria: residual sum of squares, coefficient of determination, the Student's criterion and Durbin-Watson's criterion. At the same time, graphs of all the main characteristics are plotted depending on the key parameter of the model. With the help of the developed program, a model of the unemployment rate in the Irkutsk region was construct. **Conclusion.** The model construct using the developed program turned out to be better than the traditional model of multiple linear regression. The program is universal and can be used to solve specific applied problems of data analysis.

Keywords: regression model, ordinary least squares, Leontief production function, non-elementary linear regression, coefficient of determination, unemployment

For citation: M.P. Bazilevskiy, D.V. Karbusheva. The program for estimation non-elementary linear regressions with two variables using ordinary least squares. Herald of the Daghestan State Technical University. Technical Science. 2022; 49 (3): 32-38. DOI: 10.21822 /2073-6185-2022-49-3-32-38

Введение. В настоящее время в различных областях деятельности – экономике, бизнесе, социологии, медицине и др. – накоплено огромное количество больших массивов статистических данных. Поэтому актуальной научной задачей является эффективная обработка и анализ этих данных с целью извлечения из них неизвестных ранее полезных знаний. Эффективным инструментом анализа данных является регрессионный анализ [1–3]. Его применение приводит к построению регрессионной модели влияния одной или нескольких объясняющих переменных на зависимую переменную. Решению конкретных прикладных задач с помощью регрессионного анализа посвящено множество научных работ (например, [4–8]).

Постановка задачи. Одной из главных проблем, возникающих при построении регрессионной модели, является выбор её структурной спецификации, т.е. математической формы связи между переменными. На сегодняшний день разработан весьма внушительный арсенал таких форм, подробное описание которых можно найти в [9,10]. В экономике особое внимание при этом традиционно уделяется вопросам построения производственных функций (ПФ) [11–13]. К наиболее известным относится ПФ Леонтьева вида

$$y = \min \{a_1 x_1, a_2 x_2\},$$

в которой переменная y зачастую трактуется как объемы выпускаемой продукции, x_1, x_2 – факторы производства, a_1, a_2 – определяемые технологией постоянные величины.

ПФ Леонтьева целесообразно применять тогда, когда объем выпуска продукции в моделируемой системе определяется количеством ресурса, обеспечивающего лишь наименьший возможный выпуск.

Для оценивания неизвестных параметров a_1 и a_2 ПФ Леонтьева составим регрессионную модель:

$$y_i = \min \{a_1 x_{i1}, a_2 x_{i2}\} + \varepsilon_i, \quad i = \overline{1, n}, \quad (1)$$

где n – объем выборки; $y_i, i = \overline{1, n}$ – значения объясняемой переменной y ; $x_{i1}, x_{i2}, i = \overline{1, n}$ – значения объясняющих переменных x_1 и x_2 ; a_1, a_2 – неизвестные параметры; $\varepsilon_i, i = \overline{1, n}$ – ошибки аппроксимации.

Как отмечено в [9], для непосредственного оценивания параметров модели (1) обычно применяются методы негладкой оптимизации [14], которые, как правило, являются труднореа-

лизуемыми. В [9] исследована проблема оценивания регрессии (1) с помощью метода наименьших модулей, предполагающего решение оптимизационной задачи

$$J_1 = \sum_{i=1}^n |\varepsilon_i| = \sum_{i=1}^n |y_i - \min\{a_1 x_{i1}, a_2 x_{i2}\}| \rightarrow \min.$$

В той же работе [9] установлено, что задача с функцией потерь J_1 сводится к задаче частично-булевого линейного программирования. Для её решения разработано специализированное программное обеспечение, которое находит широкое применение при моделировании реальных социально-экономических процессов и явлений (например, [15,16]).

Если регрессия (1) оценивается с помощью метода наименьших квадратов (МНК), то требуется решать оптимизационную задачу

$$J_2 = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \min\{a_1 x_{i1}, a_2 x_{i2}\})^2 \rightarrow \min.$$

В [17] предложен приближенный метод минимизации функции потерь J_2 , который подразумевает решение весьма трудоемкой вычислительной задачи. На сегодняшний день не существует специализированного программного продукта, автоматизирующего процесс её решения. Поэтому целью данной работы является разработка программы приближенного МНК-оценивания специфицированных на основе ПФ Леонтьева регрессионных моделей и её применение для моделирования уровня безработицы в Иркутской области.

Методы исследования. Алгоритм приближенного МНК-оценивания. Усложним модель (1), включив в неё свободный член a_0 :

$$y_i = a_0 + \min\{a_1 x_{i1}, a_2 x_{i2}\} + \varepsilon_i, \quad i = \overline{1, n}. \quad (2)$$

Оценка параметра a_0 показывает приближенное значение переменной y , когда $x_1 = x_2 = 0$. Таким образом, модель (2) может быть применена не только тогда, когда переменная y трактуется как выпуск продукции, а переменные x_1 и x_2 как факторы производства (в этом случае параметр a_0 справедливо был равен 0), но и при моделировании абсолютно любых по смыслу переменных.

Будем считать, что в модели (2) переменные $x_1 > 0$, $x_2 > 0$.

Очевидно, что если $a_1 > 0$, а $a_2 < 0$, то в модели (2) для любого наблюдения всегда срабатывает вторая компонента бинарной операции $\min - a_2 x_2$, а если $a_1 < 0$, а $a_2 > 0$, то, наоборот, первая компонента $- a_1 x_1$. В этих случаях оценивание регрессии (2) равносильно оцениванию обычных парных регрессий y от x_2 и y от x_1 . Поэтому будем считать, что $a_1 \cdot a_2 > 0$. Это означает, что модель (2) имеет смысл строить только тогда, когда объясняющие переменные x_1 и x_2 коррелируют с y с одинаковыми знаками.

В модели (2) вынесем параметр a_1 за знак бинарной операции \min :

$$y_i = a_0 + a_1 \min\{x_{i1}, \lambda x_{i2}\} + \varepsilon_i, \quad i = \overline{1, n}, \quad (3)$$

где $\lambda = a_2 / a_1 > 0$.

Модель (3) является частным случаем неэлементарных линейных регрессий (НЛР), подробно рассмотренных в работах [18–20].

НЛР (3) является нелинейной по параметрам. Но если придать параметру λ определенное значение, то она уже становится линейной, поэтому нетрудно получить с помощью МНК оценки параметров a_0 и a_1 .

Возникает вопрос: какое значение нужно придать параметру λ , чтобы сумма квадратов ошибок модели (3) была минимальна? При этом очевидно, что при $\lambda \rightarrow \infty$ в модели (3) всегда срабатывает первая компонента x_1 , поэтому получаем парную регрессию y от x_1 , а при $\lambda \rightarrow 0$ срабатывает вторая компонента λx_2 , поэтому получаем парную регрессию y от x_2 . В [12] доказано, что оптимальная МНК-оценка параметра λ для модели (3) лежит на отрезке

$$\lambda \in [\lambda_{\min}, \lambda_{\max}], \quad (4)$$

где $\lambda_{\min} = \min \left\{ \frac{x_{11}}{x_{12}}, \frac{x_{21}}{x_{22}}, \dots, \frac{x_{n1}}{x_{n2}} \right\}$, $\lambda_{\max} = \max \left\{ \frac{x_{11}}{x_{12}}, \frac{x_{21}}{x_{22}}, \dots, \frac{x_{n1}}{x_{n2}} \right\}$.

Тогда для приближенного МНК-оценивания параметров модели (3), можно воспользоваться следующим алгоритмом:

1. Разбить отрезок (4) точками, в каждой из которых (и на концах отрезка) найти МНК-оценки параметров a_0 и a_1 .
2. Из полученных моделей выбрать ту, для которой сумма квадратов ошибок будет минимальной.

Программа приближенного МНК-оценивания. Описанный алгоритм приближенного МНК-оценивания НЛР был реализован в виде программы с использованием среды программирования Delphi. Программа позволяет оценивать НЛР (3) как с бинарной операцией \min , так и с бинарной операцией \max .

Интерфейс программы представлен на рис. 1.

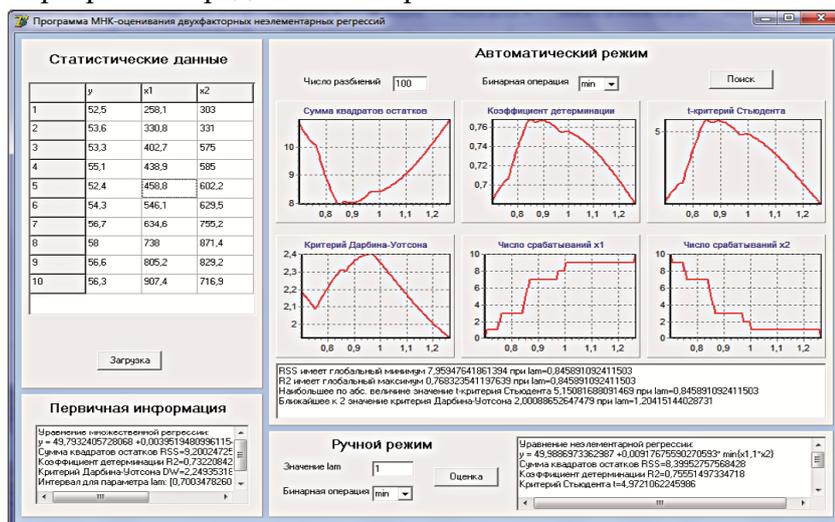


Рис. 1. Интерфейс программы
 Fig. 1. Program interface

Для работы с программой сначала нужно выбрать статистические данные (нажав кнопку «Загрузка»), которые должны храниться в текстовой файле с расширением *.txt.

В нём должны содержаться три столбца со значениями переменных, отделенные друг от друга символом табуляции. Первый столбец содержит наблюдения для переменной y . Целые и дробные части чисел отделяются друг от друга символом «,».

После выбора данных в поле «Первичная информация» автоматически выводится следующая информация: уравнение множественной линейной регрессии, сумма квадратов её остатков, коэффициент детерминации R^2 , критерий Дарбина-Уотсона, интервал для параметра λ , найденный по формуле (4). Затем нужно выбрать режим работы: ручной или автоматический. В ручном режиме нужно задать значение параметра λ , выбрать бинарную операцию (\min или \max) и нажать кнопку «Оценка». В результате в поле «Ручной режим» выводятся МНК-

для которой $\sum_{i=1}^n e_i^2 = 1931,099$, $R^2 = 0,7197$,

$$\tilde{y} = 146,72 - 3,298 \cdot 10^{-5} \max\{x_1, 39,436x_2\}, \quad (7)$$

(-6,421)

для которой $\sum_{i=1}^n e_i^2 = 2011,318$, $R^2 = 0,708$.

Как видно, НЛР (7) с бинарной операцией **max** оказалась несколько хуже линейной регрессии (5) по величине суммы квадратов остатков. При этом в ней всегда срабатывает вторая компонента бинарной операции – $39,436x_2$, из-за чего можно сделать вывод, что модель (7) есть парная линейная регрессия y от x_2 , представленная в иной математической форме. НЛР (6) с бинарной операцией **min** по величине суммы квадратов остатков оказалась лучше линейной регрессии (5) на 76,349 и лучше НЛР (7) на 80,219. При этом из-за полного отсутствия в НЛР (7) мультиколлинеарности не произошло искажения знака коэффициента при регрессоре $\min\{x_1, 28,31x_2\}$. По величине t-критерия Стьюдента этот коэффициент можно признать значимым. Представим НЛР (6) в кусочно-заданном виде:

$$\tilde{y} = \begin{cases} 146,58 - 4,633 \cdot 10^{-5} x_1, & \text{если } \frac{x_1}{x_2} \leq 28,31, \\ 146,58 - 131,16 \cdot 10^{-5} x_2, & \text{если } \frac{x_1}{x_2} > 28,31. \end{cases}$$

Тогда справедлива следующая интерпретация этой модели. Если отношение объемом ВРП x_1 к среднемесячной заработной плате x_2 не превосходит значения 28,31, то на уровень безработицы в Иркутской области влияет ВРП. При этом с ростом ВРП на 10 млрд. руб., уровень безработицы падает в среднем на 463,3 человек. А если отношение объемом ВРП x_1 к среднемесячной заработной плате x_2 больше значения 28,31, то на уровень безработицы в Иркутской области влияет заработная плата. При этом с ростом заработной платы на 1000 рублей, уровень безработицы падает в среднем на 1311,6 человек.

Вывод. В работе приводится описание алгоритма приближенного МНК-оценивания параметров двухфакторных неэлементарных линейных регрессий с бинарными операциями **min** и **max**. Получены следующие результаты.

1. На основе рассмотренного алгоритма разработана программа приближенного МНК-оценивания неэлементарных линейных регрессий с двумя переменными.

2. С помощью разработанной программы построена неэлементарная линейная регрессия уровня безработицы в Иркутской области, лишенная негативного эффекта мультиколлинеарности и оказавшаяся лучше традиционной линейной регрессии.

Библиографический список:

1. Brook R. J. Applied regression analysis and experimental design / R.J. Brook, G.C. Arnold. – CRC Press, 2018.
2. Arkes J. Regression analysis: A practical introduction / J. Arkes. – Routledge, 2019.
3. Pardoe I. Applied regression modeling / I. Pardoe. – John Wiley & Sons, 2020.
4. Boateng E. Y. A review of the logistic regression model with emphasis on medical research / E.Y. Boateng, D.A. Abaye // Journal of data analysis and information processing. – 2019. – Vol. 7. – No. 4. – pp. 190-207.
5. Parbat D. A python based support vector regression model for prediction of COVID19 cases in India / D. Parbat, M. Chakraborty // Chaos, Solitons & Fractals. – 2020. – Vol. 138. – pp. 109942.
6. Носков С.И. Дискретная модель производства алюминия в Российской Федерации / С.И. Носков // Вестник технологического университета. – 2022. – Т. 25. – № 2. – С. 80-82.
7. Werth J. Linear Regression Model Development for Analysis of Asymmetric Copper-Bisoxazoline Catalysis / J. Werth, M.S. Sigman // ACS catalysis. – 2021. – Vol. 11. – No. 7. – pp. 3916-3922.
8. Dedeturk B. K. Spam filtering using a logistic regression model trained by an artificial bee colony algorithm / B.K. Dedeturk, B. Akay // Applied Soft Computing. – 2020. – Vol. 91. – pp. 106229.

9. Носков С.И. Технология моделирования объектов с нестабильным функционированием и неопределенностью в данных / С.И. Носков. – Иркутск: Облформпечать, 1996. – 321 с.
10. Носков С.И. Построение регрессионных моделей с использованием аппарата линейно-булевого программирования / С.И. Носков, М.П. Базилевский. – Иркутск: ИрГУПС, 2018. – 176 с.
11. Клейнер Г.Б. Производственные функции: теория, методы, применение / Г.Б. Клейнер. – М.: Финансы и статистика, 1986. – 239 с.
12. Клейнер Г.Б. Экономика. Моделирование. Математика. Избранные труды. – М.: ЦЭМИ РАН, 2016. – 856 с.
13. Хацкевич Г.А. Двухфакторные производственные функции с заданной предельной нормой замещения / Г.А. Хацкевич, А.Ф. Проневич, М.В. Чайковский // Экономическая наука сегодня. – 2019. – № 10. – С. 169-181.
14. Шор Н.З. Методы минимизации недифференцируемых функций и их приложения. Киев:Наук думка,1979.200 с.
15. Носков С.И. Кусочно-линейные регрессионные модели объемов перевозки пассажиров железнодорожным транспортом / С.И. Носков, А.А. Хоняков//Модели, системы, сети в экономике, технике, природе и обществе. –2021. –№ 4 (40). – С. 80-89.
16. Носков С.И. Применение функции риска для моделирования экономических систем / С.И. Носков, А.А. Хоняков // Южно-Сибирский научный вестник. – 2020. – № 5 (33). – С. 85-92.
17. Базилевский М.П. МНК-оценивание параметров специфицированных на основе функций Леонтьева двухфакторных моделей регрессии / М.П. Базилевский // Южно-Сибирский научный вестник. 2019. № 2 (26). С. 66-70.
18. Базилевский М.П. Оценивание линейно-неэлементарных регрессионных моделей с помощью метода наименьших квадратов // Моделирование, оптимизация и информационные технологии. 2020. Т. 8. – № 4 (31).
19. Базилевский М.П. Отбор информативных операций при построении линейно-неэлементарных регрессионных моделей / М.П. Базилевский // International Journal of Open Information Technologies. 2021. Т. 9. № 5. С. 30-35.
20. Базилевский М.П. Интерпретация неэлементарных линейных регрессионных моделей / М.П. Базилевский // Информационные технологии и математическое моделирование в управлении сложными системами. – 2022. – № 1 (13). – С. 5-15.

References:

1. Brook R. J., Arnold G.C. Applied regression analysis and experimental design. CRC Press, 2018.
2. Arkes J. Regression analysis: A practical introduction. Routledge, 2019.
3. Pardoe I. Applied regression modeling. John Wiley & Sons, 2020.
4. Boateng E. Y., Abaye D.A. A review of the logistic regression model with emphasis on medical research. *Journal of data analysis and information processing*. 2019; 7(4): 190-207.
5. Parbat D., Chakraborty M. A python based support vector regression model for prediction of COVID19 cases in India. *Chaos, Solitons & Fractals*. 2020; 138: 109942.
6. Noskov S.I. Discrete model of aluminum production in the Russian Federation. [Vestnik tekhnologicheskogo universiteta] *Bulletin of the Technological University*. 2022; 25(2): 80-82. (In Russ)
7. Werth J., Sigman M.S. Linear Regression Model Development for Analysis of Asymmetric Copper-Bisoxazoline Catalysis. *ACS catalysis*. 2021; 11(7): 3916-3922.
8. Dedeturk B. K., Akay B. Spam filtering using a logistic regression model trained by an artificial bee colony algorithm. *Applied Soft Computing*. 2020; 91: 106229.
9. Noskov S.I. Technology for modeling objects with unstable operation and uncertainty in data. Irkutsk: Oblinformpechat', 1996. (In Russ)
10. Noskov S.I., Bazilevskiy M.P. Construction of regression models using the apparatus of linear-Boolean programming. Irkutsk: IrGUPS, 2018. (In Russ)
11. Kleyner G.B. Production functions: theory, methods, application. Moscow: Finansy i statistika, 1986. (In Russ)
12. Kleyner G.B. Economy. Modeling. Mathematics. Selected works. Moscow: TsEMI RAN, 2016. (In Russ)
13. Khatskevich G.A., Pronevich A.F., Chaykovskiy M.V. Two-factor production functions with given marginal rate of substitution. *Economics today*. 2019; 10: 169-181. (In Russ)
14. Shor N.Z. Methods for minimizing non-differentiable functions and their applications. Kiev: Nauk. dumka, 1979. (In Russ)
15. Noskov S.I., Khonyakov A.A. Piecewise linear regression models of passenger transportation volumes by railway. *Models, systems, networks in economics, technology, nature and society*. 2021; 4(40): 80-89. (In Russ)
16. Noskov S.I., Khonyakov A.A. Applying the risk function to model economic systems. *South Siberian Scientific Bulletin*. 2020; 5(33): 85-92. (In Russ)
17. Bazilevskiy M.P. OLS-estimation of two-factor regression models specified on Leontiev functions. *South Siberian Scientific Bulletin*. 2019; 2(26): 66-70. (In Russ)
18. Bazilevskiy M.P. Estimation linear non-elementary regression models using ordinary least squares. *Modeling, optimization and information technology*. 2020; 8(4). (In Russ)
19. Bazilevskiy M.P. Selection of informative operations in the construction of linear non-elementary regression models. *International Journal of Open Information Technologies*. 2021; 9(5): 30-35. (In Russ)
20. Bazilevskiy M.P. Interpretation of non-elementary linear regression models. *Information technology and mathematical modeling in the management of complex systems*. 2022; 1(13): 5-15. (In Russ)

Сведения об авторах:

Базилевский Михаил Павлович, кандидат технических наук, доцент, доцент кафедры математики; mik2178@yandex.ru
Карбушева Диана Витальевна, студент; karbusheva.02@mail.ru

Information about authors:

Mikhail P. Bazilevskiy, Cand. Sci. (Eng.), Assoc. Prof., Department of Mathematics; mik2178@yandex.ru
Diana V. Karbusheva, student; karbusheva.02@mail.ru

Конфликт интересов / Conflict of interest

Авторы заявляют об отсутствии конфликта интересов / The authors declare no conflict of interest

Поступила в редакцию/Received 20.08.2022.

Одобрена после рецензирования/ Revised 17.09.2022.

Принята в печать/Accepted for publication 17.09.2022.