



МАКСИМИЗАЦИЯ ЧИСЛА ДОПУСТИМЫХ ОШИБОК АППРОКСИМАЦИИ ПРИ ПОСТРОЕНИИ ЛИНЕЙНОЙ РЕГРЕССИОННОЙ МОДЕЛИ

Носков Сергей Иванович

доктор технических наук, профессор,
профессор кафедры «Информационные системы
и защита информации»,
Иркутский государственный университет
путей сообщения,
Иркутск, Россия
E-mail: sergey.noskov.57@mail.ru

Шахуров Антон Николаевич

студент кафедры «Информационные системы
и защита информации»,
Иркутский государственный университет
путей сообщения,
Иркутск, Россия
E-mail: Fantom3920@yandex.ru

Предмет исследования: ошибки аппроксимации
линейной регрессионной модели в рамках логико-алге-
браического подхода к анализу данных.

Цель исследования: разработать алгоритмический
способ решения задачи максимизации числа допусти-
мых ошибок аппроксимации с применением вычисли-
тельного аппарата линейно-булева программирования.

Методы и объекты исследования: объектом иссле-
дования является линейная регрессионная модель, ме-
тодами – линейный регрессионный анализ и аппарат ма-
тематического программирования.

Основные результаты исследования: предложен
алгоритмический способ максимизации числа допусти-
мых абсолютных и относительных ошибок аппроксима-
ции линейного регрессионного уравнения, сводящийся
к решению задач линейно-булева программирования
приемлемой для практических ситуаций размерности.
Решение сформированных задач этого типа не должно
вызывать вычислительных проблем в силу значительного
числа эффективных программных средств, например
размещенной в Интернете в свободном доступе про-
граммы LPsolve.

Ключевые слова: регрессионная модель, абсолют-
ные и относительные ошибки аппроксимации, задача
линейно-булева программирования, модель пассажиро-
оборота воздушного транспорта.

MAXIMIZING THE NUMBER OF ALLOWABLE APPROXIMATION ERRORS WHEN BUILDING A LINEAR REGRESSION MODEL

Sergey I. Noskov

Doctor of Technical Sciences, Professor,
Professor of the Department of Information Systems
and Information Security,
Irkutsk State Transport University,
Irkutsk, Russia
E-mail: sergey.noskov.57@mail.ru

Anton N. Shakhurov

student of the Department of Information Systems
and Information Security,
Irkutsk State Transport University,
Irkutsk, Russia
E-mail: Fantom3920@yandex.ru

Subject of research: errors in approximation of a linear
regression model within the framework of a logical-algebraic
approach to data analysis.

Purpose of the study: to develop an algorithmic
method for solving the problem of maximizing the number
of permissible approximation errors using a linear-Boolean
programming computer.

Methods and objects of research: the object of
research is a linear regression model, the methods are
linear regression analysis and mathematical programming
apparatus.

Main results of the study: an algorithmic method
is proposed for maximizing the number of permissible
absolute and relative errors in approximation of a linear
regression equation, which reduces to solving linear-
Boolean programming problems of a dimension acceptable
for practical situations. Solving generated problems of this
type should not cause computational problems due to a
significant number of effective software tools, for example,
the LPsolve program, which is freely available on the
Internet.

Keywords: regression model, absolute and relative
approximation errors, linear-Boolean programming problem,
air transport passenger turnover model.

ВВЕДЕНИЕ

Методы математического моделирования
являются эффективным средством исследо-
вания сложных, с множеством межфакторных
структурных взаимодействий, систем раз-
личного характера и масштаба. Эти методы
позволяют формализовывать свойственные
таким системам закономерности функциони-
рования и развития путем создания их каче-
ственных абстрактных аналогов, что открыв-
вает широкие перспективы в существенном
повышении действенности вырабатываемых
управляющих сигналов.

Регрессионные модели – весьма широ-
кий класс математических моделей, разра-

батываемых для исследования сложных объ-
ектов любой природы. Рассмотрим необходи-
мый составной элемент практически любой
модели этого типа – регрессионное уравне-
ние (зависимость) вида:

$$y_k = \sum_{i=1}^m \alpha_i x_{ki} + \varepsilon_k, k = \overline{1, n}, \quad (1)$$

где y – зависимая, а x_i – i -я независимая
переменная, α_i – i -й подлежащий оцениванию
параметр, ε_k – ошибки аппроксимации, k – но-
мер наблюдения, n – их число (длина выбор-
ки данных). Будем считать все переменные и
ошибки уравнения (1) детерминированными.

Представим уравнение (1) в векторной
форме:



$$y = X\alpha + \varepsilon,$$

где $y = (y_1, \dots, y_n)^T$, $\alpha = (\alpha_1, \dots, \alpha_m)^T$, $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T$, $X - (n \times m)$ – матрица с компонентами x_{ki} . При наличии в уравнении (1) свободного члена первый столбец матрицы X состоит из единиц.

Построение и использование регрессионных моделей в прикладных целях часто сопровождается анализом свойств ошибок аппроксимации. Так, в работе [1] представлена гладкая функция для аппроксимации функции контрольных потерь с тем, чтобы можно было использовать методы оптимизации на основе градиента для подбора модели квантильной регрессии. Обсуждаются свойства гладкого приближения, предложены два алгоритма минимизации сглаженной целевой функции. В [2] рассматривается задача выбора и оценки модели квантильной регрессии с известной структурой групп в предикторах. Для медианного случая модель оценивается путем минимизации штрафной целевой функции потерь (ошибок аппроксимации) Хубера. Статья [3] посвящена анализу существования двух непротиворечивых оценок параметров линейных предикторов в регрессии Пуассона, где ковариата измеряется с ошибками. В [4] рассматривается проблема выбора переменных в квантильной регрессии с авторегрессионными ошибками аппроксимации.

Все наиболее часто используемые в регрессионном анализе критерии адекватности моделей, в частности множественной детерминации, Фишера, Стьюдента, Дарбина – Уотсона, включают в свои расчетные формулы ошибки аппроксимации и отражают те или иные частные характеристики в качестве модельного описания сложных объектов. Так, в работе [5] критерий множественной детерминации использовался при разработке регрессионной модели индекса нормализации разницы растительности лугов и пахотных земель вдоль нарушенной полосы отвода трубопроводов Баку – Тбилиси – Джейхан и Южно-Кавказского трубопровода для целей планирования восстановления растительности. В качестве климатических факторов были выделены годовое количество осадков, годовое суммарное испарение, температура поверхности Земли, годовая минимальная и максимальная температура воздуха и солнечная радиация. Учитывались также грунтовые факторы: высота, ракурс, грунтовые воды и глубина верхнего слоя почвы. В [6] отмечается, что регрессия Стьюдента является полезным расширением нормальной модели, которую можно использовать для статистического моделирования наборов данных, включающих ошибки с тяжелыми хвостами и/или

выбросами. Обсуждается также регрессия Стьюдента с переменной дисперсией, в которой как среднее значение, так и дисперсия зависят от объясняющих переменных. Проблема, представляющая интерес, заключается в одновременном выборе значимых переменных как в модели среднего значения, так и в модели дисперсии. Описана унифицированная процедура, позволяющая выделять значимую переменную. В работе [7] предложен алгоритм применения критерия Дарбина – Уотсона для анализа автокорреляции ошибок аппроксимации. В исследовании [8] этот критерий использован при моделировании качества воздуха. В качестве независимых переменных при этом использованы концентрация твердых частиц и метеорологические параметры (температура, влажность, скорость и направление ветра) за 5 лет с трех промышленных станций мониторинга качества воздуха в Малайзии. В работе [9] предлагаются простые, основанные на критерии Дарбина – Уотсона тесты для проверки корреляции рядов, которые применимы в моделях линейной регрессии. Процедуры тестирования устойчивы при различных распределениях случайных ошибок. Асимптотические распределения предложенной статистики получены с помощью совместной центральной предельной теоремы для нескольких общих квадратичных форм и дельта-метода. В [10] с помощью критерия Фишера проведен регрессионный анализ загрязнения тяжелыми металлами почв и донных отложений прудов, находящихся в зоне с повышенной автотранспортной нагрузкой.

В работе [11] предложен метод оценивания параметров модели (1) путем минимизации средней и максимальной относительных ошибок аппроксимации.

РЕЗУЛЬТАТЫ И ОБСУЖДЕНИЕ

Наиболее популярными методами построения регрессионного уравнения (1) принято считать методы наименьших квадратов (МНК) и модулей (МНМ) (см., например, [12, 13]). Первый из них позволяет аналитически представить формулу для расчета вектора параметров:

$$\alpha = (X^T X)^{-1} X^T y.$$

Для применения МНМ следует либо воспользоваться методом вариационно-взвешенных квадратичных приближений [14], основанным на использовании специальным образом организованной итерационной процедуры, либо путем [13] сведения задачи

$$E = \sum_{k=1}^n |\varepsilon_k| \rightarrow \min \quad (2)$$

к задаче линейного программирования (ЛП) следующим образом.

Введем в рассмотрение переменные $u_k \geq 0$ (положительные части ошибок ε_k) и $v_k \geq 0$ (отрицательные части ошибок ε_k), $k = \overline{1, n}$ следующим образом:

$$u_k = \begin{cases} \varepsilon_k, & \varepsilon_k > 0 \\ 0, & \text{в противном случае,} \end{cases}$$

$$v_k = \begin{cases} -\varepsilon_k, & \varepsilon_k < 0 \\ 0, & \text{в противном случае.} \end{cases}$$

Очевидно, что при этом справедливы следующие равенства:

$$\varepsilon_k = u_k - v_k, |\varepsilon_k| = u_k + v_k, u_k v_k = 0.$$

Представим уравнение (1) в виде системы тождеств:

$$\sum_{i=1}^m \alpha_i x_{ki} + u_k - v_k = y_k, k = \overline{1, n}, \quad (3)$$

$$u_k \geq 0, v_k \geq 0, k = \overline{1, n}, \quad (4)$$

Тогда задача (2) эквивалентна задаче ЛП с ограничениями (3), (4) и целевой функцией

$$\sum_{k=1}^n (u_k + v_k) \rightarrow \min. \quad (5)$$

Предложенный в [11] метод минимизации средней относительной ошибки аппроксимации $\tilde{E} = \sum_{k=1}^n \frac{|\varepsilon_k|}{y_k} / n$ сводится к задаче ЛП с ограничениями (3), (4) и целевой функцией

$$\sum_{k=1}^n \left| \frac{1}{y_k} \right| (u_k + v_k) \rightarrow \min. \quad (6)$$

Там же предложен способ минимизации максимальной относительной ошибки аппроксимации

$$r = \max_{k=1, n} \left| \frac{\varepsilon_k}{y_k} \right| \rightarrow \min.$$

Для этого система ограничений (3), (4) дополняется неравенствами:

$$u_k + v_k - |y_k| r \leq 0, k = \overline{1, n}, \quad (7)$$

а целевая функция (6) заменяется на следующую:

$$r \rightarrow \min. \quad (8)$$

Решение трех сформированных задач ЛП (3) – (5), (3), (4), (6) и (3), (4), (7), (8) не вызывает вычислительных проблем в силу значительного числа эффективных программных средств, например размещенной в Интернете в свободном доступе программы LPsolve.

Поставим задачу максимизации числа наблюдений выборки, для которых модуль ошибки аппроксимации не превышает некоторой наперед заданной величины

d (тракуемой как допустимая абсолютная ошибка), или, формально:

$$|S| \rightarrow \max, \quad (9)$$

где

$$S = \{ k \in \{1, 2, \dots, n\} \mid |\varepsilon_k| \leq d \},$$

а через $|S|$ обозначено число элементов (мощность) множества S . Таким образом, S – множество номеров наблюдений выборки, абсолютные ошибки аппроксимации для которых допустимы.

Допустимым ошибкам можно придать и относительный характер путем постановки задачи

$$|\tilde{S}| \rightarrow \max, \quad (10)$$

где

$$\tilde{S} = \{ k \in \{1, 2, \dots, n\} \mid |\varepsilon_k| \leq \tilde{d} y_k \},$$

а \tilde{d} представляет собой долю фактических значений зависимой переменной, выраженную в процентах (например, $\tilde{d} = 0.03$ соответствует 3 %).

Введем в рассмотрение булевы переменные $\sigma_k, k = \overline{1, n}$ по правилу:

$$\sigma_k = \begin{cases} 1, & |\varepsilon_k| \leq d \\ 0, & \text{в противном случае,} \end{cases}$$

а также ограничения

$$u_k + v_k + M \sigma_k \leq M + d, \quad (11)$$

где M – наперед заданная большая положительная константа.

Тогда задача (9) эквивалентна задаче линейно-булева программирования (ЛБП) с ограничениями (3), (4), (11),

$$\sigma_k \in \{0, 1\}, k = \overline{1, n} \quad (12)$$

и целевой функцией

$$\sum_{k=1}^n \sigma_k - \delta \sum_{k=1}^n (u_k + v_k) \rightarrow \max, \quad (13)$$

где δ – наперед заданное малое положительное число, сравнимое с нулем. Присутствие в (13) второго слагаемого гарантирует выполнение условия $u_k v_k = 0$ для всех k , следующего из приведенного выше определения переменных u_k и v_k .

При решении задачи (10) необходимо ограничение (11) в задаче ЛБП (3) (4), (11) – (13) заменить на следующее:

$$u_k + v_k + M \sigma_k \leq M + \tilde{d} y_k. \quad (14)$$

Применим описанный способ идентификации параметров линейной регрессии для моделирования пассажирооборота воздушного транспорта Российской Федерации. Введем следующие обозначения:

y – пассажирооборот воздушного транспорта, млрд пасс. км;

x_1 – среднемесячная номинальная начисленная заработная плата работников организаций, руб.;

x_2 – численность трудоспособного населения в России, млн чел.

В качестве информационной базы для моделирования используем статистическую ежегодную информацию за 2002–2019 гг. [15]. С помощью МНК, МНМ и изложенного выше подхода будем строить линейную двухфакторную модель без свободного члена:

$$y_k = \alpha_1 x_{k1} + \alpha_2 x_{k2} + \varepsilon_k, k = \overline{1, 18}.$$

В результате получим:

- МНК

$$y = 0.00586x_1 + 0.36409x_2, \quad (15)$$

$$E = 173.1, \quad \tilde{E} = 0.29 \%,$$

- МНМ

$$y = 0.00577x_1 + 0.40365x_2, \quad (16)$$

$$E = 172.3, \quad \tilde{E} = 0.41 \%.$$

Значения критериев адекватности E и \tilde{E} указывают на высокую адекватность моделей (15) и (16).

В таблицах 1 и 2 отражены мощности множеств S и \tilde{S} для моделей (15) и (16).

Таблица 1. Мощность множеств S для моделей (15) и (16)

d	$ S $	
	Модель (15)	Модель (16)
1	2	2
2	4	3
3	6	3
4	7	9
5	7	9
6	7	9
7	9	9
8	10	11
9	11	11
10	11	12

Таблица 2. Мощность множеств \tilde{S} для моделей (15) и (16)

\tilde{d}	$ \tilde{S} $	
	Модель (15)	Модель (16)
0.01	2	3
0.02	6	5
0.03	6	5
0.04	8	7
0.05	9	10
0.06	10	11
0.07	10	12
0.08	10	13
0.09	13	14
0.1	14	15

Из анализа таблицы 1 следует, что при $d \leq 3$ число $|S|$ для модели (15) не меньше, чем для модели (16), при $d > 3$ – не больше.

Из таблицы же 2 следует, что такой границей является число $\tilde{d} = 0.04$.



В таблице 3 представлены результаты решения задачи ЛБП (3), (4), (11) – (13) по максимизации числа абсолютных ошибок аппроксимации для различных значений d .

Отметим, что значение $d=4$ соответствует модели (16).

В таблице 4 помещены результаты решения задачи ЛБП (3), (4), (14), (12), (13) по максимизации числа относительных ошибок аппроксимации для различных значений \tilde{d} .

Таблица 3. Результаты решения задачи (3), (4), (11) – (13) по максимизации числа абсолютных ошибок аппроксимации для различных значений d

d	α_1	α_2	$ S $	E	S
1	0.00449	0.51449	5	379.7	1,2,4,5,7
2	0.00480	0.49584	6	299.8	1,2,4,5,6,15
3	0.00481	0.47436	7	315.1	1,2,4,5,6,7,15
4	0.00577	0.40365	9	172.3	1,2,4,5,6,10,14,16,17
5	0.00474	0.51125	9	308.7	1,2,3,4,5,6,7,9,15
6	0.00576	0.42403	11	175.6	1,2,3,4,5,6,10,11,14,16,17
7	0.00576	0.41295	11	173.5	1,2,3,4,5,6,10,11,14,16,17
8	0.00590	0.39136	12	179.7	1,2,3,4,5,6,10,11,14,16,17,18
9	0.00579	0.39540	12	172.5	1,2,3,4,5,6,9,10,11,14,16,17
10	0.00597	0.32857	13	178.1	1,2,4,5,6,7,9,10,11,14,16,17,18

Таблица 4. Результаты решения задачи (3), (4), (14), (12), (13) по максимизации числа относительных ошибок аппроксимации для различных значений \tilde{d}

\tilde{d}	α_1	α_2	$ \tilde{S} $	E	S
0.01	0.00587	0.35011	6	174.8	5,6,10,14,16,17
0.02	0.00489	0.48631	7	285.8	1,2,4,5,6,9,15
0.03	0.00592	0.37499	9	178.4	4,5,6,10,11,14,16,17,18
0.04	0.00551	0.42849	10	195.6	1,2,4,5,6,9,10,14,16,17
0.05	0.00577	0.40847	11	173.0	1,2,4,5,6,10,11,14,16,17,18
0.06	0.00599	0.39015	13	189.0	1,2,4,5,6,10,11,12,13,14,16,17,18
0.07	0.00583	0.42082	14	181.0	1,2,3,4,5,6,10,11,12,13,14,16,17,18
0.08	0.00576	0.41629	15	174.4	1,2,3,4,5,6,9,10,11,12,13,14,16,17,18
0.09	0.00541	0.43853	15	208.5	1,2,3,4,5,6,7,9,10,11,14,15,16,17,18
0.1	0.00573	0.40035	16	175.9	1,2,3,4,5,6,7,9,10,11,12,13,14,16,17,18

Анализ таблиц 3, 4 позволяет сделать следующие выводы.

Поведение оценок параметров α_1 и α_2 , а также суммы модулей ошибок E не является монотонным при росте значений d и \tilde{d} . При этом мощность множеств S и \tilde{S} естественным образом возрастает.

Попадание какого-либо номера наблюдения выборки в состав множеств S или \tilde{S} для меньших значений d или \tilde{d} отнюдь не гарантирует этого для больших значений.

Разумеется, окончательный выбор конкретного значения d или \tilde{d} и, соответственно, оценок параметров приведенной выше двухфакторной модели пассажирооборота воздушного транспорта Российской Федерации остается за исследователем в зависимости от целей моделирования и его индивидуальных предпочтений относительно значений d , \tilde{d} , $|S|$, $|\tilde{S}|$.

ЗАКЛЮЧЕНИЕ И ВЫВОДЫ

В работе предложен алгоритмический способ максимизации числа допустимых абсолютных и относительных ошибок аппроксимации линейного регрессионного уравнения, сводящийся к решению задач линейно-булева программирования приемлемой для практических ситуаций размерности. Решение сформированных задач ЛП и ЛБП не должно вызывать вычислительных проблем в силу значительного числа эффективных программных средств, например размещенной в Интернете в свободном доступе программы LPsolve.

Применение МНК и МНМ позволяет построить весьма адекватные модели пассажирооборота воздушного транспорта Российской Федерации. Поведение оценок параметров модели, а также суммы модулей ошибок не является монотонным при росте уровней допустимых значений абсолютных и относительных ошибок аппроксимации. При этом мощность множеств числа допустимых ошибок естественным образом возрастает. Кроме того, попадание какого-либо номера наблюдения выборки в состав множеств номеров допустимых ошибок для их меньших значений отнюдь не гарантирует этого для больших значений.

СПИСОК ЛИТЕРАТУРЫ

- Zheng, S. Gradient descent algorithms for quantile regression with smooth approximation / S. Zheng. – DOI 10.1007/s13042-011-0031-2 // International Journal of Machine Learning and Cybernetics. – 2011. – № 2. – P. 191–207
- Sherwood, B. Quantile regression feature selection and estimation with grouped variables using Huber approximation / B. Sherwood, S. Li. – DOI 10.1007/s11222-022-10135-w // Statistics and Computing. – 2022. – № 32 (75). – URL: <https://link.springer.com/article/10.1007/s11222-022-10135-w> (date of application: 24.09.2024).
- Kukush, A. Three estimators for the poisson regression model with measurement errors / A. Kukush, H. Schneeweis, R. Wolf. – DOI 10.1007/bf02777577 // Statistical Papers. – 2004. – № 45. – P. 351–368.
- Variable selection in quantile regression when the models have autoregressive errors / W. Zhao, R. Zhang, Y. Lv, J. Liu. – DOI 10.1016/j.jkss.2012.11.003 // Oh Journal of the Korean Statistical Society. – 2014. – № 43 (3). – P. 513–530.
- Bayramov, E. Determination of main climate and ground factors controlling vegetation cover regrowth along oil and gas pipelines using multiple, spatial and geographically weighted regression procedures / E. Bayramov, M. F. Buchroithner, E. McGurty. – DOI 10.1007/s12665-011-1429-6 // Environmental Earth Sciences. – 2012. – № 66. – P. 2047–2062.
- Zhao, W. Variable selection of varying dispersion student-t regression models / W. Zhao, R. Zhang. – DOI 10.1007/s11424-014-2223-9 // Journal of Systems Science and Complexity. – 2015. – Vol. 28. – P. 961–977.
- Георгиев, Н. С. Проверка автокорреляции в остатках критерием Дарбина – Уотсона / Н. С. Георгиев, А. Д. Юрченко. – Текст : непосредственный // Аллея науки. – 2020. – № 2 (6). – С. 182–185.
- Application of Step Wise Regression Analysis in Predicting Future Particulate Matter Concentration Episode / A. Nazif, N. I. Mohammed, A. Malakahmad, M. S. Abualqumboz. – DOI 10.1007/s11270-016-2823-1 // Water, Air, & Soil Pollution. – 2016. – № 227 (117). – URL: <https://link.springer.com/article/10.1007/s11270-016-2823-1> (date of application: 24.09.2024).
- Yin, Y. Model-free tests for series correlation in multivariate linear regression / Y. Yin. – DOI 10.48550/arXiv.1901.05595 // Journal of Statistical Planning and Inference. – 2020. – № 206. – P. 179–195.
- Сорокин, А. В. Взаимосвязь накопления тяжелых металлов в донных отложениях и почве при проведении оценки загрязненности рекреационных зон автотранспортом / А. В. Сорокин, Е. В. Сотникова. – Текст : непосредственный // Технические науки – от теории к практике. – 2014. – № 36. – С. 145–151.
- Носков, С. И. Минимизация средней и максимальной относительных ошибок аппроксимации регрессионной модели / С. И. Носков. – Текст : непосредственный // Известия Тульского государственного университета. Технические науки. – 2023. – № 1. – С. 340–343.
- Дрейпер, Н. Прикладной регрессионный анализ / Н. Дрейпер, Г. Смит. – 3-е издание. – Москва : Вильямс, 2016. – 912 с. – Текст : непосредственный.
- Демиденко, Е. З. Линейная и нелинейная регрессии / Е. З. Демиденко. – Москва : Финансы и статистика, 1981. – 302 с. – Текст : непосредственный.
- Мудров, В. И. Методы обработки измерений. Квази-правдоподобные оценки / В. И. Мудров, В. А. Кушко. – 2-е издание. – Москва : Радио и связь, 1983. – 304 с. – Текст : непосредственный.
- Носков, С. И. Разработка регрессионной модели пассажирооборота воздушного транспорта Российской Федерации двумя альтернативными методами / С. И. Носков, Ю. А. Бычков, К. С. Перфильева. – Текст : электронный // Вестник кибернетики. – 2023. – № 22 (1). С. 36–42. – URL: <https://www.vestcyber.ru/jour/article/view/502> (дата обращения: 03.12.2023).