

УДК 519.254

ПОВЫШЕНИЕ ЭФФЕКТИВНОСТИ ПРОГНОЗИРОВАНИЯ ПРИ ИСПОЛЬЗОВАНИИ ЛИНЕЙНЫХ РЕГРЕССИОННЫХ МОДЕЛЕЙ

О.В.Рычка

Донецкий национальный технический университет

В работе рассматриваются сущность новых методов повышения эффективности прогнозирования, основные достоинства данных методов, в сравнении с уже известными. Описаны этапы работы программы, разработанной автором для реализации предложенных методов.

В настоящее время в различных сферах человеческой деятельности для прогнозирования часто используются регрессионные прогнозные модели. Для качественного прогноза необходимо, чтобы статистические данные были надежными и достоверными. Поэтому существует необходимость в обработке данных с целью выявления значений, которые значительно отличаются от остальных. Эти значения представляют собой аномальные измерения. Для обнаружения таких измерений существует ряд специальных критериев [1-3]. В ходе исследования были выявлены существенные недостатки данных методов [4]. Их основным недостатком является, то, что на "аномальность" исследуется по одному значению, что при большом количестве исходных статистических данных увеличивает трудоемкость использования данных критериев. В связи с выявленными недостатками существующих методов в [4] и [5] были предложены новые методы повышения качества регрессионных прогнозных моделей.

Суть данных методов заключается в том, что находятся данные, которые не попадают в прямоугольную область со сторонами $2k \cdot \sigma_e$ и $2k\sigma'_e$, где k – коэффициент, соответствующий вероятности попадания в заданную область. Данная вероятность рассчитывается по формуле (1):

$$P_0 = 2 \cdot \frac{1}{\sqrt{2\pi}} \int_0^k e^{-t^2/2} dt - 1 \quad (1)$$

Среднеквадратические отклонения невязок σ_e и σ'_e , определяются по формулам (2) и (3):

$$\sigma_e = \sqrt{\frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n-2}}, \quad (2)$$

$$\sigma'_e = \sqrt{\frac{\sum_{i=1}^n (Y_i - \hat{Y}'_i)^2}{n-2}}, \quad (3)$$

где Y_i – фактичні значення;

\hat{Y}_i – розраховані значення по вихідному рівнянню;

\hat{Y}'_i – розраховані значення по рівнянню перпендикуляра.

Основне відміння запропонованих методів між собою заключається в тому, що в одному методі дані, не попадаючі в прямокутну область відбрасуються (рис.1), а в другому – коректуються (рис.2).

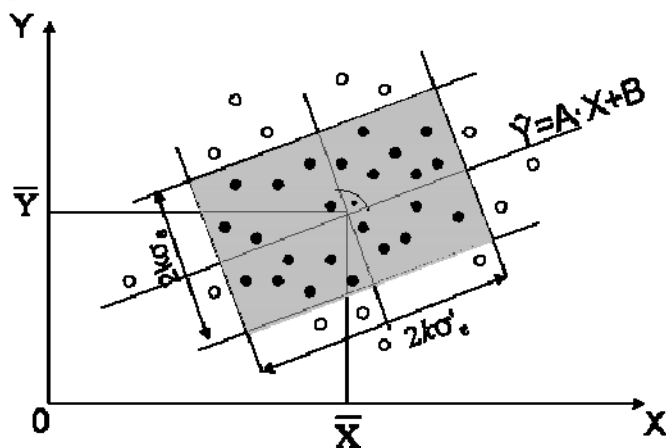


Рисунок 1 – Метод, оснований на відбрасуванні даних

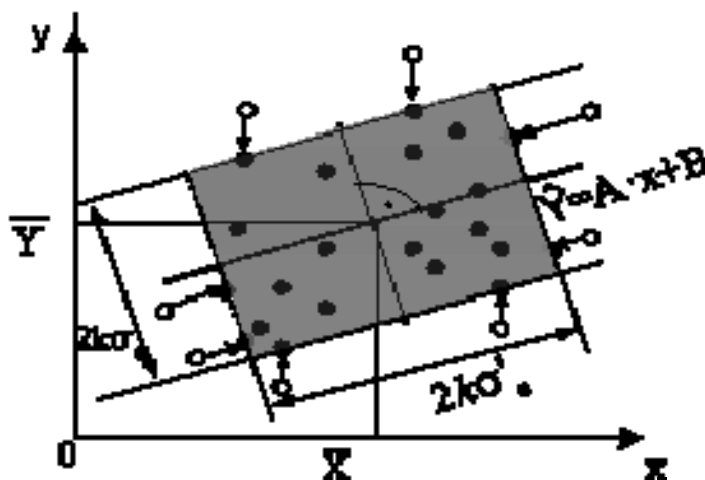


Рисунок 2 – Метод, оснований на переносі даних

Критерии оценки эффективности использования, описанных выше методов:

1. Коэффициент детерминации R^2 , который демонстрирует силу взаимосвязи между Y и X (процент случайных величин x_i , который полностью объясняет поведение случайных величин y_i данным линейным регрессионным уравнением);

2. Модуль величины смещения результата прогноза Δ (является следствием изменения положения нового регрессионного уравнения при отбрасывании части исходных измерений);

3. Доверительный интервал прогнозных значений $Y_{\text{прогн}}$ (представляет собой геометрическое место расположения прогнозных значений $Y_{\text{прогн}}$ при заданном значении $X_{\text{прогн}}$ и заданной доверительной вероятности $P_{\text{дов}}$);

4. Количество элементарных операций ЭВМ, которое необходимо для реализации отбрасывания аномальных и ненадежных измерений;

5. Точность, которая рассчитывается по формуле 4:

$$T = R^2 \cdot \frac{m}{n} \quad (4)$$

где n – исходное количество данных;

m – количество оставшихся после отбрасывания данных или количество данных, которые не подверглись изменению.

При этом наилучшим вариантом, считается вариант, при котором величина коэффициента детерминации R^2 является максимальной, при обязательном условии, что $T \geq 0.5$.

Для удобства реализации предложенных методов автором была разработана программа с использованием языка программирования Visual Basic for Applications. Основные этапы работы программы заключаются в следующем:

1. Пользователь вводит исходные данные.

2. По нажатию на кнопку, программа работает следующим образом:

а) строится вариационный ряд (по независимой переменной X);

б) находятся коэффициенты линейного регрессионного уравнения;

в) рассчитывается исходное значение коэффициента детерминации R^2 ;

г) определяется исходная величина доверительного интервала;

д) в программе предусмотрено определенное количество рабочих листов для каждой вероятности P попадания исходных статистических данных в заданную область (начиная с $P=1$ и до $P=0,5$

с шагом 0,05). На каждом из этих листов, отображаются результаты произведенных расчетов (выводятся оставшиеся или модифицированные данные, в зависимости от используемого метода, определяются новые коэффициенты уравнения, коэффициенты детерминации R^2 , величины доверительного интервала, величины смещения и величина точности);

е) формируется итоговая таблица с рассчитанными значениями всех предложенных критериев эффективности (рис.3).

	A	B	C	D	E	F	G
1		R ²	DI,%	Delta,%	Количество точек	Точность	
2	100	0,709162	20,42064		24		
3	90	0,817287	13,00199	1,814664	21		0,715126
4	85	0,814515	13,29402	2,263306	20		0,678763
5	80	0,762273	13,22544	2,355109	19		0,603466
6	75	0,838699	11,42163	1,225086	18		0,629025
7	70	0,838699	11,42163	1,225086	18		0,629025
8	65	0,871067	8,831413	0,536748	17		0,617006
9	60	0,842411	8,571503	0,173844	15		0,526507
10	50	0,780952	7,858299	0,167475	13		0,423016

Рисунок 3 – Вид итоговой таблицы

3. В программе предусмотрено отображение графиков, изображающих соответствующие области с выделением данных, которые не попадают в определенную область (рис.4).

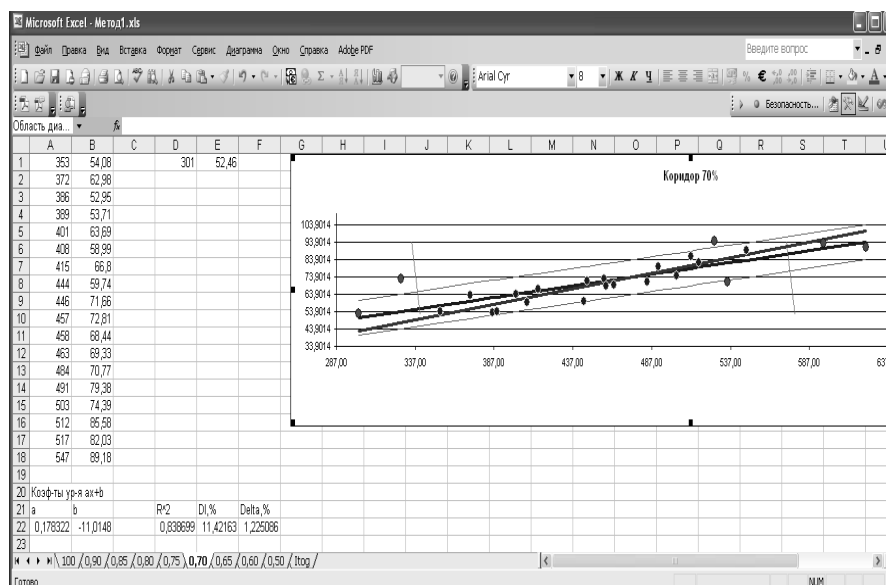


Рисунок 4 – Вид рабочего листа "0,70" после нажатия кнопки "Изобразить график"

4. По окончании работы программы рабочие листы очищаются, нажатием кнопки "Delete".

Таким образом, в данной работе были описаны методы, позволяющие повысить точность прогнозирования при использовании линейных регрессионных моделей. Основными достоинствами данных методов является:

- простота понимания и применения;
- возможность найти и обработать несколько подозрительных значений одновременно, что позволяет избежать простого перебора данных;
- хорошая формализуемость, что позволило реализовать данные методы в компьютерных технологиях.

Разработанная программа позволяет:

- сократить время реализации методов;
- определить выигрыш от использования методов;
- наглядно изобразить результаты работы методов.

Список літератури

1. Кобзарь А.И. Прикладная математическая статистика. Для инженеров и научных работников. – М.: ФИЗМАТЛИТ, 2006. – 816 с.
2. Дрейпер Н.Р., Смит Г. Прикладной регрессионный анализ. 3-е изд.: Пер. с англ. – М.: Вильямс, 2007. – 912 с.
3. Rawlings, John O. Applied regression analysis: a research tool. — 2nd ed. / John O. Rawlings, Sastry G. Pentula, David A. Dickey – USA.: Springer, 1998.
4. Смирнов А.В., Рычка О.В. Метод повышения качества прогнозных регрессионных моделей. // Наукові праці Донецького національного технічного університету. Серія "Інформатика, кібернетика та обчислювальна техніка". Випуск 12(165) – Донецьк: ДВНЗ "ДонНТУ". – 2010. – С.141-147.
5. Смирнов А.В., Рычка О.В. Новый метод улучшения качества прогнозных регрессионных моделей. // Наукові праці ДонНТУ Серія "Інформатика, кібернетика та обчислювальна техніка". Випуск 13(185) – Донецьк: ДВНЗ "ДонНТУ". – 2011. – С.168-172.

Отримано 10.07.2011