

УДК 519.233.2:519.254

**О. М. Мацуга, А. С. Черкашина***Дніпропетровський національний університет імені Олеся Гончара*

## ІНФОРМАЦІЙНА ТЕХНОЛОГІЯ РОБАСТНОГО ВІДНОВЛЕННЯ РОЗПОДІЛІВ

Розроблено інформаційну технологію робастного відновлення розподілів, яку реалізовано у вигляді програмного забезпечення «RobustProcedre». Результати тестування програми на даних імітаційного моделювання підтвердили адекватність технології.

**Ключові слова:** *інформаційна технологія, робастне оцінювання, відновлення розподілу, сплайн-розподіл.*

Разработана информационная технология робастного восстановления распределений, которую реализовано в виде программного обеспечения «RobustProcedre». Результаты тестирования программы на данных имитационного моделирования подтвердили адекватность технологии.

**Ключевые слова:** *информационная технология, робастное оценивание, восстановление распределения, сплайн-распределение.*

The information technology of distribution robust restoration is developed. It was realized as software «RobustProcedre». Software testing results on modeling data corroborate the technology.

**Key words:** *information technology, robust estimate, distribution restoration, spline-distribution.*

**Постановка проблеми.** Під час обробки і аналізу вибірових даних значної уваги потребують аномальні спостереження, значення яких різко відрізняються від інших. Такі спостереження можуть суттєво впливати на результати статистичного дослідження, зокрема результати відновлення розподілу ймовірностей. Для того щоб зменшити їх вплив, доцільно застосовувати робастні методи математичної статистики. Такі методи добре вивчені, проте їх програмна реалізація у автоматизованих системах обробки даних поки що недостатня.

**Аналіз останніх досліджень і публікацій.** Вперше теоретичний підхід до проблеми робастності в статистиці був запропонований Хьюбером [1], який ввів так звані « $M$ -оцінки», що є узагальненням оцінок максимальної правдоподібності. Суттєво інший підхід, заснований на функціях впливу, був запропонований Хампелем [2]. Стосовно задачі відновлення розподілів  $M$ -оцінки більш гнучкі, оскільки допускають пряме узагальнення на багатопараметричний випадок.

**Постановка задачі.** Нехай результати спостережень задано у вигляді вибірки  $\{x_i; i = \overline{1, N}\}$ , де  $N$  – кількість спостережень;  $x_i$  – спо-

стережуване значення в  $i$ -му експерименті. Припускається, що вибірка може містити аномальні значення. За вибірковими даними потрібно відновити розподіл ймовірностей  $F(x; \bar{\Theta})$ , де  $\bar{\Theta} = \{\theta_1, \theta_2, \dots, \theta_p\}$  – вектор параметрів розподілу, і знайти оцінку  $\hat{\bar{\Theta}} = \{\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_p\}$ , стійку відносно можливих відхилень статистичного розподілу від теоретичного.

Ставиться задача розробити інформаційну технологію робастного відновлення розподілів, яку реалізувати у вигляді сучасного програмного забезпечення. В якості моделей розподілів розглянути «чисті» та сплайн-розподіли з класу нормального, логарифмічно нормального, експоненціального та Вейбулла [3, 4].

**Основний матеріал.** Обчислювальна технологія робастного відновлення розподілів містить такі етапи:

1. Проведення первинного статистичного аналізу вибіркових даних, який передбачає побудову варіаційного ряду, розбиття ряду на класи, побудову гістограми та емпіричної функції розподілу, обчислення статистичних характеристик вибірки [3].

2. Ідентифікація розподілу з класу нормального, логарифмічно нормального, експоненціального та Вейбулла. Ідентифікацію пропонується проводити за допомогою процедури автоматизованої ідентифікації [5].

3. Робастне оцінювання параметрів ідентифікованого розподілу за нижченаведеною процедурою.

4. Перевірка вірогідності відновлення на основі будь-якого відомого критерію згоди, наприклад Колмогорова або Мізеса [3].

Реалізація нижченаведеної процедури робастного оцінювання параметрів розподілу потребує зведення функції розподілу  $F(x; \bar{\Theta})$  до лінійного вигляду  $z(t; \bar{\Theta})$ , де  $t = \phi(x)$  – перетворення над показником  $x$ ,  $z(t) = \varphi(F(x))$  – перетворення над функцією розподілу і перетворення масиву  $\{x_i, F_N(x_i); i = \overline{1, n}\}$ , де  $F_N(x)$  – значення емпіричної функції розподілу,  $n$  – кількість варіант, у масив  $\{t_i, z_i; i = \overline{1, n-1}\}$ . Для розподілів, що розглядаються у роботі, справедливі такі перетворення:

– нормального:

$$\varphi(F(x)) = \Phi^{-1}(F(x)), \quad \phi(x) = x,$$

де  $\Phi^{-1}(\square)$  – обернена до функції Лапласа;

– логарифмічно нормального:

$$\varphi(F(x)) = \Phi^{-1}(F(x)), \quad \phi(x) = \ln x;$$

– експоненціального:

$$\varphi(F(x)) = \ln \frac{1}{1 - F(x)}, \quad \phi(x) = x;$$

– Вейбулла:

$$\varphi(F(x)) = \ln \ln \frac{1}{1-F(x)}, \quad \phi(x) = \ln x.$$

«Чисті» розподіли, які описуються лише одним параметром (як експоненціальний), зводяться до лінійного вигляду

$$z(t) = \theta_1 t.$$

Лінійний вигляд сплайн-розподілів з  $k$  вузлами склеювання з класу однопараметричних розподілів можна подати виразом

$$z(t) = \begin{cases} \theta_1 t, & t \leq t^{(1)}, \\ (\theta_1 - \theta_2)t^{(1)} + \theta_2 t, & t^{(1)} \leq t \leq t^{(2)}, \\ \dots \\ \sum_{j=1}^m (\theta_j - \theta_{j+1})t^{(j)} + \theta_{m+1}t, & t^{(m)} \leq t \leq t^{(m+1)}, m = 1, \dots, k-1, \\ \dots \\ \sum_{j=1}^k (\theta_j - \theta_{j+1})t^{(j)} + \theta_{k+1}t, & t \geq t^{(k)}, \end{cases}$$

де  $t^{(m)}$  –  $m$ -й вузол склеювання.

«Чисті» двопараметричні розподіли зводяться до лінійного виду

$$z(t) = \theta_1 + \theta_2 t.$$

Сплайн-розподіли з  $k$  вузлами склеювання з класу двопараметричних розподілів зводяться до лінійного вигляду

$$z(t) = \begin{cases} \theta_1 + \theta_2 t, & t \leq t^{(1)}, \\ \theta_1 + (\theta_2 - \theta_3)t^{(1)} + \theta_3 t, & t^{(1)} \leq t \leq t^{(2)}, \\ \dots \\ \theta_1 + \sum_{j=1}^m (\theta_{j+1} - \theta_{j+2})t^{(j)} + \theta_{m+2}t, & t^{(m)} \leq t \leq t^{(m+1)}, m = 1, \dots, k-1, \\ \dots \\ \theta_1 + \sum_{j=1}^k (\theta_{j+1} - \theta_{j+2})t^{(j)} + \theta_{k+2}t, & t \geq t^{(k)}. \end{cases}$$

Процедуру знаходження робастної оцінки вектора параметрів розподілу можна подати у наступному вигляді:

1. Покладається  $q = 0$  – номер ітерації і знаходиться початкове наближення оцінки вектора параметрів розподілу  $\hat{\Theta}^{(0)}$  з умови мінімуму функціоналу залишкової дисперсії

$$S^2 = \sum_{i=1}^{n-1} w_i^{(q)} \left( z_i - z \left( t_i; \hat{\Theta}^{(q)} \right) \right)^2, \quad (1)$$

де  $w_i^{(0)} = 1$ ,  $i = \overline{1, n-1}$ ;  $z \left( t_i; \hat{\Theta} \right)$  – значення лінеаризованої функції розподілу;  $z_i$  – емпіричне значення лінеаризованої функції розподілу.

У разі відновлення «чистого» розподілу мінімізація функціоналу (1) реалізує ідею методу найменших квадратів [3]. Знаходження мінімуму функціоналу (1) для сплайн-розподілу передбачає додаткове оцінювання вузлів склеювання шляхом перебору [3]. Для цього вузли склеювання послідовно поміщаються в один з варіантів, і визначаються оцінки інших параметрів сплайн-розподілу шляхом мінімізації функціоналу (1). Серед усіх варіантів вузлів та відповідних їм оцінок обираються ті, для яких значення функціоналу (1) максимальне.

2. Збільшується номер ітерації  $q = q + 1$ .

3. Обчислюються залишки

$$\varepsilon_i = z_i - z \left( t_i; \hat{\Theta}^{(q-1)} \right), \quad i = \overline{1, n-1}.$$

4. Визначається масштабний множник

$$c = \frac{1}{0,6745} \text{median} \left| \varepsilon_i - \text{median}(\varepsilon_i) \right|.$$

5. Знаходяться значення  $y_i = \varepsilon_i / c$ ,  $i = \overline{1, n-1}$ .

6. Використовуючи поточне наближення оцінки  $\hat{\Theta}^{(q-1)}$  та обраний критерій, знаходяться ваги

$$w_i^{(q)} = \frac{\Psi(y_i)}{y_i}, \quad i = \overline{1, n-1},$$

де  $\Psi(y) = \frac{\partial \rho(y)}{\partial y}$ ;  $\rho(y)$  – функція критерію.

Під час реалізації технології використано такі запропоновані у літературі критерії [6]:

– метод найменших квадратів:

$$\rho(y) = \frac{1}{2} y^2, \quad w(y) = 1, \quad y \in (-\infty; +\infty);$$

– Хьюбера:

$$\rho(y) = \begin{cases} \frac{1}{2}y^2, y \in [-a; a] \\ a|y| - \frac{1}{2}a^2, y \notin [-a; a] \end{cases}, w(y) = \begin{cases} 1, y \in [-a; a] \\ \frac{a}{|y|}, y \notin [-a; a] \end{cases},$$

де  $a > 0$  – точка зламу;

– Рамсея:

$$\rho(y) = \frac{1}{a^2} \left( 1 - e^{-a|y|} (1 + a|y|) \right), w(y) = e^{-a|y|}, y \in (-\infty; +\infty),$$

де  $a > 0$  – точка зламу;

– хвиля Ендрюса:

$$\rho(y) = \begin{cases} a \left( 1 - \cos\left(\frac{y}{a}\right) \right), y \in [-a\pi; a\pi] \\ 2a, y \notin [-a\pi; a\pi] \end{cases}, w(y) = \begin{cases} \sin\left(\frac{y}{a}\right), y \in [-a\pi; a\pi] \\ 0, y \notin [-a\pi; a\pi] \end{cases},$$

де  $a\pi > 0$  – точка зламу;

– Тьюкі:

$$\rho(y) = \begin{cases} \frac{1}{2}y^2 - \frac{y^4}{4a^2}, y \in [-a; a] \\ \frac{1}{4}a^2, y \notin [-a; a] \end{cases}, w(y) = \begin{cases} 1 - \frac{y^2}{a^2}, y \in [-a; a] \\ 0, y \notin [-a; a] \end{cases},$$

де  $a > 0$  – точка зламу;

– Гампеля:

$$\rho(y) = \begin{cases} \frac{1}{2}y^2, y \in [-a; a] \\ a|y| - \frac{1}{2}a^2, y \in [-b; -a] \text{ або } y \in [a; b] \\ a \frac{c|y| - \frac{1}{2}y^2}{c-b} - \frac{7a^2}{6}, y \in [-c; -b] \text{ або } y \in [b; c] \\ a(b+c-a), y \in [-\infty; -c] \text{ або } y \in [c; +\infty] \end{cases},$$

$$w(y) = \begin{cases} 1, y \in [-a; a] \\ \frac{a}{|y|}, y \in [-b; -a] \text{ або } y \in [a; b] \\ a \frac{c/|y| - 1}{c - b}, y \in [-c; -b] \text{ або } y \in [b; c] \\ 0, y \in [-\infty; -c] \text{ або } y \in [c; +\infty] \end{cases},$$

де  $c > b > a > 0$  – точки зламів.

7. З умови мінімуму функціоналу (1), використовуючи ваги, знайдені на кроці 6, обчислюється наступне наближення оцінки вектора параметрів  $\hat{\Theta}^{(q)}$ .

8. Перевіряється умова

$$\max_{j=1,p} \left| \hat{\theta}_j^{(q)} - \hat{\theta}_j^{(q-1)} \right| < e.$$

Якщо вона виконується, то  $\hat{\Theta}^{(q)} = \{ \hat{\theta}_1^{(q)}, \hat{\theta}_2^{(q)}, \dots, \hat{\theta}_p^{(q)} \}$  приймаються за шукану оцінку вектора параметрів, інакше здійснюється перехід на крок 2.

Обчислювальна технологія робастного оцінювання параметрів розподілу складала ядро програмного забезпечення «RobustProcедre», яке реалізовано у середовищі Borland Delphi 7.0.

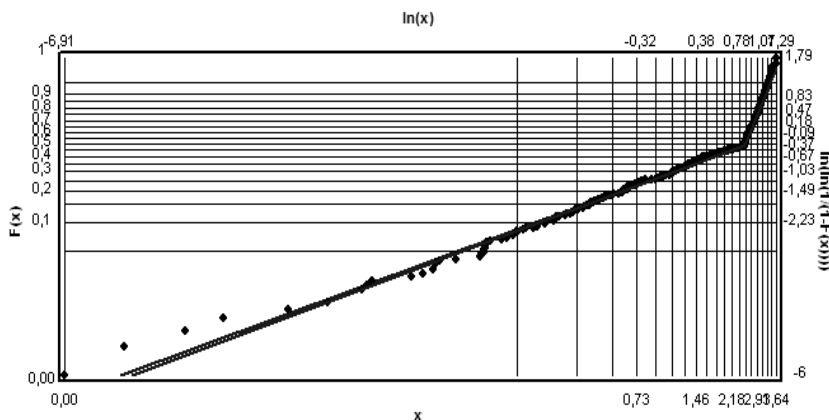
Тестування створеної інформаційної технології проведено на даних імітаційного моделювання. Експерименти полягали у моделюванні вибірки за певним законом розподілу ймовірностей, додаванні до вибірки аномальних значень із подальшим відновленням розподілу за допомогою інформаційної технології. Нижче, у якості прикладу, наведено результати одного з експериментів.

Під час експерименту моделювались дані зі сплайн-розподілу Вейбулла з одним вузлом склеювання. Параметри моделювання та знайдені оцінки параметрів подано у таблиці 1. На ймовірнісному папері розподілу Вейбулла нанесено оцінки лінеаризованої функції розподілу, знайдені методом найменших квадратів та за допомогою описаної вище технології (рис. 1). Функція, що проходить нижче на графіку, оцінена за допомогою створеної технології.

Ймовірність вірогідного відновлення розподілу методом найменших квадратів за критерієм Колмогорова складає 0,89, а робастного відновлення – 0,95.

**Параметри та знайдені оцінки параметрів  
сплайн-розподілу Вейбулла з одним вузлом**

Параметр	Значення під час моделювання	Оцінка методом найменших квадратів	Робастна оцінка
$\alpha$	3	2,93	2,92
$\beta_1$	0,8	0,79	0,80
$\beta_2$	5	5,57	5,08
$x_0$	2,5	2,53	2,51



**Рис. 1. Ймовірнісний папір розподілу Вейбулла  
з відновленою лінеаризованою функцією розподілу**

**Висновки** за результатами проведеної роботи:

– Було створено інформаційну технологію робастного відновлення розподілів з класів нормального, логарифмічно нормального, експоненціального та Вейбулла («чистих» розподілів та сплайн-розподілів з одним і двома вузлами).

– Інформаційну технологію реалізовано у вигляді програмного забезпечення «RobustProcedre».

– Проведено тестування технології на модельованих даних.

– Результати тестування свідчать про наступне: 1) оцінки параметрів, знайдені на основі запропонованої технології, більш близькі до істинних оцінок параметрів (параметрів моделювання), ніж оцінки, знайдені методом найменших квадратів; 2) робастне відновлення

більш вірогідне, ніж відновлення на основі методу найменших квадратів, за критерієм згоди Колмогорова.

### Бібліографічні посилання

1. **Хьюбер П.** Робастність в статистиці / П. Хьюбер. – М. : Мир, 1984. – 304 с.
2. Робастність в статистиці. Підхід на основі функцій впливу / Ф. Хампель, Э. Рончетти, П. Рауссеу, В. Штаэль. – М. : Мир, 1989. – 512 с.
3. Статистична обробка даних / В. П. Бабак, А. Я. Білецький, О. П. Приставка, П. О. Приставка. – К. : МІВВЦ, 2001. – 388 с.
4. **Приставка О. П.** Сплайн-розподіли у статистичному аналізі / О. П. Приставка. – Д. : Вид-во ДДУ, 1995. – 152 с.
5. **Мацуга О. М.** Обчислювальні схеми ідентифікації сплайн-розподілів за ймовірнісним папером / О. М. Мацуга, Г. С. Шубіна // Актуальні проблеми автоматизації та інформаційних технологій. – 2012. – Т. 16. – С. 112–123.
6. **Дрейпер Н.** Прикладной регрессионный анализ / Н. Дрейпер, Г. Смит. – 3-е изд. – М. : Вильямс, 2007. – 912 с.

*Надійшла до редколегії 24.06.14*