

Нечуйвігер О.П., Мотря А.І., ДКІ-ПОМ 14 пр+ПОБ14 пр ПРО ОЦІНКИ ПОХИБОК ЗАОКРУГЛЕННЯ ОСНОВНИХ АРИФМЕТИЧНИХ ОПЕРАЦІЙ

Нехай обчислення виконуються на ЕОМ у режимі з плаваючою комою (П.К.) із заокругленням результатів арифметичних операцій за класичним правилом до τ двійкових розрядів у нормалізованих мантис чисел. Похибки заокруглення в цих операціях мають вигляд

$$z = fl \begin{pmatrix} \pm \\ x \times y \\ \vdots \end{pmatrix} (1 + \varepsilon), \quad (1)$$

де ε – відносна похибка і $|\varepsilon| \leq 2^{-\tau}$. При цьому x і y вважаються заданими точно в ЕОМ, звідси випливає, що мантиси x і y мають не більш, ніж τ розрядів. Якщо вхідне число x має нормалізовану мантису з числом цифр, більшим, ніж τ (і порядок числа x може бути записаний в ЕОМ), то

$$x = x_\tau (1 + \varepsilon), \quad |\varepsilon| \leq 2^{-\tau-1}, \quad (2)$$

де x_τ – машинне заокруглене представлення числа x . Усі результати, які мають місце в режимі П.К., є прямим наслідком співвідношень (1), (2), застосування яких приводить до оцінок вигляду $(1 - 2^{-\tau})^r \leq 1 + \varepsilon \leq (1 + 2^{-\tau})^r$, які можна спростити, припустивши, що виконується умова $r \cdot 2^{-\tau} < 0,1$ (це цілком виправдано в практичних застосуваннях для будь-якого прийняттого τ). Тоді $(1 + 2^{-\tau})^r < 1 + 1,06 \cdot r \cdot 2^{-\tau}$, $(1 - 2^{-\tau})^r < 1 - 1,06 \cdot r \cdot 2^{-\tau}$, $1 - 1,06 \cdot r \cdot 2^{-\tau} < 1 + \varepsilon < 1 + 1,06 \cdot r \cdot 2^{-\tau}$, звідки $|\varepsilon| < 1,06 \cdot r \cdot 2^{-\tau}$. Останнє співвідношення використовується у всіх наступних оцінках:

$$1) \quad fl(x_1 \cdot x_2 \cdot \dots \cdot x_N) \equiv \prod_{i=1}^N x_i (1 + E), \quad \text{де } |E| < (N-1) \cdot 1,06 \cdot 2^{-\tau};$$

$$2) \quad fl(x_1 + x_2 + \dots + x_N) \equiv x_1(1 + \varepsilon_1) + x_2(1 + \varepsilon_2) + \dots + x_N(1 + \varepsilon_N), \quad \text{де } |\varepsilon_1| < (N-1) \cdot 1,06 \cdot 2^{-\tau}, \\ |\varepsilon_r| < (N-r+1) \cdot 1,06 \cdot 2^{-\tau}, \quad r = \overline{2, N}. \quad \text{Тут передбачалося, що } S_2 = fl(x_1 + x_2), \quad S_r = fl(S_{r-1} + x_r), \\ r = \overline{3, N};$$

$$3) \quad fl(x_1 \times y_1 + x_2 \times y_2 + \dots + x_N \times y_N) \equiv x_1 y_1 (1 + \varepsilon_1) + x_2 y_2 (1 + \varepsilon_2) + \dots + x_N y_N (1 + \varepsilon_N),$$

$$\text{де } |\varepsilon_1| < N \cdot 1,06 \cdot 2^{-\tau}, \quad |\varepsilon_r| < (N-r+2) \cdot 1,06 \cdot 2^{-\tau}, \quad r = \overline{2, N};$$

$$4) \quad fl \left(\frac{x_1 \cdot x_2 \cdot \dots \cdot x_m}{y_1 \cdot y_2 \cdot \dots \cdot y_n} \right) \equiv \frac{x_1 \cdot x_2 \cdot \dots \cdot x_m}{y_1 \cdot y_2 \cdot \dots \cdot y_n} (1 + E), \quad \text{де } |E| < (m+n-1) \cdot 1,06 \cdot 2^{-\tau}.$$

Для операцій додавання і множення, реалізованих на ЕОМ у режимі П.К., справедливі нерівності: $|fl(a \cdot x_1) - fl(a \cdot x_2)| \leq 2|a| \cdot |x_1 - x_2|$, $|fl(a + x_1) - fl(a + x_2)| \leq s \cdot |x_1 - x_2|$. Постійна s залежить від співвідношення порядків доданків і способу заокруглення, зафіксованого в ЕОМ. Можна підібрати способи запису й заокруглення такі, що $s = 2$.