

1 Аналіз похибок заокруглення

1.1 Види похибок

Нехай необхідно розв'язати рівняння

$$Au = f. \quad (1)$$

Неточно задані вхідні дані призводять до рівняння

$$\tilde{A}\tilde{u} = \tilde{f}. \quad (2)$$

Назвемо $\delta_1 = u - \tilde{u}$ *незсувеною похибкою*.

Застосування методу розв'язання (2) призводить до рівняння

$$\tilde{A}_h \tilde{u}_h = \tilde{f}_h, \quad (3)$$

де $h > 0$ – малий параметр. Назвемо $\delta_2 = \tilde{u} - \tilde{u}_h$ *похибкою методу*.

Реалізація методу на ЕОМ призводить до рівняння

$$\tilde{A}_h^* \tilde{u}_h^* = \tilde{f}_h^*. \quad (4)$$

Назвемо $\delta_3 = \tilde{u}_h - \tilde{u}_h^*$ *похибкою заокруглення*.

Назвемо $\delta = u - \tilde{u}_h^* = \delta_1 + \delta_2 + \delta_3$ *повною похибкою*.

Визначення 1. Кажуть, що задача (1) *коректна*, якщо

1. $\forall f \in F \exists! u \in U$;
2. задача (1) *стійка*, тобто

$$\forall \varepsilon > 0 \exists \delta > 0 \forall f : \|A - \tilde{A}\| < \delta, \|f - \tilde{f}\| < \delta \Rightarrow \|u - \tilde{u}\| < \varepsilon.$$

Якщо задача (1) некоректна, то або розв'язок її не існує, або він не єдиний, або він нестійкий, тобто

$$\exists \varepsilon > 0 \forall \delta > 0 \exists f : \|A - \tilde{A}\| < \delta, \|f - \tilde{f}\| < \delta, \|u - \tilde{u}\| > \varepsilon.$$

Абсолютна похибка $\Delta x \leq |x - x^*|$.

Відносна похибка $\delta x \leq \frac{\Delta x}{|x|}$ або $\frac{\Delta x}{|x^*|}$.

Значущими цифрами називаються всі цифри, починаючи з першої ненульової зліва.

Вірна цифра – це значуща, якщо абсолютна похибка за рахунок відкидання всіх молодших розрядів не перевищує одиниці розряду цієї цифри. Тобто, якщо $x^* = \overline{\alpha_n \dots \alpha_0}.\overline{\alpha_{-1} \dots \alpha_{-p}}$, то α_{-p} – вірна, якщо $\Delta x \leq 10^{-p}$ (інколи беруть $\Delta x \leq q \dots 10^{-p}$, $1/2 \leq w < 1$, наприклад $w = 0.55$).

1.2 Підрахунок похибок в ЕОМ

Обчислимо відносну похибку заокруглення числа x на ЕОМ з плаваючою комою. У системі числення з основою β число x представляється у вигляді

$$x = \pm (\alpha_1 \beta^{-1} + \alpha_2 \beta^{-2} + \dots + \alpha_t \beta^{-t} + \dots) \beta^p, \quad (5)$$

де $0 \leq \alpha_k < \beta$, $\alpha_1 \neq 0$, $k = 1, 2, \dots$

Якщо в ЕОМ t розрядів, то при відкиданні молодших розрядів ми оперуємо з наближеним значенням

$$x^* = \pm (\alpha_1 \beta^{-1} + \alpha_2 \beta^{-2} + \dots + \alpha_t \beta^{-t}) \beta^p$$

і, відповідно, похибка заокруглення $x - x^* = \pm \beta^p (\alpha_{t+1} \beta^{-t-1} + \dots)$. Її можна оцінити так:

$$|x - x^*| \leq \beta^{p-t-1}(\beta - 1)(1 + \beta^{-1} + \dots) \leq \beta^{p-t-1}(\beta - 1) \frac{1}{1 - \beta^{-1}} = \beta^{p-t}.$$

Враховуючи, що $\alpha_1 \neq 0$, маємо $|x| \geq \beta^p \beta^{-1} = \beta^{p-1}$. Звідси остаточно

$$\delta x \leq \frac{\beta^{p-t}}{\beta^{p-1}} = \beta^{-t+1}.$$

При точніших способах заокруглення можна отримати оцінку $\delta x \leq \beta^{-t+1}/2 = \varepsilon$. Число ε називається *машичним інсилом*. Наприклад, для $\beta = 2$, $t = 24$, $\varepsilon = 2^{-24} \approx 10^{-7}$.

1.3 Обчислення похибок обчислення значення функції

Нехай задана функція $y = f(x_1, \dots, x_n) \in C^{(1)}(\Omega)$. Необхідно обчислити її значення при наближеному значенні аргументів $\vec{x}^* = (x_1^*, \dots, x_n^*)$, де $|x_i - x_i^*| \leq \Delta x_i$ та оцінити похибку обчислення значення функції $y^* = f(x_1^*, \dots, x_n^*)$.

Маємо

$$|y - y^*| = |f(\vec{x}) - f(\vec{x}^*)| = \left| \sum_{i=1}^n \frac{\partial f}{\partial x_i}(\vec{\xi})(x_i - x_i^*) \right| \leq \sum_{i=1}^n B_i \cdots \Delta x_i,$$

де $B_i = \max_{\vec{x} \in U} \left| \frac{\partial f}{\partial x_i}(\vec{x}) \right|$, $U = \{\vec{x} = (x_1, \dots, x_n) : |x_i - x_i^*| \leq \Delta x_i, i = 1, \dots, n\} \subset \Omega$.

Отже, з точністю до величин першого порядку малості по $\Delta x = \max_i \Delta x_i$, $\Delta y = |y - y^*| \prec \sum_{i=1}^n b_i \cdots \Delta x_i$,

де $b_i = \left| \frac{\partial f}{\partial x_i}(\vec{x}^*) \right|$ і “ \prec ” означає *приблизно менше*.

Розглянемо похибки арифметичних операцій.

1. Сума: $y = x_1 + x_2$, $x_1, x_2 > 0$, $\Delta y \leq \Delta x_1 + \Delta x_2$, $\delta y = \frac{\Delta y}{y} = \frac{\Delta y}{x_1 + x_2} \leq \frac{\Delta x_1 + \Delta x_2}{x_1 + x_2} \leq \max(\delta x_1, \delta x_2)$.

2. Різниця: $y = x_1 - x_2$, $x_1 > x_2 > 0$, $\Delta y \leq \Delta x_1 + \Delta x_2$, $\delta y = \frac{\Delta y}{y} = \frac{\Delta y}{x_1 - x_2} \leq \frac{\Delta x_1 + \Delta x_2}{x_1 - x_2} = \frac{x_1 \delta x_1 + x_2 \delta x_2}{x_1 - x_2}$.

Як бачимо, при близьких аргументах зростає відносна похибка.

3. Добуток: $y = x_1 \cdot x_2$, $x_1, x_2 > 0$, $\Delta y = x_1 \Delta x_2 + x_2 \Delta x_1 + \Delta x_1 \Delta x_2 \prec x_1 \Delta x_2 + x_2 \Delta x_1$, $\delta y = \frac{\Delta y}{y} = \frac{\Delta y}{x_1 x_2} \prec \frac{x_1 \Delta x_2 + x_2 \Delta x_1}{x_1 x_2} = \frac{\Delta x_1}{x_1} + \frac{\Delta x_2}{x_2} = \delta x_1 + \delta x_2$.

4. Частка: $y = \frac{x_1}{x_2}$, $x_1, x_2 > 0$, $\Delta y = \frac{x_2 \Delta x_1 - x_1 \Delta x_2}{x_2(x_2 + \Delta x_2)} < \frac{x_2 \Delta x_1 + x_1 \Delta x_2}{x_2(x_2 + \Delta x_2)} \prec \frac{x_2 \Delta x_1 + x_1 \Delta x_2}{x_2^2}$, $\delta y = \frac{\Delta y}{y} = \frac{x_2 \Delta y}{x_1} \prec \frac{x_2 \Delta x_1 + x_1 \Delta x_2}{x_1 x_2} = \frac{\Delta x_1}{x_1} + \frac{\Delta x_2}{x_2} = \delta x_1 + \delta x_2$.

Як бачимо, при малих x_2 зростає абсолютна похибка.

Пряма задача аналізу похибок: обчислення $\Delta y, \delta y$ за заданими $\Delta x_i, i = 1, \dots, n$.

Обернена задача: знаходження $\Delta x_i, i = 1, \dots, n$ за заданими $\Delta y, \delta y$. Якщо $n > 1$ маємо одну умову $\sum_{i=1}^n b_i \Delta x_i < \varepsilon$ на багато невідомих Δx_i . Зазвичай вибирають їх із однієї з умов

$$\forall i : b_i \Delta x_i < \varepsilon/n \quad \text{або} \quad \forall i : \Delta x_i < \varepsilon/B, \quad \text{де } B = \sum_{i=1}^n b_i.$$