

# 1. Аналіз похибок заокруглення

## 1.1. Види похибок

Нехай необхідно розв'язати рівняння

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

За рахунок неточно заданих вхідних даних насправді ми маємо рівняння

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

**Означення:** Назвемо  $\delta_1 = u - \tilde{u}$  неусувною похибкою.

Застосування методу розв'язання (2) приводить до рівняння

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

де  $h > 0$  — малий параметр.

**Означення:** Назвемо  $\delta_2 = \tilde{u} - \tilde{u}_h$  похибкою методу.

Реалізація методу на ЕОМ приводить до рівняння

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

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

**Означення:** Тоді повна похибка  $\delta = u - \tilde{u}_h^* = \delta_1 + \delta_2 + \delta_3$ .

**Означення:** кажуть, що задача (1) коректна, якщо

- $\forall f \in F: \exists! u \in U$ ;
- Задача (1) стійка, тобто  $\forall \varepsilon > 0: \exists \delta > 0$ :

$$|A - \tilde{A}| < \delta, |f - \tilde{f}| < \delta \implies |u - \tilde{u}| < \varepsilon. \quad (5)$$

Якщо задача (1) некоректна, то або розв'язок її не існує, або він неєдиний, або він нестійкий, тобто  $\exists \varepsilon > 0: \forall \delta > 0$ :

$$|A - \tilde{A}| < \delta, |f - \tilde{f}| < \delta \implies |u - \tilde{u}| > \varepsilon. \quad (6)$$

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

**Означення:** Відносна похибка  $\delta x \leq \Delta x / |x|$ , або  $\Delta x / |x^*|$ .

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

**Означення:** Вірна цифра — це значуща, якщо абсолютна похибка за рахунок відкидання всіх молодших розрядів не перевищує одиниці розряду цієї цифри.

Тобто, якщо  $x^* = \alpha_n \dots \alpha_0 . \alpha_{-1} \dots \alpha_{-p} \dots$ , то  $\alpha_{-p}$  вірна, якщо  $\Delta x \leq 10^{-p}$  (інколи  $\Delta x \leq w \cdot 10^{-p}$ , де  $1/2 \leq w < 1$  наприклад,  $w = 0.55$ ).

## 1.2. Підрахунок похибок в ЕОМ

Підрахуємо відносну похибку заокруглення числа  $x$  на ЕОМ з плаваючою комою. В  $\beta$ -ічній системі числення число представляється у вигляді

$$x = \pm(\alpha_1\beta^{-1} + \alpha_2\beta^{-2} + \dots + \alpha_t\beta^{-t} + \dots) \cdot \beta^p, \quad (7)$$

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

$$x^* = \pm(\alpha_1\beta^{-1} + \alpha_2\beta^{-2} + \dots + \alpha_t\beta^{-t}) \cdot \beta^p \quad (8)$$

і відповідно похибка заокруглення

$$x - x^* = \pm\beta^p \cdot (\alpha_{t+1}\beta^{-t-1} + \dots). \quad (9)$$

Тоді її можна оцінити так

$$|x - x^*| \leq \beta^{p-t-1} \cdot (\beta - 1) \cdot (1 + \beta^{-1} + \dots) \leq \beta^{p-t-1} \cdot (\beta - 1) \cdot \frac{1}{1 - \beta^{-1}} = \beta^{p-t}. \quad (10)$$

Якщо в представленні (7) взяти  $\alpha_1 = 1$ , то  $|x| \geq \beta^p \cdot \beta^{-1}$ . Звідси остаточно

$$\delta x \leq \frac{\beta^{p-t}}{\beta^{p-1}} = \beta^{-t+1}. \quad (11)$$

При більш точних способах заокруглення можна отримати оцінку  $\delta x \leq \frac{1}{2} \cdot \beta^{-t+1} = \varepsilon$ . Число  $\varepsilon$  називається «машинним іпсилон». Наприклад, для  $\beta = 2$ ,  $t = 24$ ,  $\varepsilon = 2^{-24} \approx 10^{-7}$ .

## 1.3. Підрахунок 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

$$|y - y^*| = |f(\vec{x}) - f(\vec{x}^*)| = \left| \sum_{i=1}^n \frac{\partial f}{\partial x_i}(\vec{\xi}) \cdot (x_i - x_i^*) \right| \leq \sum_{i=1}^n B_i \cdot \Delta x_i, \quad (12)$$

$$\text{де } B_u = \max_{\vec{x} \in U} \left| \frac{\partial f}{\partial x_i}(\vec{x}) \right|.$$

Тут

$$U = \{ \vec{x} = (x_1, \dots, x_n) : |x_i - x_i^*| \leq \Delta x_i \} \subset \Omega, \quad (13)$$

для  $i = \overline{1, n}$ . Отже з точністю до величин першого порядку малості по

$$\Delta x = \max_i \Delta x_i, \quad (14)$$

$$\Delta y = |y - y^*| \prec \sum_{i=1}^n b_i \cdot \Delta x_i, \quad (15)$$

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

Розглянемо похибки арифметичних операцій.

- Сума:  $y = x_1 + x_2, x_1, x_2 > 0$ :

$$\Delta y \leq \Delta x_1 + \Delta x_2, \quad (16)$$

$$\delta y \leq \frac{\Delta x_1 + \Delta x_2}{x_1 + x_2} \leq \max(\delta x_1, \delta x_2). \quad (17)$$

- Різниця:  $y = x_1 - x_2, x_1 > x_2 > 0$ :

$$\Delta y \leq \Delta x_1 + \Delta x_2, \quad (18)$$

$$\delta y \leq \frac{x_2 \delta x_1 + x_1 \delta x_2}{x_1 - x_2}. \quad (19)$$

При близьких  $x_1, x_2$  зростає відносна похибка (за рахунок втрати вірних цифр).

- Добуток:  $y = x_1 \cdot x_2, x_1, x_2 > 0$ :

$$\Delta y \prec x_2 \Delta x_1 + x_1 \Delta x_2, \quad (20)$$

$$\delta y \leq \delta x_1 + \delta x_2. \quad (21)$$

- Частка  $y = x_1/x_2, x_1, x_2 > 0$ :

$$\Delta y \prec \frac{x_2 \Delta x_1 + x_1 \Delta x_2}{x_2^2}, \quad (22)$$

$$\delta y \leq \delta x_1 + \delta x_2. \quad (23)$$

При малих  $x_2$  зростає абсолютна похибка (за рахунок зростання результату ділення).

**Означення:** Пряма задача аналізу похибок: обчислення  $\Delta y, \delta y$  по заданих  $\Delta x_i, i = \overline{1, n}$ .

**Означення:** *Обернена задача:* знаходження  $\Delta x_i, i = \overline{1, n}$  по заданих  $\Delta y, \delta y$ . Якщо  $n > 1$ , маємо одну умову

$$\sum_{i=1}^n b_i \cdot \Delta x_i < \varepsilon \quad (24)$$

для багатьох невідомих  $\Delta x_i$ .

Вибирають їх із однієї з умов:

$$\forall i : b_i \cdot \Delta x_i < \frac{\varepsilon}{n} \quad (25)$$

або

$$\Delta x_i < \frac{\varepsilon}{\sum_{i=1}^n b_i}. \quad (26)$$

## 2. Методи розв'язання нелінійних рівнянь

*Постановка задачі.* Нехай маємо рівняння  $f(x) = 0$ ,  $\bar{x}$  — його розв'язок, тобто  $f(\bar{x}) = 0$ .

Задача розв'язання цього рівняння розпадається на етапи:

- Існування та кількість коренів.
- Відділення коренів, тобто розбиття числової вісі на інтервали, де знаходиться один корінь.
- Обчислення кореня із заданою точністю  $\varepsilon$ .

Для розв'язання перших двох задач використовуються методи математичного аналізу та алгебри, а також графічний метод. Далі розглядаються методи розв'язання третього етапу.

### 2.1. Метод ділення навпіл

Припустимо на  $[a, b]$  знаходиться лише один корінь рівняння

$$f(x) = 0 \quad (1)$$

для  $f(x) \in C[a, b]$ , який необхідно визначити. Нехай  $f(a) \cdot f(b) < 0$ .

Припустимо, що  $f(a) > 0$ ,  $f(b) < 0$ . Покладемо  $x_1 = \frac{a+b}{2}$  і підрахуємо  $f(x_1)$ . Якщо  $f(x_1) < 0$ , тоді шуканий корінь  $\bar{x}$  знаходиться на інтервалі  $(a, x_1)$ . Якщо ж  $f(x_1) > 0$ , то  $\bar{x} \in (x_1, b)$ . Далі з двох інтервалів  $(a, x_1)$  і  $(x_1, b)$  вибираємо той, на границях якого функція  $f(x)$  має різні знаки, знаходимо точку  $x_2$  — середину вибраного інтервалу, підраховуємо  $f(x_2)$  і повторюємо вказаний процес.

В результаті отримаємо послідовність інтервалів, що містять шуканий корінь  $\bar{x}$ , причому довжина кожного послідовного інтервалу вдвічі менше попереднього.

Цей процес продовжується до тих пір, поки довжина отриманого інтервалу  $(a_n, b_n)$  не стане меншою за  $b_n - a_n < 2\varepsilon$ . Тоді  $x_{n+1}$ , як середина інтервалу  $(a_n, b_n)$ , пов'язане з  $\bar{x}$  нерівністю

$$|x_{n+1} - \bar{x}| < \varepsilon. \quad (2)$$

Ця умова для деякого  $n$  буде виконуватись за теоремою Больцано-Коші. Оскільки

$$|b_{k+1} - a_{k+1}| = \frac{|b_k - a_k|}{2}, \quad (3)$$

то

$$|x_{n+1} - \bar{x}| \leq \frac{b - a}{2^{n+1}} < \varepsilon. \quad (4)$$

Звідси отримаємо нерівність для обчислення кількості ітерацій  $n$  для виконання умови (2):

$$n = n(\varepsilon) \geq \left\lceil \log \left( \frac{b - a}{\varepsilon} \right) \right\rceil + 1. \quad (5)$$

Степінь збіжності — лінійна, тобто геометричної прогресії з знаменником  $q = 1/2$ .

- **Переваги методу:** простота, надійність.

- **Недоліки методу:** низька швидкість збіжності; метод не узагальнюється на системи.

## 2.2. Метод простої ітерації

Спочатку рівняння

$$f(x) = 0 \quad (6)$$

замінюється еквівалентним

$$x = \varphi(x). \quad (7)$$

Ітераційний процес має вигляд

$$x_{n+1} = \varphi(x_n), \quad n = 0, 1, \dots \quad (8)$$

Початкове наближення  $x_0$  задається.

Для збіжності велике значення має вибір функції  $\varphi(x)$ . Перший спосіб заміни рівняння полягає в відділенні змінної з якогось члена рівняння. Більш продуктивним є перехід від рівняння (6) до (7) з функцією  $\varphi(x) = x + \tau(x) \cdot f(x)$ , де  $\tau(x)$  — знакостала функція на тому відрізку, де шукаємо корінь.

**Означення:** Кажуть, що ітераційний метод збігається, якщо  $\lim_{k \rightarrow \infty} x_k = \bar{x}$ .

Далі  $U_r = \{x : |x - a| \leq r\}$  відрізок довжини  $2r$  з серединою в точці  $a$ .

З'ясуємо умови, при яких збігається метод простої ітерації.

**Теорема 1:** Якщо

$$\max_{x \in [a, b] = U_r} |\varphi'(x)| \leq q < 1 \quad (9)$$

то метод простої ітерації збігається і має місце оцінка

$$|x_n - \bar{x}| \leq \frac{q^n}{1 - q} \cdot |x_0 - x_1| \leq \frac{q^n}{1 - q} \cdot (b - a). \quad (10)$$

**Доведення:** Нехай  $x_{k+1}, x_k \in U_r$ . Тоді

$$\begin{aligned} |x_k - x_{k-1}| &= |\varphi(x_k) - \varphi(x_{k-1})| = |\varphi'(\xi_k) \cdot (x_k - x_{k-1})| \leq |\varphi'(\xi_k)| \cdot |x_k - x_{k-1}| \leq \\ &\leq q \cdot |x_k - x_{k-1}| = \dots = q^k \cdot |x_1 - x_0|, \end{aligned} \quad (11)$$

де  $\xi_k = x_k + \theta_k \cdot (x_{k+1} - x_k)$ , а у свою чергу  $0 < \theta_k < 1$ . Далі

$$\begin{aligned} |x_{k+p} - x_k| &= |x_{k+p} - x_{k+p-1} + \dots + x_{k+1} - x_k| = |x_{k+p} - x_{k+p-1}| + \dots + |x_{k+1} - x_k| \leq \\ &\leq (q^{k+p-1} + q^{k+p-2} + \dots + q^k) \cdot |x_1 - x_0| = \frac{q^k - q^{k+p-1}}{1 - q} \cdot |x_1 - x_0| \xrightarrow[k \rightarrow \infty]{} 0. \end{aligned} \quad (12)$$

Бачимо що  $\{x_k\}$  — фундаментальна послідовність. Значить вона збіжна. При  $p \rightarrow \infty$  в (12) отримуємо (10).  $\square$

Визначимо кількість ітерацій для досягнення точності  $\varepsilon$ . З оцінки в теоремі 1 отримаємо

$$|x_n - \bar{x}| \leq \frac{q^n}{1 - q} \cdot (b - a) < \varepsilon, \quad (13)$$

звідки безпосередньо маємо

$$n(\varepsilon) = n \geq \left\lceil \frac{\ln\left(\frac{\varepsilon(1-q)}{b-a}\right)}{\ln q} \right\rceil + 1. \quad (14)$$

Практично ітераційний процес зупиняємо при:  $|x_n - x_{n-1}| < \varepsilon$ . Але ця умова не завжди гарантує, що  $|x_n - \bar{x}| < \varepsilon$ .

**Зауваження:** Умова збіжності методу може бути замінена на умову Ліпшиця

$$|\varphi(x) - \varphi(y)| \leq q \cdot |x - y|, \quad 0 < q < 1. \quad (15)$$

- **Переваги методу:** простота; при  $q < 1/2$  — швидше збігається ніж метод ділення навпіл; метод узагальнюється на системи.
- **Недоліки методу:** при  $q > 1/2$  збігається повільніше ніж метод ділення навпіл; виникають труднощі при зведенні  $f(x) = 0$  до  $x = \varphi(x)$ .

### 2.3. Метод релаксації

Якщо в методі простої ітерації для рівняння  $x = x + \tau \cdot f(x) \equiv \varphi(x)$  вибрати  $\tau(x) = \tau = \text{const}$ , то ітераційний процес приймає вигляд

$$x_{n+1} = x_n + \tau \cdot f(x_n), \quad (16)$$

де  $k = 0, 1, 2, 3 \dots$ , а  $x_0$  — задано. Метод можна записати у вигляді

$$\frac{x_{k+1} - x_k}{\tau} = f(x_k), \quad k = 0, 1, \dots \quad (17)$$

Оскільки  $\varphi'(x) = 1 + \tau \cdot f'(x)$ , то метод збігається при умові

$$|\varphi'(x)| = |1 + \tau \cdot f'(x)| \leq q < 1. \quad (18)$$

Нехай  $f'(x) < 0$ , тоді (8) запишеться у вигляді:  $-q \leq 1 + \tau \cdot f'(x) \leq q < 1$ . Звідси

$$f'(x) \leq 1 + q < 2k\tau, \quad (19)$$

і

$$0 < \tau < \frac{2}{|f'(x)|}. \quad (20)$$

Поставимо задачу знаходження  $\tau$ , для якого  $q = q(\tau) \rightarrow \min$ . Для того, щоб вибрати оптимальний параметр  $\tau$ , розглянемо рівняння для похибки  $z_k = x_k - \bar{x}$ .

Підставивши  $x_k = x + z_k$  в (16), отримаємо

$$z_{k+1} = z_k + \tau \cdot f(x + z_k). \quad (21)$$

В припущенні  $f(x) \in C^1([a, b])$  з теореми про середнє маємо

$$f(\bar{x} + z_k) = f(\bar{x}) + z_k \cdot f'(\bar{x} + \theta \cdot z_k) = z_k \cdot f'(\bar{x} + \theta \cdot z_k) = z_k \cdot f'(\xi_k), \quad (22)$$

тобто

$$z_{k+1} = z_k + \tau \cdot f'(\xi_k) \cdot z_k. \quad (23)$$

Звідси

$$|z_{k+1}| \leq |1 + \tau \cdot f'(\xi_k)| \cdot |z_k| \leq \max_U |1 + \tau \cdot f'(\xi_k)| \cdot |z_k|. \quad (24)$$

А тому

$$|z_{k+1}| \leq \max \{|1 - \tau M_1|, |1 - \tau m_1|\} \cdot |z_k|, \quad (25)$$

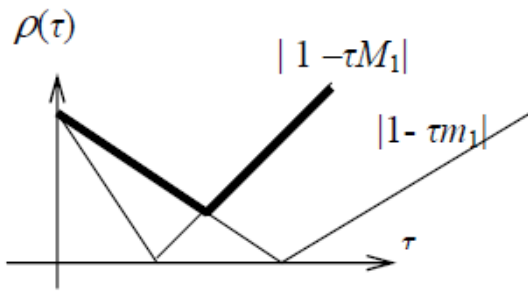
де

$$m_1 = \min_{[a,b]} |f'(x)|, \quad M_1 = \max_{[a,b]} |f'(x)| \quad (26)$$

Таким чином, задача вибору оптимального параметра зводиться до знаходження  $\tau$ , для якого функція

$$q(\tau) = \max \{|1 - \tau M_1|, |1 - \tau m_1|\} \quad (27)$$

приймає мінімальне значення:  $q(\tau) \rightarrow \min$ .



З графіка видно, що точка мінімуму визначається умовою  $|1 - \tau M_1| = |1 - \tau m_1|$ . Тому

$$1 - \tau_0 m_1 = \tau_0 M_1 - 1 \implies \tau_0 = \frac{2}{M_1 + m_1} < \frac{2}{|f'(x)|}. \quad (28)$$

При цьому значенні  $\tau$  маємо

$$q(\tau_0) = \rho_0 = \frac{M_1 - m_1}{M_1 + m_1}. \quad (29)$$

Тоді для похибки вірна оцінка

$$|x_n - \bar{x}| \leq \frac{\rho_0^n}{1 - \rho_0} \cdot (b - a) < \varepsilon. \quad (30)$$

Кількість ітерацій

$$n = n(\varepsilon) \geq \left\lceil \frac{\frac{\ln(\varepsilon(1-\rho_0))}{b-a}}{\ln \rho_0} \right\rceil + 1. \quad (31)$$

**Задача 1:** Дати геометричну інтерпретацію методу простої ітерації для випадків:



$$0 < \varphi'(x) < 1; \quad -1 < \varphi'(x) < 0; \quad \varphi'(x) < -1; \quad \varphi'(x) > 1. \quad (32)$$

**Задача 2:** Знайти оптимальне  $\tau = \tau_0$  для методу релаксації при  $f'(x) > 0$ .

## 2.4. Метод Ньютона (метод дотичних)

Припустимо, що рівняння  $f(x) = 0$  має простий дійсний корінь  $\bar{x}$ , тобто  $f(\bar{x}) = 0$ ,  $f'(\bar{x}) \neq 0$ . Нехай виконуються умови:  $f(x) \in C^1([a, b])$ ,  $f(a) \cdot f(b) < 0$ . Тоді

$$0 = f(\bar{x}) = f(x_k + \bar{x} - x_k) = f(x_k) + f'(\xi_k) \cdot (x - x_k), \quad (33)$$

де  $\xi_k = x_k + \theta_k \cdot (\bar{x} - x_k)$ ,  $0 < \theta_k < 1$ ,  $\xi_k \approx x_k$ . Тому наступне наближення виберемо з рівняння

$$f(x_k) + f'(x_k) \cdot (x_{k+1} - x_k) = 0. \quad (34)$$

Звідси маємо ітераційний процес

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}, \quad (35)$$

де  $k = 0, 1, 2, \dots$ ;  $x_0$  — задане.

Метод Ньютона ще називають методом лінеаризації або методом дотичних.

**Задача 3:** Дати геометричну інтерпретацію методу Ньютона.

Метод Ньютона можна інтерпретувати як метод простої ітерації з

$$\varphi(x) = x - \frac{f(x)}{f'(x)}, \quad (36)$$

тобто

$$\tau(x) = -\frac{1}{f'(x)}. \quad (37)$$

Тому

$$\varphi'(x) = 1 - \frac{f'(x) \cdot f'(x) - f(x) \cdot f''(x)}{(f'(x))^2} = \frac{f(x) \cdot f''(x)}{(f'(x))^2}. \quad (38)$$

Якщо  $\bar{x}$  — корінь  $f(x)$ , то  $\varphi'(\bar{x}) = 1$ . знайдеться окіл кореня, \end{equation}

$$|\varphi'(x)| = \left| \frac{f(x) \cdot f''(x)}{(f'(x))^2} \right| < 1. \quad (39)$$

Це означає, що збіжність методу Ньютона залежить від вибору  $x_0$ .

**Недолік** методу Ньютона: необхідність обчислювати на кожній ітерації не тільки значення функції, а й похідної.

Модифікований метод Ньютона позбавлений цього недоліку і має вигляд:

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_0)}, \quad k = 0, 1, 2, \dots \quad (40)$$

Цей метод має лише лінійну збіжність:  $|x_{k+1} - \bar{x}| = O(|x_k - \bar{x}|)$ .

**Задача 4:** Дати геометричну інтерпретацію модифікованого методу Ньютона.

В методі Ньютона, для якого  $f'(x_k)$  замінюється на

$$\frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}} \quad (41)$$

дає метод січних:

$$x_{k+1} = x_k - \frac{x_k - x_{k-1}}{f(x_k) - f(x_{k-1})} \cdot f(x_k), \quad (42)$$

де  $k = 1, 2, \dots$ ,  $x_0, x_1$  — задані.

**Задача 5:** Дати геометричну інтерпретацію методу січних.

## 2.5. Збіжність методу Ньютона

**Теорема 1:** Нехай  $f(x) \in C^2([a, b])$ ;  $\bar{x}$  простий дійсний корінь рівняння

$$f(x) = 0. \quad (43)$$

і  $f'(x) \neq 0$  при  $x \in U_r = \{x : |x - \bar{x}| < r\}$ . Якщо

$$q = \frac{M_2 \cdot |x_0 - \bar{x}|}{2m_1} < 1, \quad (44)$$

де

$$m_1 = \min_{U_r} |f'(x)|, \quad M_2 = \max_{U_r} |f''(x)|, \quad (45)$$

то для  $x_0 \in U_r$  метод Ньютона

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} \quad (46)$$

збігається і має місце оцінка

$$|x_n - \bar{x}| \leq q^{2^n - 1} \cdot |x_0 - \bar{x}|. \quad (47)$$

З (46) маємо

$$x_{k+1} - \bar{x} = x_k - \frac{f(x_k)}{f'(x_k)} - \bar{x} = \frac{(x_k - \bar{x}) \cdot f'(x_k) - f(x_k)}{f'(x_k)} = \frac{F(x_k)}{f'(x_k)}, \quad (48)$$

де  $F(x) = (x - \bar{x})f'(x) - f(x)$ , така, що

- $F(\bar{x}) = 0$ ;

- $F'(x) = (x - \bar{x}) \cdot f''(x).$

Тоді

$$F(x_k) = F(\bar{x}) + \int_x^{x_k} F'(t) dt = \int_x^{x_k} (t - \bar{x}) \cdot f''(t) dt. \quad (49)$$

Так як  $(t - \bar{x})$  не міняє знак на відрізку інтегрування, то скористаємося теоремою про середнє значення:

$$F(x_k) = f''(\xi_k) \int_x^{x_k} (t - \bar{x}) dt = \frac{(x_k - \bar{x})^2}{2} \cdot f''(\xi_k), \quad (50)$$

де  $\xi_k = \bar{x} + \theta_k \cdot (x_k - \bar{x})$ , де  $0 < \theta_k < 1$ . З (48), (50) маємо

$$x_{k+1} - \bar{x} = \frac{(x_k - \bar{x})^2}{2f'(\xi_k)} \cdot f''(\xi_k). \quad (51)$$

Доведемо оцінку (46) за індукцією. Так як  $x_0 \in U_r$ , то

$$|\xi_0 - \bar{x}| = |\theta_0 \cdot (x_0 - \bar{x})| < |\theta_0| \cdot |x_0 - \bar{x}| < r \quad (52)$$

звідси випливає  $\xi_0 \in U_r$ .

Тоді  $f''(\xi_0) \leq M_2$ , тому

$$|x_1 - \bar{x}| \leq \frac{(x_0 - \bar{x})^2 \cdot M_2}{2m_1} = \frac{M_2 \cdot |x_0 - \bar{x}|}{2m_1} \cdot |x_0 - \bar{x}| = q \cdot |x_0 - \bar{x}| < r, \quad (53)$$

тобто  $x_1 \in U_r$ .

Ми довели твердження (47) при  $n = 1$ . Нехай воно справджується при  $n = k$

$$|x_k - \bar{x}| \leq q^{2^k - 1} \cdot |x_0 - \bar{x}| < r, \quad (54)$$

$$|\xi_k - \bar{x}| = |\theta_k \cdot (x_k - \bar{x})| < r. \quad (55)$$

Тоді  $x_k, \xi_k \in U_r$ .

Доведемо (47) для  $n = k + 1$ . З (51) маємо

$$\begin{aligned} |x_{k+1} - \bar{x}| &\leq \frac{|x_k - \bar{x}|^2 \cdot M_2}{2m_1} \leq \left(q^{2^k - 1}\right)^2 \cdot \frac{|x_0 - \bar{x}|^2 \cdot M_2}{2m_1} = \\ &= q^{2^{k+1} - 2} \cdot \frac{|x_0 - \bar{x}| \cdot M_2}{2m_1} \cdot |x_0 - \bar{x}| = q^{2^{k+1} - 1} \cdot |x_0 - \bar{x}|. \end{aligned} \quad (56)$$

Таким чином (47) справджується для  $n = k + 1$ . Значить (47) виконується і для довільного  $n$ . Таким чином  $x_n \xrightarrow{n \rightarrow \infty} x$ .  $\square$

З (47) маємо оцінку кількості ітерацій для досягнення точності  $\varepsilon$

$$n \geq \left\lceil \log_2 \left( 1 + \frac{\ln \left( \frac{\varepsilon}{b-a} \right)}{\ln q} \right) \right\rceil + 1. \quad (57)$$

Кажуть, що ітераційний метод має *ступінь збіжності*  $m$ , якщо

$$|x_{k+1} - \bar{x}| = O(|x_k - \bar{x}|^m). \quad (58)$$

Для методу Ньютона

$$|x_{k+1} - \bar{x}| = \frac{|x_k - \bar{x}|^2 \cdot |f''(\xi_k)|}{2|f'(x_k)|}. \quad (59)$$

Звідси випливає, що

$$|x_{k+1} - \bar{x}| = O(|x_k - \bar{x}|^2). \quad (60)$$

Значить ступінь збіжності методу Ньютона  $m = 2$ . Для методу простої ітерації і ділення навпіл  $m = 1$ .

**Теорема 2:** Нехай  $f(x) \in C^2([a, b])$  та  $x$  простий корінь рівняння  $f(x) = 0$  ( $f'(x) \neq 0$ ). Якщо  $f'(x) \cdot f''(x) > 0$  ( $f'(x) \cdot f''(x) < 0$ ) то для методу Ньютона при  $x_0 = b$  послідовність наближень  $\{x_k\}$  монотонно спадає (монотонно зростає при  $x_0 = a$ ).

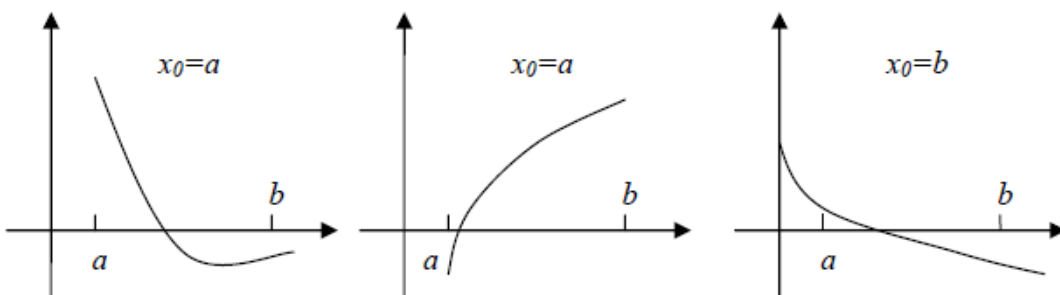
**Задача 6:** Довести [теорему 2](#) при

- $f'(x) \cdot f''(x) > 0$ ;
- $f'(x) \cdot f''(x) < 0$ .

**Задача 7:** Знайти ступінь збіжності методу січних [Калиткин Н.Н., Численные методы, с. 145–146]

Якщо  $f(a) \cdot f''(a) > 0$  та  $f''(x)$  не міняє знак, то потрібно вибирати  $x_0 = a$ ; при цьому  $\{x_k\} \uparrow \bar{x}$ .

Якщо  $f(b) \cdot f''(b) > 0$ , то  $x_0 = b$ ; маємо  $\{x_k\} \downarrow \bar{x}$ . Пояснення на рисунку 2:



**Зауваження 1:** Якщо  $\bar{x}$  —  $p$ -кратний корінь тобто

$$f^{(m)}(\bar{x}) = 0, \quad m = 0, 1, \dots, p-1; \quad f^{(p)}(\bar{x}) \neq 0, \quad (61)$$

то в методі Ньютона необхідна наступна модифікація

$$x_{k+1} = x_k - p \cdot \frac{f(x_k)}{f'(x_k)} \quad (62)$$

і

$$q = \frac{M_{p+1} \cdot |x_0 - \bar{x}|}{m_p \cdot p \cdot (p+1)} < 1. \quad (63)$$

**Зауваження 2:** Метод Ньютона можна застосовувати і для обчислення комплексного кореня

$$z_{k+1} = z_k - \frac{f(z_k)}{f'(z_k)} \quad (64)$$

В теоремі про збіжність

$$q = \frac{|x_0 - \bar{x}| M_2}{2m_1}, \quad (65)$$

де

$$m_1 = \min_{U_r} |f'(z)|, \quad M_2 = \max_{U_r} |f''(z)|. \quad (66)$$

Тут  $|z|$  — модуль комплексного числа.

**Переваги** методу Ньютона:

- висока швидкість збіжності;
- узагальнюється на системи рівнянь;
- узагальнюється на комплексні корені.

**Недоліки** методу Ньютона:

- на кожній ітерації обчислюється не тільки  $f(x_k)$ , а і похідна  $f'(x_k)$ ;
- збіжність залежить від початкового наближення  $x_0$ , оскільки від нього залежить умова збіжності

$$q = \frac{M_2 |x_0 - \bar{x}|}{2m_1} < 1; \quad (67)$$

- потрібно, щоб  $f(x) \in C^2([a, b])$ .

### 3. Методи розв'язання систем лінійних алгебраїчних рівнянь (СЛАР)

Методи розв'язування СЛАР поділяються на *прямі* та *ітераційні*. При умові точного виконання обчислень прямі методи за скінчену кількість операцій в результаті дають точний розв'язок. Використовуються вони для невеликих та середніх СЛАР  $n = 10^2 - 10^4$ . Ітераційні методи використовуються для великих СЛАР  $n > 10^5$ , як правило розріджених. В результаті отримуємо послідовність наближень, яка збігається до розв'язку.

#### 3.1. Метод Гауса

Розглянемо задачу розв'язання СЛАР

$$A\vec{x} = \vec{b}, \quad (1)$$

причому  $A = (a_{ij})_{i,j=1}^n$ ,  $\det A \neq 0$ ,  $\vec{x} = (x_i)_{i=1}^n$ ,  $\vec{b} = (b_j)_{j=1}^n$ . Метод Крамера з обчисленням визначників для такої системи має складність  $Q = O(n! \cdot n)$ .

Запишемо СЛАР у вигляді

$$\begin{cases} a_{1,1}x_1 + a_{1,2}x_2 + \dots + a_{1,n}x_n = b_1 \equiv a_{1,n+1}, \\ a_{2,1}x_1 + a_{2,2}x_2 + \dots + a_{2,n}x_n = b_2 \equiv a_{2,n+1}, \\ \dots \\ a_{n,1}x_1 + a_{n,2}x_2 + \dots + a_{n,n}x_n = b_n \equiv a_{n,n+1}. \end{cases} \quad (2)$$

Якщо  $a_{1,1} \neq 0$ , то ділимо перше рівняння на нього і виключаємо  $x_1$  з інших рівнянь:

$$\begin{cases} x_1 + a_{1,2}^{(1)}x_2 + \dots + a_{1,n}^{(1)}x_n = a_{1,n+1}^{(1)}, \\ a_{2,2}^{(1)}x_2 + \dots + a_{2,n}^{(1)}x_n = a_{2,n+1}^{(1)}, \\ \dots \\ a_{n,2}^{(1)}x_2 + \dots + a_{n,n}^{(1)}x_n = a_{n,n+1}^{(1)}. \end{cases} \quad (3)$$

Процес повторюємо для  $x_2, \dots, x_n$ . В результаті отримуємо систему з трикутною матрицею

$$\begin{cases} x_1 + a_{1,2}^{(1)}x_2 + \dots + a_{1,n}^{(1)}x_n = a_{1,n+1}^{(1)}, \\ x_2 + \dots + a_{2,n}^{(2)}x_n = a_{2,n+1}^{(2)}, \\ \dots \\ x_n = a_{n,n+1}^{(n)}. \end{cases} \quad (4)$$

Тобто

$$A^{(n)}\vec{x} = \vec{a}^{(n)}. \quad (5)$$

Це прямий хід методу Гауса. Формули прямого ходу

```

for k in range(1, n):
    for j in range(k + 1, n + 2):
        a[k, j][k] = a[k, j][k - 1] / a[k, k][k - 1]
    for i in range(k + 1, n + 1):
        a[i, j][k] = a[i, j][k - 1] - \
            a[i, j][k - 1] * a[k, j][k]

```

Звідси

$$x_n = a_{n,n+1}^{(n)}, \quad x_i = a_{i,n+1}^{(i)} - \sum_{j=i+1}^n a_{i,j}^{(n)} x_j, \quad (6)$$

для  $i = \overline{n-1, 1}$ . Це формули оберненого ходу.

Складність, тобто кількість операцій, яку необхідно виконати для реалізації методу:  $Q_f = 2/3n^2 + O(n^2)$  для прямого ходу,  $Q_b = n^2 + O(n)$  для оберненого ходу.

Умова

$$a_{k,k}^{(k-1)} \neq 0 \quad (7)$$

не суттєва, оскільки знайдеться  $m$ , для якого

$$\left| a_{m,k}^{(k-1)} \right| = \max_i \left| a_{i,k}^{(k-1)} \right| \neq 0 \quad (8)$$

(оскільки  $\det A \neq 0$ ). Тоді міняємо місцями рядки номерів  $k$  і  $m$ .

**Означення:** Елемент

$$a_{k,k}^{(k-1)} \neq 0 \quad (9)$$

називається *ведучим*.

Введемо матриці

$$M_k = \begin{pmatrix} 1 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & m_{k,k} & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & m_{n,k} & \cdots & 1 \end{pmatrix} \quad (10)$$

елементи якої обчислюється так:

$$m_{k,k} = \frac{1}{a_{k,k}^{(k-1)}}, \quad m_{i,k} = -\frac{a_{i,k}^{(k-1)}}{a_{k,k}^{(k-1)}}. \quad (11)$$

Нехай на  $k$ -му кроці  $A_{k-1}\vec{x} = \vec{b}_{k-1}$ . Множимо цю СЛАР зліва на  $M_k$ :  $M_k A_{k-1}\vec{x} = M_k \vec{b}_{k-1}$ . Позначимо  $A_k = M_k A_{k-1}$ ;  $A_0 = A$ . Тоді прямий хід методу Гауса можна записати у вигляді

$$M_n M_{n-1} \dots M_1 A \vec{x} = M_n M_{n-1} \dots M_1 \vec{b}. \quad (12)$$

Позначимо останню систему, яка співпадає з (5), так

$$U \vec{x} = \vec{c}, \quad U = (u_{i,j})_{i,j=1}^n, \quad (13)$$

причому

$$\begin{cases} u_{i,i} = 1, \\ u_{i,j} = 0, \quad i > j. \end{cases} \quad (14)$$

Таким чином  $U = M_n M_{n-1} \dots M_1 A$ . Введемо матриці

$$L_k = M_k^{-1} = \begin{pmatrix} 1 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & a_{k,k}^{(k-1)} & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & a_{n,k}^{(k-1)} & \dots & 1 \end{pmatrix} \quad (15)$$

Тоді

$$A = L_1 \dots L_n U = LU; \quad L = L_1 \dots L_n, \quad (16)$$

де  $L$  — нижня трикутна матриця,  $U$  — верхня трикутна матриця. Таким чином метод Гауса можна трактувати, як розклад матриці  $A$  в добуток двох трикутних матриць —  $LU$ -розклад.

Введемо матрицю перестановок на  $k$ -му кроці (це матриця, отримана з одиничної матриці перестановкою  $k$ -того і  $m$ -того рядка). Тоді при множенні на неї матриці  $A_{k-1}$  робимо ведучим елементом максимальний за модулем.

$$P_k = \begin{pmatrix} 1 & \dots & 0 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 1 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & 1 & \dots & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & \dots & 0 & \dots & 1 \end{pmatrix} \quad (17)$$

За допомогою цих матриць перехід до трикутної системи (13) тепер має вигляд:

$$M_n M_{n-1} P_{n-1} \dots M_1 P_1 A \vec{x} = M_n M_{n-1} P_{n-1} \dots M_1 P_1 \vec{b}. \quad (18)$$



**Твердження:** Знайдеться така матриця  $P$  перестановок, що  $PA = LU$  — розклад матриці на нижню трикутну з ненульовими діагональними елементами і верхню трикутну матрицю з одиницями на діагоналі.

Висновки про **переваги** трикутного розкладу:

- Розділення прямого і оберненого ходів дає змогу економно розв'язувати декілька систем з одноковою матрицею та різними правими частинами.
- Зберігання  $M$ , або  $L$  та  $U$  на місці  $A$ .
- Обчислюючи  $\ell$  — кількість перестановок, можна встановити знак визначника.

### 3.2. Метод квадратних коренів

Цей метод призначений для розв'язання систем рівнянь із симетричною матрицею

$$A\vec{x} = \vec{b}, \quad A^T = A. \quad (19)$$

Він оснований на розкладі матриці  $A$  в добуток:

$$A = S^T D S, \quad (20)$$

де  $S$  — верхня трикутна матриця,  $S^T$  — нижня трикутна матриця,  $D$  — діагональна матриця.

Виникає питання: як обчислити  $S$ ,  $D$  по матриці  $A$ ? Маємо

$$DS_{i,j} = \begin{cases} d_{i,i}s_{i,j}, & i \leq j, \\ 0, & i > j. \end{cases} \quad (21)$$

Далі

$$S^T D S_{i,j} = \sum_{l=1}^n s_{i,l}^T d_{l,l} s_{l,j} = \sum_{l=1}^{i-1} s_{l,i}^T s_{l,j} d_{l,l} + s_{i,i} s_{i,j} d_{i,i} + \underbrace{\sum_{l=i+1}^n s_{l,i}^T s_{l,j} d_{l,l}}_{=0} = a_{i,j}, \quad (22)$$

для  $i, j = \overline{1, n}$ .

Якщо  $i = j$ , то

$$|s_{i,i}^2| d_{i,i} = a_{i,i} - \sum_{l=1}^{i-1} |s_{l,i}^2| d_{l,l} \equiv p_i. \quad (23)$$

Тому

$$d_{i,i} = \text{sign}(p_i), \quad s_{i,i} = \sqrt{|p_i|}. \quad (24)$$

Якщо  $i < j$ , то

$$s_{i,j} = \left( a_{i,j} - \sum_{l=1}^{i-1} s_{l,i}^T d_{l,j} s_{l,j} \right) / (s_{i,i} d_{i,i}), \quad (25)$$

де  $i = \overline{1, n}$ , а  $j = \overline{i+1, n}$ .

Якщо  $A > 0$  (тобто головні мінори матриці  $A$  додатні), то всі  $d_{i,i} = +1$ .

Знайдемо розв'язок рівняння (19). Враховуючи (20), маємо:

$$S^T D \vec{y} = \vec{b} \quad (26)$$

і

$$S \vec{x} = \vec{y} \quad (27)$$

Оскільки  $S$  — верхня трикутна матриця, а  $S^T D$  — нижня трикутна матриця, то

$$y_i = \frac{b_i - \sum_{j=1}^{i-1} s_{j,i} d_{j,j} y_j}{s_{i,i} d_{i,i}}, \quad (28)$$

для  $i = \overline{1, n}$  і

$$x_i = \frac{y_i - \sum_{j=1}^{i-1} s_{i,j} x_j}{s_{i,i}}, \quad (29)$$

для  $i = \overline{n-1, 1}$ , де  $x_n = y_n / s_{n,n}$ .

Метод застосовується лише для симетричних матриць. Його складність  $Q = n^3/3 + O(n^2)$ .

**Переваги** цього методу:

- він витрачає в 2 рази менше пам'яті ніж метод Гауса для зберігання  $A^T = A$  (необхідний об'єм пам'яті  $n(n+1)/2 \sim n^2/2$ ;
- метод однорідний, без перестановок;
- якщо матриця  $A$  має багато нульових елементів, то і матриця  $S$  також.

Остання властивість дає економію в пам'яті та кількості арифметичних операцій. Наприклад, якщо  $A$  має  $m$  ненульових стрічок по діагоналі ( $m$ -діагональна), то  $Q = O(m^2 n)$ .

### 3.3. Обчислення визначника та оберненої матриці

Кількість операцій обчислення детермінанту за означенням —  $Q_{\det} = n!$ . В методі Гауса —  $PA = LU$ . Тому

$$\det P \det A = \det L \det U \quad (30)$$

звідки

$$\det A = (-1)^\ell \det L \det U = (-1)^\ell \prod_{k=1}^n a_{k,k}^{(k)}, \quad (31)$$

де  $\ell$  — кількість перестановок. Ясно, що за методом Гауса

$$Q_{\det} = \frac{2}{3} \cdot n^3 + O(n^2) \quad (32)$$

В методі квадратного кореня  $A = S^\top D S$ . Тому

$$\det A = \det S^\top \det D \det S = \prod_{k=1}^n d_{k,k} \prod_{k=1}^n s_{k,k}^2. \quad (33)$$

Тепер  $Q_{\det} = n^3/3 + O(n^2)$ .

За означенням

$$A A^{-1} = E, \quad (34)$$

де  $A^{-1}$  обернена до матриці  $A$ . Позначимо

$$A^{-1} = (\alpha_{i,j})_{i,j=1}^n. \quad (35)$$

Тоді  $\vec{\alpha}_j = (\alpha_{i,j})_{i=1}^n$  — вектор-стовпчик оберненої матриці. З (34) маємо

$$A \vec{\alpha}_j = \vec{e}_j, \quad j = \overline{1, n}. \quad (36)$$

де  $\vec{e}_j$  — стовпчики одиничної матриці:  $\vec{e}_j = (\delta_{i,j})_{i=1}^n$ ,

$$\delta_{i,j} = \begin{cases} 1, & i = j, \\ 0, & i \neq j. \end{cases} \quad (37)$$

Для знаходження  $A^{-1}$  необхідно розв'язати  $n$  систем. Для знаходження  $A^{-1}$  методом Гауса необхідна кількість операцій  $Q = 2n^3 + O(n^2)$ .

### 3.4. Метод прогонки

Це економний метод для розв'язання СЛАР з три діагональною матрицею:

$$\begin{cases} -c_0 y_0 + b_0 y_1 = -f_0, \\ a_i y_{i-1} - c_i y_i + b_i y_{i+1} = -f_i, \\ a_N y_{N-1} - c_N y_N = -f_N. \end{cases} \quad (38)$$

Матриця системи

$$A = \begin{pmatrix} -c_0 & b_0 & & 0 \\ a_0 & \ddots & \ddots & \\ & \ddots & \ddots & b_N \\ 0 & & a_N & -c_N \end{pmatrix} \quad (39)$$

тридіагональна.

Розв'язок представимо у вигляді

$$y_i = \alpha_{i+1}y_{i+1} + \beta_{i+1}, \quad i = \overline{0, N-1}. \quad (40)$$

Замінімо в (40) і  $i \mapsto i-1$  і підставимо в (33), тоді

$$(a_i\alpha_i - c_i) \cdot y_i + b_i y_{i+1} = -f_i - a_i\beta_i \quad (41)$$

Звідси

$$y_i = \frac{b_i}{c_i - a_i\alpha_i} \cdot y_{i+1} + \frac{f_i + a_i\beta_i}{c_i - a_i\alpha_i}. \quad (42)$$

Тому з (36)

$$\alpha_{i+1} = \frac{b_i}{c_i - a_i\alpha_i}, \quad \beta_{i+1} = \frac{f_i + a_i\beta_i}{c_i - a_i\alpha_i}, \quad i = \overline{1, N-1}. \quad (43)$$

Умова розв'язності (38) —  $c_i - a_i\alpha_i \neq 0$ .

Щоб знайти всі  $\alpha_i, \beta_i$ , треба задати перші значення. З (38):

$$\alpha_1 = \frac{b_0}{c_0}, \quad \beta_1 = \frac{f_0}{c_0}. \quad (44)$$

Після знаходження всіх  $\alpha_i, \beta_i$  обчислюємо  $y_N$  з системи

$$\begin{cases} a_N y_N - c_N y_N = -f_N, \\ y_{N-1} = \alpha_N y_N + \beta_N. \end{cases} \quad (45)$$

Звідси

$$y_N = \frac{f_N + a_N\beta_N}{c_N - a_N\alpha_N}. \quad (46)$$

**Алгоритм:**

```
alpha[1], beta[1] = b[0] / c[0], f[0] / c[0]

for i in range(1, N):
    z[i] = c[i] - a[i] * alpha[i]
    alpha[i + 1], beta[i + 1] = b[i] / z[i], \
        (f[i] + a[i] * beta[i]) / z[i]

y[N] = (f[N] + a[N] * beta[N]) / \
    (c[N] - a[N] * alpha[N])

for i in range(N - 1, -1, -1):
    y[i] = alpha[i + 1] * y[i + 1] + beta[i + 1]
```

Складність алгоритму  $Q = 8N - 2$ .

Метод можна застосовувати, коли  $c_i - a_i \alpha_i \neq 0, \forall i : |\alpha_i| \leq 1$ . Якщо  $|\alpha_i| \geq q > 1$  то  $\Delta y_0 \geq q^N \Delta y_N$  (тут  $\Delta y_i$  абсолютна похибка обчислення  $y_i$ ), а це приводить до експоненціального накопичення похибок заокруглення, тобто нестійкості алгоритму прогонки.

**Теорема** (про достатні умови стійкості метода прогонки): Нехай  $a_i, b_i \neq 0$ , та

$$|c_i| \geq |a_i| + |b_i|, \quad \forall i, \quad a_0 = b_N = 0, \quad (47)$$

та хоча би одна нерівність строга. Тоді  $|\alpha_i| \leq 1$  та

$$z_i = c_i - a_i \alpha_i \neq 0, \quad i = \overline{1, N}. \quad (48)$$

**Задача 8:** Довести теорему про стійкість методу прогонки.

### 3.5. Обумовленість систем лінійних алгебраїчних рівнянь

Нехай задано СЛАР

$$A\vec{x} = \vec{b}. \quad (49)$$

Припустимо, що матриця і права частина системи задані неточно і фактично розв'язуємо систему

$$B\vec{y} = \vec{h}, \quad (50)$$

де  $B = A + C, \vec{h} = \vec{b} + \vec{\eta}, \vec{y} = \vec{x} + \vec{z}$ .

Малість детермінанту  $\det A \ll 1$  не є необхідною умовою різкого збільшення похибки. Це ілюструє наступний приклад:

$$A = \text{diag}(\varepsilon), \quad a_{i,j} = \varepsilon \delta_{i,j}. \quad (51)$$

Тоді  $\det A = \varepsilon^n \ll 1$ , але  $x_i = b_i/\varepsilon$ . Тому  $\Delta x_i = \Delta b_i/\varepsilon \gg 1$ .

Оцінимо похибку розв'язку. Підставивши значення  $B, \vec{h}$ , та  $\vec{z} = \vec{y} - \vec{x}$ , отримаємо:

$$(A + C)(\vec{x} + \vec{z}) = \vec{b} + \vec{\eta}. \quad (52)$$

Віднімемо від цієї рівності (49) у вигляді  $A\vec{z} + C\vec{x} + C\vec{z} = \vec{\eta}$ . Тоді

$$A\vec{z} = \vec{\eta} - C\vec{x} - C\vec{z}, \quad \vec{z} = A^{-1}(\vec{\eta} - C\vec{x} - C\vec{z}). \quad (53)$$

**Означення:** Введемо норми векторів:  $\|\vec{z}\|$ :

$$|\vec{z}|_1 = \sum_{i=1}^n |z_i|, \quad (54)$$

$$|\vec{z}|_2 = \left( \sum_{i=1}^n |z_i|^2 \right)^{1/2}, \quad (55)$$

$$|\vec{z}|_\infty = \max_i |z_i|. \quad (56)$$

**Означення:** Норми матриці, що відповідають нормам вектора, тобто

$$|A|_m = \sup_{|\vec{x}|_m \neq 0} \frac{|A\vec{x}|_m}{|\vec{x}|_m}, \quad m = 1, 2, \infty. \quad (57)$$

такі:

$$|A|_1 = \max_j \sum_{i=1}^n |a_{i,j}|, \quad (58)$$

$$|A|_2 = \max_i \sqrt{\lambda_i(A^T A)}, \quad (59)$$

$$|A|_\infty = \max_i \sum_{j=1}^n |a_{i,j}|, \quad (60)$$

де  $\lambda_i(B)$  — власні значення матриці  $B$ .

Позначимо  $\delta(\vec{x}) = \|\vec{z}\|/\|\vec{x}\|$ ,  $\delta(\vec{b}) = \|\vec{\eta}\|/\|\vec{b}\|$ ,  $\delta(A) = \|C\|/\|A\|$  — відносні похибки  $\vec{x}$ ,  $\vec{b}$ ,  $A$ , де  $\|\cdot\|_k$  — одна з введених вище норм.

Для характеристики зв'язку між похибками правої частини та розв'язку вводять поняття обумовленості матриці системи.

**Означення:** Число обумовленості матриці  $A$  —  $\text{cond}(A) = \|A\| \cdot \|A^{-1}\|$ .

**Теорема:** Якщо  $\exists A^{-1}$  та  $\|A^{-1}\| \cdot \|C\| < 1$ , то

$$\delta(\vec{x}) \leq \frac{\text{cond}(A)}{1 - \text{cond}(A) \cdot \delta(A)} (\delta(A) + \delta(\vec{b})), \quad (61)$$

де  $\text{cond}(A)$  — число обумовленості.

*Доведення:*

$$A\vec{z} = \vec{\eta} - C\vec{x} - C\vec{z}, \quad \vec{z} = A^{-1}\vec{\eta} - A^{-1}C\vec{x} - A^{-1}C\vec{z} \quad (62)$$

$$|\vec{z}| \leq |A^{-1}\vec{\eta}| + |A^{-1}C\vec{x}| + |A^{-1}C\vec{z}| \leq |A^{-1}| \cdot |\vec{\eta}| + |A^{-1}| \cdot |C| \cdot |\vec{x}| + |A^{-1}| \cdot |C| \cdot |\vec{z}|. \quad (63)$$

$$|\vec{z}| \leq \frac{|A^{-1}| \cdot (|\vec{\eta}| + |C| \cdot |\vec{x}|)}{1 - |A^{-1}| \cdot |C|} \quad (64)$$

Оцінка похибки

$$\begin{aligned} \delta(\vec{x}) &\leq \frac{|A^{-1}|}{1 - |A^{-1}| \cdot |C|} \left( \frac{|\vec{\eta}|}{|\vec{x}|} + |C| \right) = \frac{|A^{-1}| \cdot |A|}{1 - |A^{-1}| \cdot |A| \cdot \frac{|C|}{|A|}} \left( \frac{|\vec{\eta}|}{|A| \cdot |\vec{x}|} + \delta(A) \right) \leq \\ &\leq \frac{\text{cond}(A)}{1 - \text{cond}(A) \cdot \delta(A)} \left( \frac{|\vec{\eta}|}{|\vec{x}|} + \delta(A) \right). \quad \square \end{aligned} \quad (65)$$

**Наслідок:** Якщо  $C \equiv 0$ , то  $\delta(\vec{x}) \leq \text{cond}(A) \cdot \delta(\vec{b})$ .

**Властивості**  $\text{cond}(A)$ :

- $\text{cond}(A) \geq 1$ ;
- $\text{cond}(A) \geq \max_i |\lambda_i(A)| / \min_i |\lambda_i(A)|$ ;
- $\text{cond}(AB) \leq \text{cond}(A) \cdot \text{cond}(B)$ ;
- $A^\top = A^{-1} \implies \text{cond}(A) = 1$ .

Друга властивість має місце оскільки довільна норма матриці не менше її найбільшого за модулем власного значення. Значить  $\|A\| \geq \max |\lambda_A|$ . Оскільки власні значення матриць  $A^{-1}$  та  $A$  взаємно обернені, то

$$|A^{-1}| \geq \max \frac{1}{|\lambda_A|} = \frac{1}{\min |\lambda_A|}. \quad (66)$$

Якщо  $1 \ll \text{cond}(A)$ , то система називається *погано обумовленою*.

Оцінка впливу похибок заокруглення при обчисленні розв'язку СЛАР така (Дж. Уілкінсон):

$\delta(A) = O(n\beta^{-t})$ ,  $\delta(\vec{b}) = O(\beta^{-t})$ , де  $\beta$  — розрядність ЕОМ,  $t$  — кількість розрядів, що відводиться під мантису числа. З оцінки (61) витікає:  $\delta(\vec{x}) = \text{cond}(A) \cdot O(n\beta^{-t})$ . Висновок: найпростіший спосіб підвищити точність обчислення розв'язку погано обумовленої СЛАР — збільшити розрядність ЕОМ при обчисленнях. Інші способи пов'язані з розглядом цієї СЛАР як некоректної задачі із застосуванням відповідних методів її розв'язання.

**Приклад** погано обумовленої системи — системи з матрицею Гільберта

$$H_n = \left( \frac{1}{i+j-1} \right)_{i,j=1}^n, \quad (67)$$

наприклад  $\text{cond}(H_8) \approx 10^9$ .

## 4. Ітераційні методи для систем

### 4.1. Ітераційні методи розв'язання СЛАР

Систему

$$A\vec{x} = \vec{b} \quad (1)$$

зводимо до вигляду

$$\vec{x} = B\vec{x} + \vec{f}. \quad (2)$$

Будь яка система

$$\vec{x} = \vec{x} - C \cdot (A\vec{x} - \vec{b}) \quad (3)$$

має вигляд (2) і при  $\det C \neq 0$  еквівалентна системі (1). Наприклад, для  $C = \tau \cdot E$ :

$$\vec{x} = \vec{x} - \tau \cdot (A\vec{x} - \vec{b}). \quad (3')$$

#### 4.1.1. Метод простої ітерації

Цей метод застосовується до рівняння (2)

$$\vec{x}^{(k+1)} = B\vec{x}^{(k)} + \vec{f}, \quad (4)$$

де  $\vec{x}^{(0)}$  — початкове наближення, задано.

**Теорема:** Ітераційний процес збігається, тобто

$$\left| \vec{x}^{(k)} - \vec{x} \right| \xrightarrow[k \rightarrow \infty]{} 0, \quad (5)$$

якщо

$$|B| \leq q < 1. \quad (6)$$

При цьому має місце оцінка

$$\left| \vec{x}^{(n)} - \vec{x} \right| \leq \frac{q^n}{1 - q} \cdot \left| \vec{x}^{(1)} - \vec{x}^{(0)} \right|. \quad (7)$$

#### 4.1.2. Метод Якобі



Припустимо  $\forall i: a_{i,i} \neq 0$ . Зведемо систему (1) до вигляду

$$x_i = - \sum_{j=1}^{i-1} \frac{a_{i,j}}{a_{i,i}} \cdot x_j - \sum_{j=i+1}^n \frac{a_{i,j}}{a_{i,i}} \cdot x_j + \frac{b_i}{a_{i,i}}, \quad (8)$$

де  $i = \overline{1, n}$ .

Ітераційний процес запишемо у вигляді

$$x_i^{(k+1)} = - \sum_{j=1}^{i-1} \frac{a_{i,j}}{a_{i,i}} \cdot x_j^{(k)} - \sum_{j=i+1}^n \frac{a_{i,j}}{a_{i,i}} \cdot x_j^{(k)} + \frac{b_i}{a_{i,i}}, \quad (9)$$

де  $k = 0, 1, \dots$ , а  $i = \overline{1, n}$ .

**Теорема:** Ітераційний процес збігається до розв'язку, якщо виконується умова

$$\forall i : \sum_{\substack{j=1 \\ i \neq j}}^n |a_{i,j}| \leq |a_{i,i}|. \quad (10)$$

Це умова діагональної переваги матриці  $A$ .

**Теорема:** Якщо ж

$$\forall i : \sum_{\substack{j=1 \\ i \neq j}}^n |a_{i,j}| \leq q \cdot |a_{i,i}|, \quad 0 \leq q < 1. \quad (11)$$

то має місце оцінка точності:

$$|\vec{x}^{(n)} - \vec{x}| \leq \frac{q^n}{1 - q} \cdot |\vec{x}^{(0)} - \vec{x}|. \quad (12)$$

#### 4.1.3. Метод Зейделя

В компонентному вигляді ітераційний метод Зейделя записується так:

$$x_i^{(k+1)} = - \sum_{j=1}^{i-1} \frac{a_{i,j}}{a_{i,i}} \cdot x_j^{(k+1)} - \sum_{j=i+1}^n \frac{a_{i,j}}{a_{i,i}} \cdot x_j^{(k)} + \frac{b_i}{a_{i,i}}, \quad (13)$$

де  $k = 0, 1, \dots$ , а  $i = \overline{1, n}$ .

На відміну від методу Якобі на  $k$ -му-кроці попередні компоненти розв'язку беруться з  $(k + 1)$ -ої ітерації.

**Теорема:** Достатня умова збіжності методу Зейделя —  $A^T = A > 0$ .

#### 4.1.4. Матрична інтерпретація методів Якобі і Зейделя

Подано матрицю  $A$  у вигляді

$$A = A_1 + D + A_2, \quad (14)$$

де  $A_1$  — нижній трикутник матриці  $A$ ,  $A_2$  — верхній трикутник матриці  $A$ ,  $D$  — її діагональ. Тоді систему (1) запишемо у вигляді

$$D\vec{x} = A_1\vec{x} + A_2\vec{x} + \vec{b}, \quad (15)$$

або

$$\vec{x} = D^{-1}A_1\vec{x} + D^{-1}A_2\vec{x} + D^{-1}\vec{b}, \quad (16)$$

Матричний запис методу Якобі:

$$\vec{x}^{(k+1)} = D^{-1}A_1\vec{x}^{(k)} + D^{-1}A_2\vec{x}^{(k)} + D^{-1}\vec{b}, \quad (17)$$

методу Зейделя:

$$\vec{x}^{(k+1)} = D^{-1}A_1\vec{x}^{(k+1)} + D^{-1}A_2\vec{x}^{(k)} + D^{-1}\vec{b}, \quad (18)$$

**Теорема:** Необхідна і достатня умова збіжності методу Якобі: всі корені рівняння

$$\det(D + \lambda(A_1 + A_2)) = 0 \quad (19)$$

по модулю більше 1.

**Теорема:** Необхідна і достатня умова збіжності методу Зейделя: всі корені рівняння

$$\det(A_1 + D + \lambda A_2) = 0 \quad (20)$$

по модулю більше 1.

#### 4.1.5. Однокрокові (двошарові) ітераційні методи

Канонічною формою однокрокового ітераційного методу розв'язку СЛАР є його запис у вигляді

$$B_k \frac{\vec{x}^{(k+1)} - \vec{x}^{(k)}}{\tau_{k+1}} + A\vec{x}^{(k)} = \vec{b}, \quad (21)$$

Тут  $\{B_k\}$  — послідовність матриць (пере-обумовлюючі матриці), що задають ітераційний метод на кожному кроці;  $\{\tau_{k+1}\}$  — ітераційні параметри.

**Означення:** Якщо  $B_k = E$ , то ітераційний процес називається *явним*

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} - \tau_{k+1} (A\vec{x}^{(k)} + \vec{b}). \quad (22)$$

**Означення:** Якщо  $B_k \neq E$ , то ітераційний процес називається *неявним*

$$B_k \vec{x}^{(k+1)} = F^k. \quad (23)$$

У цьому випадку на кожній ітерації необхідно розв'язувати СЛАР.

**Означення:** Якщо  $\tau_{k+1} \equiv \tau$ ,  $B_k \equiv B$ , то ітераційний процес називається *стаціонарним*; інакше — *нестабілізованим*.

Методам, що розглянуті вище відповідають:

- методу простої ітерації:  $B_k = E$ ,  $\tau_{k+1} = \tau$ ;
- методу Якобі:  $B_k = D$ ,  $\tau_{k+1} = 1$ ;
- методу Зейделя:  $B_k = D + A_1$ ,  $\tau_{k+1} = 1$ .

#### 4.1.6. Збіжності стаціонарних ітераційних процесів у випадку симетричних матриць

Розглянемо випадок симетричних матриць  $A^T = A$  і стаціонарний ітераційний процес  $B_k \equiv E$ ,  $\tau_{k+1} \equiv \tau$ .

Нехай для  $A$  справедливі нерівності

$$\gamma_1 E \leq A \leq \gamma_2 E, \quad \gamma_1, \gamma_2 > 0. \quad (24)$$

Тоді при виборі  $\tau = \tau_0 = \frac{2}{\gamma_1 + \gamma_2}$  ітераційний процес збігається. Найбільш точним значенням  $\gamma_1, \gamma_2$  при яких виконуються обмеження (24) є  $\gamma_1 = \min \lambda_i(A)$ ,  $\gamma_2 = \max \lambda_i(A)$ . Тоді

$$q = q_0 = \frac{\gamma_2 - \gamma_1}{\gamma_2 + \gamma_1} = \frac{1 - \xi}{1 + \xi}, \quad \xi = \frac{\gamma_1}{\gamma_2}. \quad (25)$$

і справедлива оцінка

$$\left| \vec{x}^{(n)} - \vec{x} \right| \leq \frac{q^n}{1-q} \cdot \left| \vec{x}^{(0)} - \vec{x} \right|. \quad (26)$$

**Зауваження:** аналогічно обчислюється  $q$  і для методу релаксації розв'язання нелінійних рівнянь, де  $\gamma_1 = m = \min |f'(x)|$ ,  $\gamma_2 = M_1 = \max |f'(x)|$ .

Явний метод з багатьма параметрами  $\{\tau_k\}$ :

$$B \equiv E, \quad \tau_k : \min_{\tau} q(\tau), \quad n = n(\varepsilon) \rightarrow \min, \quad (27)$$

які обчислюються за допомогою нулів багаточлена Чебишова, називаються ітераційним методом з чебишевським набором параметрів.

#### 4.1.7. Метод верхньої релаксації

Узагальненням методу Зейделя є метод верхньої релаксації:

$$(D + \omega A_1) \cdot \frac{\vec{x}^{(k+1)} + \vec{x}^{(k)}}{\omega} + A\vec{x}^{(k)} = \vec{b}, \quad (28)$$

де  $D$  — діагональна матриця з елементами  $a_{i,i}$  по діагоналі.  $\omega > 0$  — заданий числовий параметр.

Тепер  $B = D + \omega A_1$ ,  $\tau = \omega$ . Якщо  $A^T = A > 0$ , то метод верхньої релаксації збігається при умові  $0 < \omega < 2$ . Параметр підбирається експериментально з умови мінімальної кількості ітерацій.

#### 4.1.8. Методи варіаційного типу

До цих методів відносяться: метод мінімальних нев'язок, метод мінімальних поправок, метод найшвидшого спуску, метод спряжених градієнтів. Вони дозволяють обчислювати наближення без використання апіорної інформації про  $\gamma_1$ ,  $\gamma_2$  в (24).

Нехай  $B = E$ . Для методу мінімальних нев'язок параметри  $\tau_{k+1}$  обчислюються з умови

$$\left| \vec{r}^{(k+1)} \right|^2 = \left| \vec{r}^{(k)} \right|^2 - 2\tau_{k+1} \cdot \left\langle \vec{r}^{(k)}, A\vec{r}^{(k)} \right\rangle + \tau_{k+1}^2 \cdot \left| A\vec{r}^{(k)} \right|^2 \rightarrow \min. \quad (29)$$

Тому

$$\tau_{k+1} = \frac{\left\langle A\vec{r}^{(k)}, \vec{r}^{(k)} \right\rangle}{\left| \vec{r}^{(k)} \right|^2}, \quad (30)$$

де  $\vec{r}^{(k)} = A\vec{x}^{(k)} - \vec{b}$  — нев'язка.

Умова для завершення ітераційного процесу:

$$\left| \vec{r}^{(n)} \right| < \varepsilon. \quad (31)$$

Швидкість збіжності цього методу співпадає із швидкістю методу простої ітерації з одним оптимальним параметром  $\tau_0 = \frac{2}{\gamma_1 + \gamma_2}$ .

Аналогічно будуються методи з  $B \neq E$ . Матриця  $B$  називається переобумовлювачем і дозволяє підвищити швидкість збіжності ітераційного процесу. Його вибирають з умов

- легко розв'язувати СЛАР  $B\vec{x}^{(k)} = F_k$  (діагональний, трикутний, добуток трикутних та інше);
- зменшення числа обумовленості матриці  $B^{-1}A$  у порівнянні з  $A$ .

## 4.2. Методи розв'язання нелінійних систем

Розглянемо систему рівнянь

$$\begin{cases} f_1(x_1, \dots, x_n) = 0, \\ \dots \\ f_n(x_1, \dots, x_n) = 0. \end{cases} \quad (32)$$

Перепишемо її у векторному вигляді:

$$\vec{f}(\vec{x}) = 0. \quad (33)$$

### 4.2.1. Метод простої ітерації

В цьому методі рівняння (33) зводиться до еквівалентного вигляду

$$\vec{x} = \vec{\Phi}(\vec{x}). \quad (34)$$

Ітераційний процес представимо у вигляді:

$$\vec{x}^{(k+1)} = \vec{\Phi}(\vec{x}^{(k)}). \quad (35)$$

початкове наближення  $\vec{x}^{(0)}$  — задано.

Нехай оператор  $\vec{\Phi}$  визначений на множині  $H$ . За теоремою про стискуючі відображення ітераційний процес (35) сходиться, якщо виконується умова

$$\left| \vec{\Phi}(\vec{x}) - \vec{\Phi}(\vec{y}) \right| \leq q \cdot |\vec{x} - \vec{y}|, \quad 0 < q < 1, \quad (36)$$

або

$$\left| \vec{\Phi}'(\vec{x}) \right| \leq q < 1, \quad (37)$$

де  $\vec{x} \in U_r$ ,  $\vec{\Phi}'(\vec{x}) = \left( \frac{\partial \varphi_i}{\partial x_j} \right)_{i,j=1}^n$ . Для похибки справедлива оцінка

$$\left| \vec{x}^{(m)} - \vec{x} \right| \leq \frac{q^n}{1-q} \cdot \left| \vec{x}^{(0)} - \vec{x} \right|. \quad (38)$$

Частинним випадком методу простої ітерації є метод релаксації для рівняння (33):

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} - \tau \cdot \vec{F}(\vec{x}^{(k)}), \quad (39)$$

де  $\tau < 2 / \left\| \vec{F}'(\vec{x}) \right\|$ .

#### 4.2.2. Метод Ньютона

Розглянемо рівняння

$$\vec{F}(\vec{x}) = 0. \quad (40)$$

Представимо його у вигляді

$$\vec{F}(\vec{x}^{(k)}) + \vec{F}'(\vec{\xi}^{(k)}) \cdot (\vec{x} - \vec{x}^{(k)}) = 0, \quad (41)$$

де

$$\vec{\xi}^{(k)} = \vec{x}^{(k)} + \theta_k \cdot (\vec{x}^{(k)} - \vec{x}), \quad (42)$$

де  $0 < \theta_k < 1$ . Тут  $\vec{F}'(\vec{x}) = \left( \frac{\partial f_i}{\partial x_j} \right)_{i,j=1}^n$  — матриця Якобі для  $\vec{F}(\vec{x})$ . Можемо наближено вважати  $\vec{\xi}^{(k)} \approx \vec{x}^{(k)}$ . Тоді з (41) матимемо

$$\vec{F}(\vec{x}^{(k)}) + \vec{F}'(\vec{x}^{(k)}) \cdot (\vec{x}^{(k+1)} - \vec{x}^{(k)}) = 0. \quad (43)$$

Ітераційний процес представимо у вигляді:

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} - \vec{F}'(\vec{x}^{(k)})^{-1} \cdot \vec{F}(\vec{x}^{(k)}). \quad (44)$$

Для реалізації методу Ньютона потрібно, щоб існувала обернена матриця

$$\vec{F}'(\vec{x}^{(k)})^{-1}. \quad (45)$$

Можна не шукати обернену матрицю, а розв'язувати на кожній ітерації СЛАР

$$\begin{aligned} A_k \vec{z}^{(k)} &= \vec{F}(\vec{x}^{(k)}), \\ \vec{x}^{(k+1)} &= \vec{x}^{(k)} - \vec{z}^{(k)}, \end{aligned} \quad (46)$$

де  $k = 0, 1, 2, \dots$ , і  $\vec{x}^{(0)}$  — задано, а матриця  $A_k = \vec{F}'(\vec{x}^{(k)})$ .

Метод має квадратичну збіжність, якщо добре вибрано початкове наближення. Складність методу (при умові використання методу Гаусса розв'язання СЛАР (46) на кожній ітерації  $Q_n = \frac{2}{3}n^3 + O(n^2)$ , де  $n$  — розмірність системи (33).

#### 4.2.3. Модифікований метод Ньютона

Ітераційний процес має вигляд:

$$\vec{x}^{(k+1)} = \vec{x}^{(k)} - \vec{F}'(\vec{x}^{(0)})^{-1} \cdot \vec{F}(\vec{x}^{(k)}). \quad (47)$$

Тепер обернена матриця обчислюється тільки на нульовій ітерації. На інших — обчислення нового наближення зводиться до множення матриці  $A_0 = \vec{F}'(\vec{x}^{(0)})^{-1}$  на вектор  $\vec{F}(\vec{x}^{(k)})$  та додавання до  $\vec{x}^{(k)}$ .

Запишемо метод у вигляді системи лінійних рівнянь (аналог (46))

$$\begin{aligned} A_0 \vec{z}^{(k)} &= \vec{F}(\vec{x}^{(k)}), \\ \vec{x}^{(k+1)} &= \vec{x}^{(k)} - \vec{z}^{(k)}, \end{aligned} \quad (48)$$

де  $k = 0, 1, 2, \dots$

Оскільки матриця  $A_0$  розкладається на трикутні (або обертається) один раз, то складність цього методу на одній ітерації (окрім нульової)  $Q_n = O(n^2)$ . Але цей метод має лінійну швидкість збіжності.

Можливе циклічне застосування модифікованого методу Ньютона, тобто коли обернену матрицю похідних шукаємо та обертаємо через певне число кроків ітераційного процесу.

**Задача 9:** Побудувати аналог методу січних для систем нелінійних рівнянь.

## 5. Алгебраїчна проблема власних значень

Нехай задано матрицю  $A \in \mathbb{R}^{n \times n}$ . Тоді задача на власні значення ставиться так: знайти число  $\lambda$  та вектор  $x \neq 0$ , що задовольняють рівнянню

$$Ax = \lambda x. \quad (1)$$

**Означення:**  $\lambda$  називається *власним значенням*  $A$ , а  $x$  — *власним вектором*.

З (1)

$$\det(A - \lambda E) = P_n(\lambda) = (-1)^n \lambda^n + a_n \lambda^{n-1} + \dots + a_0 = 0. \quad (2)$$

Тут  $P_n(\lambda)$  — характеристичний багаточлен.

Для розв'язання багатьох задач механіки, фізики, хімії потрібне знаходження всіх власних значень  $\lambda_i, i = \overline{1, n}$ , а іноді й всіх власних векторів  $x_i$ , що відповідають  $\lambda_i$ .

**Означення:** Цю задачу називають *повною проблемою власних значень*.

В багатьох випадках потрібно знайти лише максимальне або мінімальне за модулем власне значення матриці. При дослідженні стійкості коливальних процесів іноді потрібно знайти два максимальних за модулем власних значення матриці.

**Означення:** Останні дві задачі називають *частковими проблемами власних значень*.

### 5.1. Степеневий метод

1. Знаходження  $\lambda_{\max}$ :  $\lambda_{\max} \equiv |\lambda_1| \geq |\lambda_2| \geq |\lambda_3| \geq \dots$

Нехай  $x^{(0)}$  — заданий вектор, будемо послідовно обчислювати вектори

$$x^{(k+1)} = Ax^{(k)}, \quad k = 0, 1, \dots \quad (3)$$

Тоді  $x^{(k)} = A^k x^{(0)}$ . Нехай  $\{e_i\}_{i=1}^n$  — система власних векторів. Представимо  $x^{(0)}$  у вигляді:

$$x^{(0)} = \sum_{i=1}^n c_i e_i. \quad (4)$$

Оскільки  $Ae_i = \lambda_i e_i$ , то  $x^{(k)} = \sum_{i=1}^n c_i \lambda_i^k e_i$ . При великих  $k$ :  $x^{(k)} \approx c_1 \lambda_1^k e_1$ . Тому



$$\mu_1^{(k)} = \frac{x_m^{(k+1)}}{x_m^{(k)}} = \lambda_1 + O\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right). \quad (5)$$

Значить  $\mu_1^{(k)} \xrightarrow[k \rightarrow \infty]{} \lambda_1$ .

Якщо матриця  $A = A^T$  симетрична, то існує ортонормована система векторів  $\langle e_i, e_j \rangle = \delta_{ij}$ . Тому

$$\begin{aligned} \mu_1^{(k)} &= \frac{\langle x^{(k+1)}, x^{(k)} \rangle}{\langle x^{(k)}, x^{(k)} \rangle} = \frac{\left\langle \sum_i c_i \lambda_i^{k+1} e_i, \sum_j c_j \lambda_j^k e_j \right\rangle}{\left\langle \sum_i c_i \lambda_i^k e_i, \sum_j c_j \lambda_j^k e_j \right\rangle} = \frac{\sum_i c_i^2 \lambda_i^{2k+1}}{\sum_i c_i^2 \lambda_i^{2k}} = \\ &= \frac{c_1^2 \lambda_1^{2k+1} + c_2^2 \lambda_2^{2k+1} + \dots}{c_1^2 \lambda_1^{2k} + c_2^2 \lambda_2^{2k} + \dots} = \lambda_1 + O\left(\left|\frac{\lambda_2}{\lambda_1}\right|^{2k}\right) \xrightarrow[k \rightarrow \infty]{} \lambda_1. \end{aligned} \quad (6)$$

Це означає збіжність до максимального за модулем власного значення з квадратичною швидкістю.

Якщо  $|\lambda_1| > 1$ , то при проведенні ітерацій відбувається зріст компонент вектора  $x^{(k)}$ , що приводить до «переповнення» (overflow). Якщо ж  $|\lambda_1| < 1$ , то це приводить до зменшення компонент (underflow). Позбутися негативу такого явища можна нормуючи вектори  $x^{(k)}$  на кожній ітерації.

Алгоритм степеневого методу знаходження максимального за модулем власного значення з точністю  $\varepsilon$  виглядає так:

```
e[0] = x[0] / norm(x[0])

k = 0
while True:
    k += 1

    x[k + 1] = A * x[k]
    μ[k][1] = scalar_product(x[k + 1], e[k])
    e[k + 1] = x[k + 1] / norm(x[k + 1])

    if abs(μ[k + 1][1] - μ[k][1]) < ε:
        break

λ[1] = μ[k + 1][1]
```

За цим алгоритмом для симетричної матриці  $A^T = A$  швидкість прямування  $\mu_1^{(k)}$  до  $\lambda_{\max}$  — квадратична.

2. Знаходження  $\lambda_2$ :  $|\lambda_1| \geq |\lambda_2| \geq |\lambda_3| \geq \dots$ . Нехай  $\lambda_1, e_1$  відомі.

**Задача 10:** Довести, що якщо  $|\lambda_1| \geq |\lambda_2| \geq |\lambda_3| \geq \dots$  то

$$\mu_2^{(k)} = \frac{x_m^{(k+1)} - \lambda_1 x_m^{(k)}}{x_m^{(k)} - \lambda_1 x_m^{(k-1)}} \xrightarrow{k \rightarrow \infty} \lambda_2, \quad (7)$$

де  $x^{(k+1)} = Ax^{(k)}$ ,  $x_m^{(k)}$  —  $m$ -та компонента  $x^{(k)}$ .

**Задача 11:** Побудувати алгоритм обчислення  $\lambda_2, e_2$ , використовуючи нормування векторів та скалярні добутки для обчислення  $\mu_2^{(k)}$ .

3. Знаходження мінімального власного числа  $\lambda_{\min}(A) = \min_i |\lambda_i(A)|$ .

Припустимо, що  $\lambda_i(a) > 0$  то відоме  $\lambda_{\max}$ . Розглянемо матрицю  $B = \lambda_{\max}E - A$ .  
Маємо

$$\forall i: \quad \lambda_i(B) = \lambda_{\max} - \lambda_i(A). \quad (8)$$

Тому  $\lambda_{\max}(B) = \lambda_{\max}(A) - \lambda_{\min}(A)$ . Звідси  $\lambda_{\min}(A) = \lambda_{\max}(A) - \lambda_{\max}(B)$ .

Якщо  $\exists i: \lambda_i(A) < 0$ , то будемо матрицю  $\bar{A} = \sigma E + A$ ,  $\sigma > 0$ :  $\bar{A} > 0$  і для неї попередній розгляд дає необхідний результат. Замість  $\lambda_{\max}$  в матриці  $B$  можна використовувати  $\|A\|$ .

Ще один спосіб обчислення мінімального власного значення полягає в використанні обернених ітерацій:

$$Ax^{(k+1)} = x^{(k)}, \quad k = 0, 1, \dots \quad (9)$$

Але цей метод вимагає більшої кількості арифметичних операцій: складність методу на основі формули (3):  $Q = O(n^2)$ , а на основі (9) —  $Q = O(n^3)$ , оскільки треба розв'язувати СЛАР, але збігається метод (9) швидше.

## 5.2. Ітераційний метод обертання

Цей метод розв'язання повної проблеми власних значень для симетричних матриць  $A^T = A$ . Існує матриця  $U$ , що приводить матрицю  $A$  до діагонального виду:

$$A = U\Lambda U^T, \quad (10)$$

де  $\Lambda$  — діагональна матриця, по діагоналі якої стоять власні значення  $\lambda_i$ ;  $U$  — унітарна матриця, тобто:  $U^{-1} = U^T$ .

З (1) маємо

$$\Lambda = U^T A U. \quad (11)$$

Нехай  $\exists \tilde{U}$  — матриця, така що  $\tilde{\Lambda} = \tilde{U}^T A \tilde{U}$  і  $\tilde{\Lambda} = \left( \tilde{\lambda}_{ij} \right)_{i,j=1}^n$ ,  $|\tilde{\lambda}_{ij}| < \delta \ll 1, i \neq j$ .

Тоді діагональні елементи мало відрізняються від власних значень

$$|\tilde{\lambda}_{ij} - \lambda_i(A)| < \varepsilon = \varepsilon(\delta). \quad (12)$$

Введемо

$$t(A) = \sum_{\substack{i,j=1 \\ i \neq j}}^n a_{ij}^2. \quad (13)$$

З малості величини  $t(A)$  випливає, що діагональні елементи малі. По  $A = A_0$  за допомогою матриць обертання  $U_k$ :

$$U_k = \begin{pmatrix} 1 & \cdots & 0 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & \cos \phi & \cdots & -\sin \phi & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & \sin \phi & \cdots & \cos \phi & \cdots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & 0 & \cdots & 1 \end{pmatrix}. \quad (14)$$

що повертають систему векторів на кут  $\varphi$ , побудуємо послідовність  $\{A_k\}$  таку, що  $A_k \xrightarrow[k \rightarrow \infty]{} \Lambda$ .

**Задача 12:** Показати, що матриця обертання  $U_k$  є унітарною:  $U_k^{-1} = U_k^T$ .

Послідовно будуюмо:

$$A_{k+1} = U_k^T A_k U_k, \quad (15)$$

**Означення:** Процес (15) називається *монотонним*, якщо:  $t(A_{k+1}) < t(A_k)$ .

**Задача 13:** Довести, що для матриці (15) виконується:

$$a_{i,j}^{(k+1)} = a_{i,j}^{(k)} \cos(2\varphi) + \frac{1}{2} \left( a_{j,j}^{(k)} - a_{i,i}^{(k)} \right) \sin(2\varphi), \quad (16)$$

Показати, що

$$t(A_{k+1}) = t(A_k) - 2\left(a_{i,j}^{(k)}\right)^2 \quad (17)$$

якщо вибрати  $\varphi$  з умови  $a_{i,j}^{(k+1)} = 0$ .

Звідси

$$\varphi = \varphi_k = \frac{1}{2} \arctan\left(p^{(k)}\right), \quad (18)$$

тобто

$$p^{(k)} = \frac{2a_{i,j}^{(k)}}{a_{i,i}^{(k)} - a_{j,j}^{(k)}}, \quad (19)$$

де

$$\left|a_{i,j}^{(k)}\right| = \max_{\substack{m,l \\ m \neq l}} \left|a_{m,l}^{(k)}\right|. \quad (20)$$

Тоді  $t(A_k) \xrightarrow[k \rightarrow \infty]{} 0$ . Чим більше  $n$  тим більше ітерацій необхідно для зведення  $A$  до  $\Lambda$ .

Якщо матриця несиметрична, то застосовують QR- або QL-методи.

## 6. Інтерполявання функцій

### 6.1. Постановка задачі інтерполявання

Нехай функція  $f(x) \in C([a, b])$  задана своїми значеннями  $y_i = f(x_i)$ ,  $x_i \in [a, b]$ ,  $i = \overline{0, n}$ , причому  $x_i \neq x_j$  для  $i \neq j$ .

**Означення:** Функція  $\Phi(x)$  називається *інтерполуючою* для  $f(x)$  на сітці  $\{x_i\}_{i=0}^n$ , якщо  $\Phi(x_i) = y_i$ ,  $i = \overline{0, n}$ .

Задача інтерполявання функції має не єдиний розв'язок.

**Означення:** Виберемо систему лінійно незалежних функцій  $\{\varphi_k(x)\}_{k=0}^n$ ,  $\varphi_k(x) \in C([a, b])$  і побудуємо лінійну комбінацію

$$\Phi(x) = \Phi_n(x) = \sum_{k=0}^n c_k \cdot \varphi_k(x), \quad (1)$$

яка називається *узагальненим багаточленом*.

Умови інтерполявання дають СЛАР

$$\sum_{k=0}^n c_k \cdot \varphi_k(x_i) = y_i, \quad i = \overline{1, n} \quad (2)$$

розв'язком якої є  $\vec{c} = (c_0, \dots, c_n)$ .

Якщо

$$D(x_0, \dots, x_n) = \begin{vmatrix} \varphi_0(x_0) & \cdots & \varphi_n(x_0) \\ \vdots & \ddots & \vdots \\ \varphi_0(x_n) & \cdots & \varphi_n(x_n) \end{vmatrix} \neq 0, \quad (3)$$

то система (2) має єдиний розв'язок.

**Означення:** Система функцій  $\{\varphi_k(x)\}_{k=0}^n$  називається *системою Чебишова*, якщо  $\forall \{x_i\}_{i=0}^n$  таких, що  $x_i \in [a, b]$  і  $x_i \neq x_j$  при  $i \neq j$  виконується  $D(x_0, \dots, x_n) \neq 0$ .

**Приклади** систем Чебишова:

1.  $\varphi_k(x) = x^k$  — алгебраїчна система.

Визначник  $D(x_0, \dots, x_n) \neq 0$  є визначником Вандермонда:

$$\begin{aligned} D(x_0, \dots, x_n) &= \begin{vmatrix} 1 & x_0 & \cdots & x_0^n \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & \cdots & x_n^n \end{vmatrix} = \\ &= \prod_{0 \leq k < m \leq n} (x_k - x_m) \neq 0, \end{aligned} \quad (4)$$

2.  $\varphi_k(x) = L_k(x)$  — ортогональні багаточлени Лежандра;

3.  $\varphi_k(x) = T_k(x)$  — ортогональні багаточлени Чебишова.

4.  $\varphi_k(x)$ :  $1, \cos(x), \sin(x), \dots, \cos(nx), \sin(nx)$ .

Тоді

$$\begin{aligned}\Phi_n(x) &= T_n(x) = \\ &= a_0 + \sum_{k=1}^n (a_k \cdot \cos(kx) + b_k \cdot \sin(kx))\end{aligned}\quad (5)$$

— тригонометричний багаточлен.

## 6.2. Інтерполяційна формула Лагранжа

Якщо  $\varphi_k(x) = x^k$ , то

$$\Phi_n(x) = P_n(x) = \sum_{k=0}^n c_k \cdot x^k. \quad (6)$$

Задача інтерполювання функції  $f(x)$  алгебраїчним, багаточленом полягає в знаходженні коефіцієнтів  $c_k$ ,  $k = \overline{0, n}$  для яких виконується умова  $f(x_i) = \varphi(x_i)$ ,  $i = \overline{0, n}$ .

Представимо інтерполяційний багаточлен у вигляді

$$P_n(x) = L_n(x) = \sum_{k=0}^n f(x_k) \cdot \Phi_k^{(n)}(x). \quad (7)$$

**Означення:** Тут  $L_n(x)$  — інтерполяційний поліном,  $\Phi_k^{(n)}(x)$  — поліноми  $n$ -го степеня, які називають множниками Лагранжа.

З умови  $L_n(x_i) = f(x_i)$  випливає, що множник Лагранжа повинен задовольняти умови

$$\Phi_k^{(n)}(x_i) = \delta_{i,k}. \quad (8)$$

Оскільки  $\Phi_k^{(n)}(x)$  — багаточлен степеня  $n$ , то він має вигляд

$$\Phi_k^{(n)}(x) = A_k(x - x_0) \dots (x - x_{k-1})(x - x_{k+1}) \dots (x - x_n), \quad (9)$$

де  $A_k$  — число.

Знайдемо його з умови  $\Phi_k^{(n)}(x_k) = 1$ :

$$A_k = \frac{1}{(x_k - x_0) \dots (x_k - x_{k-1})(x_k - x_{k+1}) \dots (x_k - x_n)}. \quad (10)$$

Таким чином багаточлен  $\Phi_k^{(n)}(x)$  мають вигляд:

$$\Phi_k^{(n)}(x) = \frac{(x - x_0) \dots (x - x_{k-1})(x - x_{k+1}) \dots (x - x_n)}{(x_k - x_0) \dots (x_k - x_{k-1})(x_k - x_{k+1}) \dots (x_k - x_n)} \quad (11)$$

Позначивши

$$\omega_n(x) = \prod_{i=0}^n (x - x_i), \quad (12)$$

маємо

$$\Phi_k^{(n)}(x) = \frac{\omega_n(x)}{(x - x_k) \cdot \omega'_n(x_k)}. \quad (13)$$

Остаточна формула Лагранжа має вигляд

$$L_n(x) = \sum_{k=0}^n f(x_k) \cdot \frac{\omega_n(x)}{(x - x_k) \cdot \omega'_n(x_k)} \quad (14)$$

### 6.3. Залишковий член інтерполяційного полінома

В заданих точках (точки інтерполявання) значення функції та полінома співпадають, але в інших точках в загальному випадку не співпадають. Отже доцільно розглянути питання про похибку інтерполявання.

**Означення:** Замінюючи функцію  $f(x)$  на  $L_n(x)$  ми допускаємо похибку  $r_n(x) = f(x) - L_n(x)$ . Це *залишковий член* інтерполявання.

З означення випливає, що,  $r_n(x_i) = 0$ ,  $x_i \in [a, b]$ . Оцінимо похибку у кожній точці  $x \in [a, b]$ . Введемо допоміжну функцію:

$$g(t) = f(t) - L_n(t) - K \cdot \omega_n(t), \quad (15)$$

де  $t \in [a, b]$ , і  $g(x_i) = 0$  для  $i = \overline{0, n}$ .

Знайдемо таке  $K$ , щоб  $g(x) = 0$ , в деякій точці  $x \in [a, b]$ ,  $x \neq x_i$ ,  $i = \overline{0, n}$ . Легко бачити, що

$$K = \frac{f(x) - L_n(x)}{\omega_n(x)}. \quad (16)$$

Припустимо що  $f(x) \in C^{(n+1)}([a, b])$ , тоді  $g(t) \in C^{(n+1)}([a, b])$ . Функція  $g(t) = 0$  в  $(n+2)$ -х точках, а саме  $t = x$ ,  $t = x_i$ ,  $i = \overline{0, n}$ . З теореми Ролля випливає, що існує  $(n+1)$ -а точка, де  $g'(t_i) = 0$ ,  $i = \overline{0, n}$ . Продовжуючи цей процес, отримуємо, що існує хоча б одна  $\xi \in [a, b]$  така, що  $g^{(n+1)}(\xi) = 0$ . Оскільки

$$g^{(n+1)}(t) = f^{(n+1)}(t) - 0 - K \cdot (n+1)!, \quad (17)$$

то  $\exists \xi$ , що

$$g^{(n+1)}(\xi) = f^{(n+1)}(\xi) - (n+1)! \cdot \frac{f(x) - L_n(x)}{\omega_n(x)} = 0. \quad (18)$$

Звідси

$$r_n(x) = f(x) - L_n(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \cdot \omega_n(x). \quad (19)$$

Оскільки  $\xi$  невідомо, то використовують оцінку залишкового члена:

$$|r_n(x)| = |f(x) - L_n(x)| \leq \frac{M_{n+1}}{(n+1)!} \cdot |\omega_n(x)|, \quad (20)$$

де  $M_{n+1} = \sup_{x \in [a, b]} |f^{(n+1)}(x)|$ .

## 6.4. Багаточлени Чебишова. Мінімізація залишкового члена інтерполяційного полінома

Як вибрати вузли інтерполяції щоб похибка інтерполювання була мінімальною? Спочатку обґрунтуємо теоретичний апарат, завдяки якому будемо досліджувати це питання.

**Означення:** Багаточленом Чебишова ( $n$ -того степеня, 1-го роду) називається поліном, який задається такими рекурентними співвідношеннями

$$T_{n+1}(x) - 2x \cdot T_n(x) + T_{n-1}(x) = 0, \quad (21)$$

де початкові значення

$$T_0(x) = 1, \quad T_1(x) = x. \quad (22)$$

Знайдемо явний вигляд багаточлена Чебишова. Будемо шукати розв'язок рівняння (21) у вигляді  $T_n(x) = q^n$ , де  $q = q(x)$ . Підставивши в (21), отримуємо характеристичне рівняння  $q^2 - 2xq + 1 = 0$ . Тоді при  $|x| \geq 1 \implies q_{1,2} = x \pm \sqrt{x^2 - 1}$ , а при  $|x| < 1 \implies q_{1,2} = \cos(\varphi) \pm i \sin(\varphi)$ ,  $\varphi = \arccos(x)$ .

Розглянемо обидва випадки детальніше:

1. при  $|x| \leq 1$ :  $T_n(x) = A \cdot \cos(n\varphi) + B \cdot \sin(n\varphi)$ . З (21) випливає  $A = 1$ ,  $B = 0$  і тому

$$T_n(x) = \cos(n \arccos(x)). \quad (23)$$

2. при  $|x| > 1$ :

$$T_n(x) = \frac{1}{2} \left( \left( x + \sqrt{x^2 - 1} \right)^n + \left( x - \sqrt{x^2 - 1} \right)^n \right). \quad (24)$$

Знайдемо нулі та екстремуми багаточлена Чебишова:  $T_n(x) = 0$ ,  $x \in [-1, 1]$ ,  $\cos(n \arccos(x)) = 0$ ,  $\arccos(x) = \frac{2k+1}{2n} \pi$ ,  $k = \overline{0, n-1}$ .

Отже нулі багаточлена Чебишова:

$$x_k = \cos\left(\frac{(2k+1)\pi}{2n}\right) \in [-1, 1], \quad k = \overline{0, n-1}. \quad (25)$$

Локальні екстремуми багаточлена Чебишова на  $x \in [-1, 1]$ :

$$x'_k = \cos\left(\frac{k\pi}{2n}\right), \quad k = \overline{0, n}. \quad (26)$$

Коефіцієнт при старшому члені багаточлена дорівнює  $2^{n-1}$ .

**Означення:** Введемо нормований багаточлен Чебишова  $\bar{T}_n(x) = 2^{1-n} T_n(x) = x^n + \dots$

Тоді

$$\left| \bar{T}_n \right|_{C([-1,1])} = \max_{x \in [-1,1]} |T_n(x)| = 2^{1-n}. \quad (27)$$

**Означення:** Відхиленням двох функцій  $f(x)$  та  $\Phi(x)$  називається величина

$$\Delta(f, \Phi) = |f(x) - \Phi(x)|_{C([a,b])}. \quad (28)$$



**Теорема (Чебишова):** Серед усіх багаточленів  $n$ -го степеня з коефіцієнтом 1 при старшому степені  $\overline{T}_n(x)$  найменше відхиляється від 0 на  $[-1, 1]$ , тобто

$$\left| \overline{T}_n(x) - 0 \right|_{C([-1,1])} = \inf_{\overline{P}_n(x)} \left| \overline{P}_n(x) \right|_{C([-1,1])} = 2^{1-n}. \quad (29)$$

*Доведення.* Будемо доводити від супротивного: припустимо, що існує багаточлен, такий, що

$$\overline{Q}_n(x) < 2^{1-n}. \quad (30)$$

Тоді  $Q_{n-1}(x) = \overline{T}_n(x) - \overline{Q}_n(x)$  — поліном степеня не вище  $n - 1$  і не рівний тотожно нулю. Дослідимо його знаки:

$$\operatorname{sgn}(Q_{n-1}(x'_k)) = \operatorname{sgn}(\overline{T}_n(x'_k) - \overline{Q}_n(x'_k)) = \operatorname{sgn}(\overline{T}_n(x'_k)) = \alpha \cdot (-1)^k, \quad (31)$$

де  $\alpha = \pm 1$ .

Значить  $\exists z_k, k = \overline{0, n-1}$  таке, що  $Q_{n-1}(z_k) = 0$ . Це протиріччя, бо  $Q_{n-1}(x)$  — поліном степеня  $\leq n - 1$ .  $\square$

Тепер узагальнимо наш багаточлен Чебишова на довільний проміжок. Нагадаємо  $T_n(t) = \cos(n \arccos t)$ ,  $-1 \leq t \leq 1$ . Від змінної  $t \in [-1, 1]$  перейдемо до  $x \in [a, b]$ . Запровадимо заміну

$$t = \frac{2x}{b-a} - \frac{b+a}{b-a}, \quad x = \frac{b+a}{2} + \frac{b-a}{2}t. \quad (32)$$

Тоді

$$T_n^{[a,b]}(t) = \overline{T}_n\left(\frac{2x}{b-a} - \frac{b+a}{b-a}\right) = 2^{1-n} \cos\left(n \arccos\left(\frac{2x - (b+a)}{b-a}\right)\right). \quad (33)$$

Побудований нами багаточлен Чебишова на  $[a, b]$  не є нормованим.

*Нормований багаточлен Чебишова на  $[a, b]$ :*

$$\overline{T}_n^{[a,b]}(x) = \frac{(b-a)^n}{2^{2n-1}} \cos\left(n \arccos\left(\frac{2x - (b+a)}{b-a}\right)\right). \quad (34)$$

Відповідно його нулі

$$x_k = \frac{a+b}{2} - \frac{b-a}{2} \cdot t_k, \quad t_k = \cos\left(\frac{(2k+1)\pi}{2n}\right), \quad (35)$$

де  $k = \overline{0, n-1}$ , а точки екстремуму

$$x'_k = \frac{a+b}{2} - \frac{b-a}{2} \cdot t'_k, \quad t'_k = \cos\left(\frac{k\pi}{n}\right), \quad k = \overline{0, n}. \quad (36)$$

Теорема Чебишова вірна і у випадку  $[a, b]$ . Тепер

$$\left| \overline{T}_n^{[a,b]} \right|_{C([a,b])} = \frac{(b-a)^n}{2^{2n-1}}. \quad (37)$$

Перейдемо до питання мінімізації залишкового члена. Нагадаємо, що

$$|r_n(x)| = |f(x) - L_n(x)| \leq \frac{M_{n+1}}{(n+1)!} \cdot |\omega_n(x)|, \quad (38)$$

$$\text{де } M_{n+1} = \max_{x \in [a,b]} |f^{(n+1)}(x)|, \omega_n(x) = \prod_{i=0}^n (x - x_i) = x^{n+1} + \dots$$

Поставимо задачу

$$\inf_{\bar{P}_n(x)} \max_{x \in [a,b]} |\omega_n(x)|. \quad (39)$$

З теоремою Чебишова  $\omega_n(x) = \bar{T}_{n+1}^{[a,b]}(x)$  поліном Чебишова. Якщо співпадають поліноми, то співпадають їх нулі. Отже:  $x_k$  — вузли інтерполяції співпадають з нулями багаточлена Чебишова

$$x_k = \frac{a+b}{2} - \frac{b-a}{2} \cdot t_k, \quad t_k = \cos\left(\frac{(2k+1)\pi}{2(n+1)}\right), \quad (40)$$

де  $k = \overline{0, n}$ .

В цьому випадку

$$|r_n(x)| = |f(x) - L_n(x)| \leq \frac{M_{n+1}}{(n+1)!} \cdot \frac{(b-a)^{n+1}}{2^{2n+1}}. \quad (41)$$

Цю оцінку не можна покращити! Так для  $f(x) = \bar{P}_{n+1}^{[a,b]}(x) = x^{n+1} + \dots$  її  $(n+1)$  похідна дорівнює  $(n+1)!$ , тому  $M_{n+1} = (n+1)!$ . Різниця  $f(x) - L_n(x) = \bar{T}_{n+1}^{[a,b]}(x)$ , отже

$$\max_{x \in [a,b]} |f(x) - L_n(x)| = \frac{(b-a)^{n+1}}{2^{2n+1}}. \quad (42)$$

## 6.5. Розділені різниці

Розділені різниці є аналогом похідної для функції, що задана таблицею.

**Означення:** Розділеною різницею 1-го порядку для функції  $f(x)$  називатимемо

$$f(x_i; x_j) = \frac{f(x_i) - f(x_j)}{x_i - x_j}. \quad (43)$$

Розділеною різницею 2-го порядку для функції  $f(x)$  називатимемо

$$f(x_i; x_j; x_k) = \frac{f(x_i; x_j) - f(x_j; x_k)}{x_i - x_k}. \quad (44)$$

Аналогічно визначаються розділені різниці довільного порядку.

Наведемо деякі властивості розділених різниць:

$$1. \quad f(x_0; \dots; x_n) = \sum_{i=0}^n \frac{f(x_i)}{\prod_{i \neq j} (x_i - x_j)}. \quad (45)$$

2. Розділена різниця — лінійний функціонал:

$$(\alpha_1 f_1 + \alpha_2 f_2)(x_0; x_1) = \alpha_1 f_1(x_0; x_1) + \alpha_2 f_2(x_0; x_1). \quad (46)$$

3. Розділена різниця — симетричний функціонал:

$$f(x_1; \dots; x_i; \dots; x_j; \dots; x_n) = f(x_1; \dots; x_j; \dots; x_i; \dots; x_n). \quad (47)$$

$$4. \forall f(x) \in C^{(n)}([a, b]): \exists \xi \in [a, b]: f(x_0; x_1; \dots; x_n) = \frac{f^{(n)}(\xi)}{n!}.$$

**Задача 14:** Довести першу властивість розділених різниць.

Таблиця розділених різниць має вигляд:

$x_i$	$f_i$	р.р.1	р.р.2	...	р.р.n
$x_0$	$f(x_0)$				
		$f(x_0; x_1)$			
$x_1$	$f(x_1)$		$f(x_0; x_1; x_2)$		
		$f(x_1; x_2)$			
$x_2$	$f(x_2)$				
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
...	...	...	...	...	$f(x_0; \dots; x_n)$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	
$x_{n-2}$	$f(x_{n-2})$				
		$f(x_{n-2}; x_{n-1})$			
$x_{n-1}$	$f(x_{n-1})$		$f(x_{n-2}; x_{n-1}; x_n)$		
		$f(x_{n-1}; x_n)$			
$x_n$	$f(x_n)$				

## 6.6. Інтерполяційна формула Ньютона

Запишемо формулу Лагранжа інтерполяційного багаточлена

$$L_n(x) = \sum_{i=0}^n f(x_i) \cdot \frac{\omega_n(x)}{(x - x_i) \cdot \omega'_n(x_i)}, \quad (48)$$

$$\text{де } \omega_n(x) = \prod_{j=0}^n (x - x_j).$$

Позначимо  $\Phi_j(x) = L_j(x) - L_{j-1}(x)$ . Тоді, оскільки

$$L_n(x) = L_0(x) + (L_1(x) - L_0(x)) + \dots + (L_n(x) - L_{n-1}(x)), \quad (49)$$

і

$$L_j(x_i) = L_{j-1}(x_i) = f(x_i), \quad i = \overline{0, j-1}, \quad (50)$$

то

$$\Phi_j(x_i) = A_j \cdot (x - x_0) \cdot \dots \cdot (x - x_{j-1}), \quad (51)$$

де  $A_j$  визначається з умови  $\Phi_j(x_j) = L_j(x_j) - L_{j-1}(x_j) = f(x_j) - L_{j-1}(x_j)$ . Звідси

$$\Phi_j(x) = \frac{f(x_j) - L_{j-1}(x_j)}{(x_j - x_0) \dots (x_j - x_{j-1})} \cdot (x - x_0) \dots (x - x_{j-1}). \quad (52)$$

Тоді

$$\begin{aligned} A_j &= \frac{f(x_j) - L_{j-1}(x_j)}{(x_j - x_0) \dots (x_j - x_{j-1})} = \frac{f(x_j)}{(x_j - x_0) \dots (x_j - x_{j-1})} - \\ &\quad - \sum_{i=0}^{j-1} \frac{f(x_i)}{(x_i - x_0) \dots (x_i - x_{i-1})(x_i - x_{i+1}) \dots (x_j - x_i)} = \\ &= \frac{f(x_j)}{(x_j - x_0) \dots (x_j - x_{j-1})} + \sum_{i=0}^{j-1} \frac{f(x_i)}{(x_i - x_0) \dots (x_i - x_j)} = \\ &= \sum_{i=0}^j \frac{f(x_i)}{(x_i - x_0) \dots (x_i - x_{i-1})(x_i - x_{i+1}) \dots (x_i - x_j)} = f(x_0; \dots; x_j). \end{aligned} \quad (53)$$

Звідси маємо інтерполяційну формулу Ньютона вперед ( $x_0 \rightarrow x_n$ ):

$$L_n(x) = f(x_0) + f(x_0; x_1)(x - x_0) + \dots + f(x_0; \dots; x_n)(x - x_0) \dots (x - x_{n-1}). \quad (54)$$

Аналогічно, інтерполяційна формула Ньютона назад ( $x_n \rightarrow x_0$ ):

$$L_n(x) = f(x_n) + f(x_n; x_{n-1})(x - x_n) + \dots + f(x_n; \dots; x_0)(x - x_n) \dots (x - x_1). \quad (55)$$

Маємо рекурсію за степенем багаточлена

$$L_n(x) = L_{n-1}(x) + f(x_0; \dots; x_n)(x - x_0) \dots (x - x_1). \quad (56)$$

Звідси

$$L_n(x) = f(x_0) + (x - x_0)(f(x_0; x_1) + (x - x_1)(\dots + (x - x_{n-1})f(x_0; x_1; \dots; x_n) \dots)) \quad (57)$$

і цю формулу розкриваємо починаючи з середини (це аналог схеми Горнера обчислення значення багаточлена).

Введемо нову формулу для похибки інтерполявання. Для  $x \neq x_i, i = \overline{0, n}$  розглянемо розділену різницю

$$f(x; x_0; \dots; x_n) = \frac{f(x)}{(x - x_0) \dots (x - x_n)} + \sum_{k=0}^n \frac{f(x_k)}{\prod_{i \neq k} (x - x_i)}. \quad (58)$$

Звідси

$$\begin{aligned} f(x) &= f(x_0) \cdot \frac{(x - x_1) \dots (x - x_n)}{(x_0 - x_1) \dots (x_0 - x_n)} + \dots + f(x_n) \cdot \frac{(x - x_1) \dots (x - x_{n-1})}{(x_n - x_1) \dots (x_n - x_{n-1})} + \\ &\quad + f(x; x_0; \dots; x_n)(x - x_0) \dots (x - x_n) = L_n(x) + f(x; x_0; \dots; x_n) \cdot \omega_n(x). \end{aligned} \quad (59)$$

Тоді похибка має вигляд

$$r_n(x) = f(x) - L_n(x) = f(x; x_0; \dots; x_n) \cdot \omega_n(x). \quad (60)$$

Це нова форма для залишкового члена.

Порівнюючи з формулою залишкового члена в пункті 6.3, маємо

$$f(x; x_0; \dots; x_n) = \frac{f^{(n+1)}(\xi)}{(n+1)!}, \quad (61)$$

що доводить останню властивість розділених різниць.

Нехай маємо сітку рівновіддалених вузлів:  $x_i = a + ih$ ,  $h = \frac{b-a}{n}$ ,  $i = \overline{0, n}$ ,  $x_0 = a$ ,  $x_n = b$ . Розначимо  $\Delta f_0 = f_1 - f_0$ ,  $\Delta^2 f_0 = \Delta f_1 - \Delta f_0 = f_2 - 2f_1 + f_0$ ,  $\dots$  — скінченні різниці.

Запишемо формули Ньютона у нових позначеннях:

$$L_n(x) = L_n(x_0 + th) = f_0 + t\Delta f_0 + \dots + \frac{t(t-1)\dots(t-n+1)}{n!} \cdot \Delta^n f_0, \quad (62)$$

де  $t = \frac{x-x_0}{h}$ .

Це інтерполяційна формула Ньютона вперед по рівновіддалених вузлах.

**Задача 15:** Побудувати інтерполяційну формулу Ньютона назад по рівновіддалених вузлах.

## 6.7. Інтерполювання з кратними вузлами

Нехай  $f(x)$  задана таблицею значень  $f^{(j)}(x_i)$ ,  $i = \overline{0, n}$ ,  $j = \overline{0, k_i - 1}$ ,  $k_i$  — кратності відповідних вузлів. Побудуємо  $H_m^{(i)}(x_j) = f^{(i)}(x_j)$  — інтерполяційний багаточлен Ерміта по кратним вузлах, де

$$m = \sum_{i=1}^n k_i - 1. \quad (63)$$

Якщо  $k_i = 1$ , то  $H_m(x) = L_n(x)$ .

Для побудови  $H_m(x)$  в загальному випадку для кожної точки  $x_i$  введемо  $k_i$  точок  $x_{i,j}^\varepsilon = x_i + j\varepsilon$ ,  $i = \overline{0, n}$ ,  $j = \overline{0, k_i - 1}$ . З умови  $\forall i: x_{i,k_i-1}^\varepsilon = x_i + \varepsilon(k_i - 1) < x_{i+1}$  можна вибрати  $\varepsilon$ .

Нехай  $f(x) \in C^{(m)}([a, b])$ . Запишемо інтерполяційну формулу Ньютона:

$$\begin{aligned} L_m^\varepsilon = & f(x_{0,0}^\varepsilon) + f(x_{0,0}^\varepsilon; x_{0,1}^\varepsilon) (x - x_{0,0}^\varepsilon) + \dots \\ & + f(x_{0,0}^\varepsilon; \dots; x_{n,k_n-1}^\varepsilon) (x - x_{0,0}^\varepsilon) \dots (x - x_{n,k_n-1}^\varepsilon). \end{aligned} \quad (64)$$

При  $\varepsilon \rightarrow 0$  маємо  $x_{i,j}^\varepsilon \rightarrow x_i$ . Крім того

$$f(x_{i,0}^\varepsilon; \dots; x_{i,k_i-1}^\varepsilon) = f(x_i; \dots; x_i) = \frac{f^{(k_i)}(x_i)}{k_i!}. \quad (65)$$

Тому  $L_m^\varepsilon(x) \rightarrow H_m(x)$  та

$$R_m(x) = f(x) - H_m(x) = \frac{f^{(m+1)}(\xi)}{(m+1)!} \cdot \Omega_m(x), \quad (66)$$

де  $\Omega_m(x) = (x - x_0)^{k_0} \dots (x - x_n)^{k_n}$ .

## 6.8. Збіжність процесу інтерполювання

Виникає питання, чи буде прямувати до нуля похибка інтерполювання  $f(x) - L_n(x)$ , якщо число вузлів  $n$  збільшувати?

Введемо норму

$$|f(x) - L_n|_{C([a,b])} = \max_{x \in [a,b]} |f(x) - L_n(x)|. \quad (67)$$

Тоді для довільної  $f(x) \in C^{(n+1)}([a, b])$  справджується оцінка

$$|f(x) - L_n(x)|_{C([a,b])} \leq \frac{M_{n+1}}{(n+1)!} |\omega_n(x)|_{C([a,b])}, \quad (68)$$

де  $M_{n+1} = \max_{x \in [a,b]} |f^{(n+1)}(x)|$ ,  $\omega_n(x) = \prod_{i=0}^n (x - x_i)$ .

А яка оцінка буде для довільної неперервної функції?

**Означення:** Кажуть, що інтерполяційний процес для функції  $f(x)$  збігається в точці  $x \in [a, b]$ , якщо

$$\forall \{x_i\}_{i=1}^n : h = \max_{i=1,n} \rightarrow 0 : \lim_{n \rightarrow \infty} L_n(x) = f(x), \quad (69)$$

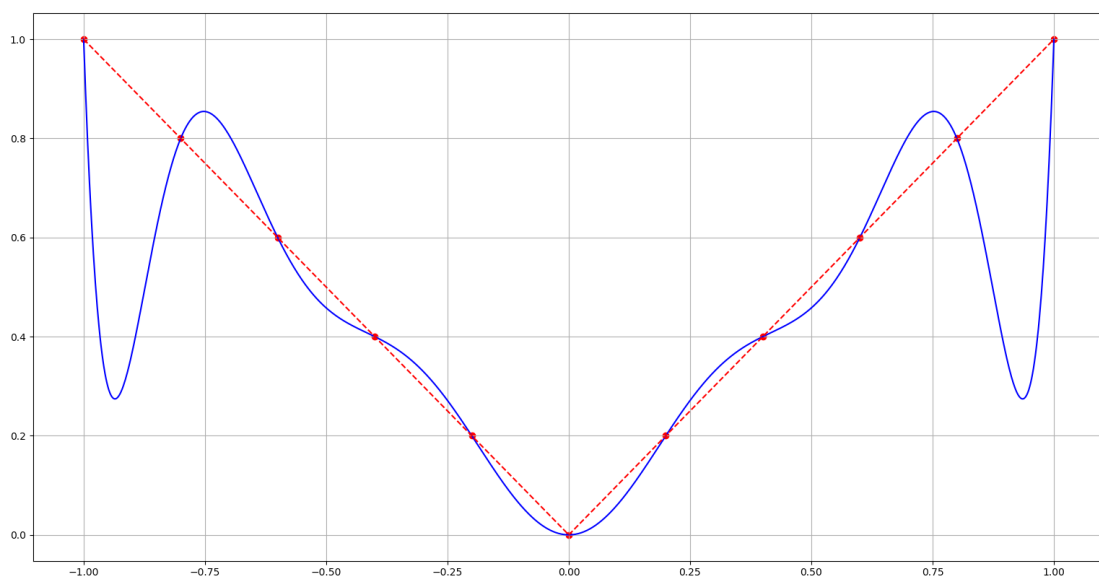
де, як завжди,  $h_i = x_i - x_{i-1}$ .

**Означення:** Якщо  $\|f(x) - L_n(x)\|_{C([a,b])} \xrightarrow{n \rightarrow \infty} 0$ , то інтерполяційний процес збігається рівномірно.

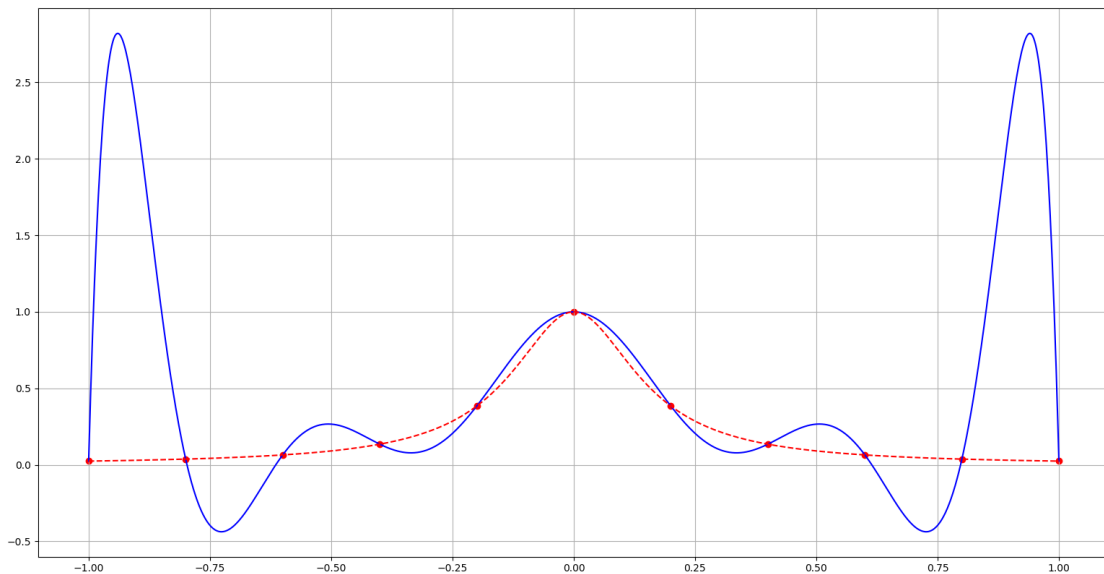
Розглянемо приклади поведінки інтерполяційних багаточленів при  $n \rightarrow \infty$  для деяких функцій.

**Приклад 1:** Послідовність інтерполяційних багаточленів (сітка рівномірна), побудованих для неперервної функції  $f(x) = |x|$ ,  $-1 \leq x \leq 1$  (функція неперервна, але негладка), не збігається на  $x \in [-1, 1]$ , крім точок  $x = -1, 0, 1$ .

На рисунку дано графіки самої функції (штрихова лінія) та інтерполяційного багаточлена (суцільна лінія) на рівномірній сітці  $x_i = -1 + ih$ ,  $h = 2/n$ ,  $i = 0, n$  для  $n = 10$ :

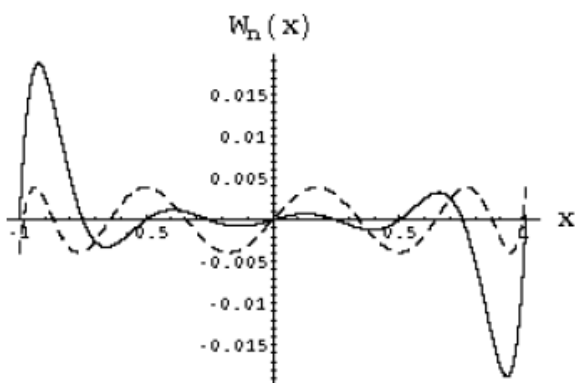


**Приклад 2:** Функція Рунге  $f(x) = \frac{1}{1+40x^2}$ ,  $-1 \leq x \leq 1$  (функція аналітична!). Для рівномірної сітки  $x_i = -1 + ih$ ,  $h = 2/n$ ,  $i = \overline{0, n}$  маємо графіки: суцільна лінія — інтерполяційного багаточлена; пунктирна — самої функції для  $n = 10$ :



Пояснити чому рівномірна сітка дає великі похибки інтерполювання біля кінців інтервалу інтерполювання допомагає наступний рисунок. На цьому рисунку суцільною лінією представлено графік функції

$\omega_n = \prod_{i=0}^n (x - x_i)$  ( $n = 8$ ) для рівномірної сітки. Як бачимо максимальні за модулем значення цієї функції припадають на кінці інтервалу.



Для порівняння на цьому ж рисунку (штрихова лінія) побудовано графік  $\omega_n = \prod_{i=0}^n (x - x_i)$ , що відповідає чебишовським вузлам, які мінімізують похибку інтерполювання. Тепер відхилення  $\omega_n(x)$  розподілено рівномірно по всьому проміжку інтерполювання.

**Теорема (Фабера):**  $\forall \{x_i\}_{i=0}^n$  існує  $f(x) \in C([a, b])$ , для якої інтерполяційний процес не збігається рівномірно.

**Теорема (Марцинкевича):**  $\forall f(x) \in C([a, b])$  існують  $\{x_i\}_{i=0}^n$  такі, що послідовність  $\{L_n(x)\}$  збігається рівномірно до  $f(x)$ .

**Теорема:** Стала Лебега

$$|P_n| = \max_j \sum_{j=0}^n |\varphi_j^{(n)}(x)|, \quad (70)$$

де

$$\varphi_j^{(n)}(x) = \frac{\omega_n(x)}{(x - x_j) \cdot \omega'_n(x_j)}. \quad (71)$$

**Теорема:** Для  $f(x) \in C([a, b])$ :

$$|f(x) - L_n(x)|_{C([a,b])} \leq (1 + |P_n|) \cdot E_n(f), \quad (72)$$

де

$$E_n(f) = \inf_{Q_n(x)} |f(x) - Q_n(x)|_{C([a,b])} \quad (73)$$

— відхилення багаточлена  $n$ -го степеня найкращого рівномірного наближення від  $f(x)$ .

**Теорема:** Нехай  $P_n^E$  — оператор інтерполяції на рівномірній сітці,  $P_n^T$  — оператор інтерполяції на чебишовській сітці. Тоді на  $[-1, 1]$  маємо наближені оцінки:

$$|P_n^E| \approx C_1 \cdot 2^n, \quad |P_n^T| \approx C_2 \cdot \ln(n). \quad (74)$$

Останні оцінки пояснюють розбіжність процесу інтерполювання при великих  $n$ .

## 6.9. Кусково-лінійна інтерполяція

Інтерполяція багаточленом Лагранжа або Ньютона на відрізок  $[a, b]$  з використанням великої кількості вузлів інтерполяції часто приводить до поганого наближення через розбіжність процесу інтерполювання. Для того щоб уникнути великої похибки, весь відрізок  $[a, b]$  розбивають на частинні відрізки  $[x_{i-1}, x_i]$  і на кожному з частинних відрізків замінюють функцію  $f(x)$  багаточленом невисокого степеню. В цьому і полягає кусково-поліноміальна інтерполяція.

Розглянемо найпростішу таку інтерполяцію — лінійну. Нехай задана  $f(x)$  значеннями  $f(x_i)$ ,  $i = \overline{0, n}$ . Побудуємо функцію  $\Phi_1(x)$  — лінійну на  $x \in [x_{i-1}, x_i]$ , що інтерполює ці значення:

$$\Phi_1(x) = L_1^i(x) = f(x_{i-1}) \cdot \frac{x - x_{i-1}}{x_i - x_{i-1}} + f(x_i) \cdot \frac{x_i - x}{x_i - x_{i-1}}, \quad (75)$$

де  $x \in [x_{i-1}, x_i]$ .

Представимо її у вигляді

$$\Phi_1(x) = \sum_{i=0}^n f(x_i) \cdot \varphi_i(x). \quad (76)$$

З умов інтерполювання маємо

$$\Phi_1(x_j) = \sum_{i=0}^n f(x_i) \cdot \varphi_i(x_j) = f(x_j). \quad (77)$$



Звідси

$$\varphi_i(x_j) = \delta_{i,j} = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}. \quad (78)$$

Значить

$$\varphi_i(x) = \begin{cases} 0, & a \leq x \leq x_{i-1} \\ \frac{x-x_{i-1}}{x_i-x_{i-1}}, & x_{i-1} \leq x \leq x_i \\ \frac{x_{i+1}-x}{x_{i+1}-x_i}, & x_i \leq x \leq x_{i+1} \\ 0, & x_{i+1} \leq x \leq b \end{cases} \quad (79)$$

**Теорема:** Для довільної  $f(x) \in C^{(2)}([a, b])$  справедлива оцінка

$$|f(x) - \Phi_1(x)|_{C([a,b])} \leq \frac{M_2}{8} \cdot |h|^2, \quad (80)$$

де  $\Phi_1(x)$  — кусково-лінійна функція, побудована по значеннях  $f(x_i)$ ,  $i = \overline{0, n}$ ,  $|h| = \max_i h_i$ ,  
 $h_i = x_i - x_{i-1}$ .

*Доведення.* Маємо для  $x \in [x_{i-1}, x_i]$ :

$$z(x) = f(x) - \Phi_1(x) = f(x) - L_1^i(x) = \frac{f''(\xi_i)}{2!} \cdot (x - x_{i-1}) \cdot (x - x_i). \quad (81)$$

Звідси

$$|f(x) - \Phi_1(x)| \leq \frac{M_2^i}{2} \cdot |(x - x_{i-1})(x - x_i)| \leq \frac{M_2^i \cdot h_i^2}{8}, \quad (82)$$

де

$$M_2^i = \max_{x \in [x_{i-1}, x_i]} |f''(x)|.$$

Остання оцінка отримана з нерівності

$$\max_{[x_{i-1}, x_i]} |(x - x_{i-1}) \cdot (x - x_i)| = \frac{h_i^2}{4}. \quad (83)$$

Тоді

$$\max_{i=\overline{1,n}} \max_{x \in [x_{i-1}, x_i]} |z(x)| \leq \frac{h^2 M_2}{8}, \quad (84)$$

де  $M_2 = \max_{x \in [a,b]} |f''(x)|$ ,  $h_i = \max_i h_i$ , що доводить (80).  $\square$

**Задача 16:** Довести оцінку  $|f'(x) - \Phi'_1(x)| \leq |h| \cdot M_2$ .

Отже маємо збіжність процесу інтерполявання за допомогою кусково-лінійної функції

$$\left| f(x) - \Phi_1^{(n)}(x) \right|_{C([a,b])} \xrightarrow{h \rightarrow 0, n \rightarrow \infty} 0, \quad (85)$$

тобто

$$\left\{ \Phi_1^{(n)}(x) \right\} \implies f(x). \quad (86)$$

Розглянемо деякі простори:

1.  $H_0 = L_2([a, b])$  — гільбертів простір, в якому скалярний добуток визначається так:

$$\langle u, v \rangle = \int_a^b (u(x) \cdot v(x)) \, dx \quad (87)$$

а норма  $\|u\|_0 = \sqrt{\langle u, u \rangle}$ .

2.  $H_k = W_2^k([a, b])$ . Тепер скалярний добуток

$$\langle u, v \rangle_k = \sum_{m=0}^k \int_a^b (u^{(m)}(x) \cdot v^{(m)}(x)) \, dx, \quad (88)$$

а норма  $\|u\|_k = \sqrt{\|u^{(0)}\|^2 + \dots + \|u^{(k)}\|^2}$ .

**Теорема:** Нехай  $f(x) \in H_2 = W_2^2([a, b])$ . Тоді

$$\left| f^{(k)} - \Phi_1^{(k)} \right|_0 \leq |h|^{2-k} \cdot |f|_2, \quad (89)$$

де  $k = 1, 2$ .

Зауважимо, що кусково-лінійна інтерполяція негладка, тому на практиці застосовують квадратичні, а найчастіше — кубічні поліноми на кожному проміжку.

## 6.10. Кусково-кубічна ермітова інтерполяція

Нехай деяка функція  $f(x)$  задана в точках  $x_i$  своїми значеннями та значеннями похідної:  $y_i = f(x_i)$ ,  $y'_i = f'(x_i)$ ,  $i = 0, n$ . Позначимо через  $\Phi_3(x)$  функцію, яка буде інтерполювати задану. Тоді

$$\Phi_3(x) = H_3^i(x), \quad x \in [x_{i-1}, x_i]. \quad (90)$$

Неважко написати явний вигляд цього багаточлена  $H_3^i(x)$  на проміжку:

$x_i$	$y_i$			
		$y'_i$		
$x_i$	$y_i$		$\frac{y_{i-1,i} - y'_i}{h_i}$	
		$y_{i-1,i}$		$\frac{y'_i - 2y_{i-1,i} + y'_{i-1}}{h_i^2}$
$x_{i-1}$	$y_{i-1}$		$\frac{y'_{i-1} - y_{i-1,i}}{h_i}$	
		$y'_{i-1}$		
$x_{i-1}$	$y_{i-1}$			

$$H_3^i(x) = y_i + y_i'(x - x_i) + \frac{y_{i-1,i} - y_i'}{h_i} \cdot (x - x_i)^2 + \frac{y_i' - 2y_{i-1,i} + y_{i-1}'}{h_i^2} \cdot (x - x_i)^2 \cdot (x - x_{i-1}) \quad (91)$$

Можна представити кусково-кубічну функцію і в такому вигляді:

$$\Phi_3(x) = \sum_{i=0}^n (y_i \cdot \varphi_i^0(x) + y_i' \cdot \varphi_i^1(x)). \quad (92)$$

Умови інтерполявання:  $\Phi_3(x_i) = y_i$ ,  $\Phi_3'(x_i) = y_i'$ ,  $i = \overline{0, n}$ . Якщо ці умови підставити в \u{yiref{eq:6.10.2}}, то отримаємо умови на базисні функції:

$$\varphi_i^0(x_j) = \delta_{i,j}, \quad (93)$$

$$(\varphi_i^0)'(x_j) = 0, \quad (94)$$

$$\varphi_i^1(x_j) = 0, \quad (95)$$

$$(\varphi_i^1)'(x_j) = \delta_{i,j}. \quad (96)$$

для  $i, j = \overline{0, n}$ .

Ці функції кусково-кубічні, тобто  $\varphi_i^k(x) \in \pi_3$ ,  $x \in [x_{i-1}, x_{i+1}]$ ,  $k = 0, 1$  ( $\pi_3$  — множина багаточленів третього степеня), на всіх інших проміжках вони нульові. Нехай  $h_i \equiv h$ , і позначимо  $s = \frac{x - x_i}{h}$ ,  $x \in [x_{i-1}, x_i] \implies s \in [-1, 0]$ .

1. введемо  $\bar{\varphi}_1^0(s) = \varphi_i^0(x)$ ,  $x \in [x_{i-1}, x_{i+1}]$ ,  $x \in [0, 1]$ . Побудуємо цю функцію. Вона задовольняє умовам:

$$\bar{\varphi}_1^0(0) = 1, \quad (97)$$

$$\bar{\varphi}_1^0(1) = 0, \quad (98)$$

$$(\bar{\varphi}_1^0)'(0) = 0, \quad (99)$$

$$(\bar{\varphi}_1^0)'(1) = 0. \quad (100)$$

Її явний вигляд отримаємо за допомогою таблиці розділених різниць:

0	1			
		0		
0	1		-1	
		-1		2
1	0		1	
		0		
1	0			

$$H_3(s) = 1 + 0 \cdot s - 1 \cdot s^2 + 2s^2(s - 1) = 2s^3 - 3s^2 + 1 \equiv \bar{\varphi}_1^0(s). \quad (101)$$

Аналогічно

$$2. \bar{\varphi}_2^0(s) = -2s^3 - 3s^2 + 1, \varphi_i^0(x) = \bar{\varphi}_2^0(s), x \in [x_{i-1}, x_i], s \in [-1, 0];$$

$$3. \bar{\varphi}_1^1(s) = s(s - 1)^2, \varphi_i^1(x) = h\bar{\varphi}_1^1(s), x \in [x_i, x_{i+1}], s \in [0, 1];$$

$$4. \bar{\varphi}_2^1(s) = s(s+1)^2, \varphi_i^0(x) = h\bar{\varphi}_2^1(s), x \in [x_{i-1}, x_i], s \in [-1, 0].$$

А тепер будуюмо явний вигляд функцій  $\varphi_i^k(x)$  для довільного проміжку  $x \in [x_{i-1}, x_{i+1}]$ :

$$\varphi_i^0(x) = \begin{cases} 0, & a \leq x \leq x_{i-1}, \\ -2s^3 - 3s^2 + 1, & x_{i-1} \leq x \leq x_i, \\ 2s^3 - 3s^2 + 1, & x_i \leq x \leq x_{i+1}, \\ 0, & x_{i+1} \leq x \leq b, \end{cases} \quad (102)$$

i

$$\varphi_i^1(x) = \begin{cases} 0, & a \leq x \leq x_{i-1}, \\ hs(s+1)^2, & x_{i-1} \leq x \leq x_i, \\ hs(s-1)^2, & x_i \leq x \leq x_{i+1}, \\ 0, & x_{i+1} \leq x \leq b, \end{cases} \quad (103)$$

де  $s = \frac{x-x_i}{h}$  (якщо сітка нерівномірна, то в формулах замість  $h$ , буде  $h_i$  або  $h_{i+1}$  на відповідних інтервалах).

Оцінімо  $\|f(x) - \Phi_3(x)\|_{C([a,b])}$ . Розглянемо для  $x \in [x_{i-1}, x_i]$ :

$$f(x) - \Phi_3(x) = f(x) - H_3^i(x) = \frac{f^{(4)}(\xi)}{4!} \cdot (x - x_{i-1})^2(x - x_i)^2. \quad (104)$$

Зразу потрібно зробити припущення, що  $f(x) \in C^4([a, b])$ . З тих же міркувань, що і для кусково-лінійної функції, максимум знаходиться в точці  $\bar{x}_i = \frac{x_i + x_{i-1}}{2}$  тому для модуля похибки маємо:

$$|f(x) - \Phi_3(x)| \leq \frac{M_4^i}{24} \left(\frac{h^2}{4}\right)^2 = \frac{M_4^i h^4}{384}, \quad (105)$$

$$|f(x) - \Phi_3(x)|_{C([a,b])} \leq \frac{M_4 h^4}{384}. \quad (106)$$

Звідси отримаємо теорему:

**Теорема:** Якщо функція  $f(x) \in C^4([a, b])$  задана в точках  $x_i$  своїми значеннями  $y_i = f(x_i)$ ,  $y'_i = f'(x_i)$ ,  $i = \overline{0, n}$ , то для кусково-кубічної ермітової інтерполяції

$$\Phi_3(x) = \sum_{i=0}^n \left( y_i \varphi_i^{(0)}(x) + y'_i \varphi_i^{(1)}(x) \right) \quad (107)$$

має місце оцінка

$$|f(x) - \Phi_3(x)|_{C([a,b])} \leq \frac{M_4 h^4}{384}. \quad (108)$$

А для похідної

$$|f'(x) - \Phi'_3(x)|_{C([a,b])} \leq M \cdot M_4 h^3, \quad (109)$$

де  $M$  — стала незалежна від  $h$ .

**Задача 17:** Довести оцінку для  $\|f'(x) - \Phi'_3(x)\|_{C([a,b])}$ .

Порівняємо кусково-лінійну  $\Phi_1(x)$  та кусково-кубічну інтерполяцію  $\Phi_3(x)$ : при згущенні сітки у 2 рази точність лінійної підвищується в 4 рази, а кубічної — у 16 разів, але треба задавати більше даних.

## 6.11. Кубічні інтерполяційні сплайни

Сплайн (spline) в перекладі означає рейка, якою користувалися креслярі при проведенні гладкої кривої, що з'єднувала задані точки на площині.

**Означення:** Функція  $s(x)$  називається *сплайном степеня  $m$  і дефекту  $k$* , якщо

1.  $s(x) \in \pi_m$  (множина поліномів степеня  $m$ ) для  $x \in [x_{i-1}, x_i]$ ,  $i = \overline{1, n}$ .
2.  $s(x) \in C^{(m-k)}([a, b])$ .

**Приклади:**

1.  $\Phi_1(x)$ :  $m = 1, k = 1$ ;
2.  $\Phi_3(x)$ :  $m = 3, k = 2$ ;

Зараз ми побудуємо сплайн, для якого  $m = 3, k = 1$ .

**Означення:** Функція  $s_3(x) = s(x)$  називається *кубічним інтерполяційним природнім сплайном*, якщо

1. Кубічність:

$$s(x) \in \pi_3, \quad x \in [x_{i-1}, x_i], \quad i = \overline{1, n} \quad (110)$$

2. Дефект 1:

$$s(x) \in C^{(2)}([a, b]) \quad (111)$$

3. Інтерполуює  $f(x)$ :

$$s(x_i) = f(x_i), \quad i = \overline{0, n} \quad (112)$$

4. Природній:

$$s''(a) = s''(b) = 0. \quad (113)$$

Умови (113), так звані *умови природності*, необхідні, щоб разом було  $4n$  умов для знаходження  $4n$  коефіцієнтів сплайну. Замість них можуть бути такі умови:

$$s''(a) = A, \quad s''(b) = B \quad (4.a)$$

$$s'(a) = A, \quad s'(b) = B \quad (4.b)$$

$$s(a) = s(b), \quad s'(a) = s'(b), \quad s''(a) = s''(b) \quad (4.c)$$

Умови (4.c) — це так звані умови періодичності.

Побудуємо природній сплайн. З (110) та (111) маємо

$$s''(x) = m_{i-1} \cdot \frac{x_i - x}{h_i} + m_i \cdot \frac{x - x_{i-1}}{h_i}, \quad (114)$$

де  $m_i = s''(x_i)$  і вони є невідомими коефіцієнтами:  $h_i = x_i - x_{i-1}$ .

Двічі інтегруючи (114), маємо

$$s(x) = m_{i-1} \cdot \frac{(x_i - x)^3}{6h_i} + m_i \cdot \frac{(x - x_{i-1})^3}{6h_i} + \left(f_{i-1} - \frac{m_{i-1}h_i^2}{6}\right) \cdot \frac{x_i - x}{h_i} + \left(f_i - \frac{m_i h_i^2}{6}\right) \cdot \frac{x - x_{i-1}}{h_i}, \quad (115)$$

для  $x \in [x_{i-1}, x_i]$ .

З (113) маємо  $m_0 = m_n = 0$ .

Враховуючи, що  $s'(x_i - 0) = s'(x_i + 0)$  отримаємо СЛАР для знаходження всіх  $m_i = s''(x_i)$ :

$$\begin{cases} \frac{h_i m_{i-1}}{6} + \frac{(h_i + h_{i+1})m_i}{3} + \frac{h_{i+1}m_{i+1}}{6} = \\ = \frac{f_{i+1} - f_i}{h_i} - \frac{f_i - f_{i-1}}{h_i}, \quad i = \overline{1, n-1}, \\ m_0 = m_n = 0. \end{cases} \quad (116)$$

Це тридіагональна СЛАР; її можна розв'язати методом прогонки за  $Q = O(N)$  арифметичних операцій.

**Задача 18:** Написати СЛАР для кубічного інтерполяційного сплайну, якщо  $s'(a) = A$ ,  $s'(b) = B$  та розробити алгоритм її розв'язання (тобто написати формули методу прогонки).

**Теорема:** Нехай  $f(x) \in C^{(4)}([a, b])$ , тоді має місце оцінка

$$|f^{(k)}(x) - s^{(k)}(x)|_{C([a,b])} \leq M_4 |h|^{4-k}, \quad (117)$$

де  $k = 0, 1, 2$  і  $M_4 = \max_{[a,b]} |f^{(4)}(x)|$ ,  $|h| = \max_i h_i$ .

Введемо клас функцій

$$U = \left\{ u(x) : u(x) \in W_2^2([a, b]), u(x_i) = f_i, i = \overline{0, n} \right\} \quad (118)$$

— це функції досить гладкі і приймають задані значення. Якщо ввести такий функціонал

$$\Phi(u) = \int_a^b (u''(x))^2 dx, \quad (119)$$

то

$$\Phi(s) = \inf_{u \in U} \Phi(u), \quad (120)$$

де  $s(x)$  — кубічний природний інтерполяційний сплайн.

Оскільки кривизна графіка кривої  $u(x)$  пропорційна  $u''(x)$ , то це фактично означає, що сплайн має в середньоквадратичному розумінні найменшу кривизну серед всіх функцій  $u(x) \in W_2^2([a, b])$ , що інтерполують значення  $f(x_i)$ .

Для того, щоб не розв'язувати СЛАР (116) інколи будують наближення до сплайну  $\tilde{s}(x)$ , яке отримується заміною  $m_i = s''(x_i)$  на

$$f_{\tilde{x}, \hat{x}, i} \equiv \frac{1}{\bar{h}_i} \left( \frac{f_{i+1} - f_i}{h_{i+1}} - \frac{f_i - f_{i-1}}{h_i} \right) \approx f''(x_i) \approx s''(x_i), \quad (121)$$

де  $\bar{h}_i = \frac{h_i + h_{i+1}}{2}$ , причому  $f''(x_i) - f_{\bar{x}\bar{x},i} = O(|h|^2)$ ; При цьому  $\tilde{s}(x) - s(x) = O(h^4)$ . Відмітимо, що  $\tilde{s}(x)$  не є сплайном дефекту 1.

**Зауваження 1:** Складемо матрицю  $A$  розмірності  $(n-1) \times (n-1)$ :

$$A = \begin{pmatrix} \frac{h_1+h_2}{3} & \frac{h_2}{6} & 0 & \dots & 0 \\ \frac{h_2}{6} & \frac{h_2+h_3}{3} & \frac{h_3}{6} & \ddots & \vdots \\ 0 & \frac{h_3}{6} & \frac{h_3+h_4}{3} & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots \\ 0 & \dots & 0 & \frac{h_{n-1}}{6} & \frac{h_{n-1}+h_n}{3} \end{pmatrix} \quad (122)$$

і матрицю  $H$  розмірності  $(n+1) \times (n-1)$ :

$$H = \begin{pmatrix} \frac{1}{h_1} & -\left(\frac{1}{h_1} + \frac{1}{h_2}\right) & \frac{1}{h_2} & 0 & \dots & 0 \\ 0 & \frac{1}{h_2} & -\left(\frac{1}{h_2} + \frac{1}{h_3}\right) & \frac{1}{h_3} & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \frac{1}{h_{n-1}} & -\left(\frac{1}{h_{n-1}} + \frac{1}{h_n}\right) & \frac{1}{h_n} \end{pmatrix} \quad (123)$$

Тоді можна записати СЛАР (116) відносно моментів  $\vec{m} = (m_1, m_2, \dots, m_{n-1})$  вигляді:

$$A\vec{m} = F\vec{f}, \quad (124)$$

де

$$\vec{f} = (f_0, f_1, \dots, f_n)^T \quad (125)$$

**Зауваження 2:** Нагадаємо формулу для інтерполяційного багаточлена Лагранжа

$$L_n(x) = \sum_{i=0}^n f(x_i) \Phi_i^{(n)}(x), \quad (126)$$

де  $\Phi_i^{(n)}$  — множники Лагранжа. Це представлення інтерполяційного багаточлена Лагранжа по системі функцій  $\{\Phi_i^{(n)}\}$ . Для

$$\Phi_1 = \sum_{i=1}^n f(x_i) \varphi_i(x) \quad (127)$$

маємо представлення по системі кусково-лінійних функцій  $\{\varphi_i(x)\}$ . Для

$$\Phi_3(x) = \sum_{i=1}^n (f(x_i) \varphi_i^0(x) + f'(x_i) (\varphi_i^1)'(x)) \quad (128)$$

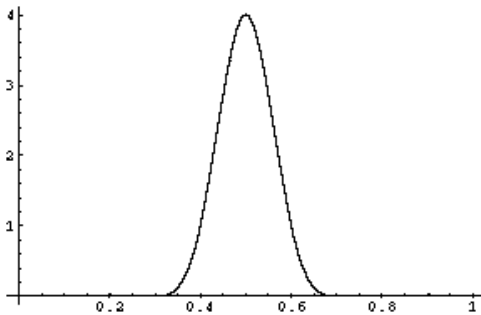
— представлення по системі  $\{\varphi_i^0, (\varphi_i^1)'\}$ .

Аналогічно, якщо представити кубічний сплайн у вигляді

$$s_3(x) = \sum_{i=0}^n c_i B_3^i(x), \quad (129)$$

то відповідна система для кубічного сплайну буде  $\{B_3^i(x)\}_{i=1}^n$ . Тут  $B_3^i(x)$  — так званий кубічний  $B_3$ -сплайн. Формула дається, а графік представлено на рис.:

$$B_3^i(z) = \frac{1}{6h} \begin{cases} \left(\frac{z-x_{i-2}}{h}\right)^3, & z \in [x_{i-2}, x_{i-1}]; \\ -3\left(\frac{z-x_{i-1}}{h}\right)^3 + 3\left(\frac{z-x_{i-1}}{h}\right)^2 + 3\left(\frac{z-x_{i-1}}{h}\right) + 1, & z \in [x_{i-1}, x_i]; \\ -3\left(\frac{x_{i+1}-z}{h}\right)^3 + 3\left(\frac{x_{i+1}-z}{h}\right)^2 + 3\left(\frac{x_{i+1}-z}{h}\right) + 1, & z \in [x_i, x_{i+1}]; \\ \left(\frac{x_{i+2}-z}{h}\right)^3, & z \in [x_{i+1}, x_{i+2}]; \\ 0, & z < x_{i-2} \vee x_{i+2} < z. \end{cases} \quad (130)$$



**Задача 19:** Показати, що  $B_3^i$  є кубічним сплайном дефекту 1.

Для знаходження коефіцієнтів  $c_i$  записується СЛАР з умов інтерполявання.

## 6.12. Параметричні сплайни

На практиці часто виникає задача побудови кривої по заданим точкам  $(x_i, y_i)_{i=1}^n$ . В цьому випадку використовують сплайни. Якщо відповідна функція  $y = f(x)$  однозначна, то сплайн будується за алгоритмом з попереднього пункту.

Окремо розглянемо випадок, коли точки  $(x_i, y_i)_{i=1}^n$  в площині  $(x, y)$  розташовані у довільний спосіб:

В цьому випадку відповідна функція задається параметрично

$$x = x(t), \quad y = y(t), \quad t \in [A, B]. \quad (131)$$

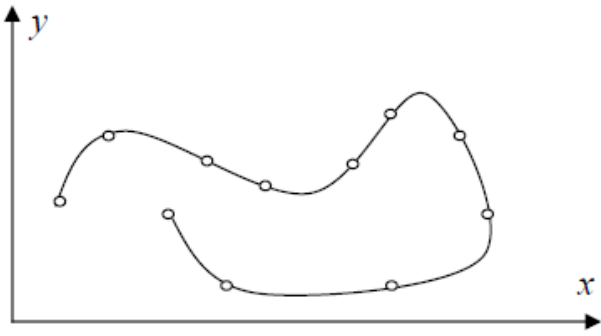
Для значень  $x_i, i = \overline{1, n}$  побудуємо кубічний сплайн  $s_x(t)$  такий, що  $s_x(t_i) = x_i, i = \overline{1, n}$ , а для  $y_i, i = \overline{1, n}$  будуємо сплайн  $s_y(t)$ , для якого  $s_y(t_i) = y_i, i = \overline{1, n}$ .

**Означення:** Тоді параметрична функція

$$(s_x(t), s_y(t)), \quad t \in [A, B]. \quad (132)$$

називається *параметричним сплайном* для функції (132).

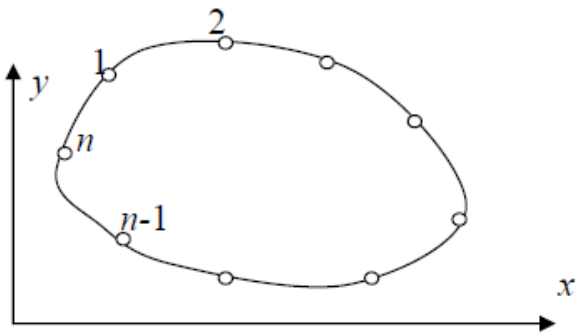




Стає питання про вибір параметру  $t$ . Нехай  $t_i = i, i = \overline{1, n}$ , тобто для табличних даних  $(x_i, y_i)_{i=1}^n$  параметром виступає номер точки в площині  $(x, y)$ . Тоді для параметричного сплайну неперервний параметр  $t$  змінюється на інтервалі  $t \in [1, n]$ .

Побудова сплайнів  $s_x(t)$  та  $s_y(t)$  здійснюється за алгоритмом наведеним в попередньому пункті по значенням  $f_i = x_i, i = \overline{1, n}$  та  $f_i = y_i, i = \overline{1, n}$ .

Розглянемо тепер побудову замкненої гладкої кривої:



Параметризуємо її як в попередньому випадку. Відмінність полягає в тому, що тепер функції  $x = x(t)$  та  $y = y(t)$  періодичні з періодом  $T = n$ , тобто

$$\forall t: \quad x(t) = x(t + n), \quad y(t) = y(t + n). \quad (133)$$

Наприклад, для значень в точках маємо:

$$x_1 = x_{n+1}, \quad y_1 = y_{n+1}; \quad x_0 = x_n, \quad y_0 = y_n. \quad (134)$$

Побудуємо алгоритм реалізації періодичного параметричного кубічного сплайну. Як і для звичайного сплайну на інтервалі  $t \in [t_i, t_{i+1}]$  маємо:

$$\begin{aligned} s(t) = & \frac{m_{i-1}(t_i - t)^3}{6h_i} + \frac{m_i(t - t_{i-1})^3}{6h_i} + \\ & + \left( f_{i-1} - \frac{m_{i-1}h_i^2}{6} \right) \frac{t_i - t}{h_i} + \\ & + \left( f_i - \frac{m_ih_i^2}{6} \right) \frac{t - t_{i-1}}{h_i}, \end{aligned} \quad (135)$$

де  $s(t)$  — одна з функцій  $s_x(t)$  або  $s_y(t)$ ;  $f_i = x_i, i = \overline{1, n}$  або  $f_i = y_i, i = \overline{1, n}$ ;  $h_i = t_{i+1} - t_i = 1$ . Для знаходження коефіцієнтів сплайну  $m_i = s''(t_i)$  з умови неперервності першої похідної сплайна маємо СЛАР:

$$\begin{cases} \frac{h_i m_{i-1}}{6} + \frac{(h_i + h_{i+1})m_i}{3} + \frac{h_{i+1}m_{i+1}}{6} = \\ = \frac{f_{i+1} - f_i}{h_i} - \frac{f_i - f_{i-1}}{h_i}, \quad i = \overline{1, n}, \\ m_0 = m_n, \quad m_1 = m_{n+1}, \end{cases} \quad (136)$$

Додаткові умови на коефіцієнти  $m_i$  впливають з періодичності сплайну та його других похідних.

Системі (74) відповідає матриця розмірності  $(n \times n)$ :

$$A = \begin{pmatrix} \frac{h_1+h_2}{3} & \frac{h_2}{6} & 0 & \cdots & 0 & \left\langle \frac{h_1}{6} \right\rangle \\ \frac{h_2}{6} & \frac{h_2+h_3}{3} & \frac{h_3}{6} & \ddots & & 0 \\ 0 & \frac{h_3}{6} & \frac{h_3+h_4}{3} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 \\ 0 & & \ddots & \ddots & \ddots & \frac{h_n}{6} \\ \left\langle \frac{h_1}{6} \right\rangle & 0 & \cdots & 0 & \frac{h_n}{6} & \frac{h_n+h_1}{3} \end{pmatrix} \quad (137)$$

яка є майже тридіагональною: «заважають» два елементи матриці, що виділені кутовими дужками.

Для розв'язання таких систем застосовують метод циклічної прогонки.

Розглянемо алгоритм цього методу для більш загальної системи:

$$\begin{cases} a_i y_{i-1} - c_i y_i + b_i y_{i+1} = -f_i, \quad i = \overline{1, n}, \\ y_0 = y_n, \quad y_{n+1} = y_1, \quad a_1 = a_n, \quad b_{n+1} = b_1, \end{cases} \quad (138)$$

Формули методу [ЛМС, стор. 391–392]:

1.  $\alpha_2 = b_1/c_1, \beta_2 = f_1/c_1, \gamma_2 = a_1/c_1;$
2.  $z_i = c_i - a_i \alpha_i; \alpha_{i+1} = b_i/z_i; \beta_{i+1} = (f_i + a_i \beta_i)/z_i; \gamma_{i+1} = a_i \gamma_i/z_i, i = \overline{2, n};$
3.  $p_{n-1} = \beta_n; q_{n-1} = \alpha_n + \gamma_n;$
4.  $p_i = \alpha_{i+1} p_{i+1} + \beta_{i+1}; q_i = \alpha_{i+1} q_{i+1} + \gamma_{i+1}, i = \overline{n-2, 1};$
5.  $y_n = (\beta_{n+1} + a_{n+1} p_1)/(1 - \alpha_{n+1} q_1 - \gamma_{n+1});$
6.  $y_i = p_i + y_n q_i, i = \overline{1, n-1}$

Метод стійкий ( $|\alpha_i| < 1, 1 - \alpha_{n+1} \alpha_1 - \gamma_{n+1} \neq 0$ ), якщо  $a_i, b_i > 0, c_i > b_i + a_i$ . Для системи (74) ці умови виконані.

Метод економний, оскільки кількість арифметичних операцій, що витрачається на його реалізацію,  $Q = O(n)$ .

Розглянуті в цьому пункті параметричні сплайни мають хороші апроксимативні та екстремальні властивості, тому побудовані по ним криві добре відновлюють задані як при малій, так досить великій кількості точок інтерполювання

### 6.13. Застосування інтерполювання

1. Складання таблиць. Нехай  $r_1^i(x)$  залишковий член лінійної інтерполяції по двох сусідніх точках  $x_{i-1}, x_i$ :

$$r_1^i(x) = f(x) - L_1^i(x) = \frac{f''(\xi)}{2!} \cdot (x - x_{i-1})(x - x_i). \quad (139)$$

Тоді

$$|r_1^i(x)| \leq \frac{M_2^i}{2} \cdot |(x - x_{i-1})(x - x_i)| \leq \frac{M_2^i h^2}{8}, \quad x \in [x_{i-1}, x_i]. \quad (140)$$

Таким чином

$$|f(x) - L_1^i(x)|_{C([a,b])} \leq \frac{M_2 h^2}{8}. \quad (141)$$

Ця оцінка може бути використана при складанні таблиць функцій, які при відновлення проміжних значень лінійною інтерполяцією сусідніх значень забезпечують точність  $\varepsilon$ .

Для того, щоб похибка була меншою за  $\varepsilon$  потрібно вибрати

$$h \leq \sqrt{\frac{8\varepsilon}{M_2}}. \quad (142)$$

Аналогічно, для квадратичного інтерполювання маємо

$$|f(x) - L_2^i(x)|_{C([a,b])} \leq \frac{M_3 h^3}{9\sqrt{3}} < \varepsilon. \quad (143)$$

Звідси

$$h \leq \sqrt[3]{\frac{9\sqrt{3}\varepsilon}{M_3}}. \quad (144)$$

2. Розв'язування рівнянь. Нехай необхідно розв'язати рівняння

$$f(x) = \bar{y}. \quad (145)$$

При  $\bar{y} = 0$  маємо рівняння  $f(x) = 0$ . Нехай  $\bar{x}$  корінь рівняння (145).

1. Обернене інтерполювання. Якщо відома обернена функція  $x = x(y)$ , то  $\bar{x} = x(\bar{y})$ . Нехай функція  $f(x)$  задана значеннями  $y_i = f(x_i)$ ,  $x_i \in [a, b]$ . Побудуємо інтерполяційний багаточлен  $L_n(y)$  по значеннях  $\{y_i, x_i\}_{i=0}^n$  де  $y_i$  вважаються значеннями аргументу, а  $x_i$  — значеннями оберненої функції. Тоді наближення до кореня є  $x^* = L_n(y)$ .

Оцінимо похибку:

$$|\bar{x} - x^*| = |x(\bar{y}) - L_n(\bar{y})| \leq \frac{M_{n+1}}{(n+1)!} |\omega_n(\bar{y})|, \quad (146)$$

$$\text{де } M_{n+1} = \max_{y_{\min} \leq y \leq y_{\max}} \left| \frac{d^{n+1}}{dy^{n+1}} x(y) \right|, \quad |\omega_n(y)| = (y - y_0) \dots (y - y_n).$$

**Недоліком** методу є складність обчислення похідних старших порядків оберненої функції.

2. Пряме інтерполювання. Нехай знову відомі  $y_i = f(x_i)$ ,  $x_i \in [a, b]$ . Тоді замість рівняння (145) розв'язуємо рівняння

$$L_n(x^*) = y, \quad (147)$$

де  $L_n(x)$  інтерполяційний багаточлен по значенням  $\{x_i, y_i\}_{i=0}^n$ .

**Недоліками** методу є необхідність розв'язування алгебраїчних рівнянь  $n$ -го степеня та необхідність вибирати шуканий розв'язок серед  $n$  коренів багаточлена степеня  $n$ .

Але **позитивним** є те, що функція є все таки алгебраїчною (а саме багаточленом).

Оцінимо похибку такого способу обчислення кореня. Маємо:

$$f(x^*) - L_n(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \cdot \omega_n(x). \quad (148)$$

Далі  $f(x^*) - y = f(x^*) - f(\bar{x})$ , звідки

$$|f(x^*) - f(\bar{x})| \leq \frac{M_{n+1}}{(n+1)!} \cdot |\omega_n(x)|. \quad (149)$$

Тут тепер  $M_{n+1} = \max_{x \in [a, b]} |f^{(n+1)}(x)|$ .

За теоремою Лагранжа  $f(x^*) - f(\bar{x}) = f'(\eta)(x^* - \bar{x})$ .

Припустимо, що  $f'(x) \neq 0$ . Це означає, що на проміжку  $[a, b]$  функція  $f(x)$  монотонна. Звідси

$$|x^* - \bar{x}| \leq \frac{|f(x^*) - f(\bar{x})|}{\min_{x \in [a, b]} |f'(x)|} \leq \frac{M_{n+1}}{\min_{x \in [a, b]} |f'(x)|} \cdot \frac{|\omega_n(x)|}{(n+1)!}. \quad (150)$$

### 3. Метод інтерполювання побудови характеристичного багаточлена.

Одним з найпростіших методів побудови характеристичного багаточлена є наступний. Відомо, що багаточлен  $n$ -го степеня однозначно визначається своїми значеннями в  $(n+1)$ -й точці. Тому для побудови  $P_n(\lambda) = \det(A - \lambda E)$  виберемо на проміжку де знаходяться власні значення (наприклад,  $\lambda \in [-\|A\|_k, \|A\|_k]$ , де  $k = 1$  або  $k = \infty$ ) деякі точки  $\lambda_i, i = \overline{0, n}$ . За допомогою методу Гауса для несиметричних матриць або методу квадратних коренів для симетричних матриць обчислимо  $P_n(\lambda_i) = \det(A - \lambda_i E)$  і по цих значення за формулою, наприклад, Ньютона побудуємо інтерполяційний багаточлен, який співпадатиме з характеристичним.

Далі розв'язується рівняння  $P_n(\lambda) = 0$  одним з відомих методів для нелінійного рівняння. Характерно, що часто для цього використовують метод парабол (обернене інтерполювання по трьох точках, або заміна рівняння  $n$ -го степеня в околі кореня на квадратне рівняння за допомогою інтерполяційного багаточлена другого степеня).

Зауважимо, що знаходження власних значень за допомогою характеристичного багаточлена пов'язана з проблемою нестійкості коренів характеристичного багаточлена відносно похибок обчислення коефіцієнтів цього багаточлена. Тому застосовують цей метод для невеликих розмірностей  $n \leq 10$  матриці  $A$ .

## 6.14. Тригонометрична інтерполяція

Інтерполяція відбувається за системою функцій

$$1, \sin(x), \cos(x), \sin(2x), \cos(2x), \dots, \sin(kx), \cos(kx), \dots \quad (151)$$

що є відрізком тригонометричного ряду Фур'є. Щоб знайти  $T_n(x)$  потрібно визначити  $2n+1$  коефіцієнт, а значить задати  $(2n+1)$  значень періодичної з періодом  $2\pi$  функції  $y_i = f_i(x), i = \overline{0, 2n}$ .

Покажемо, що

$$T_n(x) = \sum_{i=0}^{2n} f(x_i) \Phi_i(x), \quad (152)$$

де

$$\Phi_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^{2n} \frac{\sin\left(\frac{x-x_j}{2}\right)}{\sin\left(\frac{x_i-x_j}{2}\right)}, \quad (153)$$

тобто  $T_n(x_k) = f(x_k)$ , та  $\Phi_i(x_k) = \delta_{ik}$ . Дійсно

$$\Phi_i(x_k) = \prod_{\substack{j=0 \\ j \neq i}}^{2n} \frac{\sin\left(\frac{x_k-x_j}{2}\right)}{\sin\left(\frac{x_i-x_j}{2}\right)}, \quad (154)$$

для  $i \neq k$ , а

$$\Phi_i(x_i) = \prod_{\substack{j=0 \\ j \neq i}}^{2n} \frac{\sin\left(\frac{x_i-x_j}{2}\right)}{\sin\left(\frac{x_i-x_j}{2}\right)} = 1. \quad (155)$$

Таким чином за допомогою формули (152) ми уникли необхідності підраховувати коефіцієнти Фур'є  $a_k, b_k$ .

Нехай функція  $f(x)$  є парною та неперервною на проміжку  $[-\pi, \pi]$ . Тоді по значенням в  $(n+1)$ -й точці,  $y_i = f_i(x), i = \overline{0, n}, x_i \in [0, \pi]$  можна побудувати парний тригонометричний багаточлен:

$$T_n(x) = \sum_{i=0}^n f(x_i) \prod_{\substack{j=0 \\ j \neq i}}^n \frac{\cos(x) - \cos(x_i)}{\cos(x_i) - \cos(x_j)}. \quad (156)$$

Якщо ж функція є непарною на проміжку  $[-\pi, \pi]$ , то по значенням в  $n$  точках  $y_i = f_i(x), i = \overline{1, n}, x_i \in [0, \pi]$  можна побудувати непарний інтерполяційний багаточлен:

$$T_n(x) = \sum_{i=1}^n f(x_i) \frac{\sin(x)}{\sin(x_i)} \prod_{\substack{j=1 \\ j \neq i}}^n \frac{\cos(x) - \cos(x_i)}{\cos(x_i) - \cos(x_j)}. \quad (157)$$

**Задача 20:** Показати, що тригонометричні багаточлени (156), (157) є інтерполуючими для функції  $f(x)$ . Яке значення функції інтерполуює (157) при  $x = 0$ ? Чому?

## 6.15. Двовимірна інтерполяція

Побудова багаточлена для функції від двох змінних  $z = f(x, y)$ , що інтерполуює значення  $z_i = f(x_i, y_i)$  в точках  $A_i(x_i, y_i)$ , пов'язана з такими труднощами

1. Якщо в одновимірному випадку кількість вузлів та степінь багаточлена пов'язані простою залежністю:  $n+1$  точка  $x_i$  дозволяють побудувати багаточлен  $n$ -го степеня  $L_n(x)$ , то в двовимірному випадку багаточлен  $n$ -го степеня від двох змінних

$$P_n(x, y) = \sum_{0 \leq k+m \leq n} a_{k,m} x^k y^m, \quad (158)$$

має  $N = \frac{(n+1)(n+2)}{2}$  коефіцієнтів  $a_{k,m}$ . Тому необхідно задати значення в точках  $A_i(x_i, y_i)_{i=1}^N$ .

2. Не всяке розташування вузлів допустиме. Якщо розглянути умови інтерполювання

$$P_n(x_i, y_i) = \sum_{0 \leq k+m \leq n} a_{k,m} x_i^k y_i^m = z_i, \quad (159)$$

то для розв'язності цієї СЛАР необхідно виконання умови  $\det B \neq 0$ , де матриця  $B$  має коефіцієнти:

$\begin{equation}$

$\end{equation}$

Ця умова, наприклад, для лінійної інтерполяції  $n = 1$  та  $N = 3$  вимагає, щоб вузли  $A_i(x_i, y_i)$  не лежали на одній прямій. Якщо  $n = 2$ , то  $N = 6$  і необхідно розглядати точки, які не лежать на деякій кривій другого порядку і т. д.

Розглянемо випадки, коли можна записати багаточлен для двовимірної інтерполяції в явному вигляді.

Нехай область, в якій інтерполюється функція є прямокутником:

$$\bar{\Omega} = (x, y) : 0 \leq x \leq L_1, 0 \leq y \leq L_2. \quad (160)$$

Введемо сітку

$$x_i = ih_x, \quad h_x = L_1/N, \quad i = \overline{0, N}; \quad y_j = jh_y, \quad h_y = L_2/M, \quad j = \overline{0, M} \quad (161)$$

Тоді інтерполяційний багаточлен має вигляд

$$P(x, y) = \sum_{i=0}^N \sum_{j=0}^M f(x_i, y_j) \prod_{\substack{p=0 \\ p \neq i}}^N \prod_{\substack{q=0 \\ q \neq j}}^M \left( \frac{x - x_i}{x_p - x_i} \cdot \frac{y - y_j}{y_q - y_j} \right) \quad (162)$$

Розглянемо випадок, коли  $N = M = 1$ . Тоді

$$\begin{aligned} P_{1,1}(x, y) &= f(x_0, y_0) \cdot \frac{x - x_1}{x_0 - x_1} \cdot \frac{y - y_1}{y_0 - y_1} + f(x_0, y_1) \cdot \frac{x - x_1}{x_0 - x_1} \cdot \frac{y - y_0}{y_1 - y_0} + \\ &+ f(x_1, y_0) \cdot \frac{x - x_0}{x_1 - x_0} \cdot \frac{y - y_1}{y_0 - y_1} + f(x_1, y_1) \cdot \frac{x - x_0}{x_1 - x_0} \cdot \frac{y - y_0}{y_1 - y_0} = \\ &= a_0 + a_1 x + a_2 y + a_{1,2} xy. \end{aligned} \quad (163)$$

Це так звана білінійна інтерполяція, тобто лінійна по кожній окремій змінній.

Формула (162) являє собою приклад інтерполювання на всій області. В одновимірному випадку при великих степенях багаточлена отримують погане наближення через розбіжність процесу інтерполювання. Так ж картина має місце і в двовимірному випадку. Тому найчастіше застосовують кусково-поліноміальну апроксимацію.

Коротко наведемо деякі відомості про кусково-поліноміальне інтерполювання з теорії методу скінчених елементів розв'язання крайових задач для диференціальних рівнянь в частинних похідних.

Нехай область  $\Omega \subset \mathbb{R}^2$  — багатокутник в площині. Представимо її у вигляді

$$\Omega = \bigsqcup_{i=1}^n K_i, \quad (164)$$

де  $K_i \in T_h$ .

**Означення:**  $T_h$  називається «триангуляцією» області  $\Omega$ ,

а  $K_i$  — багатокутники з непорожньою внутрішністю що не мають спільних внутрішніх точок, причому  $\text{diam } K_i \leq h$ , де  $h$  — характеристика щільності розбиття.

Найчастіше  $K_i$  це трикутники або прямокутники.

Нехай  $v \in X$  — функція, яку ми будемо інтерполювати. Позначимо  $X_h$  простір, що апроксимує  $X$ , а його елементи  $c_h \in X_h$ . Причому звуження цієї функції на область  $K_i$ , тобто  $v_h|_{K_i}$  є поліномом.

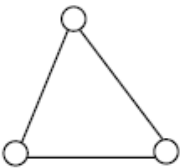
Позначимо  $\Pi_k, k \geq 0$  — простір багаточленів степеня  $k$  по сукупності змінних  $x, y$ ; його розмірність  $\dim \Pi_k = \frac{(k+1)(k+2)}{2}$ . Нехай  $\Theta_k, k \geq 0$  — простір багаточленів степеня по кожній окремій змінній  $x, y$ ; його розмірність  $\dim \Theta_k = (k+1)^2$ .

**Наприклад,**  $P_1(x, y) = a_0 + a_1x + a_2y \in \Pi_1$  — поліном степеня 1 по  $x, y$ , а  $Q_1(x, y) = a_0 + a_1x + a_2y + a_{1,2}xy \in \Theta_1$  лінійна по кожній окремій змінній.

Позначимо через  $X_h^k = \{v_h \in C^0(\Omega) : v_h|_k, \forall k \in T_k\}$  — простір інтерполантів при розбитті області на трикутники, а  $Y_h^k = \{v_h \in C^0(\Omega) : v_h|_k \in \Theta_k, \forall k \in T_h\}$  — при розбитті на прямокутники.

**Приклад 1:** Утворимо  $X_h^1, k = 1$ .

Будуємо багаточлен 1-го степеня по двох змінних. Оскільки  $\dim \Pi_1 = 3$ , то для цього треба задати значення функції в трьох точках. Точки, які задано —  $A_i, i = \overline{1,3}$  вибираємо вершинами трикутника, як на малюнку:



Тоді поліном першого степеня  $z = P_1(x, y)$  є розв'язком такого рівняння відносно  $z$ :

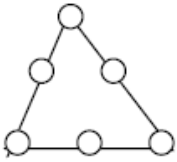
$$\begin{vmatrix} 1 & x & y & z \\ 1 & x_1 & y_1 & z_1 \\ 1 & x_2 & y_2 & z_2 \\ 1 & x_3 & y_3 & z_3 \end{vmatrix} = 0. \quad (165)$$

Тут  $f_i = f(x_i, y_i), i = \overline{1,3}$ .

**Задача 21:** Знайти явний вигляд  $z = P_1(x, y)$  — інтерполяційного багаточлена по значенням в точках  $A_i = (x_i, y_i), i = 1, 2, 3$ .

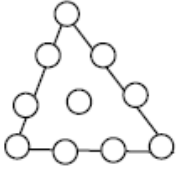
**Приклад 2:** Для  $X_h^2, k = 2, \dim \Pi_2 = 6$ .

Треба задати 6 значень, щоб забезпечити однозначність наближення. Тому вибираємо точки інтерполювання так:



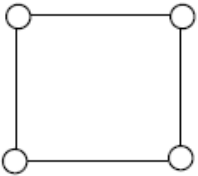
**Приклад 3:**  $X_h^2$ ,  $k = 3$ ,  $\dim \Pi_3 = 10$ .

Потрібно задати 10 точок, як на наступному малюнку:

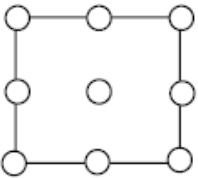


**Приклад 4:**  $Y_h^1$ ,  $k = 1$ ,  $\dim \Theta_1 = 4$ .

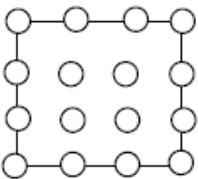
Формула для  $Q_1(x, y)$  наведена в (??). Точки:



**Приклад 5:**  $Y_h^2$ ,  $k = 2$ ,  $\dim \Theta_2 = 9$ .



**Приклад 6:**  $Y_h^3$ ,  $k = 3$ ,  $\dim \Theta_4 = 16$ .



Нехай  $X = W_2^m(\Omega) = H^m(\Omega)$  — це простір з нормою

$$\|v\|_m^2 = \sum_{k=0}^m \|v_k\|^2, \quad (166)$$

де

$$|v_k|^2 = \int_{\Omega} (D^m v)^2 d\Omega = \int_{\Omega} \left( (D_{x^k}^k)^2 + (D_{x^{k-1}y}^k)^2 + \dots + (D_{y^k}^k)^2 \right) d\Omega \quad (167)$$



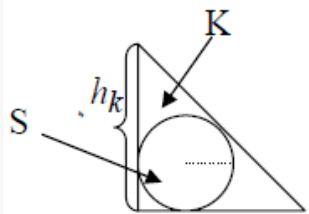
$$D_{x^k}^k v = \frac{\partial^k v}{\partial x^k}; \quad D_{x^{k-1}y}^k v = \frac{\partial^k v}{\partial x^{k-1} \partial y}; \quad \dots \quad (168)$$

Якщо  $\|v\|_m \leq M < \infty$ , то  $v \in W_2^m(\Omega)$ , класу функцій інтегрованих з квадратом до  $m$ -ї похідної.

Розглянемо розбиття на трикутники. Накладемо обмеження на них.

**Означення:** Розбиття  $T_h$  називається *регулярним*, якщо  $\exists \sigma \geq 1$  таке, що

$$\max_{k \in T_h} \frac{h_k}{\rho_k} \leq \sigma, \quad h_k = \text{diam } K, \quad S \subset K, \quad \rho_k = \mu(S) : \quad (169)$$



Якщо  $h_k/\rho_k \gg 1$ , то  $K$  вироджується в пряму і це погано.

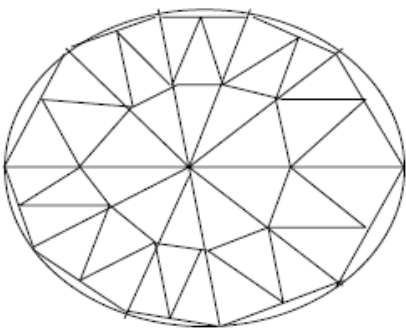
**Теорема:** Нехай  $v \in W_2^{l+1}(\Omega)$ ,  $1 \leq l \leq k$ ,  $T_h$  — регулярна триангуляція  $v_h \in X_h = \{v_h : v_h|_K = P_k\}$ . Тоді

$$|v - v_h|_m \leq Ch^{l+1-m} |v|_{l+1}, \quad m = 0, 1, \quad k \geq 1. \quad (170)$$

**Наприклад:** для  $k = 1, m = 0: l = 1, \|v - v_h\|_{L_2(\Omega)} = \|v - v_h\|_0 \leq Ch^2 |v|_2$ . Якщо ж  $k = 3, l = 3$ , то  $\|v - v_h\|_0 \leq h^4 |v|_4$ .

Ця теорема дозволяє стверджувати збіжність процесу інтерполювання. І чим більше степінь полінома на кожному елементі тим вища швидкість збіжності.

Узагальнимо результат теореми на область з гладкою границею:



Для цього вибираємо точки на границі і будуємо вписаний багатогранник. Його триангулюємо. Далі на кожному трикутнику будуємо інтерполант. В результаті отримуємо  $v_h \in X_h^k$ .

Тоді для  $k = 1, l = 1: \|v - v_h\|_0 \leq ch^{3/2} |v|_2$ .

## 7. Чисельне диференціювання

### 7.1. Побудова формул чисельного диференціювання

Задача чисельного диференціювання виникає у випадку коли необхідно обчислити похідну функції, значення якої задані таблицею. Нехай задано

$$f_i = f(x), \quad i = \overline{0, n}, \quad x_i \in [a, b]. \quad (1)$$

Проінтерполюємо ці значення. Тоді

$$f(x) = L_n(x) + r_n(x), \quad (2)$$

де залишковий член у формі Ньютона має вигляд:

$$r_n(x) = f(x; x_0, \dots, x_n) \omega_n(x), \quad (3)$$

де

$$\omega_n(x) = \prod_{i=0}^n (x - x_i). \quad (4)$$

Звідси

$$f^{(k)}(x) = L_n^{(k)}(x) + r_n^{(k)}(x). \quad (5)$$

За наближене значення похідної в точці  $x$  беремо  $f^{(k)}(x) \approx L_n^{(k)}(x)$ ,  $x \in [a, b]$ .

Оцінимо похибку наближення —  $r^{(k)}(x)$ . За формулою Лейбніца:

$$r_n^{(k)}(x) = \sum_{j=0}^k C_k^j f^{(j)}(x; x_0, \dots, x_n) \omega_n^{(k-j)}(x). \quad (6)$$

З властивості розділених різниць маємо для  $f(x) \in C^{n+k+1}[a, b]$ :

$$f^{(j)}(x; x_0, \dots, x_n) = j! \cdot \underbrace{f(x, \dots, x; x_0, \dots, x_n)}_{j+1} = \frac{j!}{(n+j+1)!} \cdot f^{(n+j+1)}(\xi_j), \quad (7)$$

де всі  $x_j \in [a, b]$ .

Остаточно вираз для похибки наближення похідної має вигляд:

$$r_n^{(k)}(x) = \sum_{j=0}^k \frac{k!}{(k-j)!(n+j+1)!} \cdot f^{(n+j+1)}(\xi_j) \omega_n^{(k-j)}(x). \quad (8)$$

Оцінка похибки матиме вигляд:

$$\left| f^{(k)}(x) - L_n^{(k)}(x) \right| \leq M \sum_{j=0}^k \frac{k!}{(k-j)!(n+j+1)!} \cdot \left| \omega_n^{(k-j)}(x) \right|, \quad (9)$$

де  $M = \max_{0 \leq j \leq k} \max_{x \in [a, b]} \left| f^{(n+j+1)}(x) \right|$ .

Нагадаємо, що процес інтерполювання розбіжний. Крім того, якщо  $k > n$ , то  $L_n^{(k)}(x) \equiv 0$ . Тому не можна брати великими значення  $n$  та  $k$ . Як правило  $k = 1, 2$ , іноді  $k = 3, 4$ . Відповідно,  $n = k$ , або  $n = k + 1$ , або  $n = k + 2$ .

Подивимося як залежить порядок збіжності процесу чисельного диференціювання від кроку. Нехай  $x_i = x_0 + ih$ ,  $h > 0$  — крок. Тоді за умови  $x_n - x_0 = O(h)$ :

$$\omega_n(x) = (x - x_0) \cdot \dots \cdot (x - x_n) = O(h^{n+1}), \quad (10)$$

де  $x \in [x_0, x_n]$ .

Перша похідна від  $\omega_n(x)$  має порядок на одиницю менше, тобто

$$\omega'_n(x) = O(h^n). \quad (11)$$

Далі

$$r_n^{(k)}(x) = O(h^{n+1-k}), \quad (12)$$

тому

$$f^{(k)}(x) - L_n^{(k)} = O(h^{n+1-k}). \quad (13)$$

При умові  $n \geq k$  останній вираз збігається до нуля, тобто

$$f^{(k)}(x) - L_n^{(k)}(x) \xrightarrow{h \rightarrow 0} 0. \quad (14)$$

Далі

$$r_n^{(k)}(x) = \underbrace{f(x; x_0, \dots, x_n) \omega_n^{(k)}(x)}_{O(h^{n+1-k})} + \underbrace{\sum_{j=1}^k C_k^j f^{(j)}(x; x_0, \dots, x_n) \omega_n^{(k-j)}(x)}_{O(h^{n+2-k})}. \quad (15)$$

Якщо

$$\omega_n^{(k)}(\bar{x}) = 0, \quad (16)$$

то

$$r_n^{(k)}(\bar{x}) = O(h^{n+2-k}). \quad (17)$$

Точки  $x = \bar{x}$  називаються *точками підвищеної точності формул чисельного диференціювання*.

**Приклад 1:** Виведемо формули чисельного диференціювання для  $k = 1$ ,  $n = 1$ .

Виберемо точки  $x_0, x_1 = x_0 + h$  і інтерполяційний багаточлен має вигляд:

$$L_1(x) = f_0 + (x - x_0) \cdot \frac{f_1 - f_0}{h}. \quad (18)$$

Для похідної отримаємо вираз:

$$f'(x) \approx L'_1(x) = \frac{f_1 - f_0}{h}, \quad x \in [x_0, x_1]. \quad (19)$$

Розписавши за формулою Тейлора, отримаємо вираз для похибки:

$$r'_1(x) = \frac{f^{(3)}(\xi_1)}{3!} \cdot (x - x_0)(x - x_1) + \frac{f^{(2)}(\xi_0)}{2!} \cdot (2x - x_1 - x_0) = O(h). \quad (20)$$

Якщо  $2\bar{x} - x_1 - x_0 = 0$ , то  $r'_1(\bar{x}) = O(h^2)$ . Тобто  $\bar{x} = \frac{x_1+x_0}{2}$  — точка підвищеної точності. Більш точно (див. приклад 3):

$$|r'_1(\bar{x})| \leq \frac{h^2 M_3}{24}, \quad (21)$$

де  $M_3 = \max_{x \in [a,b]} |f^{(3)}(x)|$ .

**Приклад 2:** Аналогічно виведемо формули чисельного диференціювання для  $k = 1, n = 2$ .

Виберемо точки  $x_0, x_1 = x_0 + h, x_2 = x_0 + 2h$ . Інтерполяційний поліном має вигляд:

$$L_2(x) = f_0 + (x - x_0) \cdot \frac{f_1 - f_0}{h} + (x - x_0)(x - x_1) \cdot \frac{f_2 - 2f_1 + f_0}{2h^2}. \quad (22)$$

Тоді замінімо  $f'(x) \approx L_2'(x) = \frac{f_1 - f_0}{h} + (2x - x_0 - x_1) \cdot \frac{f_2 - 2f_1 + f_0}{2h^2}, x \in [x_0, x_2]$ .

Якщо сюди підставити  $x = x_0$ , то отримаємо  $f'(x_0) \approx \frac{-f_2 + 4f_1 - 3f_0}{2h}$ . Для точки  $x = x_1$  маємо  $f'(x_1) \approx \frac{f_2 - f_0}{2h} = f_{x,1}^0$ . Для точки  $x = x_2$  маємо  $f'(x_2) \approx \frac{f_0 - 4f_1 + 3f_2}{2h}$ . Для похибки маємо оцінку  $r'_2(x) = O(h^2)$ .

Позначимо

- для  $x \in [x_i, x_{i+1}]$ ,  $f_{x,i} = \frac{f_{i+1} - f_i}{h} \approx f'(x)$ , (різницева похідна вперед);
- для  $x \in [x_{i-1}, x_i]$  —  $f_{\bar{x},i} = \frac{f_i - f_{i-1}}{h} \approx f'(x)$ , (різницева похідна назад);
- для  $x \in [x_{i-1}, x_{i+1}]$  —  $f_{x,i}^0 \approx f'(x)$  (центральна різницева похідна).

Замість  $f'(x_i)$  можна взяти будь-яке із значень:  $f_{x,i}$ ,  $f_{\bar{x},i}$  або  $f_{x,i}^0$ .

**Задача 21:** Знайти точки підвищеної точності формул чисельного диференціювання для  $k = 1, n = 2$  і оцінити похибку в цих точках.

**Приклад 3:** При  $n = 1, k = 1$  оцінимо точність формул чисельного диференціювання за формулою Тейлора.

1. Нехай  $f(x) \in C^2([a, b])$ . Тоді

$$f'(x_0) - \frac{f(x_0 + h) - f(x_0)}{h} = f'(x_0) - \frac{1}{h} \left( f_0 + hf'_0 + \frac{h^2}{2} f''(\xi) - f_0 \right) = -\frac{h}{2} \cdot f''(\xi), \quad (23)$$

$$\left| f'(x_0) - \frac{f_1 - f_0}{h} \right| \leq \frac{M_2 h}{2}, \quad (24)$$

де  $M_2 = \max_{[x_0, x_1]} |f''(\xi)|$ .

2. Нехай  $f(x) \in C^3([a, b])$ . Тоді, розписавши розклад по формулі Тейлора до третьої похідної, маємо оцінку:

$$\begin{aligned}
f'(\bar{x}) - \frac{f(x_0 + h) - f(x_0)}{h} &= \\
&= f'(\bar{x}) - \frac{1}{h} \left( f'(\bar{x}) + \frac{h}{2} \cdot f'(\bar{x}) + \frac{h^2}{8} \cdot f''(\bar{x}) + \frac{h^3}{48} \cdot f'''(\xi) \right) - \\
&\quad - f(\bar{x}) + \frac{f}{2} \cdot f'(\bar{x}) - \frac{h^2}{8} \cdot f''(\bar{x}) + \frac{h^3}{48} \cdot f'''(\eta) \Big) = -\frac{h^2}{24} \cdot f'''(\zeta).
\end{aligned} \tag{25}$$

$$\left| f'(\bar{x}) - \frac{f_1 - f_0}{h} \right| \leq \frac{h^2 M_3}{24}, \tag{26}$$

$$\text{де } \bar{x} = \frac{x_1 + x_0}{2}.$$

**Задача 22:** Показати, що якщо  $f(x) \in C^3([a, b])$ , то  $\left| f'(x_1) - \frac{f_2 - f_0}{2h} \right| \leq \frac{M_3 h^2}{6}$ .

**Приклад 4:** При  $n = 2, k = 2$  маємо:

$$\begin{aligned}
L_2(x) &= f_{i-1} + \frac{f_i - f_{i-1}}{h} \cdot (x - x_{i-1}) + \\
&\quad + \frac{f_{i+1} - 2f_i + f_{i-1}}{h^2} \cdot (x - x_{i-1}) \cdot (x - x_i)
\end{aligned} \tag{27}$$

$$L_2''(x) = \frac{f_{i+1} - 2f_i + f_{i-1}}{h^2}. \tag{28}$$

Для  $f(x) \in C^4([a, b])$  оцінимо точність формул чисельного диференціювання за формулою Тейлора:

$$\begin{aligned}
f''(x_1) - \frac{f_2 - 2f_1 + f_0}{h^2} &= \\
&= f''(x_1) - \frac{f(x_1 + h) - 2f(x_1) + f(x_1 - h)}{h^2} = \\
&= f''(x_1) - \frac{1}{h^2} \left( f_1 + hf_1'' + \frac{h^2}{2} \cdot f_1'' + \frac{h^3}{6} \cdot f_1''' + \frac{h^4}{24} \cdot f^{(4)}(\xi_1) - 2f_1 + \right. \\
&\quad \left. + f_1 - hf_1'' + \frac{h^2}{2} \cdot f_1'' - \frac{h^3}{6} \cdot f_1''' + \frac{h^4}{24} \cdot f^{(4)}(\xi_2) \right) = \frac{h^2}{12} \cdot f^{(4)}(\xi),
\end{aligned} \tag{29}$$

де  $\xi_1, \xi_2, \xi \in [x_0, x_2]$ .

Отже,

$$\left| f_1'' - \frac{f_2 - 2f_1 + f_0}{h^2} \right| \leq \frac{M_4 h^2}{12}. \tag{30}$$

**Задача 23:** Побудувати формулу чисельного диференціювання  $k = 2, n = 2$  у випадку нерівновіддалених вузлів:  $x_0, x_1 = x_0 + h_1, x_2 = x_1 + h_2$ . Оцінити точність формули. Знайти точки підвищеної точності оцінити похибку.

Крім інтерполяційних формул для чисельного диференціювання можна застосовувати сплайни. Нехай  $f_i = f(x_i)$ . Побудуємо інтерполяційний сплайн першого степеня  $s_1(x)$ , для якого має місце оцінка  $\left| f^{(k)}(x) - s_1^{(k)}(x) \right| = O(h^{2-k}), k = 0, 1$ . Звідси при  $k = 1$  маємо  $f'(x) - s_1'(x) = O(h)$ .

Для кубічного інтерполяційного сплайну  $s_3(x)$  маємо для першої та другої похідних:

$$\left| f^{(k)}(x) - s_3^{(k)}(x) \right| = O(h^{4-k}), \quad k = 1, 2. \quad (31)$$

## 7.2. Про обчислювальну похибку чисельного диференціювання

Нехай значення функції обчислені з деякою похибкою. Постає питання про вплив цих похибок на значення похідних обчислених за формулами чисельного диференціювання.

Перед цим зробимо зауваження про вплив збурення функції на значення звичайних похідних.

Нехай  $f(x) \in C^1([a, b])$  і її збурення має вигляд:

$$\tilde{f}(x) = f(x) + \frac{\sin(\omega x)}{n}. \quad (32)$$

При  $n \rightarrow \infty$  маємо  $\|f(x) - \tilde{f}(x)\|_{C([a,b])} = \frac{1}{n} \rightarrow 0$ , звідси  $\tilde{f}(x) \xrightarrow{n \rightarrow \infty} f(x)$ . Таким чином це малі збурення. Маємо  $\tilde{f}'(x) = f'(x) + \frac{\omega}{n} \cdot \cos(\omega x)$ . Нехай  $\omega = n^2$ , тоді

$$\left| f(x) - \tilde{f}(x) \right|_{C([a,b])} = \frac{|\omega|}{n} = n \xrightarrow{n \rightarrow \infty} \infty. \quad (33)$$

Цей приклад ілюструє нестійкість оператора диференціювання. Є сподівання, що ця нестійкість має місце і для чисельного диференціювання.

Нехай  $\tilde{f}_i = f_i + \delta_i$ ,  $f_i = f(x_i)$ ,  $i = \overline{0, n}$ ,  $|\delta_i| \leq \delta$ . Розглянемо вплив похибок  $\delta_i$  на конкретних формулах чисельного диференціювання.

**Приклад 1:** Оцінімо вплив збурень на похибку обчислення першої похідної  $n = 1$ ,  $k = 1$ .

$$f'_i = \frac{\tilde{f}_i - \tilde{f}_{i-1}}{h} = f'_i - \frac{f_i - f_{i-1}}{h} - \frac{\delta_i - \delta_{i-1}}{h}, \quad (34)$$

$$\left| f'_i - \frac{\tilde{f}_i - \tilde{f}_{i-1}}{h} \right| \leq \left| f'_i - \frac{f_i - f_{i-1}}{h} \right| + \left| \frac{\delta_i - \delta_{i-1}}{h} \right| \leq \frac{M_2 h}{2} + \frac{2\delta}{h} \xrightarrow{h \rightarrow 0} \infty \quad (35)$$

Таким чином, як і для аналітичного диференціювання, маємо некоректність: при малих збуреннях  $|\delta_i| \ll \delta_i$  можуть бути як завгодно великі похибки, якщо  $\frac{\delta}{h} \rightarrow \infty$  при  $h \rightarrow 0$ .

Мінімізуємо вплив цих збурень. Позначимо

$$\varphi(h) = \frac{M_2 h}{2} + \frac{2\delta}{h}. \quad (36)$$

Тоді мінімум цієї функції досягається для таких  $h$ :

$$\varphi'(h) = \frac{M_2}{2} - \frac{2\delta}{h^2} = 0, \quad (37)$$

звідки  $h_0 = 2\sqrt{\frac{\delta}{M_2}}$ . При такому значенні  $h$  оцінка похибки (35) така:

$$\varphi(h_0) = 2\sqrt{M_2 \delta} = O(\sqrt{\delta}) \xrightarrow{\delta \rightarrow 0} 0. \quad (38)$$

**Приклад 2:** Подивимось на вплив збурень на похибку обчислення першої похідної при використанні центральної різницевої похідної.

$$\left| f'_i - \frac{\tilde{f}_{i+1} - \tilde{f}_{i-1}}{2h} \right| \leq \left| f'_i - \frac{f_{i+1} - f_{i-1}}{2h} \right| + \left| \frac{\delta_{i+1} - \delta_{i-1}}{2h} \right| \leq \frac{M_3 h^2}{6} + \frac{\delta}{h} = \varphi(h). \quad (39)$$

З рівняння  $\varphi'(h) = \frac{M_3 h}{3} - \frac{\delta}{h^2} = 0$  маємо:  $h_0^3 = \frac{3\delta}{M_3}$ ,  $h_0 = \sqrt[3]{\frac{3\delta}{M_3}}$ . Отже,

$$\varphi(h_0) = \frac{M_3}{6} \sqrt[3]{\frac{9\delta^2}{M_3^2}} + \frac{\delta}{\sqrt[3]{\frac{3\delta}{M_3}}} = \frac{1}{2} \sqrt[3]{\frac{M_3 \delta^2}{3}} + \sqrt[3]{\frac{M_3 \delta^2}{3}} = \frac{3}{2} \sqrt[3]{\frac{M_3 \delta^2}{3}} = O\left(\sqrt[3]{\delta^2}\right). \quad (40)$$

Таким чином швидкість збіжності при  $\delta \rightarrow 0$  похибки формули чисельного диференціювання центральною похідною вища ніж для формули з [прикладу 1](#) (похідна вперед або назад).

**Задача 24:** Дослідити похибку чисельного диференціювання для  $n = 2$ ,  $k = 2$ , вибрати оптимальний крок  $h_0$ , дати оцінку  $\varphi(h_0)$ .

## 8. Апроксимування функцій

### 8.1. Постановка задачі апроксимації

Наближення функцій застосовують у випадках, якщо

- функція складна (трансцендентна або є розв'язком складної задачі) і її замінюють функцією, яка легко обчислюється (найчастіше, поліномом);
- необхідно побудувати функцію неперервного аргументу для функції, яка задана своїми значеннями (таблична);
- таблична функція наближається табличною ж функцією (згладжування).

Інтерполювання не кращий спосіб наближення функцій через розбіжність цього процесу для поліномів. Тим більше доцільність застосування інтерполювання сумнівна, якщо функція таблична, а її значення неточні. Потрібно будувати апроксимуючу функцію з інших міркувань.

Найбільш загальний принцип: наблизити  $f(x)$  функцією  $\Phi(x)$  так, щоб досягалася деяка задана точність  $\varepsilon$ :

$$|f(x) - \Phi(x)| < \varepsilon \quad (1)$$

Але розв'язок в такій постановці може не існувати або бути не єдиним.

Загальна постановка задачі наближення така. Нехай маємо елемент  $f$  лінійного нормованого простору  $R$ . Побудуємо підпростір  $M_n$ , в якому елементи є лінійною комбінацією:

$$\Phi = \sum_{i=0}^n c_i \varphi_i \in M_n \subset R \quad (2)$$

по елементах лінійно незалежної системи

$$\{\varphi_i\}_{i=0}^{\infty}, \quad \varphi_i \in R. \quad (3)$$

Відхилення  $\Phi \in M_n$  від  $f \in R$  є число

$$\Delta(f, \Phi) = |f - \Phi|. \quad (4)$$

Позначимо

$$\inf_{\Phi \in M_n} |f - \Phi| = \Delta(f). \quad (5)$$

**Означення.** Елемент  $\Phi_0$  такий, що

$$\Delta(f, \Phi_0) = |f - \Phi_0| = \inf_{\Phi \in M_n} |f - \Phi| = \Delta(f), \quad (6)$$

називається *елементом найкращого наближення* (ЕНН).

Ясно, що умову точності треба перевіряти на цьому елементі. У випадку її невиконання треба збільшувати кількість елементів  $n$  в (2).



**Теорема 1:** Для будь-якого лінійного нормованого простору  $R$  існує елемент найкращого наближення  $\Phi_0 \in M_n$ .

*Доведення:* Введемо

$$F(\vec{c}) = F(c_0, c_1, \dots, c_n) = |f - \Phi| = \left| f - \sum_{i=0}^n c_i \varphi_i \right|. \quad (7)$$

Це неперервна функція аргументів  $\vec{c} = (c_0, c_1, \dots, c_n)$ . Для елементів, які задовольняють умові

$$|\Phi| > 2|f|, \quad f \in R_1, \quad \Phi \in M_n, \quad (8)$$

маємо

$$F(\vec{c}) = |f - \Phi| \geq |\Phi| - |f| > 2|f| - |f| = |f| > \Delta(f). \quad (9)$$

Значить ЕНН  $\Phi_0 \in \{\Phi : \|\Phi\| \leq 2\|f\|\} = \overline{U} \subset M_n$ . За теоремою Кантора  $\exists \Phi_0$ , де  $F(\vec{c})$  досягає мінімуму. Причому  $\|f - \Phi_0\| \leq \|f - \Phi\|$ .  $\square$

Елементів найкращого наближення в лінійному нормованому просторі може бути і декілька.

**Означення:** Простір  $R$  називається *строго нормованим*, якщо з умови

$$|f + g| = |f| + |g|, \quad |f| \neq 0, \quad |g| \neq 0 \quad (10)$$

випливає, що  $\exists \lambda \neq 0$  таке, що

$$g = \lambda f. \quad (11)$$

**Теорема 2:** Якщо простір  $R$  строго нормований, то елемент найкращого наближення  $\Phi_0$  єдиний.

*Доведення:* від супротивного. Нехай існують  $\Phi_0^{(1)} \neq \Phi_0^{(2)}$  — два елементи найкращого наближення. Візьмемо  $\alpha \in [0, 1]$ , тоді

$$\begin{aligned} \Delta(f) &\leq \|f - \alpha \Phi_0^{(1)} - (1 - \alpha) \Phi_0^{(2)}\| = \|\alpha (f - \Phi_0^{(1)}) + (1 - \alpha) (f - \Phi_0^{(2)})\| \leq \\ &\leq \alpha \|f - \Phi_0^{(1)}\| + (1 - \alpha) \|f - \Phi_0^{(2)}\| = \alpha \Delta(f) + (1 - \alpha) \Delta(f) = \Delta(f). \end{aligned} \quad (12)$$

Тобто всі « $\leq$ » можна замінити на « $=$ » Отримаємо

$$|\alpha (f - \Phi_0^{(1)}) + (1 - \alpha) (f - \Phi_0^{(2)})| = \alpha |f - \Phi_0^{(1)}| + (1 - \alpha) |f - \Phi_0^{(2)}|. \quad (13)$$

За припущенням  $\exists \lambda$  таке, що

$$\alpha (f - \Phi_0^{(1)}) + \lambda(1 - \alpha) (f - \Phi_0^{(2)}) = 0. \quad (14)$$

Виберемо  $\alpha = 1/2$ . Тоді

$$f - \Phi_0^{(1)} = \lambda (f - \Phi_0^{(2)}). \quad (15)$$

Оскільки

$$|f - \Phi_0^{(1)}| = |f - \Phi_0^{(2)}| = \Delta(f), \quad (16)$$

то остання рівність має місце тільки для  $\lambda = 1$ . Звідси

$$f - \Phi_0^{(1)} = f - \Phi_0^{(2)} \implies \Phi_0^{(1)} = \Phi_0^{(2)}. \quad (17)$$

Отже, ми отримали протиріччя з припущенням, що і доводить існування єдиного елемента найкращого наближення.  $\square$

**Теорема 3:** Гільбертів простір  $H$  — строго нормований.

*Доведення:* Нехай

$$|f + g| = |f| + |g|, \quad (18)$$

$$|f + g|^2 = |f|^2 + 2|f| \cdot |g| + |g|^2. \quad (19)$$

З іншого боку

$$|f + g|^2 = \langle f + g, f + g \rangle = |f|^2 + 2\langle f, g \rangle + |g|^2. \quad (20)$$

Звідси  $\|f\| \cdot \|g\| = \langle f, g \rangle$ . Для довільного гільбертового простору  $\langle f, g \rangle \leq \|f\| \cdot \|g\|$ .

Таким чином на елементах (18) нерівність Коші-Буняковського перетворюється в рівність. Розглянемо

$$|f - \lambda g|^2 = |f|^2 - 2\lambda \langle f, g \rangle + \lambda^2 |g|^2 = |f|^2 - \lambda |f| \cdot |g| + \lambda^2 |g|^2 = (|f| - \lambda |g|)^2. \quad (21)$$

Тоді для  $\lambda = \|f\|/\|g\|$  маємо  $\|f - \lambda g\| = 0$ . Звідси  $\exists \lambda: f = \lambda g$ , тобто  $H$  — строго нормований.  $\square$

**Наслідок:**  $R = H \implies \exists! \Phi_0 \in M_n$ .

**Приклади** (строго нормованих просторів):

$$1. L_2([a, b]) \text{ з нормою } \|u\| = \sqrt{\int_a^b u^2 dx}.$$

$$2. L_p([a, b]) \text{ з нормою } \|u\| = \left( \int_a^b u^p dx \right)^{1/p}, p > 1.$$

Простір  $C([a, b])$  не є строго нормованим, але в ньому існує єдиний елемент найкращого наближення (про цей факт в наступному пункті).

## 8.2. Найкраще рівномірне наближення

**Означення:** Найкраще рівномірне наближення — це наближення в просторі  $R = C([a, b])$ , де  $\|f\|_{C([a, b])} = \max_{x \in [a, b]} |f|$  — рівномірна метрика.

**Теорема 1 (Хаара):** Для того, щоб  $\forall f \in C([a, b])$  існував єдиний елемент найкращого рівномірного наближення необхідно і достатньо, щоб система  $\{\varphi_i\}_{i=0}^\infty$  була системою Чебишова.

**Означення:** Система  $\{\varphi\}_{i=1}^{\infty}$  називається *системою Чебишова*, якщо елемент  $\Phi_n(x) = \sum_{i=0}^n c_i \varphi_i(x)$  має не більше  $n$  нулів, причому  $\sum_{i=1}^n c_i^2 \neq 0$ .

**Наприклад,** системою Чебишова є поліноміальна система  $\{x^i\}_{i=0}^{\infty}$ .

**Означення:** Позначимо  $Q_n^0(x)$  — *багаточлен найкращого рівномірного наближення* (далі — БНРН.).

Його відхилення від  $f$ :

$$\Delta(f) = \|Q_n^0(x) - f(x)\|_C = \inf_{Q_n(x)} \|Q_n(x) - f(x)\|. \quad (22)$$

**Теорема 2 (Чебишова):**  $Q_n^0(x)$  — БНРН неперервної функції  $f(x)$  тоді та тільки тоді, якщо на відрізку  $[a, b]$  існує хоча б  $(n+2)$ -а точки  $a \leq x_0 \leq \dots \leq x_m \leq b, m \geq n+1$  такі, що

$$f(x_i) - Q_n^0(x_i) = \alpha(-1)^i \Delta(f), \quad (23)$$

де  $i = \overline{0, m}, \alpha = \pm 1$ .

**Означення:** Точки  $\{x_i\}_{i=0}^m$ , які задовольняють умовам теореми Чебишова, називаються *точками чебишовського альтернансу*.

**Теорема 3:**  $Q_n^0(x)$  — БНРН для неперервної функції єдиний.

*Доведення.* Припустимо, існують два БНРН степеня  $n$ :  $Q_n^{(1)}(x) \neq Q_n^{(2)}(x)$ :

$$\Delta(f) = \|f - Q_n^{(1)}\|_C = \|f - Q_n^{(2)}\|_C. \quad (24)$$

Звідси випливає, що

$$\left| f - \frac{Q_n^{(1)}(x) + Q_n^{(2)}(x)}{2} \right| \leq \left| \frac{f - Q_n^{(2)}(x)}{2} \right| = \Delta(x), \quad (25)$$

тобто багаточлен

$$\frac{Q_n^{(1)}(x) + Q_n^{(2)}(x)}{2} \quad (26)$$

також є БНРН. Нехай  $x_0, x_1, \dots, x_m$  — відповідні йому точки чебишовського альтернансу.

Це означає, що

$$\left| \frac{Q_n^{(1)}(x_i) + Q_n^{(2)}(x_i)}{2} - f(x_i) \right| = \Delta(f), \quad (27)$$

або

$$\left(Q_n^{(1)}(x_i) - f(x_i)\right) + \left(Q_n^{(2)}(x_i) - f(x_i)\right) = 2\Delta(f). \quad (28)$$

Оскільки  $\left|Q_n^{(k)}(x_i) - f(x_i)\right| \leq \Delta(f)$ ,  $k = 1, 2$ , то (28) можливе лише у тому випадку, коли

$$Q_n^{(1)}(x_i) - f(x_i) = Q_n^{(2)}(x_i) - f(x_i), \quad (29)$$

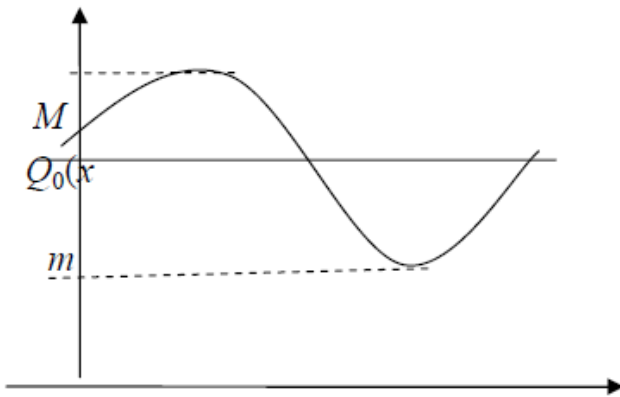
для усіх  $i = \overline{0, n+1}$ .

Звідки випливає, що  $Q_n^{(1)}(x) = Q_n^{(2)}(x)$ , а це суперечить початковому припущенню.  $\square$

### 8.3. Приклади побудови БНРН

Скінченного алгоритму побудови БНРН для довільної функції не існує. Є ітераційний [ЛМС, 73–79]. Але в деяких випадках можна побудувати БНРН за теоремою Чебишова.

1. Потрібно наблизити багаточленом нульового степеня.



Нехай  $M = \max_{[a,b]} f(x) = f(x_0)$ ,  $m = \min_{[a,b]} f(x) = f(x_1)$ , тоді  $Q_0(x)$  — БНРН має вигляд (див. рис. 11):

$$Q_0(x) = \frac{M+m}{2}, \quad (30)$$

де  $\Delta(x_0) = \frac{M+m}{2}$ , а  $x_0, x_1$  — точки чебишовського альтернансу.

2. Опукла функція  $f(x) \in C([a, b])$  наближається багаточленом першого степеня

$$Q_1(x) = c_0 + c_1x. \quad (31)$$

Оскільки  $f(x)$  опукла, то різниця  $f(x) - (c_0 + c_1x)$  може мати лише одну внутрішню точку екстремуму. Тому точки  $a, b$  є точками чебишовського альтернансу. Нехай  $\xi$  третя — точка чебишовського альтернансу. Згідно з теоремою Чебишова, маємо систему:

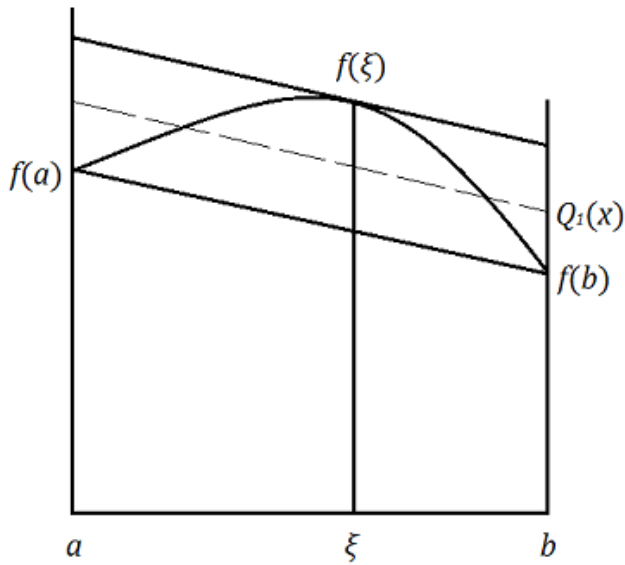
$$\begin{cases} f(a) - c_0 - c_1a = \alpha\Delta(f), \\ f(\xi) - c_0 - c_1\xi = -\alpha\Delta(f), \\ f(b) - c_0 - c_1b = \alpha\Delta(f). \end{cases} \quad (32)$$

Звідси  $f(b) - f(a) = c_1(b - a)$  та  $c_1 = \frac{f(b)-f(a)}{b-a}$ .

Цю систему треба замкнути, використавши ще одне рівняння з умови: точка  $\xi$  є точкою екстремуму різниці  $f(x) - (c_0 + c_1x)$ . Тому для диференційованої функції  $f(x)$  для визначення  $\xi$  маємо

рівняння (дотична і січна паралельні):

$$f'(\xi) = c_1 = \frac{f(b) - f(a)}{b - a} \quad (33)$$



Геометрично ця процедура виглядає наступним чином (див. [рис. 12](#)). Проводимо січну через точки  $(a, f(a))$ ,  $(b, f(b))$ . Для неї тангенс кута дорівнює  $c_1$ . Проводимо паралельну їй дотичну до кривої  $y = f(x)$ , а потім пряму, рівновіддалену від січної та дотичної, яка і буде графіком  $Q_1(x)$ . При цьому  $x_0 = a$ ,  $x_1 = \xi$ ,  $x_2 = b$ .

3. Потрібно наблизити  $f(x) = x^{n+1}$ ,  $x \in [-1, 1]$  багаточленом степеня  $n$ :  $Q_n^0(x)$ . Введемо

$$\bar{P}_{n+1}(x) = x^{n+1} - Q_n(x) = x^{n+1} - a_1 x^n - \dots \quad (34)$$

Далі

$$\Delta(f) = \inf_{Q_n(x)} \|x^{n+1} - Q_n^0(x)\|_C = \inf_{\bar{P}_{n+1}} \|\bar{P}_{n+1} - 0\|_C = \|\bar{T}_{n+1}(x)\|. \quad (35)$$

Звідси

$$x^{n+1} - Q_n^0(x) = \bar{T}_{n+1}(x), \quad (36)$$

або

$$Q_n^0(x) = x^{n+1} - \bar{T}_{n+1}(x). \quad (37)$$

**Задача 25:** Для прикладу 3 вказати точки чебишовського альтернансу  $\{x_i\}$ ,  $i = \overline{0, n+1}$ .

4. Потрібно наблизити  $f(x) = P_{n+1}(x) = a_0 + \dots + a_{n+1}x^{n+1}$ ,  $a_{n+1} \neq 0$ ,  $x \in [a, b]$  БНРН степеня  $n$ . Запишімо його у вигляді:

$$Q_n^0(x) = P_{n+1}(x) - a_{n+1}(x) \bar{T}_{n+1}^{[a,b]} \quad (38)$$

де  $\bar{T}_{n+1}^{[a,b]}(x)$  — нормований багаточлен Чебишова на проміжку  $x \in [a, b]$ .

Дійсно це БНРН: вираз у правій частині є багаточленом степеня  $n$ , оскільки коефіцієнт при  $x^{n+1}$  дорівнює нулю, а його нулі

$$x_k = \frac{b+a}{2} + \frac{b-a}{2} \cdot t_k, \quad (39)$$

де

$$t_k = \cos\left(\frac{(2k+1)\pi}{2(n+1)}\right), \quad (40)$$

для  $k = \overline{0, n}$  є точками чебишевського альтернату для  $Q_n^0(x)$ .

**Задача 26:** Показати, що для  $f(x)$  парної (непарної) функції БНРН це багаточлен по парних (непарних) степенях  $x$ .

5. Телескопічний метод. Дуже часто БНРН точно знайти не вдається. В таких випадках шукається багаточлен, близький до нього. Бажано щоб цей багаточлен був невисокого степеня (менше арифметичних операцій на його обчислення) Спочатку будують такий багаточлен

$$P_n(x) = \sum_{j=0}^n a_j x^j, \quad (41)$$

щоб відхилення від  $f(x)$  була достатньо малою. (наприклад меншою за  $\varepsilon/2$ ).

Можна це зробити, наприклад, за формулою Тейлора. Потім наближають багаточлен  $P_n(x)$  багаточленом найкращого рівномірного наближення  $P_{n-1}(x)$  (за алгоритмом п. 4; для простоти  $x \in [-1, 1]$ ):

$$P_{n-1}(x) = P_n(x) - a_n T_n(x) 2^{1-n}. \quad (42)$$

Оскільки  $|T_n(x)| \leq 1$  на відрізку  $[-1, 1]$ , то

$$|P_{n-1}(x) - P_n(x)| \leq |a_n| \cdot 2^{1-n}. \quad (43)$$

Далі наближають багаточлен  $P_{n-1}(x)$  багаточленом найкращого рівномірного наближення  $P_{n-2}(x)$  і т. д. Пониження степеня продовжується до тих пір, поки сумарна похибка від таких послідовних апроксимацій залишається меншою за задане мале число  $\varepsilon$ .

## 8.4. Найкраще середньоквадратичне наближення

Наблизимо функцію  $f(x) \in H$  з гільбертового простору  $H$  функціями із скінченно-вимірного підпростору  $M_n$  простору  $H$ . Тут  $H$  — гільбертів простір із скалярним добутком  $\langle u, v \rangle$ , норма і відстань для якого визначаються формулами:

$$|u| = \sqrt{\langle u, u \rangle}, \quad \Delta(u, v) = |u - v|. \quad (44)$$

Побудуємо

$$u = \sum_{i=0}^n c_i \varphi_i \in M_n \subset H, \quad (45)$$

де  $\{\varphi_i\}_{i=0}^\infty$  — лінійно-незалежна система елементів з  $H$ .

**Означення:** Елементом найкращого середньоквадратичного наближення (в подальшому ЕНСКН) називатимемо  $\Phi_0$  такий, що

$$|f - \Phi_0| = \sqrt{\langle f - \Phi_0, f - \Phi_0 \rangle} = \inf_{\Phi \in M_n} |f - \Phi|. \quad (46)$$

**Теорема 1:** Нехай  $f \in H$ ,  $\Phi_0 \in M_n$  — елемент найкращого середньоквадратичного наближення, тобто

$$|f - \Phi_0| = \inf_{\Phi \in M_n} |f - \Phi|, \quad (47)$$

тоді

$$\forall \Phi \in M_n : \quad \langle f - \Phi_0, \Phi \rangle = 0. \quad (48)$$

*Доведення.*

Нехай (48) не виконується, тобто  $\exists \Phi_1 \in M_n$ :

$$\langle f - \Phi_0, \phi_1 \rangle = \alpha \neq 0. \quad (49)$$

Без обмеження загальності можемо вважати, що  $\|\Phi_1\| = 1$ .

Побудуємо  $\Phi_2 = \Phi_0 + \alpha \Phi_1$ , тоді

$$|f - \Phi_2|^2 = \langle f - \Phi_2, f - \Phi_2 \rangle = |f - \Phi_0|^2 - \alpha^2 < |f - \Phi_0|^2. \quad (50)$$

Отже, елемент  $\Phi_2$  кращий за елемент найкращого середньоквадратичного наближення  $\Phi_0$ . А це суперечність.  $\square$

**Наслідок:**  $f = \Phi_0 + \nu$ , де  $\Phi_0 \in M_n$ , а  $\nu \perp M_n$  (поправка  $\nu$  — з ортогонального доповнення до  $M_n$ ).

Знайти ЕНСН

$$\Phi_0 = \sum_{i=0}^n c_i \varphi_i \quad (51)$$

означає знайти коефіцієнти  $c_i$ .

Для виконання (48) достатньо, щоб

$$\langle f - \Phi_0, \varphi_k \rangle = 0, \quad k = \overline{0, n}. \quad (52)$$

Підставимо (51) у формулу (52):

$$\langle f - \sum_{i=0}^n c_i \varphi_i, \varphi_k \rangle = 0. \quad (53)$$

Таким чином маємо СЛАР для  $c_i$ :

$$\sum_{i=0}^n c_i \langle \varphi_i, \varphi_k \rangle = \langle f, \varphi_k \rangle, \quad k = \overline{0, n}. \quad (54)$$

З теорему 1 витікає лише достатність умов (54) для знаходження коефіцієнтів  $c_i$ . Розглянемо задачу

$$|f - \Phi_0| = \inf_{\Phi \in M_n} |f - \Phi|, \quad (55)$$

як задачу мінімізації функції багатьох змінних:

$$F(a_0, \dots, a_n) = |f - \Phi|^2 = \left| f - \sum_{i=0}^n a_i \varphi_i \right|^2 \rightarrow \min. \quad (56)$$

Умови мінімуму цієї функції приводять до (54).

**Задача 27:** Показати, що для коефіцієнтів  $c_i$  елемента найкращого середньо-квадратичного наближення умови (54) є необхідними та достатніми.

Матриця СЛАР (54) складається з елементів  $g_{ik} = \langle \varphi_i, \varphi_k \rangle$ , тобто це матриця Грамма:  $G = (g_{ik})_{i,k=0}^n$ . Оскільки це матриця Грамма лінійно-незалежної системи, то  $\det G \neq 0$ , що ще раз доводить існування та єдиність ЕНСН. Оскільки  $G^T = G$ , то для розв'язку цієї системи використовують метод квадратних коренів.

Якщо взяти  $x \in [0, 1]$  та  $\varphi_i(x) = x^i, i = \overline{0, n}, H = L_2(0, 1)$ , то

$$g_{ik} = \int_0^1 x^i x^k dx = \frac{1}{i+k+1}, \quad i, k = \overline{0, n}. \quad (57)$$

Це матриця Гілберта, яка є погано обумовленою:  $\text{cond} G \approx 10^7, n = 6$ . Праві частини

$$f_k = \langle f, \varphi_k \rangle = \int_0^1 f(x) x^k dx \quad (58)$$

як правило, обчислюються наближено, тому похибки обчислення  $c_i$  можуть бути великими.

Що робити? Якщо вибирати систему  $\{\varphi_i\}_{i=0}^\infty$  ортонормованою, тобто

$$\langle \varphi_i, \varphi_k \rangle = \delta_{ik} = \begin{cases} 1, & i = k, \\ 0, & i \neq k, \end{cases} \quad (59)$$

то система (54) має явний розв'язок

$$\Phi_0 = \sum_{i=0}^n \langle f, \varphi_i \rangle \cdot \varphi_i. \quad (60)$$

Якщо  $\{\varphi_i\}$  — повна ортонормована система, то довільну функцію можна представити у вигляді ряду Фур'є:

$$f = \sum_{i=0}^{\infty} \langle f, \varphi_i \rangle \cdot \varphi_i, \quad (61)$$

i

$$f - \Phi_0 = \sum_{i=n+1}^{\infty} c_i \varphi_i = \nu \quad (62)$$

— залишок (похибка). Таким чином ЕНСН є відрізком ряду Фур'є. Далі



$$\begin{aligned}
|f - \Phi_0|^2 &= \langle f - \Phi_0, f - \Phi_0 \rangle = |f|^2 - 2\langle f, \Phi_0 \rangle + |\Phi_0|^2 = \\
&= |f|^2 - 2|\Phi_0|^2 - 2\langle \nu, \Phi_0 \rangle + |\Phi_0|^2 = |f|^2 - |\Phi_0|^2 = \\
&= \sum_{i=0}^{\infty} c_i^2 - \sum_{i=0}^n c_i^2 = \sum_{i=n+1}^{\infty} c_i^2 \xrightarrow{n \rightarrow \infty} 0.
\end{aligned} \tag{63}$$

Останнє випливає з відповідної теореми математичного аналізу. Таким чином, якщо  $\{\varphi_i\}_{i=0}^{\infty}$  — повна ортонормована система, то

$$\sum_{i=n+1}^{\infty} c_i^2 \xrightarrow{n \rightarrow \infty} 0, \tag{64}$$

$$\Phi_0^{(n)} \xrightarrow{n \rightarrow \infty} f, \tag{65}$$

Значить вірна

**Теорема 2:** В гільбертовому просторі  $H$  послідовність ЕНСН  $\{\Phi_0^{(n)}\}$  по повній ортонормованій системі  $\{\varphi_i\}_{i=0}^{\infty}$  збігається до  $f$ .

**Зауваження 1:** Відхилення можна обчислити за формулою:

$$\Delta^2(f) = |f - \Phi_0|^2 = |f|^2 - 2\langle f, \Phi_0 \rangle + |\Phi_0|^2 = |f|^2 - |\Phi_0|^2 = |f|^2 - \sum_{i=0}^n c_i^2. \tag{66}$$

Якщо  $\{\varphi_i\}_{i=0}^{\infty}$  — ортогональна система, але не нормована, тобто  $\langle \varphi_i, \varphi_k \rangle = \delta_{ik} \|\varphi_i\|^2$ , то

$$c_i = \frac{\langle f, \varphi_i \rangle}{\|\varphi_i\|^2}, \tag{67}$$

$$\Phi_0 = \sum_{i=0}^n \frac{\langle f, \varphi_i \rangle}{\|\varphi_i\|^2} \cdot \varphi_i, \tag{68}$$

$$|f - \Phi_0|^2 = |f|^2 - \sum_{i=0}^n \frac{c_i^2}{\|\varphi_i\|^2}. \tag{69}$$

Для функції  $f(x)$ , щоб побудувати ЕНСН покладемо  $H = L_{2,\alpha}([a, b])$ , в якому скалярний добуток виберемо наступним чином

$$\langle u, v \rangle = \int_a^b u(x)v(x) d\alpha(x), \tag{70}$$

де  $\alpha(x)$  — зростаюча функція. Можливі випадки:

1.  $\alpha(x) \in C^1([a, b])$ , тоді  $d\alpha(x) = \rho(x) dx > 0$  та

$$\langle u, v \rangle = \int_a^b \rho(x)u(x)v(x) dx. \tag{71}$$

2.  $\alpha(x)$  — функція стрибків,  $\alpha(x) = \alpha(x_k - 0)$ , де  $x_{k-1} \leq x \leq x_k$ ,  $k = \overline{1, n}$ . Якщо ввести  $\rho_k = \alpha(x_k + 0) - \alpha(x_k - 0)$ , то

$$\langle u, v \rangle = \sum_{k=1}^n \rho_k u(x_k) v(x_k). \quad (72)$$

Перший вибір  $\alpha(x)$  використовується при апроксимації функцій неперервного аргументу, а другий — для табличних функцій.

## 8.5. Системи ортогональних функцій

Як вибрати ортонормальну або ортогональну систему функцій  $\{\varphi_i\}_{i=0}^{\infty}$ ?

Розглянемо деякі з найбільш вживаних таких систем.

1. Якщо  $H = L_2([-1, 1])$ ;  $\rho \equiv 1$  (ваговий множник), то  $\varphi_i(x) = L_i(x)$  — система *багаточленів Лежандра*, які мають вигляд

$$L_n(x) = \frac{1}{2^n n!} \cdot \frac{d^n}{dx^n} \cdot (x^2 - 1)^n. \quad (73)$$

Використовують також рекурентні формули

$$(n+1)L_{n+1}(x) = (2n+1)xL_n(x) - nL_{n-1}(x), \quad (74)$$

до яких додаємо умови

$$L_0(x) = 1, \quad L_1(x) = x. \quad (75)$$

Це ортогональна система в тому сенсі, що

$$\langle L_i, L_k \rangle = \int_{-1}^1 L_i(x) L_k(x) dx = \delta_{ik} |L_i(x)|^2, \quad (76)$$

де  $\|L_i(x)\|^2 = \frac{2}{2i+1}$  і тому  $c_i = \frac{\langle f, L_i \rangle}{\|L_i\|^2} = \frac{2i+1}{2} \langle f, L_i \rangle$ .

**Зауваження:** Якщо потрібно побудувати наближення на довільному проміжку  $(a, b)$ , то бажано перейти до проміжку  $(-1, 1)$ , тобто по  $f(x)$  на  $[a, b]$  побудувати  $f(t)$  з  $t \in [-1, 1]$  заміною  $x = At + B$ ,  $t = \alpha x + \beta$  та для побудови багаточлена НСКН для  $f(t)$  використати багаточлени Лежандра  $L_i(t)$ .

Можна робити навпаки — систему багаточленів перевести з  $[a, b]$  на  $[-1, 1]$ , але це вимагає більше обчислень і процес побудови ЕНСН складніше.

2. Якщо  $H = L_{2,\rho}([-1, 1])$ ,  $\rho(x) = 1/\sqrt{1-x^2}$ , скалярний добуток

$$\langle u, v \rangle = \int_{-1}^1 \frac{u(x)v(x)}{\sqrt{1-x^2}} dx \quad (77)$$

(це невластні інтеграли другого роду), то  $\varphi_i(x) = T_i(x)$ , де  $\{T_i(x)\}$  — система ортогональних *багаточленів Чебишова* 1-го роду, які мають вигляд

$$T_n(x) = \cos(n \arccos(x)). \quad (78)$$

Рекурентна формула для цих багаточленів:

$$T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x), \quad (79)$$

до якої додаємо умови  $T_0 = 1, T_1 = x$ .

Для цієї системи

$$|T_n|^2 = \begin{cases} \pi, & n = 0, \\ \pi/2, & n = 1, 2, \dots \end{cases} \quad (80)$$

3.  $H$  гільбертів простір з ваговим множником  $\rho(x) = (1-x)^\alpha(1+x)^\beta$ . Система  $\varphi_i(x) = P_n^{(\alpha,\beta)}(x)$  — багаточленів Якобі,  $\alpha, \beta > -1$  ( $\alpha, \beta$  — числові параметри) ортогональна в сенсі скалярного добутку

$$\langle u, v \rangle = \int_{-1}^1 (1-x)^\alpha(1+x)^\beta u(x)v(x) dx. \quad (81)$$

Ця система є узагальненням випадків 1. та 2.

Диференціальна формула для багаточленів:

$$P_n^{(\alpha,\beta)}(x) = \frac{(-1)^n}{2^n n!} (1-x)^{-\alpha}(1+x)^{-\beta} \frac{d^n}{dx^n} \left( (1-x)^{n+\alpha}(1+x)^{n+\beta} \right). \quad (82)$$

Рекурентна формула:

$$2(n+1)(n+\alpha+\beta+1)(2n+\alpha+\beta)P_{n+1}^{(\alpha,\beta)}(x) = \quad (83)$$

$$= (2n+\alpha+\beta+1)((2b+\alpha+\beta)(2n+\alpha+\beta+2)x + \alpha^2 - \beta^2)P_n^{(\alpha,\beta)}(x) - \quad (84)$$

$$- 2(n+\alpha)(n+\beta)(2n+\alpha+\beta+2)P_{n-1}^{(\alpha,\beta)}(x), \quad (85)$$

де  $P_0^{(\alpha,\beta)} = 1, P_{-1}^{(\alpha,\beta)} = 0$ , і

$$|P_n^{(\alpha,\beta)}|^2 = \frac{2^{\alpha+\beta+1}\Gamma(\alpha+n+1)\Gamma(\beta+n+1)}{n!(\alpha+\beta+2n+1)\Gamma(\alpha+\beta+n+1)}, \quad (86)$$

та

$$\Gamma(z) = \int_0^\infty e^{-t} t^{z-1} dt, \quad (87)$$

а  $\Gamma(z+1) = z \cdot \Gamma(z), \Gamma(n+1) = n!$ .

Коли  $\alpha = \beta = 0$ :  $P_n^{(0,0)}(x) = L_n(x)$ , а для  $\alpha = \beta = -1/2$ :  $P_n^{(-1/2,-1/2)}(x) = T_n(x)$ .

4.  $H = L_{2,\rho}([0, \infty)), \rho(x) = x^\alpha e^{-x}, \alpha > -1$ .

Цьому ваговому множнику відповідає система багаточленів Лагерра  $\varphi_i(x) = L_i^\alpha(x)$ , які задаються диференціальною формулою:

$$L_n^\alpha(x) = (-1)^n x^{-\alpha} e^x \frac{d^n}{dx^n} (x^{\alpha+n} e^{-x}) \quad (88)$$

або в рекурентній формі

$$(n+1)L_{n+1}^\alpha = (2n+\alpha+1-x)L_n^\alpha - (n+\alpha)L_{n-1}^\alpha, \quad (89)$$

де  $L_0^\alpha = 1$ ,  $L_{-1}^\alpha = 0$  та з нормою  $\|L_n^\alpha\|^2 = n! \cdot \Gamma(\alpha + n + 1)$ .

5.  $H = L_{2,\alpha}((-\infty, \infty))$ ,  $\rho(x) = e^{-x^2}$ . Систему ортогональних функцій вибираємо як систему багаточленів Ерміта  $\varphi_i(x) = H_i(x)$ , які задаються диференціальною формулою:

$$H_n(x) = (-1)^n e^{x^2} \frac{d^n}{dx^n} e^{-x^2}, \quad (90)$$

або в рекурентній формі

$$H_{n+1}(x) = 2xH_n(x) - 2nH_{n-1}(x) \quad (91)$$

де  $H_0 = 1$ ,  $H_{-1} = 0$  та  $\|H_n\|^2 = 2^n n! \sqrt{\pi}$ .

6.  $H = L_2([0, 2\pi])$ ,  $\rho(x) \equiv 1$ ,  $f(x) = f(x + 2\pi)$ .  $f(x)$  —  $2\pi$ -періодичні функції. За систему ортонормованих функцій вибираємо *тригонометричну систему*

$$\varphi_0(x) = \frac{1}{\sqrt{2\pi}}, \quad (92)$$

$$\varphi_{2k-1}(x) = \frac{\cos(kx)}{\sqrt{\pi}}, \quad (93)$$

$$\varphi_{2k}(x) = \frac{\sin(kx)}{\sqrt{\pi}}. \quad (94)$$

Елемент найкращого середньоквадратичного наближення представляє собою тригонометричний багаточлен

$$\Phi_0(x) \equiv T_n(x) = \frac{a_0}{2} + \sum_{k=1}^n (a_k \cos(kx) + b_k \sin(kx)), \quad (95)$$

формули для обчислення цих коефіцієнтів наведені в наступному пункті.

7. Якщо потрібно апроксимувати табличну функцію, то  $H = \ell_2$ ,  $x_i = i$ ,  $i = \overline{0, N}$ ,

$$\langle u, v \rangle = \frac{1}{N+1} \sum_{i=0}^N u_i v_i, \quad (96)$$

і за систему ортогональних функцій вибираємо наступну систему багаточленів:  $\varphi_k(x) = p_k^{(N)}(x)$ ,  $k = \overline{0, m}$  ( $m \leq N$ ) — систему багаточленів Чебишова дискретного аргументу, які задається формулою

$$p_k^{(N)}(x) = \sum_{j=0}^k \frac{(-1)^j C_k^j C_{k+j}^j}{N^{(j)}} \cdot x^{(j)} \quad (97)$$

де  $x^{(j)} = x(x-1) \dots (x-j+1)$  — факторіальний багаточлен;  $C_k^j$  — число сполук.

Рекурентна формула:

$$\frac{(m+1)(N-m)}{2(2m+1)} \cdot p_{m+1}^{(N)} = \left( \frac{N}{2} - x \right) p_m^{(N)} - \frac{m(N+m+1)}{2(2m+1)} \cdot p_{m-1}^{(N)}, \quad (98)$$

з початковими значеннями  $p_0^{(N)} = 1$ ,  $p_{-1}^{(N)} = 0$ .

Наприклад  $p_1^{(N)} = 1 - \frac{2x}{N}$ ,  $p_2^{(N)} = 1 - \frac{6x}{N} + \frac{6x^2}{N(N-1)}$ .

У випадку, якщо задані вузли  $t_i = t_0 + ih$ ,  $i = \overline{0, N}$ , то робимо заміну  $x_i = \frac{t_i - t_0}{h} = i$ .

## 8.6. Середньоквадратичне наближення періодичних функцій

Нехай маємо періодичну функцію  $f(x)$  неперервного аргументу, з періодом  $T = 2\pi$ , тобто  $f(x + 2\pi) = f(x)$ . В просторі  $H_2 = L_2([0, 2\pi])$  визначений скалярний добуток:

$$\langle u, v \rangle = \int_0^{2\pi} u(x)v(x) dx \quad (99)$$

В якості системи лінійно-незалежних функцій  $\{\varphi_i\}$  виберемо тригонометричну систему функцій:

$$\varphi_0(x) = 1; \quad \varphi_{2k-1}(x) = \cos(kx); \quad \varphi_{2k}(x) = \sin(kx), \quad (100)$$

для  $k = 1, 2, \dots$ , яка є повною нормованою системою в  $L_2([0, 2\pi])$ .

Будемо шукати  $\Phi(x)$  у вигляді тригонометричного багаточлена

$$\Phi_0(x) \equiv T_n(x) = \frac{a_0}{2} + \sum_{k=1}^n (a_k \cos(kx) + b_k \sin(kx)), \quad (101)$$

За теорією найкращого середньоквадратичного наближення коефіцієнти обчислюємо за формулами:

$$\begin{cases} a_0 = \langle f, \varphi_0 \rangle = \frac{1}{2\pi} \int_0^{2\pi} f(x) dx, \\ a_k = \langle f, \varphi_{2k-1} \rangle = \frac{1}{\pi} \int_0^{2\pi} f(x) \cos(kx) dx, \\ b_k = \langle f, \varphi_{2k} \rangle = \frac{1}{\pi} \int_0^{2\pi} f(x) \sin(kx) dx. \end{cases} \quad (102)$$

Відхилення:

$$\Delta^2(f) = |f|^2 - \left( 2\pi a_0^2 + \sum_{k=1}^n \pi(a_k^2 + b_k^2) \right). \quad (103)$$

Тепер нехай функція  $f(x)$  задана таблично:  $f_i = f(x_i)$ ,  $i = \overline{1, N}$ . Тригонометрична система  $\varphi_0(x)$ ,  $\varphi_{2k-1}(x)$ ,  $\varphi_{2k}(x)$  — ортогональна в  $H = L_2(\omega)$  для  $\omega = \{x_i = \pi i/N, i = \overline{1, n}\}$  в сенсі скалярного добутку

$$\langle u, v \rangle = \frac{1}{N} \sum_{i=1}^N u_i v_i, \quad u_i = u(x_i). \quad (104)$$

Тоді

$$\begin{cases} a_0 = \frac{1}{N} \sum_{i=1}^N f_i, \\ a_k = \frac{2}{N} \sum_{i=1}^N f_i \cos(kx_i), \\ b_k = \frac{2}{N} \sum_{i=1}^N f_i \sin(kx_i), \end{cases} \quad (105)$$

Це *формули Бесселя*. В формулі (101):  $\Phi(x) \equiv T_n(x)$  (тобто багаточлен той же), але коефіцієнти визначаємо за формулою (105).

**Зауваження:** Як правило кількість даних значень  $N \gg 2n + 1$ . Але якщо  $N = 2n + 1$ , то  $n = \frac{N-1}{2}$  і  $N$  — непарне. При цьому  $T_{\frac{N-1}{2}}(x)$  — БНСКН і звідси

$$\Delta^2(f) = \left| f(x) - T_{\frac{N-1}{2}}(x) \right|^2 = \frac{1}{N} \sum_{i=1}^N \left( f(x_i) - T_{\frac{N-1}{2}}(x) \right)^2 \rightarrow \inf_{a_k, b_k}. \quad (106)$$

Оскільки найменше значення відхилення  $\Delta^2(f) = 0$ , то тригонометричний багаточлен найкращого середньоквадратичного наближення співпадає з інтерполяційним тригонометричним багаточленом і

$$T_{\frac{N-1}{2}}(x) = f(x_i). \quad (107)$$

Для визначення коефіцієнтів  $a_i, b_i$  за формулою Бесселя (105) необхідна кількість операцій  $Q = O(N^2)$ . Існують алгоритми, які дозволяють обчислити за  $Q = O(N \cdot \log(N))$  операцій. Це так званий алгоритм швидкого перетворення Фур'є. Якщо в (105) існує група доданків, які рівні між собою, тобто число  $N$  можна представити як  $N = p_1 p_2$ , то можна так вибрати сітку, що  $Q = O(N \cdot \max(p_1, p_2))$ . Якщо ж  $N = n^m$ , то  $Q = O(Nm) = O(N \log_2(N))$ .

## 8.7. Метод найменших квадратів (МНК)

Нехай в результаті вимірювань функції  $f(x)$  маємо таблицю значень:

$$y_i \approx f(x_i), \quad i = \overline{1, N}, \quad x_i \in [a, b]. \quad (108)$$

За даними цієї таблиці треба побудувати аналітичну формулу  $\Phi(x; a_0, a_1, \dots, a_n)$  таку, що

$$\Phi(x_i; a_0, a_1, \dots, a_n) \approx y_i, \quad i = \overline{1, N}. \quad (109)$$

Виконувати це інтерполюванням тобто задавати

$$\Phi(x_i; a_0, a_1, \dots, a_n) = y_i, \quad i = \overline{1, N} \quad (110)$$

нераціонально, бо  $N \gg n$  і система перевизначена; її розв'язки як правило не існують. Вигляд функції  $\Phi(x; a_0, a_1, \dots, a_n)$  і число параметрів  $a_i$  у деяких випадках відомі. В інших випадках вони визначаються за графіком, побудованим за відомими значеннями  $f(x_i)$  так, щоб залежність (109) була досить простою і добре відображала результати спостережень. Але такі міркування не дають змогу побудувати єдиний елемент та й ще найкращого наближення.

Тому визначають параметри  $a_0, \dots, a_n$  так, щоб у деякому розумінні всі рівняння системи (109) одночасно задовольнялись з найменшою похибкою, наприклад, щоб виконувалося:

$$I(a_0, \dots, a_n) = \sum_{i=1}^N (y_i - \Phi(x_i; a_0, \dots, a_n))^2 \rightarrow \min. \quad (111)$$

Такий метод розв'язання системи (109) і називають методом найменших квадратів, оскільки мінімізується сума квадратів відхилення  $\Phi(x; a_0, a_1, \dots, a_n)$  від значень  $f(x_i)$ .

Для реалізації мінімуму необхідно та достатньо виконання умов:

$$\frac{\partial I}{\partial a_i} = 0, \quad i = \overline{0, n}, \quad (112)$$

Якщо  $\Phi(x_i; a_0, \dots, a_n)$  лінійно залежить від параметрів  $a_0, \dots, a_n$ , тобто

$$\Phi(x; a_0, \dots, a_n) = \sum_{k=0}^n a_k \varphi_k(x), \quad (113)$$

то з (110) маємо СЛАР:

$$\sum_{k=0}^n a_k \varphi_k(x_i) = y_i, \quad i = \overline{1, N}, \quad (114)$$

яку називають *системою умовних рівнянь*. Позначивши

$$C = (\varphi_k(x_i))_{\substack{i=\overline{1, N} \\ j=\overline{0, n}}}, \quad (115)$$

$$\vec{a} = (a_0, \dots, a_n)^T, \quad (116)$$

$$\vec{y} = (y_1, \dots, y_N)^T, \quad (117)$$

маємо матричний запис СЛАР (114):

$$C\vec{a} = \vec{y}. \quad (118)$$

Помноживши систему умовних рівнянь (118) зліва на транспоновану до  $C$  матрицю  $C^T$  отримаємо систему нормальних рівнянь

$$C^T C \vec{a} = C^T \vec{y}, \quad (119)$$

де  $G = A = C^T C$ ,  $\dim G = n + 1$ ,  $G = (g_{ik})_{i,k=0}^n$ , а самі

$$g_{ik} = \sum_{j=1}^N c_{ij}^T c_{jk} = \sum_{j=1}^N c_{ji} c_{jk} = \sum_{j=1}^N \varphi_k(x_j) \varphi_i(x_j), \quad (120)$$

а

$$C^T \vec{y} = \left( \sum_{i=1}^N c_{ik} y_i \right)_{k=0}^n \quad (121)$$

з якої власно і обчислюють невідомі коефіцієнти.

Покажемо, що МНК є методом знаходження ЕНСКН, якщо визначити скалярний добуток

$$\langle u, v \rangle = \sum_{i=1}^N u(x_i) v(x_i). \quad (122)$$

Поставимо задачу знаходження ЕНСН:

$$\Delta(f, \Phi) = |f - \Phi|^2 = \langle f - \Phi, f - \Phi \rangle = \sum_{i=1}^N (y_i - \Phi(x_i, \vec{a}))^2 \rightarrow \inf. \quad (123)$$

За теорією середньоквадратичного наближення для цього необхідно, щоб коефіцієнти  $a_0, \dots, a_n$  знаходилися з системи:

$$\sum_{j=0}^n a_k \langle \varphi_k, \varphi_j \rangle = \langle \varphi_k, f \rangle, \quad (124)$$

де  $k = \overline{0, n}$ , а це співпадає з (119).

Якщо відома інформація про обчислювальну похибку для значень  $f(x_i)$ :  $|f(x_i) - y_i| < \varepsilon_i$ , то вибирають такий скалярний добуток

$$\langle u, v \rangle = \sum_{i=1}^N \rho_i u(x_i) v(x_i), \quad (125)$$

де  $\rho_i = 1/\varepsilon_i^2$ .

Нехай тепер  $\Phi(x; a_0, \dots, a_n)$  — нелінійна функція параметрів  $a_0, \dots, a_n$ , наприклад:

$$\Phi = a_0 e^{a_1 x} + a_2 e^{a_3 x} + \dots, \quad (126)$$

або

$$\Phi = a_0 \cos(a_1 x) + a_2 \sin(a_3 x) + \dots \quad (127)$$

Складемо функціонал:

$$S(a_0, \dots, a_n) = \sum_{i=1}^N \rho_i (y_i - \Phi(x, \vec{a}))^2 \rightarrow \int_a. \quad (128)$$

Оскільки тепер  $\Phi(x; a_0, \dots, a_n)$  нелінійна, то застосуємо метод лінеаризації.

Нехай відомі наближені значення  $\vec{a}^{(0)} = (a_0^{(0)}, a_1^{(0)}, \dots, a_n^{(0)})$ . Розкладемо  $\Phi(x, \vec{a})$  в околі  $a^{(0)}$ . Тоді отримаємо лінійне наближення до  $\Phi(x, \vec{a})$ :

$$\Phi(x, \vec{a}) \approx \Phi(x, \vec{a}^{(0)}) + \sum_{k=0}^n \frac{\partial \Phi}{\partial a_k} \langle x, \vec{a}^{(0)} \rangle (a_k - a_k^{(0)}). \quad (129)$$

Якщо ввести позначення

$$\vec{x} = \vec{a} - \vec{a}^{(0)}, \quad (130)$$

$$y_i^* = y_i - \Phi(x, \vec{a}^{(0)}), \quad (131)$$

$$c_{i,k} = \Phi'_{a_k}(x_i, \vec{a}^{(0)}), \quad (132)$$

то отримаємо систему умовних рівнянь відносно поправок до  $\vec{a}^{(0)}$ :

$$C\vec{z} = \vec{y}^*. \quad (133)$$



Замінімо її на систему нормальних рівнянь

$$C^T C \vec{z} = c^T \vec{y}^*. \quad (134)$$

Знайшовши  $\vec{z}$ , обчислюємо наступне наближення:  $\vec{a}^{(1)} = \vec{a}^{(0)} + \vec{z}$ . Цей процес можна продовжувати: на кожній ітерації знаходимо  $\vec{z}^{(m)}$ ,  $m = 0, 1, \dots$  і уточнюємо наближення до  $\vec{a}$ :  $\vec{a}^{(m)} = \vec{a}^{(m-1)} + \vec{z}^{(m-1)}$ .

Умова припинення ітерацій

$$|\vec{z}^{(m)}| = \sqrt{\sum_{k=0}^n \left(z_k^{(m)}\right)^2} < \varepsilon. \quad (135)$$

Важливим є вибір початкового наближення  $\vec{a}^{(0)}$ . З системи умовних рівнянь (нелінійної) виберемо деякі  $n + 1$ . Розв'язок цієї системи і дасть початкове наближення.

Для деяких простих нелінійних залежностей від невеликої кількості параметрів задачу можна ліанеризувати аналітично. Наприклад, розглянемо наближення даних алометричним законом

$$y_i \approx f(x_i), \quad \Phi(x, A, \alpha) = Ax^\alpha. \quad (136)$$

Система умовних рівнянь має вигляд:

$$\Phi(x_i) = Ax_i^\alpha = y_i, \quad i = \overline{1, N}. \quad (137)$$

Прологарифмуємо її:

$$\psi(x_i) = \ln(\Phi(x_i)) = \ln(A) + \alpha \ln(x_i) = \ln(y_i), \quad i = \overline{1, N}. \quad (138)$$

Введемо  $a = \ln(A)$ . Тепер функція  $\psi(x, a, \alpha)$  лінійна. Система умовних рівнянь відносно параметрів  $a$  та  $\alpha$  має вигляд:

$$C\vec{z} = \vec{b}, \quad (139)$$

де  $\vec{z} = (a, \alpha)$ ,  $\vec{b} = (\ln(y_i))_{i=1}^N$ , а

$$C = \begin{pmatrix} 1 & \ln(x_1) \\ \vdots & \vdots \\ 1 & \ln(x_N) \end{pmatrix} \quad (140)$$

Запишемо систему нормальних рівнянь для методу найменших квадратів:

$$G = C^T C = \begin{pmatrix} N & \sum_{i=1}^N \ln(x_i) \\ \sum_{i=1}^N \ln(x_i) & \sum_{i=1}^N (\ln(x_i))^2 \end{pmatrix}, \quad (141)$$

$$C^T \vec{b} = \begin{pmatrix} \sum_{i=1}^N \ln(y_i) \\ \sum_{i=1}^N \ln(x_i) \ln(y_i) \end{pmatrix}. \quad (142)$$

Розв'язавши систему (141)–(142), знаходимо  $\alpha$ , та  $A = \exp(\alpha)$ .

## 8.8. Згладжуючі сплайни

Якщо значення в точках  $x_i$  неточно  $\tilde{f}_i = f_i + \varepsilon_i$ , то застосовують згладжування. Для цього треба побудувати нову таблицю із згладженими значеннями  $\bar{f}_i$ .

Наведемо деякі прості формули згладжування:

1.  $m = 1$ :

- $\bar{f}_i = \frac{1}{3} (\tilde{f}_{i-1} + \tilde{f}_i + \tilde{f}_{i+1}), N = 3;$
- $\bar{f}_i = \frac{1}{5} (\tilde{f}_{i-2} + \dots + \tilde{f}_{i+2}), N = 5;$
- $\bar{f}_i = \frac{1}{2k} (\tilde{f}_{i-k} + \dots + \tilde{f}_{i+k}).$

2.  $m = 3$ :

- $\bar{f}_i = \frac{1}{3 \cdot 5} (-3\tilde{f}_{i-2} + 12\tilde{f}_{i-1} + 17\tilde{f}_i + 12\tilde{f}_{i+1} - 3\tilde{f}_{i+2}), N = 5.$

Їх отримуємо в такий спосіб: до  $\tilde{f}_i$  застосовуємо апроксимацію, будуємо багаточлен НСКН

$$Q_m(x) = \sum_{k=0}^m c_k p_k^N(x), \quad (143)$$

де  $p_k^N$  — система багаточленів Чебишова дискретного аргументу. Беремо значення  $\bar{f}_i = Q_m(x_i)$ , які приводять до наведених вище формул.

Але ці формули не дають гарантію, що в результаті ми отримаємо функцію, яка задовольняє умові  $|\tilde{f}_i - f_i| < \varepsilon_i$ .

Згладжуючі сплайни дають можливість побудувати наближення з заданою точністю. Нагадаємо деякі відомості про сплайни. Явний вигляд кубічного сплайна:

$$\begin{aligned} s(x) = & m_{i-1} \cdot \frac{(x_i - x)^3}{6h_i} + m_i \cdot \frac{(x - x_{i-1})^3}{6h_i} + \\ & + \left( f_{i-1} - \frac{m_{i-1}h_i^2}{6} \right) \cdot \frac{x_i - x}{h_i} + \\ & + \left( f_i - \frac{m_i h_i^2}{6} \right) \cdot \frac{x - x_{i-1}}{h_i}, \end{aligned} \quad (144)$$

для  $x \in [x_{i-1}, x_i]$ , де  $h_i = x_i - x_{i-1}$ .

Тут  $s(x_i) = f_i, i = \overline{0, n}$ , а  $m_i = s''(x_i)$  задовольняють систему:

$$\begin{cases} \frac{h_i m_{i-1}}{6} + \frac{(h_i + h_{i+1}) m_i}{3} + \frac{h_{i+1} m_{i+1}}{6} = \\ = \frac{f_{i+1} - f_i}{h_{i+1}} - \frac{f_i - f_{i-1}}{h_i}, \quad i = \overline{1, n-1} \\ m_0 = m_n = 0. \end{cases} \quad (145)$$

В матричній формі ця система має вигляд

$$A \vec{m} = H \vec{f}. \quad (146)$$

Тут

$$\vec{m} = (m_1, \dots, m_{n-1})^\top, \quad \vec{f} = (f_0, \dots, f_n)^\top, \quad (147)$$

а матриці

$$A = \begin{pmatrix} \frac{h_1+h_2}{3} & \frac{h_2}{6} & 0 & \dots & 0 & 0 \\ \frac{h_2}{6} & \frac{h_2+h_3}{3} & \frac{h_3}{6} & \dots & 0 & 0 \\ \ddots & \ddots & \ddots & \ddots & \vdots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & 0 & \frac{h_{n-2}}{6} & \frac{h_{n-2}+h_{n-1}}{3} & \frac{h_{n-1}}{6} \\ 0 & 0 & \dots & 0 & \frac{h_{n-1}}{6} & \frac{h_{n-1}+h_n}{3} \\ 0 & 0 & 0 & \dots & 0 & 1 \end{pmatrix} \quad (148)$$

і

$$H = \begin{pmatrix} \frac{1}{h_1} & -\left(\frac{1}{h_1} + \frac{1}{h_2}\right) & \frac{1}{h_2} & 0 & \dots & 0 & 0 \\ 0 & \frac{1}{h_2} & -\left(\frac{1}{h_2} + \frac{1}{h_3}\right) & \frac{1}{h_3} & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \ddots & \ddots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & \frac{1}{h_{n-1}} & -\left(\frac{1}{h_{n-1}} + \frac{1}{h_n}\right) & \frac{1}{h_n} \\ 0 & 0 & 0 & 0 & \dots & 0 & 0 \end{pmatrix} \quad (149)$$

Кубічний інтерполяційний сплайн мінімізує функціонал:

$$\Phi(u) = \int_a^n (u''(x))^2 dx : \quad (150)$$

$$\Phi(s) = \int_{u \in U} \Phi(u), \quad (151)$$

де

$$U = \left\{ u(x) : u(x_i) = f_i, i = \overline{0, n}, u(x) \in W_2^2([a, b]) \right\}. \quad (152)$$

Введемо функціонал

$$\Phi_1(u) = \Phi(u) + \sum_{i=0}^n \rho_i \left( \tilde{f}_i - u(x_i) \right)^2. \quad (153)$$

**Означення:** Згладжуючим сплайном назвемо функцію  $g$ , яка є розв'язком задачі:

$$\Phi_1(g) = \inf_{u \in W_2^2([a, b])} \Phi_1(u). \quad (154)$$

Перший доданок в  $\Phi_1(u)$  дає мінімум «згину», другий — середньоквадратичне наближення до значень  $\tilde{f}_i$ . Покажемо, що  $g$  є сплайном.

Нехай існує функція  $g(x)$ . Побудуємо кубічний сплайн такий, що  $s(x_i) = g(x_i)$ . З того, що  $g(x)$  є розв'язком задачі (154), маємо  $\Phi_1(s) \geq \Phi_1(g)$ , а тоді

$$\int_a^b (s''(x))^2 dx + \sum_{i=0}^n \rho_i \left( \tilde{f}_i - s(x_i) \right)^2 \geq \int_a^b (g''(x))^2 dx + \sum_{i=0}^n \rho_i \left( \tilde{f}_i - g(x_i) \right)^2 \quad (155)$$

Звідси  $\Phi(s) \geq \Phi(g)$ .

Оскільки кубічний інтерполяційний сплайн  $s(x)$  мінімізує функціонал (150), то  $\Phi(s) \leq \Phi(g)$ . Тому  $\Phi(s) = \Phi(g)$ . Звідки  $s = g$ .

Позначимо

$$\mu_i = g(x_i), \quad i = \overline{0, n}, \quad (156)$$

Якби значення  $\mu_i$  були б відомі, то для побудови  $g$  достатньо було б розв'язати систему

$$A\vec{m} = H\vec{\mu} \quad (157)$$

Підставимо (144) та (156) в  $\Phi_1(g)$ :

$$\Phi_1(g) = \sum_{i=1}^n \int_{x_{i-1}}^{x_i} \left( m_{i-1} \cdot \frac{x_i - x}{h_i} + m_i \cdot \frac{x - x_{i-1}}{h_i} \right)^2 dx + \sum_{i=0}^n \rho_i \left( \tilde{f}_i - \mu_i \right)^2 = \inf \Phi_1(u). \quad (158)$$

Після перетворень маємо:

$$\begin{aligned} \Phi_1(g) &= \sum_{i=1}^n \int_{x_{i-1}}^{x_i} \left( m_{i-1} \cdot \frac{x_i - x}{h_i} + m_i \cdot \frac{x - x_{i-1}}{h_i} \right)^2 dx + \sum_{i=0}^n \rho_i \left( \tilde{f}_i - \mu_i \right)^2 = \\ &= \sum_{i=1}^n m_i \left( \frac{h_i m_{i-1}}{6} + \frac{(h_i + h_{i+1}) m_i}{3} + \frac{h_{i+1} m_{i+1}}{6} \right) dx + \sum_{i=0}^n \rho_i \left( \tilde{f}_i - \mu_i \right)^2 = \\ &= \langle A\vec{m}, \vec{m} \rangle + \sum_{i=0}^n \rho_i \left( \tilde{f}_i - \mu_i \right)^2. \end{aligned} \quad (159)$$

**Задача 28:** Показати, що для кубічного згладжуючого сплайну  $g$  мають місце формули вище.

Оскільки  $\Phi_1(g)$  представляє собою квадратичну функція відносно  $\vec{m} = (m_0, \dots, m_N)$ , то необхідною і достатньою умовою мінімуму є

$$\frac{\partial \Phi_1}{\partial \mu_j} = 0, \quad j = \overline{0, n}. \quad (160)$$

Знаходимо:

$$\begin{aligned} \frac{\partial \Phi_1}{\partial \mu_j} &= \frac{\partial}{\partial \mu_j} \langle A\vec{m}, \vec{m} \rangle + 2\rho_j (\mu_j - \tilde{f}_j) = \\ &= 2 \left\langle \frac{\partial}{\partial \mu_j} (A\vec{m}), \vec{m} \right\rangle + 2\rho_j (\mu_j - \tilde{f}_j) = \\ &= 2 \left\langle \frac{\partial}{\partial m_j} (H\vec{\mu}), \vec{m} \right\rangle + 2\rho_j (\mu_j - \tilde{f}_j) = \\ &= 2 \left\langle \frac{\partial \vec{m}}{\partial m_j}, H^\top \vec{m} \right\rangle + 2\rho_j (\mu_j - \tilde{f}_j) = \\ &= 2 (H^\top \vec{m})_j + 2\rho_j (\mu_j - \tilde{f}_j) = 0. \end{aligned} \quad (161)$$

Отже, з умови мінімізації функціоналу

$$\Phi_1(u) + \int_a^b (u''(x))^2 dx + \sum_{i=0}^n \rho_i (\tilde{f}_i - u(x_i))^2. \quad (162)$$

ми отримали таку систему рівнянь :

$$2 (H^\top \vec{m})_i + 2\rho_i (\mu_i - \tilde{f}_i) = 0, \quad (163)$$

де, як і раніше і  $\mu_i$  — це невідомі значення згладжуючого сплайну:

$$\mu_i = s(x_i), \quad m_i = s''(x_i). \quad (164)$$

Можна записати (163) у матричному вигляді, якщо ввести матрицю  $R = \text{diag } \rho_i$ :

$$H^\top \vec{m} + R\vec{\mu} = R\vec{f}. \quad (165)$$

Тут  $\vec{f}$  — вектор заданих значень функції.

Таким чином маємо для  $\vec{m}$  та  $\vec{\mu}$  дві системи (157) і (165). Виключаючи  $\vec{\mu}$  отримаємо таку систему лінійних рівнянь

$$(A + HR^{-1}H^\top) \vec{m} = H\vec{f}. \quad (166)$$

Розв'язавши її, можемо обчислити

$$\vec{\mu} = \vec{f} - R^{-1}H^\top \vec{m} \quad (167)$$

і підставити знайдені значення  $\mu_i$  та  $m_i$  в формулу для сплайну

$$\begin{aligned}
 g(x) = & m_{i-1} \cdot \frac{(x_i - x)^3}{6h_i} + m_i \cdot \frac{(x - x_{i-1})^3}{6h_i} + \\
 & + \left( \mu_{i-1} - \frac{m_{i-1} \cdot h_i^2}{6} \right) \cdot \frac{x_i - x}{h_i} + \\
 & + \left( \mu_i - \frac{m_i \cdot h_i^2}{6} \right) \cdot \frac{x - x_{i-1}}{h_i},
 \end{aligned} \tag{168}$$

Тепер звернемо увагу на матрицю системи (165) :

$$A' = A + HR^{-1}H^T. \tag{169}$$

Оскільки матриці  $H, H^T$  — трьохдіагональні, то матриця  $HR^{-1}H^T$  буде п'ятидіагональною, а тому п'ятидіагональною буде й  $A'$ .

Розв'язують зазвичай системи з такими матрицями наступним чином:

1. або методом квадратних коренів; для матриць із такою структурою цей метод має складність  $Q = O(nm) = O(2n) = O(n)$ , оскільки в нашому випадку півширина діагональної смуги  $m = 2$ .
2. або методом п'ятидіагональної прогонки [Самарский А. А., Николаев С. Н., «Методы решения сеточных уравнений»], що також має складність  $O(n)$ .

**Зауваження:**  $\rho_i$  вибирають так:  $\rho_i = 1/\varepsilon_i^2$ .

## 9. Чисельне інтегрування

### 9.1. Постановка задачі чисельного інтегрування

Нехай потрібно знайти

$$I(f) = \int_a^b \rho(x)f(x) \, dx, \quad (1)$$

де  $f$  — задана функція,  $\rho(x) > 0$  — деякий ваговий множник. Ця задача часто вимагає чисельного вирішення, оскільки

- значна кількість інтегралів типу (1) не можуть бути обчислені аналітично;
- інформація про функцію  $f$  може бути задана у вигляді таблиці.

Нагадаємо, що за означенням

$$I(f) = \lim_{\Delta \rightarrow 0} \sum_{i=1}^n \rho(\xi_i)f(\xi_i)\Delta x_i, \quad (2)$$

де  $\Delta x_i = x_i - x_{i-1}$ , а  $\{x_i\}_{i=0}^n$  — розбиття проміжку  $[a, b]$ ,  $x_i \in [a, b]$ ,  $\xi_i \in [x_{i-1}, x_i]$ .

**Означення:** Тому візьмемо як наближення таку суму, яка називається *квадратурною формулою*:

$$I_n(f) = \sum_{k=0}^n c_k f(x_k), \quad (3)$$

де  $x_k$  — вузли квадратурної формули, а  $c_k$  — її вагові множники.

Задача полягає в тим, щоб вибрати  $\{x_n, c_k\}_{k=0}^n$ , так щоб похибка була найменша:

$$R_n(f) = I(f) - I_n(f) \rightarrow \min \quad (4)$$

**Означення:** Квадратурну формулу (3) називають *квадратурною формулою замкнутого типу*, якщо  $x_0 = a$  та  $x_n = b$ , і *відкритого типу*, якщо  $x_0 > a$  та  $x_n < b$ .

**Означення:** Кажуть, що квадратурна формула (3) має  *$m$ -ий степінь алгебраїчної точності*, якщо

$$\forall f \in \pi_m : R_n(f) = 0, \quad (5)$$

де  $\pi_m$  — множина поліномів  $m$ -го степеня, і

$$\exists P_{m+1}(x) \in \pi_{m+1} : R_n(P_{m+1}) \neq 0. \quad (6)$$

Цю умову можна замінити умовою

$$R_n(x^\alpha) = 0, \quad \alpha = \overline{0, m}, \quad R_n(x^{m+1}) \neq 0, \quad (7)$$

вона більш зручна для перевірки.

Розглянемо деякі підходи до побудови квадратурних формул:

1. *Інтерполяційний*. Він приводить до квадратурних формул інтерполяційного типу. В інтегралі (1) покладають  $f(x) \approx L_n(x)$  по деяких вузлах  $\{x_k\}_{k=0}^n$  (вузли фіксовані). Тоді:

$$I_n(f) \approx I(L_n(x)) = \int_a^b \rho(x) \sum_{k=0}^m \frac{f(x_k) \omega_n(x)}{(x - x_k) \omega'_n(x_k)} dx = \sum_{k=0}^n f(x_k) \int_a^b \rho(x) \frac{\omega_n(x)}{(x - x_k) \omega'_n(x_k)} dx. \quad (8)$$

Отже вузлами цієї квадратурної формули є вузли інтерполяційного багаточлена, а вагові множники

$$c_j = \int_a^b \rho(x) \frac{\omega_n(x)}{(x - x_k) \omega'_n(x_k)} dx \quad (9)$$

2. *Найвищого алгебраїчного степеня точності*. Вибираємо одночасно  $x_k$  і  $c_k$  з умови  $R_n(x) = 0$ ,  $\alpha = 0, m$ , щоб  $m$  було максимальним. Отримуємо систему нелінійних алгебраїчних рівнянь, розв'язавши яку отримуємо квадратурні формули найвищого алгебраїчного степеня точності.
3. *Складені квадратурні формули*. Проміжок  $[a, b]$  розбиваємо на підпроміжки (наприклад однокової довжини), а потім на кожному проміжку використовуємо, з невеликим ступенем, формули з пункту 1 або 2. Наприклад, для формул інтерполяційного типу:

$$I(f) = \sum_{i=1}^N \int_{x_{i-1}}^{x_i} \rho(x) f(x) dx \approx \sum_{i=1}^N \sum_{k=0}^n c_k^i f(x_k^i) = I_h(f). \quad (10)$$

**Означення:** Кажуть, що квадратурна формула складеного типу  $I_h$  має *порядок (ступінь) точності*  $p$  по кроку  $h$ , якщо  $R_h(f) = I(f) - I_h(f) = O(h^p)$ .

4. *Квадратурні формули оптимальні на класі функцій*. Вибираємо  $\{x_k, c_k\}$  так, щоб досягався

$$\inf_{x_k, c_k} \sup_{f \in F} R_n(f). \quad (11)$$

Це ми можемо робити, коли знаємо з яким класом функцій маємо справу.

**Зауваження 1** (про квадратурні формули інтерполяційного типу): При підвищенні степеня інтерполяції погіршується якість наближення функції внаслідок розбіжності процесу інтерполяції:

$$|f - L_n|_C \not\rightarrow 0. \quad (12)$$

Але

$$R_n \xrightarrow{n \rightarrow \infty} 0. \quad (13)$$

наприклад, для  $f \in C([a, b])$ .

І все ж таки розбіжність процесу інтерполювання дає взнаки:

$$\max_k |c_k| \xrightarrow{n \rightarrow \infty} \infty. \quad (14)$$

і це приводить до поганих наслідків чисельного інтегрування. Дійсно, розглянемо випадок, коли функція задана неточними значеннями:



$$\tilde{f}(x_k) = f(x_k) + \delta_k, \quad |\delta_k| < \delta. \quad (15)$$

Тоді

$$\delta I_n(f) = I_n(\tilde{f}) - I_n(f) = \sum_{k=0}^n c_k \delta_k. \quad (16)$$

Якщо всі  $c_k > 0$ , то

$$|\delta I| = \sum_{k=0}^n c_k |\delta_k| \leq \delta \sum_{k=0}^n c_k = \delta(b-a). \quad (17)$$

При  $\rho \equiv 1$ , якщо підставити  $f \equiv 1$ , то отримаємо

$$b-a = \int_a^b dx = \sum_{k=0}^n c_k. \quad (18)$$

При  $\rho \neq 1$ :

$$\sum_{k=0}^n c_k = \int_a^b \rho(x) dx, \quad (19)$$

бо хоча б нульовий степінь точності будь-яка квадратурна формула повинна мати.

Нагадаємо, що

$$\max_k |c_k| \xrightarrow{n \rightarrow \infty} \infty, \quad (20)$$

а оскільки  $\sum c_k > 0$ , то  $\exists c_k > 0$  і  $\exists c_k < 0$ , тому з ростом  $n$  зростає  $|c_k|$ , а відповідно і вплив похибки на результат. Тому не можна використовувати великі степені і використовують складені квадратурні формули.

**Зауваження 2:** Ясно, що квадратурні формули інтерполяційного типу мають алгебраїчний степінь точності принаймні  $m = n$ , бо ми заміняємо  $f \mapsto L_n$ , а якщо  $f \in \pi_n$ , то  $f \equiv L_n$ . Але виявляється, що для парних  $n$  та симетричному розташуванні вузлів інтегрування,  $m = n + 1$ , тобто алгебраїчний степінь точності на одиницю вищий степеня інтерполяції.

## 9.2. Квадратурні формули прямокутників

Припустимо, що  $\rho \equiv 1$ . Тоді можна побудувати такі квадратурні формули інтерполяційного типу при  $n = 0$ :

1. лівих прямокутників:  $x_0 = a: I_0^L = (b-a) \cdot f(a)$ ;
2. правих прямокутників:  $I_0^R = (b-a) \cdot f(b), x_0 = b$ ;
3. середніх прямокутників:

$$I_0 = (b-a) \cdot f(x_0), \quad x_0 = \frac{a+b}{2} \quad (21)$$

Знайдемо тепер алгебраїчну степінь точності цих квадратурних формул. Для лівих прямокутників:

$$I_0^L(1) = b - a = I(1), \quad (22)$$

$$I_0^L(x) = (b - a) \cdot a \neq \frac{b^2 - a^2}{2} = \int_a^b x \, dx = I(x), \quad (23)$$

отже степінь точності  $m = 0$ . Така ж вона буде і для  $I_0^R$ . А для середніх прямокутників

$$I_0(x) = (b - a) \cdot \frac{a + b}{2} = I(x), \quad (24)$$

$$I_0(x^2) = (b - a) \cdot \left(\frac{a + b}{2}\right)^2 \neq \frac{b^3 - a^3}{3} = \int_a^b x^2 \, dx = I(x^2), \quad (25)$$

тому  $m = 1$ . Отож нею і будемо користуватися.

Оцінимо для неї похибку. Взагалі для формули інтерполяційного типу:

$$R_n(f) = I(f) - I_n(f) = I(f) - I(L_n) = I(f - L_n) = I(r_n) = \int_a^b r_n(x) \, dx, \quad (26)$$

де  $r_n(x)$  — залишковий член інтерполяції. Далі

$$|R_n(f)| \leq (b - a) \cdot \max_x |r_n(x)| \leq (b - a) \cdot \frac{M_{n+1}}{n + 1} \cdot \max_x |\omega_n(x)|. \quad (27)$$

Для  $I_0$ :

$$\begin{aligned} |R_0(f)| &= \left| \int_a^b r_0(x) \, dx \right| \leq \int_a^b |r_0(x)| \, dx \leq \int_a^b \frac{M_1}{1!} |x - x_0| \, dx = \\ &= M_1 \int_a^b |x - x_0| \, dx \leq M_1 \cdot \frac{b^2 - a^2}{4}. \end{aligned} \quad (28)$$

Але це погана оцінка, вона не використовує той факт, що квадратурна формула має степінь точності на одиницю вищу. Отримаємо кращу оцінку. Маємо при  $f \in C^2([a, b])$ :

$$f(x) = f(x_0) + (x - x_0) \cdot f'(x_0) + \frac{(x - x_0)^2}{2} \cdot f''(\xi), \quad (29)$$

де  $x_0 \equiv \frac{a+b}{2}$ , а  $\xi \in [a, b]$ . Тоді

$$\begin{aligned} R_0(f) &= \int_a^b f(x) \, dx - \int_a^b L_0(x) \, dx = \int_a^b (f(x) - f(x_0)) \, dx = \\ &= \int_a^b \left( (x - x_0) f'(x_0) + \frac{(x - x_0)^2}{2} f''(\xi) \right) \, dx = \\ &= \int_a^b \frac{(x - x_0)^2}{2} \cdot f''(\xi) \, dx = f''(\eta) \int_a^b \frac{(x - x_0)^2}{2} \, dx = \frac{f''(\eta)}{24} \cdot (b - a)^3. \end{aligned} \quad (30)$$

Таким чином

$$|R_0(f)| \leq \frac{M_2}{24} \cdot (b - a)^3 \quad (31)$$

Але тут у нас немає впливу на точність (величину похибки). Тому використовують формулу складеного типу. Якщо сітка рівномірна, то складена квадратурна формула прямокутників має вигляд

$$I(f) = \sum_{i=1}^n \int_{x_{i-1}}^{x_i} f(x) dx \approx \sum_{i=1}^N h \cdot f(\bar{x}_i) = I_h(f), \quad (32)$$

де  $\bar{x}_i = x_{i-1/2} = x_i - h/2$ .

Оцінімо похибку цієї квадратурної формули:

$$R_h(f) = I(f) - I_h(f) = \sum_{i=1}^N \int_{x_{i-1}}^{x_i} (f(x) - f(\bar{x}_i)) dx = \sum_{i=1}^N f''(\eta_i) \cdot \frac{h^3}{24}, \quad (33)$$

$$|R_h(f)| \leq \frac{M_2}{24} n h^3 = \frac{M_2 h^2 (b - a)}{24}. \quad (34)$$

Тобто ця формула має степінь точності  $p = 2$  по кроку  $h$ . (Не слід плутати з алгебраїчним степенем точності  $m = 1$  для цієї формули).

Якщо  $f(x) \in C^4([a, b])$ , то

$$\begin{aligned} f(x) - f(\bar{x}_i) &= f(\bar{x}_i) + (x - \bar{x}_i)f'(\bar{x}_i) + \frac{(x - \bar{x}_i)^2}{2} \cdot f''(\bar{x}_i) + \\ &+ \frac{(x - \bar{x}_i)^3}{6} \cdot f'''(\bar{x}_i) + \frac{(x - \bar{x}_i)^4}{24} \cdot f^{(4)}(\xi_i) - f(\bar{x}_i). \end{aligned} \quad (35)$$

При непарних степенях інтеграли пропадуть і тому:

$$\begin{aligned} R_h(f) &= \sum_{i=1}^N \int_{x_{i-1}}^{x_i} \frac{(x - \bar{x}_i)^2}{2} \cdot f''(\bar{x}_i) dx + \sum_{i=1}^N \int_{x_{i-1}}^{x_i} \frac{(x - \bar{x}_i)^4}{24} \cdot f^{(4)}(\xi_i) dx = \\ &= \frac{h^2}{24} \sum_{i=1}^N h \cdot f''(\bar{x}_i) + \sum_{i=1}^N \frac{h^5 \cdot f^{(4)}(\eta_i)}{1920}. \end{aligned} \quad (36)$$

Оскільки

$$\sum_{i=1}^N h \cdot f''(\bar{x}_i) \quad (37)$$

це квадратурна формула середніх прямокутників для  $f''(x)$  з похибкою  $O(h^2)$ , то

$$R_h(f) = \frac{h^2}{24} \int_a^b f''(x) dx + O(h^4) = O(h^4), \quad (38)$$

і

$$R_h(f) = \overset{\circ}{R}_h(f) + \alpha(h), \quad (39)$$

де

$$\overset{\circ}{R}_h(f) = \frac{h^2}{24} \int_a^b f''(x) dx, \quad (40)$$

де  $\alpha(h) = O(h^4)$ .

Ця формула використовується для побудови програм, що автоматично вибирають крок інтегрування.

### 9.3. Формула трапеції

Нехай  $x_0 = a$ ,  $x_1 = b$ ,  $L_1(x) = f(x)$ . Тоді отримаємо формулу:

$$I_1(f) = \frac{b-a}{2} \cdot (f(a) + f(b)) \quad (41)$$

Формула має алгебраїчний степінь точності  $m = 1$ , оскільки  $I(x^2) \neq I_1(x^2)$ . Це формула замкнутого типу. Залишковий член:

$$R_1(f) = \int_a^b \frac{f''(\xi)(x-a)(x-b)}{2} dx = -\frac{(b-a)^3}{12} \cdot f''(\xi). \quad (42)$$

Оцінка залишкового члена:

$$|R_1(f)| \leq M_2 \cdot \frac{(b-a)^3}{12}. \quad (43)$$

З геометричної точки зору замінюється площа криволінійної трапеції площею звичайної трапеції.

Складена квадратурна формула трапецій:

$$\begin{aligned} I_h(f) &= \sum_{i=1}^N \frac{h}{2} \cdot (f(x_{i-1}) + f(x_i)) = \\ &= \frac{h}{2} \cdot f(a) + \sum_{i=1}^{N-1} h \cdot f(x_i) + \frac{h}{2} \cdot f(b), \end{aligned} \quad (44)$$

де  $x_i = a + ih$ ,  $h = \frac{b-a}{N}$ ,  $i = \overline{0, N}$  та

$$|R_h(f)| \leq M_2 \cdot \frac{b-a}{12} \cdot h^2, \quad f \in C^2([a, b]). \quad (45)$$

Якщо  $f \in C^4([a, b])$ , то

$$R_h(f) = \overset{\circ}{R}_h(f) + \alpha(h), \quad (46)$$

де

$$\overset{\circ}{R}_h(f) = -\frac{h^2}{12} \int_a^b f^{(4)}(x) dx, \quad (47)$$

а  $\alpha(h) = O(h^4)$ .

**Задача 29:** Використовуючи явний вигляд головних членів похибки складених квадратурних формул прямокутників та трапецій, побудувати лінійною комбінацією цих двох формул квадратурну формулу

четвертого степеня точності за кроком  $h$ .

## 9.4. Квадратурна формула Сімпсона

Нехай  $x_0 = a$ ,  $x_1 = \frac{a+b}{2}$ ,  $x_2 = b$ . Замість  $f$  використовуємо  $L_2(x)$ . Тоді отримуємо квадратурну формулу:

$$I_2(f) = \frac{b-a}{6} \left( f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right). \quad (48)$$

Це *квадратурна формула Сімпсона*.

**Задача 30:** Довести, що алгебраїчна степінь точності квадратурної формули Сімпсона  $m = 3$ .

**Задача 31:** Довести, що для  $f \in C^4([a, b])$  залишковий член квадратурної формули Сімпсона має місце представлення:

$$R_2(f) = \frac{1}{24} \int_a^b (x-a) \left( x - \frac{a+b}{2} \right)^2 (x-b) f^{(4)}(\xi) dx = \frac{f^{(4)}(\xi)}{2880} \cdot (b-a)^5, \quad (49)$$

та вірна оцінка:

$$|R_2(f)| \leq \frac{M_4}{2880} \cdot (b-a)^5. \quad (50)$$

Складена квадратурна формула Сімпсона має вигляд:

$$\begin{aligned} I_h(f) &= \sum_{i=1}^N \frac{h}{6} (f(x_{i-1}) + 4f(x_{i-1/2}) + f(x_i)) = \\ &= \frac{h}{6} (f(x_0) + 4f(x_{1/2}) + 2f(x_1) + \dots + 2f(x_{N-1}) + 4f(x_{N-1/2}) + f(x_N)). \end{aligned} \quad (51)$$

Якщо  $f \in C^4([a, b])$ , то має місце оцінка:

$$|R_h(f)| \leq \frac{M_4}{2880} \cdot (b-a) \cdot h^4, \quad p = 4. \quad (52)$$

Якщо  $f \in C^6([a, b])$ , то

$$R_h(f) = \overset{\circ}{R}_h(f) + \alpha(h), \quad (53)$$

де

$$\overset{\circ}{R}_h(f) = \frac{h^4}{2880} \int_a^b f^{(4)}(x) dx, \quad (54)$$

**Задача 32:** Побудувати інтерполяційну квадратурну формулу для  $n = 3$ ,  $x_0 = a$ ,  $x_1 = \frac{3a+b}{4}$ ,  $x_2 = \frac{a+3b}{4}$ ,  $x_3 = b$ . Який алгебраїчний степінь точності вона має?

## 9.5. Принцип Рунге

Нехай задана деяка величина  $I$  (сіткова функція, інтеграл, неперервна функція). Нехай  $I_h \approx I$  та  $I_n \rightarrow I$  при  $h \rightarrow 0$ . Нехай похибка послідовності  $I_h$  представляється у вигляді:

$$R_h = I - I_h = \overset{\circ}{R}_h + \alpha(h), \quad (55)$$

де  $\overset{\circ}{R}_h = C \cdot h^m$  — головний член похибки,  $C$  не залежить від  $h$ ,  $\alpha(h) = o(h^m)$ . Обчислимо  $I_{h/2}$ . З (55) випливає, що

$$I = I_h + Ch^m + \alpha(h), \quad (56)$$

$$I = I_{h/2} + C \cdot \frac{h^m}{2^m} + \alpha(h). \quad (57)$$

Звідси

$$I_{h/2} - I_h = \frac{Ch^m}{2^m} \cdot (2^m - 1) + \alpha(h). \quad (58)$$

З (55):

$$\overset{\circ}{R}_{h/2} = \frac{Ch^m}{2^m} = \frac{I_{h/2} - I_h}{2^m - 1}, \quad (59)$$

та

$$\overset{\circ}{R}_h = \frac{2^m}{2^m - 1} \cdot (I_{h/2} - I_h). \quad (60)$$

**Означення:** Формула (59) носить назву *апостеріорної оцінки* похибки обчислення  $I$  за допомогою наближення  $I_{h/2}$ . (Апріорні оцінки це оцінки отримані до обчислення величини  $I_h$ , апостеріорні оцінки — під час її обчислення).

З формули (59) впливає такий алгоритм обчислення інтегралу із заданою точністю  $\varepsilon$ :

- обчислюємо  $I_h, I_{h/2}, \overset{\circ}{R}_{h/2}$ ;
- перевіряємо чи  $\left| \overset{\circ}{R}_{h/2} \right| < \varepsilon$ .
  - Якщо так, то  $I \approx I_{h/2}$ ;
  - Якщо ж ні, то:
    - обчислюємо  $I_{h/2}, I_{h/4}, \overset{\circ}{R}_{h/4}$ ;
    - перевіряємо  $\left| \overset{\circ}{R}_{h/4} \right| < \varepsilon$  і т. д.
- Процес продовжуємо поки не буде виконана умова  $\left| \overset{\circ}{R}_{h/2^k} \right| < \varepsilon, k = 1, 2, \dots$

**Зауваження:** Ми даємо оцінку не похибки, а її головного члена з точністю  $\alpha(h)$ , тому такий метод може давати збої, якщо не виконана умова

$$|\alpha(h)| \ll \left| \overset{\circ}{R}_{h/2^k} \right|. \quad (61)$$

За допомогою головного члена похибки можна отримати краще значення для  $I$ :

$$\tilde{I}_{h/2} = I_{h/2}^{(1)} = I_{h/2} + \overset{\circ}{R}_{h/2} = \frac{2^m}{2^m - 1} \cdot I_{h/2} - \frac{1}{2^m - 1} \cdot I_h. \quad (62)$$

Це екстраполяційна формула Річардсона:  $I_h - \tilde{I}_{h/2} = \alpha(h)$ .

Для квадратурної формули трапецій  $p = 2$  і

$$I - I_h = Ch^2 + O(h^4), \quad (63)$$

$$\overset{\circ}{R}_{h/2} = \frac{I_{h/2} - I_h}{3}. \quad (64)$$

Маємо

$$R_h = -\frac{h^2}{12} \int_a^b f''(x) dx + O(h^4) = O(h^2). \quad (65)$$

Отже, якщо застосовувати екстраполяційну формулу Річардсона, то

$$\tilde{I}_{h/2} = \frac{4}{3} \cdot I_{h/2} - \frac{1}{3} \cdot I_h, \quad (66)$$

і  $I_h - \tilde{I}_{h/2} = O(h^4)$ .

**Задача 33:** Написати явний вигляд квадратурної формули, яка отримується екстраполяцією Річардсона з квадратурної формули трапецій.

Якщо невідомо  $m$ , то можна використати принцип Рунге для його знаходження. Для цього використаємо  $I_h, I_{h/2}, I_{h/4}$ :

$$I_{h/2} - I_h = \frac{Ch^m}{2^m} \cdot (2^m - 1) + \alpha(h), \quad (67)$$

$$I_{h/4} - I_{h/2} = \frac{Ch^m}{2^{2m}} \cdot (2^m - 1) + \alpha(h), \quad (68)$$


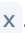
З точністю  $\alpha(h)$  маємо

$$2^m = \frac{I_{h/2} - I_h}{I_{h/4} - I_{h/2}}. \quad (69)$$

Звідки

$$m = \log_2 \left( \frac{I_{h/2} - I_h}{I_{h/4} - I_{h/2}} \right). \quad (70)$$

Оцінка  $\left| \overset{\circ}{R}_{h/4} \right| < \varepsilon$  — найбільш точна, тому  $I \approx I_{h/4}$ .

Покажемо чому формулу прямокутників рідко використовують з принципом Рунге. Нехай точки, в яких обчислюється значення функції позначаються: в  $I_h$  — , в  $I_{h/2}$  — .

Для формули трапецій використовуються такі точки:

$I_h$ :    o \_ \_ \_ o \_ \_ \_ o \_ \_ \_ o

$I_{h/2}$ : o \_ x \_ o \_ x \_ o \_ x \_ o

Для формули прямокутників:

$I_h$ :    \_ \_ o \_ \_ \_ o \_ \_ \_ o \_ \_

$I_{h/2}$ : \_ x \_ x \_ x \_ x \_ x \_ x \_

Як бачимо для формули трапецій необхідно підраховувати нові значення в  $N$  точках, а для формули прямокутників в  $2N$ .

Для економного використання обчислених значень функції на сітці з кроком  $h$  для формули трапецій запишемо:

$$I_h = \frac{h}{2} \left( f(a) + 2 \sum_{i=1}^{N-1} f(x_i) + f(b) \right), \quad (71)$$

$$I_{h/2} = \frac{h}{4} \left( f(a) + 2 \sum_{i=1}^{N-1} f(x_i) + 2 \sum_{i=1}^{N-1} f(x_{i-1/2}) + f(b) \right) = \frac{1}{2} \cdot I_h + \frac{h}{2} \sum_{i=1}^{N-1} f(x_{i-1/2}). \quad (72)$$

Отже, на одному кроці принципу Рунге кількість обчислень  $Q_t = O(N)$ , а для  $Q_r = O(2N)$ .

Цей принцип застосовується і для формули Сімпсона  $m = 4$ . Головна частина залишкового члена для цієї формули:

$$\overset{\circ}{R}_{h/2} = \frac{I_{h/2} - I_h}{15}. \quad (73)$$

$$\tilde{I}_{h/2} = \frac{16}{15} \cdot I_{h/2} - \frac{1}{15} \cdot I_h, \quad (74)$$

$$I_h - \tilde{I}_{h/2} = O(h^6). \quad (75)$$

Розглянемо використання так званих *адаптивних квадратурних формул*, в яких змінний крок вибирається за принципом Рунге. Для цього запишемо формулу трапецій із змінним кроком:

$$I_h(f) = \sum_{i=1}^N \frac{h_i}{2} \cdot (f(x_{i-1}) + f(x_i)), \quad (76)$$

де  $h_i = x_i - x_{i-1}$ .

Оцінімо похибку на кожному інтервалі:

$$R_{h_i} = I_i - I_{h_i} = \int_{x_{i-1}}^{x_i} f(x) dx - \frac{h_i}{2} (f(x_{i-1}) + f(x_i)) = -\frac{h_i^3}{6} \cdot f''(x_{i-1/2}) + O(h_i^5). \quad (77)$$

Таким чином  $p = 3$  і головний член похибки:

$$\overset{\circ}{R}_{h/2} = \frac{I_{h_i/2} - I_{h_i}}{7}. \quad (78)$$



Умова припинення ділення навпіл проміжку  $[x_{i-1}, x_i]$ :

$$\left| \overset{\circ}{R}_{h_i/2} \right| \leq \frac{\varepsilon \cdot h_i}{b-a}. \quad (79)$$

Це забезпечує точність на всьому інтервалі

$$\left| \overset{\circ}{R}_{h/2} \right| = \left| \sum_{i=1}^N R_{h_i/2} \right| \leq \sum_{i=1}^N \frac{\varepsilon \cdot h_i}{b-a} = \varepsilon \cdot \frac{b-a}{b-a} = \varepsilon. \quad (80)$$

Ще одне застосування принципу Рунге — високоточне обчислення інтегралу від достатньо гладкої функції за допомогою *таблиці Ромберга*. Для побудови цієї таблиці обчислимо за допомоги складеної квадратурної формули трапецій із сталим кроком  $h$  послідовність значень  $I_h = I_h^{(0)}$ ,  $I_{h/2} = I_{h/2}^{(0)}$ ,  $I_{h/4} = I_{h/4}^{(0)}$ ,  $I_{h/8} = I_{h/8}^{(0)}$ , ... які мають похибку  $O(h^2)$ . За допомогою екстраполяції Річардсона (62) з коефіцієнтами лінійної комбінації  $(4/3, -1/3)$  уточнимо ці значення (див. також формулу (66)). Отримаємо  $I_{h/2}^{(1)}$ ,  $I_{h/4}^{(1)}$ ,  $I_{h/8}^{(1)}$ . Вони мають похибку  $O(h^4)$ . Знову використовуємо екстраполяцію Річардсона з коефіцієнтами лінійної комбінації  $(16/15, -1/15)$ . Отримаємо  $I_{h/4}^{(2)}$ ,  $I_{h/8}^{(2)}$ , які мають точність  $O(h^6)$  і т. д.

**Означення:** Отримані значення можна розмістити в такій *таблиці Ромберга*:

	0	1	2	3
$h$	$I_h^{(0)}$			
$h/2$	$I_{h/2}^{(0)}$	$I_{h/2}^{(1)}$		
$h/4$	$I_{h/4}^{(0)}$	$I_{h/4}^{(1)}$	$I_{h/4}^{(2)}$	
$h/8$	$I_{h/8}^{(0)}$	$I_{h/8}^{(1)}$	$I_{h/8}^{(2)}$	$I_{h/8}^{(3)}$

Всі значення крім останнього  $I_{h/8}^{(3)}$  можна оцінити за принципом Рунге (див. формулу (59)). Використання формули (72) для обчислення  $I_{h/2^k}^{(0)}$  та лінійні комбінації (59) дають простий та економічний алгоритм обчислення  $I$ . Початкове значення  $h$  можна брати рівним  $b-a$ , або  $\frac{b-a}{n}$ , де  $n$  ціле.

## 9.6. Квадратурні формули найвищого алгебраїчного степеня точності

Розглянемо інтеграл:

$$I(f) = \int_a^b \rho(x) f(x) dx, \quad (81)$$

де

$$\rho(x) > 0, \quad x \in [a, b], \quad \left| \int_a^b \rho(x) x^i dx \right| < \infty \quad (82)$$

Розглянемо задачу побудови квадратурної формули:

$$I_n(f) = \sum_{k=1}^n c_k f(x_k), \quad (83)$$

яка при заданому  $n$  була б точною для алгебраїчного багаточлена можливо більшого степеня.

**Означення:** Такі квадратурні формули існують, вони називаються квадратурні формули *найвищого алгебраїчного степеня точності* або формули Гауса (або Гауса-Крістофеля).

В (83) невідомими є  $c_k, x_k, k = \overline{1, n}$ . Їх обирають з умови, що (83) точна для будь-якого багаточлена степеня  $p$ , а це еквівалентно умові, щоб формула була точною для функції  $f(x) = x^\alpha, \alpha = \overline{0, p}$ . Звідси отримуємо умови:

$$I_n(x^\alpha) = \int_a^b \rho(x) x^\alpha dx = \sum_{k=1}^n c_k x_k^\alpha, \quad \alpha = \overline{0, p}. \quad (84)$$

Ми хочемо отримати формули для  $m \rightarrow \max$ . Щоб кількість рівнянь була рівною кількості невідомих нам потрібно, щоб  $p + 1 = 2n$ .

**Задача 34:** Побудувати квадратурну формулу найвищого степеня точності (розв'язати систему рівнянь (84) для  $a = -1, b = 1, \rho(x) = 1$ ).

**Теорема (Гаусса):** Квадратурна формула (83) буде точною для будь-якого багаточлена степеня  $p = 2n - 1$ , тобто  $\forall f(x) \in \pi_{2n-1}$  тоді і тільки тоді, коли виконуються умови:

1. поліном  $\omega(x) = (x - x_1)(x - x_2) \dots (x - x_n)$  ортогональний з вагою  $\rho(x)$  до будь-якого багаточлена степеня менше  $n, Q_{n-1}$ :

$$\int_a^b \omega(x) Q_{n-1}(x) \rho(x) dx = 0; \quad (85)$$

2. формула (83) є квадратурною формулою інтерполяційного типу, тобто коефіцієнти обчислюються за формулою:

$$c_k = \int_a^b \rho(x) \frac{\omega(x)}{(x - x_k) \omega'(x_k)} dx. \quad (86)$$

*Доведення:*

1. **Необхідність:** Нехай формула (83) точна для багаточлена степеня  $p = 2n - 1$ , тобто  $I(f) = I_n(f), \forall f(x) \in \pi_{2n-1}$ . Тоді

$$I(f) = \int_a^b \rho(x) \omega(x) Q_{n-1}(x) dx = \sum_{k=1}^n c_k \omega(x_k) Q_{n-1}(x_k) = 0. \quad (87)$$

Тобто виконується (85). Тепер покладемо

$$f(x) = \frac{\omega(x)}{(x - x_j) \omega'(x_j)} \in \pi_{n-1} \subset \pi_{2n-1}. \quad (88)$$

Отримаємо

$$\begin{aligned}\int_a^b \rho(x) f(x) dx &= \int_a^b \rho(x) \frac{\omega(x)}{(x - x_j) \omega'(x_j)} dx = \\ &= \sum_{k=1}^n c_k \frac{\omega(x_k)}{(x_k - x_j) \omega'(x_j)} = \sum_{k=1}^n c_k \delta_{kj} = c_j.\end{aligned}\quad (89)$$

тобто виконується і умова (86).

2. *Достатність*: Нехай виконується (85) і (86). Подамо  $\forall f(x) \in \pi_{2n-1}$  у вигляді

$$f(x) - \omega(x) Q_{n-1}(x) + R_{n-1}(x). \quad (90)$$

Розглянемо

$$\begin{aligned}I(f) &= \int_a^b \rho(x) f(x) dx = \int_a^b \rho(x) (\omega(x) Q_{n-1}(x) + R_{n-1}(x)) dx = \\ &= \sum_{k=1}^n c_k \omega(x_k) Q_{n-1}(x_k) + \sum_{k=1}^n c_k R_{n-1}(x_k).\end{aligned}\quad (91)$$

Оскільки  $R_{n-1}(x_k) = f(x_k) - \omega(x_k) Q_{n-1}(x_k) = f(x_k)$ , то

$$I(f) = \sum_{k=1}^n c_k f(x_k) = I_n(f). \quad (92)$$

Тобто формула (83) є точною для будь-якого багаточлена степеня  $2n - 1$ .  $\square$

Отже, з точністю до сталого множника багаточлени  $\omega(x)$  співпадають з багаточленами  $n$ -того степеня ортогональної системи багаточленів. Ця система ортогональна на проміжку  $[a, b]$  з вагою  $\rho(x)$ .

Вивчимо деякі властивості квадратурних формул Гауса:

1. Покажемо, що  $c_k, x_k$  визначаються однозначно.

Представимо багаточлен  $\omega(x)$  у вигляді

$$\omega(x) = a_0 + a_1 x + \dots + a_{n-1} x^{n-1} + x^n. \quad (93)$$

Умови ортогональності (85) приймуть вигляд

$$\int_a^b \rho(x) \omega(x) x^\alpha dx = \int_a^b \rho(x) (a_0 + \dots + x^n) x^\alpha dx = 0, \quad (94)$$

де  $\alpha = \overline{0, n-1}$ , або

$$\int_a^b \rho(x) (a_0 + a_1 x + \dots + a_{n-1} x^{n-1}) x^\alpha dx = - \int_a^b \rho(x) x^n x^\alpha dx. \quad (95)$$

Покажемо, що відповідна однорідна система рівнянь

$$\int_a^b \rho(x)(a_0 + a_1x + \dots + a_{n-1}x^{n-1})x^\alpha dx = 0, \quad (96)$$

де  $\alpha = \overline{0, n-1}$ , має єдиний розв'язок  $a_0 = a_1 = \dots = a_{n-1} = 0$ .

Помножимо систему (96) на  $a_\alpha$  і просумуємо по всіх  $\alpha = \overline{0, n-1}$ :

$$\begin{aligned} \sum_{\alpha=0}^{n-1} a_\alpha \int_a^b \rho(x)(a_0 + a_1x + \dots + a_{n-1}x^{n-1})x^\alpha dx = \\ = \int_a^b \rho(x) \sum_{\alpha=0}^{n-1} \sum_{j=0}^{n-1} a_j a_\alpha x^\alpha x^j dx = \int_a^b \rho(x) \left( \sum_{j=0}^{n-1} a_j x^j \right)^2 dx = 0. \end{aligned} \quad (97)$$

Звідси і з умови  $\rho(x) > 0$  випливає, що  $a_0 = a_1 = \dots = a_{n-1} = 0$ . Тому і відповідна неоднорідна система має єдиний розв'язок. Отже існує єдиний багаточлен  $\omega(x)$  степеня  $n$ , який ортогональний з вагою  $\rho(x)$  до будь-якого багаточлена степеня  $n-1$ .

2. Покажемо, що найвищий степінь точності формули Гаусса  $p = 2n - 1$ . З теореми випливає, що  $p \geq 2n - 1$ . Покажемо, що існує багаточлен степеня  $2n$ , для якого формула не є точною. Для цього введемо функцію  $f(x) = \omega^2(x) \in \pi_{2n}$ . Маємо

$$I(f) = \int_a^b \rho(x) \omega^2(x) dx > 0, \quad (98)$$

але

$$I_n(f) = \sum_{k=1}^n c_k \omega^2(x_k) = 0. \quad (99)$$

Отже,  $I(f) \neq I_n(f)$ . Звідси  $p \leq 2n - 1$ , тобто  $p = 2n - 1$ .

3. Коефіцієнти формул Гаусса додатні, тобто  $c_k > 0$ . Розглянемо багаточлени

$$\varphi_j = \left( \frac{\omega(x)}{(x - x_j)\omega'(x_j)} \right)^2, \quad (100)$$

які мають степінь  $2n - 2$  і властивості:

$$\varphi_i(x_k) = \delta_{ik}, \quad I(\varphi_j) = \int_a^b \rho(x) \varphi_j(x) dx > 0. \quad (101)$$

Оскільки для цих багаточленів справедлива теорема Гауса, то

$$I(\varphi_j) = I_n(\varphi_j) = \sum_{k=1}^n c_k \varphi_j(x_k) = \sum_{k=1}^n c_k \left( \frac{\omega(x)}{(x - x_j)\omega'(x_j)} \right)^2 = \sum_{k=1}^n c_k \delta_{jk}^2 = c_j. \quad (102)$$

Звідси випливає, що  $c_j > 0$ ,  $j = \overline{1, n}$ .

**Теорема:** Нехай вагова функція  $\rho(x) > 0$ ,  $x \in [a, b]$ ,  $f(x) \in C^{2n}([a, b])$ . Тоді існує точка  $\xi \in [a, b]$  така, що

$$R_n(f) = \frac{f^{(2n)}(\xi)}{(2n)!} \int_a^b \rho(x) \omega^2(x) dx. \quad (103)$$

*Доведення:* Розглянемо інтерполяційний багаточлен Ерміта з двократними вузлами

$$H_{2n-1}(x) : f(x_i) = H_{2n-1}(x_i) \wedge f'(x_i) = H'_{2n-1}(x_i), \quad (104)$$

де  $i = \overline{1, m}$ . Для нього

$$r_{2n-1}(x) = f(x) - H_{2n-1}(x) = \frac{f^{(2n)}(\xi)}{(2n)!} \cdot \omega^2(x). \quad (105)$$

Звідси

$$R_{2n-1}(x) = \int_a^b \rho(x) r_{2n-1}(x) dx = \frac{f^{(2n)}(\xi)}{(2n)!} \int_a^b \rho(x) \omega^2(x) dx. \quad \square \quad (106)$$

## 9.7. Частинні випадки квадратурної формули Гауса

1. Розглянемо відрізок  $[-1, 1]$  і ваговий множник  $\rho(x) = 1$ , тобто виведемо формули Гауса для обчислення інтегралу

$$\int_{-1}^1 f(x) dx. \quad (107)$$

Щоб знайти вузли квадратурна формули розглянемо багаточлени Лежандра:

$$(n+1)L_{n+1}(x) = (2n+1)xL_n(x) - nL_{n-1}(x), \quad (108)$$

де  $L_0 = 1, L_1(x) = x$ .

Багаточлени Лежандра задовольняють умовам [теорему 1 \(пункт 1\)](#), тому  $\omega(x) = L_n(x)$  і вузлами квадратурної формули є корені цього багаточлена. Вагові множники цієї формули обчислюються за формулою

$$c_k = \int_{-1}^1 \frac{\omega(x)}{(x - x_k)\omega'(x_k)} dx, \quad (109)$$

а залишковий член

$$R_n(f) = \frac{2^{2n+1}}{(2n+1)!(2n)!} \cdot \frac{(n!)^2}{(2n)!} \cdot f^{(2n)}(\xi). \quad (110)$$

**Приклад:** Побудуємо квадратурну формулу

$$\int_0^1 f(x) dx \approx \sum_{k=1}^n c_k f(x_k). \quad (111)$$

При  $n = 2$  потрібно знайти  $c_0, c_1, x_0, x_1$ .

**Розв'язок:** Заміною  $t = 2x - 1$  переведемо  $x \in [0, 1]$  на проміжок  $t \in [-1, 1]$ . Запишемо  $L_2(t) = 3t^2/2 - 1/2$ . Звідси

$$L_2(x) = \frac{3(2x - 1) - 1}{2} = \frac{12x^2 - 12x + 2}{6} = 6x^2 - 6x + 1 = 0. \quad (112)$$

Звідси  $x_1 = \frac{3-\sqrt{3}}{6}$ ,  $x_2 = \frac{3+\sqrt{3}}{6}$ . За формулою (86) знайдемо

$$c_1 = \int_0^1 \frac{x - x_2}{x_1 - x_2} dx = \frac{1}{2}, \quad (113)$$

$$c_2 = \int_0^1 \frac{x - x_1}{x_2 - x_1} dx = \frac{1}{2}. \quad (114)$$

Тобто квадратурна формула має вигляд

$$\int_0^1 f(x) dx = \frac{1}{2} \left( f\left(\frac{3-\sqrt{3}}{6}\right) + f\left(\frac{3+\sqrt{3}}{6}\right) \right). \quad (115)$$

2. Розглянемо відрізок  $[-1, 1]$  і вагу  $\rho(x) = 1/\sqrt{1-x^2}$ , тобто виведемо формулу Гауса для обчислення інтегралу

$$I(f) = \int_{-1}^1 \frac{f(x) dx}{\sqrt{1-x^2}}. \quad (116)$$

Багаточлени Чебишова задовольняють умовам **теорему 1 (п.1)**, тому

$$\omega(x) = \overline{T}_n(x) = \frac{\cos(n \arccos(x))}{2^{n-1}}. \quad (117)$$

Вузлами квадратурної формули є корені цього багаточлена Чебишова, тобто корені рівняння  $\cos(n \arccos(x)) = 0$ . Звідси

$$x_k = \cos\left(\frac{(2k-1)\pi}{2n}\right), \quad k = \overline{1, n}. \quad (118)$$

Відповідні коефіцієнти обчислюються за формулами

$$c_k = \int_{-1}^1 \frac{\overline{T}_n(x) dx}{\sqrt{1-x^2} T'_n(x)(x-x_k)} = \frac{\pi}{n}, \quad (119)$$

для  $k = \overline{1, n}$ .

**Означення** Отже, квадратурні формули найвищого степеня точності (ці формули називають формулами Ерміта) мають вигляд

$$I_n(f) = \frac{\pi}{n} \sum_{k=1}^n f(x_k), \quad (120)$$

де  $x_k$  — корені багаточлена Чебишова.

Залишковий член має вигляд

$$R_n(f) = \frac{\pi}{2^{2n-1}(2n)!} \cdot f^{(2n)}(\xi). \quad (121)$$

3. Розглянемо проміжок  $(-\infty, \infty)$  і вагу  $\rho(x) = \exp\{-x^2\}$ , тобто виведемо формулу Гауса для обчислення інтегралу

$$\int_{-\infty}^{\infty} \exp\{-x^2\} f(x) dx. \quad (122)$$

За теорією

$$\omega(x) = H_n(x) = (-1)^n e^{x^2} \frac{d^n}{dx^n} e^{-x^2}, \quad (123)$$

де  $H_n(x)$  — багаточлени Ерміта. Багаточлени Ерміта обчислюються також за рекурентними формулами

$$H_{n+1}(x) = 2xH_n(x) - 2nH_{n-1}(x), \quad (124)$$

з початковими умовами  $H_{-1} = 0, H_0 = 1$ .

Коефіцієнти квадратурної формули обчислюються за формулами

$$c_k = \frac{2^{n+1} n! \sqrt{\pi}}{(H'_n(x_k))^2} \quad (125)$$

Залишковий член

$$R_n(f) = \frac{n! \sqrt{\pi}}{2^n (2n)!} \cdot f^{(2n)}(\xi). \quad (126)$$

**Наприклад:** при  $n = 1, H_1(x) = 2x$ . Корінь  $x_0 = 0$ ,

$$c_0 = \int_{-\infty}^{\infty} \exp\{-x^2\} dx = \sqrt{\pi}. \quad (127)$$

Квадратурна формула має вигляд:

$$I_1(x) = \sqrt{\pi} f(0). \quad (128)$$

4. Розглянемо відрізок  $[0, \infty]$  і ваговий множник  $\rho(x) = x^\alpha e^{-x}$ , тобто виведемо формулу Гауса для обчислення інтегралу

$$\int_0^{\infty} x^\alpha e^{-x} f(x) dx. \quad (129)$$

За теорією

$$\omega(x) = L_n^\alpha(x) = (-1)^n x^{-\alpha} e^x \frac{d^n}{dx^n} (x^{\alpha+n} e^{-x}), \quad (130)$$

де  $L_n^\alpha(x)$  — багаточлени Лагера. Коефіцієнти обчислюються за формулами

$$c_k = \frac{P(n+1)P(n+\alpha+1)}{x_k(L_n^\alpha(x_k))^2}. \quad (131)$$

Залишковий член при  $\alpha = 0$  рівний

$$R_n(f) = \frac{(n!)^2}{(2n)!} \cdot f^{(2n)}(\xi). \quad (132)$$

## 5. Інтегрування швидко осцилюючих функцій.

Розглянемо інтеграл

$$I(f) = \int_a^b f(x)e^{j\omega x} dx, \quad j^2 = -1. \quad (133)$$

При великих  $\omega$  застосування складених квадратурних формул інтерполяційного типу приводить до великої похибки і при малих кроках  $h$ . Розглянемо  $e^{j\omega x}$  як ваговий коефіцієнт, тобто  $\rho(x) = e^{j\omega x}$ . Замінімо  $[a, b]$  на  $[-1, 1]$ :

$$x_i = \frac{a+b}{2} + \frac{b-a}{2} \cdot d_j, \quad (134)$$

де  $d_j \in [-1, 1]$ , а  $i = \overline{1, n}$ .

**Зауваження:** вузли можуть бути не рівновіддалені, якщо рівновіддалені, то  $d_i = -1 + 2i/2n$ ,  $i = \overline{1, n}$ .

Замінімо  $f(x)$  на інтерполяційний багаточлен Лагранжа  $L_{n-1}(x)$  з вузлами  $x_i$  і отримаємо формулу

$$I_n(f) = \int_a^b L_{n-1}(x)e^{j\omega x} dx, \quad (135)$$

яка буде точною для всіх багаточленів не вище  $n - 1$  степеня. Тобто, якщо в (135) підставити багаточлен Лагранжа, то можна обчислити інтеграл і отримати квадратурну формулу

$$S_n^\omega(f) = \frac{b-a}{2} \cdot \exp\left\{j\omega \cdot \frac{a+b}{2}\right\} \cdot \sum_{i=1}^n D_i\left(\omega \cdot \frac{b-a}{2}\right) f(x_i), \quad (136)$$

де

$$D_i(p) = \int_{-1}^1 \left( \prod_{\substack{k=1 \\ k \neq i}}^n \frac{\xi - d_k}{d_i - d_k} \right) \exp\{jp\xi\} d\xi. \quad (137)$$

При  $n = 3$ ,  $d_1 = -1$ ,  $d_2 = 0$ ,  $d_3 = 1$  це формула Філона. Можна брати і більше точок, наприклад,  $n = 5$ ,  $d_1 = -1$ ,  $d_2 = -1/2$ ,  $d_3 = 0$ ,  $d_4 = 1/2$ ,  $d_5 = 1$ .

Ці формули не потрібно застосовувати, коли немає швидко осцилюючого множника.

## 9.8. Обчислення невластних інтегралів

Розглянемо обчислення інтегралів з такими особливостями:



## 1. інтеграли другого роду, тобто

$$I = \int_a^b F(x) dx, \quad F(x) \xrightarrow{x \rightarrow a \vee x \rightarrow b} \infty; \quad (138)$$

## 2. інтеграли першого роду

$$I = \int_a^\infty F(x) dx. \quad (139)$$

Розглянемо спочатку інтеграли другого роду, тобто:

$$I = \int_a^b F(x) dx, \quad F(x) \xrightarrow{x \rightarrow a \vee x \rightarrow b} \infty. \quad (140)$$

1. *Мультиплікативний спосіб*. Представимо підінтегральну функцію у вигляді  $F(x) = \rho(x)f(x)$ , причому  $\rho(x)$  — особлива, а  $f(x)$  — гладка. Далі для обчислення інтегралу

$$I = \tilde{I}(f) = \int_a^b \rho(x)f(x) dx \quad (141)$$

використовуємо відповідні квадратурні формули Гаусса.

**Приклад 1:** Потрібно обчислити інтеграл

$$I = \int_{-1}^1 \frac{dx}{\sqrt{1-x^4}}. \quad (142)$$

*Розв'язок.* Точки  $x = \pm 1$  є особливими.

Представимо підінтегральну функцію у вигляді:

$$F(x) = \underbrace{\frac{1}{\sqrt{1-x^2}}}_{\rho(x)} \cdot \underbrace{\frac{1}{\sqrt{1+x^2}}}_{f(x)} \quad (143)$$

отримаємо інтеграл вигляду

$$I = \int_0^1 \frac{1}{\sqrt{1+x^2}} \frac{dx}{\sqrt{1-x^2}}, \quad (144)$$

де  $\rho(x) = 1/\sqrt{1-x^2}$ .

Далі використовуємо квадратурну формулу Ерміта (120) з попереднього пункту і обчислюємо інтеграл.

**Приклад 2:** Обчислити інтеграл

$$I = \int_0^{\pi} \ln(\sin(x)) \, dx. \quad (145)$$

Розв'язок. Особливі точки  $x = 0, x = \pi$ .

Зведемо цю особливість до степеневій:

$$\rho(x) = \frac{1}{\sqrt{x(\pi - x)}}, \quad (146)$$

тоді

$$f(x) = \sqrt{x(\pi - x)} \cdot \ln(\sin(x)) \xrightarrow{x \rightarrow 0, \pi} 0. \quad (147)$$

Для знаходження інтегралу з таким  $\rho(x)$  застосовуємо квадратурні формули Чебишова. Неприємності виникають, оскільки  $f'(x) \xrightarrow{x \rightarrow 0, \pi} \infty$  (хоча квадратурні формули даватимуть наближене значення).

Тому застосовують другий спосіб розв'язання проблеми:

2. **Адитивний спосіб.** Представимо підінтегральну функцію у вигляді  $F(x) = f(x) + \psi(x)$ , причому  $\psi(x)$  — особлива,  $f(x)$  — гладка. Розбиваємо інтеграл на два:  $I = I_1 + I_2$ .

1.  $I_1 = \int_a^b f(x) \, dx$  — обчислюють чисельно (наприклад, формули Сімпсона чи трапецій),

2.  $I_2 = \int_a^b \psi(x) \, dx$  — пробують обчислити аналітично (можливо апроксимувати функцію  $\psi(x)$ , наприклад, рядом).

**Приклад 3:** Обчислити інтеграл з прикладу 2.

$$I = \int_0^{\pi} \ln(\sin(x)) \, dx, \quad (148)$$

$$I = 2 \int_0^{\pi/2} \ln(\sin(x)) \, dx. \quad (149)$$

Представимо

$$\ln(\sin(x)) = \ln\left(\frac{\sin x}{x}\right) + \ln(x). \quad (150)$$

Отримаємо два інтеграли:

1.  $I_1 = \int_a^b \ln\left(\frac{\sin x}{x}\right) \, dx$  обчислюємо чисельно,

$$2. I_2 = \int_0^{\pi/2} \ln(x) dx = \frac{\pi}{2} \left( \ln\left(\frac{\pi}{2}\right) - 1 \right).$$

Розглянемо тепер інтеграли першого роду

$$I = \int_a^{\infty} F(x) dx. \quad (151)$$

### 1. Заміни.

- Нехай  $a > 0$ . Зробимо заміну  $t = \frac{x-a}{x}$ ,  $x = \frac{a}{1-t}$ . Тоді

$$I = a \int_0^1 F\left(\frac{a}{1-t}\right) \frac{dt}{(1-t)^2}, \quad (152)$$

а це інтеграл другого роду.

- Якщо  $a = 0$ , то робимо заміну  $t = e^{-x}$ ,  $x = -\ln(t)$ , тоді

$$I = \int_0^1 F(-\ln(t)) \cdot \frac{dt}{t}. \quad (153)$$

Знову отримуємо інтеграл другого роду.

- Якщо  $a < 0$  (не можна зробити заміну  $t = \frac{x-a}{x}$ , тому що виникає особливість в точці  $x = 0$ ), розбиваємо інтеграл на два:

$$I = \int_a^0 F(x) dx + \int_0^{\infty} F(x) dx \quad (154)$$

і обчислюємо за допомогою попередніх пунктів.

### 2. Мультиплікативний спосіб обчислення інтегралів першого роду ґрунтується на представленні підінтегральної функції у вигляді

$$F(x) = \rho(x) f(x), \quad (155)$$

де, наприклад,

$$\rho(x) = x^\alpha \cdot e^{-x}, \quad x \in [0, \infty). \quad (156)$$

Такий ваговий коефіцієнт відповідає багаточленам Лагера. При  $x \in (-\infty, \infty)$ ,  $\rho(x) = \exp\{-x^2\}$  приходимо до багаточленів Ерміта.

### 3. Обрізання границі. Ще один спосіб ґрунтується на обрізанні верхньої границі. Для цього інтеграл запишемо у вигляді

$$I = \int_a^{\infty} F(x) dx = \int_a^b F(x) dx + \int_b^{\infty} F(x) dx. \quad (157)$$

Верхня границя  $b$  обчислюють з умови

$$\left| \int_b^{\infty} F(x) dx \right| < \frac{\varepsilon}{2}, \quad (158)$$

де  $\varepsilon$  — задана точність. Для обчислення  $\int_a^b F(x) dx$  використовують квадратурні формули складеного типу.

## 9.9. Обчислення кратних інтегралів

Розглянемо інтеграл

$$I = \int_a^b \int_c^d f(x, y) dx dy. \quad (159)$$

Цей інтеграл зводиться до повторного, якщо ввести

$$F(x) = \int_c^d f(x, y) dy. \quad (160)$$

Тоді

$$I = \int_a^b F(x) dx. \quad (161)$$

До однократних інтегралів можна застосувати квадратурну формулу середніх прямокутників. Тоді

$$I \approx I_0 = (b - a)F(\bar{x}) \int_c^d f(\bar{x}, y) dy \approx (b - a)(d - c)f(\bar{x}, \bar{y}), \quad (162)$$

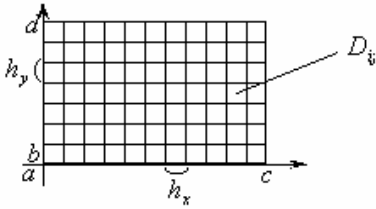
де

$$\bar{x} = \frac{a + b}{2}, \quad \bar{y} = \frac{c + d}{2}. \quad (163)$$

*Кубатурна формула* (це формула наближеного обчислення інтегралу (159)), якщо використовується формула трапеції має вигляд

$$I \approx I_1 = \frac{(b - a)(d - c)}{4} \cdot (f(a, c) + f(b, c) + f(a, d) + f(b, d)). \quad (164)$$

Точність залежить від поведінки функції та від довжини інтервалів  $[a, b]$ ,  $[c, d]$ . Аналог формул складеного типу для (159) будується таким чином: розбиваємо  $D$  на комірки:



Тоді

$$D = \bigsqcup_{i,j} D_{i,j}, \quad (165)$$

де  $D_{i,j} = \{x_{i-1} \leq x \leq x_i, y_{j-1} \leq y \leq y_j\}$ , а у свою чергу  $x_i = a + ih_x$  для  $i = \overline{0, N_x}$ , і  $y_j = c + jh_y$  для  $j = \overline{0, N_y}$ , і нарешті  $h_x = \frac{b-a}{N_x}$  і  $h_y = \frac{d-c}{N_y}$ .

Тоді для кожного інтегралу по комірці застосовуємо кубатурну формулу прямокутників (164):

$$I = \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} \iint_{D_{i,j}} f(x, y) dx dy \approx I_{0,h} = \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} h_x h_y f(\bar{x}_i, \bar{y}_j), \quad (166)$$

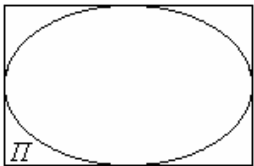
де  $\bar{x}_i = x_i - h_x/2$ , а  $\bar{y}_j = y_j - h_y/2$ .

Якщо  $f(x, y) \in C^{(2)}(D)$ , то  $I - I_{0,h} = O(h_x^2 + h_y^2)$ . Степінь точності по крокам сітки — 2. Складність методу пропорційна кількості комірок:  $Q = O(N_x N_y) = O(N^2)$ , якщо  $N \approx N_x \approx N_y$ . В 3-вимірному просторі  $f = f(x, y, z)$  складність —  $Q = O(N^3)$ .

Якщо  $D$  не прямокутник, то замість  $f$  введемо

$$\bar{f}(x, y) = \begin{cases} f(x, y), & x \in D, \\ 0, & x \in \Pi \setminus D, \end{cases} \quad (167)$$

де  $\Pi$  — найменший охоплюючий  $D$  прямокутник  $D \subset \Pi$ :



Тоді

$$I = \iint_{\Pi} \bar{f}(x, y) dx dy, \quad (168)$$

що розглядався вище.

Недолік такого підходу:  $f(x, y)$  може бути розривною функцією і через це низька точність обчислення інтегралу.

Наступний підхід базується на відповідній заміні змінних

$$x = x(\xi, \eta), \quad y = y(\xi, \eta), \quad (169)$$

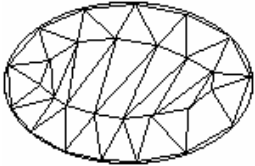
такій, що  $D \mapsto \Pi$ , тоді

$$I = \iint_{\Pi} f(x(\xi, \eta), y(\xi, \eta)) J(\xi, \eta) d\xi d\eta, \quad (170)$$

де  $\Pi$  — прямокутник в площині  $(\xi, \eta)$ ;  $J(\xi, \eta)$  — Якобів переходин. Якщо границя області  $D$  гладка, то якобів буде мати особливості, що також знижує швидкість збіжності.

Ще один підхід в обчисленні інтегралу по довільній області  $D$  базується на триангулюванні області. Якщо область довільного вигляду, то її можна розбити на трикутники таким чином:

$$D = \bigsqcup_{k=1}^N D_k. \quad (171)$$



Тоді

$$I = \sum_{k=1}^N I_k, \quad I_k = \iint_{D_k} f(x, y) dx dy \quad (172)$$

Застосуємо кубатурні формули до кожного  $I_k \approx I_k^h$ . Для цього замінимо функцію поліномом першого степеня

$$f(x, y) \approx L_1(x, y) = A + Bx + Cy. \quad (173)$$

**Задача 35:** Побудувати явний вигляд кубатурної формули, яка дозволяє наближено обчислити  $I_k$  по трикутнику  $D_k$ , якщо замінити  $f(x, y) \approx L_1(x, y)$  на інтерполяційний багаторилен 1-го степеня.

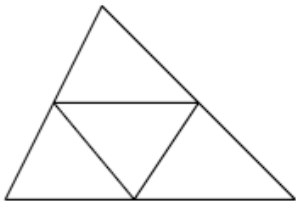
Точність такого підходу

$$I - I^h = I - \sum_{k=1}^N I_k^h = O(h^2), \quad (174)$$

де  $h = \max_k \text{diam } D_k$ .

Складність обчислення інтеграла:  $Q = O(h^{-2})$  для  $D \subset \mathbb{R}^2$ ;  $Q = O(h^{-3})$  для  $D \subseteq \mathbb{R}^3$ .

Можна згустити сітку, поділивши один трикутник  $D_k$  на чотири:



І, нарешті, розглянемо ідею методу статистичних випробувань (метод Монте-Карло) для обчислення інтегралу

$$I = \iint_D f(x, y) \, dx \, dy. \quad (175)$$

Замінімо наближено

$$I = \iint_{\Pi} \bar{f}(x, y) \, dx \, dy \approx \frac{1}{N} \sum_{i=1}^N \bar{f}(\xi_i, \eta_i) \cdot \mu(\Pi_i), \quad (176)$$

де  $\mu$  — міра нп відповідному просторі,  $\Pi$  — найменший охоплюючий  $D$  прямокутник  $D \subset \Pi$ ;  $f(x, y)$  — продовження функції  $f$ ;  $\xi_i, \eta_i$  — незалежні реалізації рівномірно розподілених на  $[a, b]$  та  $[c, d]$  випадкових величин  $\xi$  та  $\eta$ . Складність цього методу  $Q = O(N)$ . Оцінка точності —  $I - I_N = O\left(1/\sqrt{N}\right)$  носить ймовірністний характер.

- **Позитивна** сторона методу — незалежна від розмірності складність;
- **негативна** — низька точність.

## 10. Чисельні методи розв'язання задачі Коші для звичайних диференціальних рівнянь

Постановка задачі: нехай потрібно знайти розв'язок диференціального рівняння з початковими умовами

$$\frac{du}{dt} = f(t, u), \quad t > t_0, \quad u(t_0) = u_0, \quad (1)$$

де  $u = (u_1, \dots, u_m)^\top$ ,  $f = (f_1, \dots, f_m)^\top$ ,  $u_k = u_k(t)$ ,  $f_k(t; u_1, \dots, u_m)$ .

Якщо довільна функція  $f_k$  неперервна по кожній своїй змінній та по  $u_j$  вони Ліпшиць-неперервні, тобто

$$|f_k(t; \dots, u_j, \dots) - f_k(t; \dots, v_j, \dots)| \leq L_j |u_j - v_j|, \quad \forall j, k, \quad (2)$$

то розв'язок задачі (1) існує і єдиний.

Нехай задано рівняння  $m$ -ого порядку та початкові умови:

$$\begin{cases} v^{(m)}(t) = F(t; v, v', \dots, v^{(m-1)}), \\ t = t_0 : v = v_1, v' = v_2, \dots, v^{(m-1)} = v_m. \end{cases} \quad (3)$$

Введемо компоненти вектора  $\vec{u}$ :  $u_k(t) = v^{(k-1)}(t)$ . Тоді задача (2) записується у вигляді системи (1):

$$\begin{cases} \frac{du_1}{dt} = u_2, & u_1(t_0) = v_1, \\ \dots\dots\dots, \\ \frac{du_{m-1}}{dt} = u_m, & u_{m-1}(t_0) = v_{m-1}, \\ \frac{du_m}{dt} = F(t; u_1, \dots, u_m), & u_m(t_0) = v_m. \end{cases} \quad (4)$$

Тому далі, в основному, розглядаються методи розв'язання задачі (1).

### 10.1. Наближені аналітичні методи

1. Метод послідовних наближень (метод Пікара):

Потрібно розв'язати диференціальне рівняння з відповідними початковими умовами:

$$\frac{du}{dt} = f(x, u), \quad u(t_0) = u_0. \quad (5)$$

Проінтегруємо (5)

$$u(t) = u(t_0) + \int_{t_0}^t f(\xi, u(\xi)) d\xi. \quad (6)$$

Задаємо  $u^{(0)}(t)$  і запишемо ітераційний процес

$$u^{(k+1)}(t) = u_0 + \int_{t_0}^t f(\xi, u^{(k)}(\xi)) d\xi, \quad k = 0, 1, \dots \quad (7)$$



Існує  $T = T(u_0, L)$  ( $L = \max_j L_j$  — стала Ліпшиця) така, що

$$u^{(n)}(t) \xrightarrow{n \rightarrow \infty} u(t), \quad t \in [t_0, T]. \quad (8)$$

Тому  $u(t) \approx u^{(n)}(t)$ .

**Недолік** методу: необхідно проведення аналітичного інтегрування.

## 2. Метод рядів Тейлора:

Нехай розв'язок задачі (5) можна представити у вигляді ряду

$$u(t) = \sum_{k=0}^{\infty} \frac{u^{(k)}(t_0)}{k!} \cdot (t - t_0)^k. \quad (9)$$

Будемо шукати наближення у вигляді скінченної суми:

$$u(t) \approx u^N(t) = \sum_{k=0}^N \frac{u^{(k)}(t_0)}{k!} \cdot (t - t_0)^k, \quad t \in [t_0, t_1]. \quad (10)$$

Для визначення  $u^{(k)}(t_0)$  диференціюємо рівняння (5) по  $t$ :

$$u^{(0)}(t_0) = u_0, \quad u^{(1)}(t_0) = f(t_0, u_0), \quad (11)$$

$$u^{(2)}(t_0) = f_{t,0} + f_{u,0}f_0, \quad u^{(3)}(t_0) = f_{tt,0} + 2f_{tu,0}f_0 + f_{uu,0}f_0^2, \quad (12)$$

і так далі.

Якщо  $\tau = t_1 - t_0$  малий параметр, то:

$$|u(t) - u^N(t)| = O(\tau^{N+1}). \quad (13)$$

**Недоліки** методу:

- зростання кількості доданків при обчисленні  $u^{(n)}(t_0)$ ;
- необхідно аналітичного диференціювання.

## 10.2. Методи типу Ейлера

Розглянемо задачу Коші:

$$\frac{du}{dt} = f(t, u), \quad u(t_0) = u_0. \quad (14)$$

Використаємо перше наближення за допомогою рядів Тейлора на проміжку  $t_0 \leq t \leq t_1$ :

$$u(t) \approx u^1(t) = u(t_0) + (t - t_0)f(t_0, u_0). \quad (15)$$

Обчислимо наближене значення в точці  $t_1$ :

$$u(t_1) = u_1 \approx u^1(t_1) = u_0 + \tau f(t_0, u_0) \quad (16)$$

де  $\tau = t_1 - t_0$  — деякий крок. Якщо позначити  $y_1 = u(t_1)$ ,  $y_0 = u(t_0)$ , то маємо формулу

$$y_1 = y_0 + \tau f(t_0, y_0). \quad (17)$$

Застосовуючи такий підхід для  $t_n \leq t \leq t_{n+1}$ , отримаємо рекурентну формулу

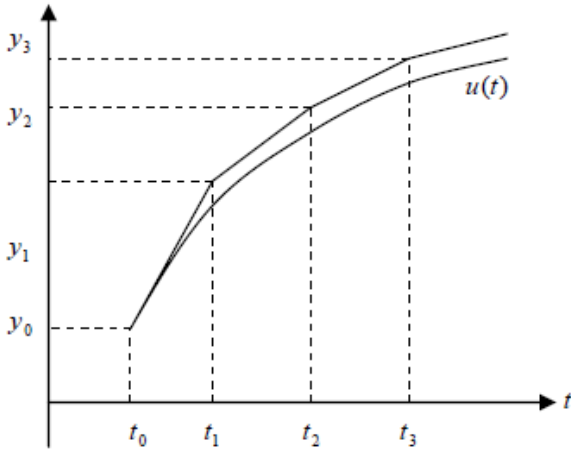
$$y_{n+1} = y_n + \tau f(t_n, y_n), \quad (18)$$

для  $n = 0, 1, \dots$  з початковими умовами  $y_0 = u_0$ .

Це формула методу Ейлера. Крок інтегрування може змінюватися:

$$\tau = \tau_n = t_{n+1} - t_n. \quad (19)$$

Геометрична інтерпретація методу Ейлера представлена на рис:



Його друга назва — метод ламаних.

Цей метод відноситься до однокрокових, тобто розв'язок на наступному кроці обчислюється тільки по одному значенню на попередньому кроці.

Загальна формула однокрокових методів:

$$y_{n+1} = y_n + \tilde{\Phi}(t_n, y_n, f). \quad (20)$$

Для методу Ейлера  $\tilde{\Phi} = \tau f(t_n, y_n)$ .

**Означення:** Величина  $R(\tau) = y(t_{n+1}) - u(t_{n+1})$ , де  $y_{n+1}$  обчислюється за формулою (20), причому  $y_n = u(t_n)$  називається *похибка методу (20) на одному кроці*.

Загальна похибка на  $(n + 1)$ -му кроці складається з похибок на всіх попередніх кроках.

**Означення:** Величина  $z_{n+1} = y_{n+1} - u(t_{n+1})$ , де  $y_{n+1}$  і всі попередні  $y_k$ ,  $k = 1, 2, \dots$  також наближені, називається *похибкою методу (20)*.

Якщо виразити  $y_{n+1} = u(t_{n+1}) + z_{n+1}$  і підставити в (20), то

$$u(t_{n+1}) + z_{n+1} = u(t_n) + z_n + \tau \tilde{\Phi}(t_n, u_n + z_n). \quad (21)$$

Тоді можна записати рівняння для  $z_n$ :

$$z_{n+1} = z_n + \tau \left( -\frac{u_{n+1} - u_n}{\tau} + \tilde{\Phi}(t_n, u_n) \right) + \tau (\tilde{\Phi}(t_n, u_n + z_n) - \tilde{\Phi}(t_n, u_n)). \quad (22)$$

**Означення:** Величина  $\psi_n = -\frac{u_{n+1} - u_n}{\tau} + \tilde{\Phi}(t_n, u_n)$  називається *похибкою апроксимації методу (20)*.

Для методу Ейлера

$$\begin{aligned}\psi_n &= -\frac{u_{n+1} - u_n}{\tau} + f(t_n, u_n) = \\ &= -\frac{1}{\tau} \left( u_n + \tau \cdot \frac{du}{dt}(t_n) + O(\tau^2) - u_n \right) + f(t_n, u_n) = \\ &= -\frac{du}{dt}(t_n) + f(t_n, u_n) + O(\tau) = O(\tau) \xrightarrow{\tau \rightarrow 0} 0.\end{aligned}\quad (23)$$

Похибка апроксимації це нев'язка, коли замість  $y$  в різницеве рівняння (20) підставляємо точний розв'язок задачі Коші  $u$ .

**Означення:** Метод (20) має похибку на одному кроці степеня  $m$ , якщо  $|R_n(\tau)| = O(\tau^{m+1})$ .

**Означення:** Кажуть, що чисельний метод має похибку апроксимації степеня  $p$ , якщо  $|\psi_n| = O(\tau^p)$ .

Для методу Ейлера похибка на одному кроці

$$\begin{aligned}R_n(\tau) &= y_{n+1} - u(t_{n+1}) = u(t_n) + \tau f(t_n, u(t_n)) - u(t_n + \tau) = \\ &= u(t_n) + \tau f(t_n, u(t_n)) - u(t_n) - \tau y'(t_n) + O(\tau^2) = O(\tau^2).\end{aligned}\quad (24)$$

Отже  $m = 1$ . Маємо зв'язок похибок:  $\psi(\tau) = R(\tau)/\tau$  і тому  $p = 1$ .

**Означення:** Кажуть, що метод (20) має степінь точності  $m$ , якщо

$$\forall n : x_n = y_n - u(t_n) = O(\tau^m).\quad (25)$$

**Теорема:** Нехай  $f(y, u) \in C^{(1)}(\overline{D}_T)$ , та  $|f_u(t, u)| \leq L$ , де  $\overline{D}_T = \{t_0 < t \leq T, |u(t)| \leq K\}$ . Тоді метод Ейлера (20) має степінь точності  $m = 1$  і для нього має місце оцінка

$$|z_n| \leq M \cdot \max_j |\psi_j| = O(\tau),\quad (26)$$

де  $M = M(T, L)$ .

**Доведення:** Для  $z$  маємо задачу

$$z_{n+1} = z_n + \tau \psi_n + \tau(f(t_n, u_n + z_n) - f(t_n, u_n)).\quad (27)$$

Оцінімо  $z_{n+1}$ :

$$\begin{aligned}|z_{n+1}| &\leq |z_n| + \tau |\psi_n| + \tau L |z_n| + \tau |\psi_n| \leq \dots \leq (1 + \tau L)^n |z_0| + \sum_{j=0}^n \tau (1 + \tau L)^{n-j} |\psi_j| \leq \\ &\leq (1 + \tau L)^n \max_j |\psi_j| \cdot \sum_{j=0}^n \tau \leq T e^{L\tau} \max_j |\psi_j| \leq M \max_j |\psi_j| = O(\tau),\end{aligned}\quad (28)$$

адже

$$(1 + \tau L)^n = (1 + \tau L)^{n\tau/\tau} = (1 + \tau L)^{t_n/\tau} \leq (1 + \tau L)^{T/\tau} \leq e^{LT}.\quad (29)$$

Позначимо  $T e^{LT} = M$  і отримуємо бажану оцінку.  $\square$

Таким чином порядок точності методу Ейлера  $m = 1$ .

Метод Ейлера можна вивести із таких міркувань. Інтегруємо (14) по  $t$ :  $t_n < t < t_{n+1}$ :

$$u(t_{n+1}) = u(t_n) + \int_{t_n}^{t_{n+1}} f(t, u(t)) dt. \quad (30)$$

Застосуємо формулу лівих прямокутників для інтегралу:

$$\int_{t_n}^{t_{n+1}} f(t, u(t)) dt \approx \tau f(t_n, u(t_n)) \quad (31)$$

і підставимо в (30). Отримаємо формулу для явного методу Ейлера:

$$y_{n+1} = y_n + \tau f(t_n, u(t_n)) \quad (32)$$

Застосуємо формулу правих прямокутників для інтегрування

$$\int_{t_n}^{t_{n+1}} f(t, u(t)) dt \approx \tau f(t_{n+1}, u(t_{n+1})) \quad (33)$$

і підставимо в (30). Отримаємо формулу для неявного методу Ейлера:

$$y_{n+1} = y_n + \tau f(t_{n+1}, y_{n+1}). \quad (34)$$

Ці формули 1-го степеня точності по кроку  $\tau$ . Заміняючи інтеграл за квадратурною формулою трапеції:

$$\int_{t_n}^{t_{n+1}} f(t, u(t)) dt \approx \frac{\tau}{2} (f(t_n, u(t_n)) + f(t_{n+1}, u(t_{n+1}))) \quad (35)$$

ми отримаємо формулу методу трапеції інтегрування задачі Коші:

$$y_{n+1} = y_n + \frac{\tau}{2} (f(t_n, u(t_n)) + f(t_{n+1}, u(t_{n+1}))). \quad (36)$$

Це неявний метод.

Обчислимо похибку апроксимації цього методу:

$$\begin{aligned} \psi_n &= -\frac{u_{n+1} - u_n}{\tau} + \frac{1}{2} (f(t_n, u_n) + f(t_{n+1}, u_{n+1})) = \\ &= -\frac{1}{\tau} \left( u_{n+1/2} + \frac{\tau}{2} \cdot \dot{u}_{n+1/2} + \frac{1}{2} \left( \frac{\tau}{2} \right)^2 \ddot{u}_{n+1/2} + O(\tau^3) - \right. \\ &\quad \left. - u_{n+1/2} + \frac{\tau}{2} \cdot \dot{u}_{n+1/2} - \frac{1}{2} \left( \frac{\tau}{2} \right)^2 \ddot{u}_{n+1/2} + O(\tau^3) \right) + \\ &\quad + \frac{1}{2} \left( f_{n+1/2} + \frac{\tau}{2} \cdot \dot{f}_{n+1/2} + O(\tau^2) + f_{n+1/2} - \frac{\tau}{2} \cdot \dot{f}_{n+1/2} + O(\tau^2) \right) = \\ &= -\dot{u}_{n+1/2} + O(\tau^2) + f_{n+1/2} = O(\tau^2). \end{aligned} \quad (37)$$

Таким чином метод трапецій має другий порядок апроксимації.

**Задача 36:** Показати, що похибка на одному кроці методу трапецій є величина порядку  $O(\tau^3)$ .

Таким чином ми отримали більш точний метод. Недолік методу трапецій — неявність (треба розв'язувати нелінійне рівняння). Для розв'язання проблеми неявності застосуємо метод предиктор-корректор (предиктор — попередник, коректор — уточнювач).

Обчислимо попередньо за явним методом Ейлера

$$\bar{y}_{n+1} = y_n + \tau f(t_n, y_n). \quad (38)$$

Уточнимо за методом трапецій

$$y_{n+1} = y_n + \frac{\tau}{2} (f(t_n, y_n) + f(t_{n+1}, \bar{y}_{n+1})). \quad (39)$$

Формули (38) та (39) утворюють *метод Ейлера-Коші*.

Оцінимо похибку на одному кроці

$$R(\tau) = y_1 + u(t_0 + \tau) - u_0 - \frac{\tau}{2} (f(t_0, y_0) + f(t_0 + \tau, u_0 + \tau f(t_0, u_0))) - u(t_0 + \tau). \quad (40)$$

Звідси  $R(0) = 0$ . Далі

$$\begin{aligned} R'(\tau) = & \frac{1}{2} (f(t_0, y_0) + f(t_0 + \tau, y_0 + \tau f_0)) - \\ & - \frac{du}{dt}(t_0 + \tau) + \frac{\tau}{2} \cdot f_t(t_0 + \tau, u_0 + \tau f(t_0, u_0)) + \\ & + \frac{\tau}{2} \cdot f_{tt}(t_0 + \tau, u_0 + \tau f_0(t_0, u_0)) \cdot f(t_0, u_0), \end{aligned} \quad (41)$$

$$\text{тому } R'(0) = \frac{f_0 + f_0}{2} - \frac{du}{dt}(t_0) = 0.$$

Наступна похідна

$$R''(\tau) = f_t(t_0 + \tau, u_0 + \tau f_0) + f_{tt}(t_0 + \tau, u_0 + \tau f_0) \cdot f_0 + \tau(\dots) - u''(t_0), \quad (42)$$

$$\text{тому } R''(0) = f_{t,0} + f_{u,0} \cdot f_0 - \frac{d^2u}{dt^2}(t_0) = 0,$$

$$\text{оскільки } \frac{d^2u}{dt^2} = \frac{d}{dt} f(t, u) = f_t(\dots) + f_{tt}(\dots) \cdot f_0.$$

І, нарешті,  $R'''(0) = \frac{3}{2} (f_{tt,0} + 2f_{tu,0} + f_{uu} \cdot f_0^2) - u'''(t_0) \neq 0$ . Таким чином похибка на одному кроці має порядок  $p = 2$ .

Ще один метод цього типу:

$$y_{n+1/2} = y_n + \frac{\tau}{2} \cdot f(t_n, y_n), \quad (43)$$

$$y_{n+1} = y_n + \tau f(t_{n+1/2}, y_{n+1/2}). \quad (44)$$

Це модифікований метод Ейлера. На етапі коректор (44) використовуємо формулу середніх прямокутників, а предиктор — це метод Ейлера на інтервалі  $(t_n, t_{n+1/2})$ .

**Задача 37:** Показати, що для модифікованого методу Ейлера  $R(\tau) = O(\tau^3)$ , тобто  $m = 2$ .

### 10.3. Методи типу Рунге-Кутта

Методи типу Ейлера мають низьку точність ( $m = 1, 2$ ). Рунге запропонував метод третього, а Кутта розвинув його ідею та отримав методи четвертого порядку точності.

Розглянемо задачу Коші

$$\frac{du}{dt} = f(t, u), \quad u(t_0) = u_0. \quad (45)$$

Явний  $m$ -етапний (стадійний) метод Рунге-Кутта полягає в наступному.

Нехай розв'язок  $y_n = y(t_n)$  вже відомий. Задаються числові коефіцієнти  $\alpha_i, \beta_{i,j}, i = \overline{2, m}, j = \overline{1, i-1}$  та  $p_i, i = \overline{1, m}$  і послідовно обчислюються за формулами:

$$k_1 = \tau f(t_n, y_n), \quad (46)$$

$$k_2 = \tau f(t_n + \alpha_2 \tau, y_n + \beta_{2,1} \tau k_1), \quad (47)$$

$$k_3 = \tau f(t_n + \alpha_3 \tau, y_n + \beta_{3,1} \tau k_1 + \beta_{3,2} \tau k_2), \quad (48)$$

$$\dots\dots\dots \quad (49)$$

$$k_m = \tau f \left( t_n + \alpha_m \tau, y_n + \sum_{j=1}^{m-1} \beta_{m,j} \tau k_j \right). \quad (50)$$

Потім з формули

$$y_{n+1} = y_n + \sum_{i=1}^m p_i k_i \quad (51)$$

знаходимо нове значення  $y_{n+1} = y(t_{n+1}) \approx u(t_{n+1})$ .

Коефіцієнти  $\alpha_i, \beta_{i,j}, p_i$  вибираємо з міркувань точності. Наприклад, для того, щоб рівняння (51) апроксимувало рівняння (45), необхідно, щоб

$$\sum_{i=1}^m p_i = 1. \quad (52)$$

Інформація про  $m$ -стадійний метод записується в таблиці Батчера:

$\alpha_1 = 0$				
$\alpha_2$	$\beta_{2,1}$			
$\vdots$	$\vdots$	$\ddots$		
$\alpha_m$	$\beta_{m,1}$	$\dots$	$\beta_{m,m-1}$	
	$p_1$	$\dots$	$p_{m-1}$	$p_m$

Похибка на одному кроці:

$$\begin{aligned}
 R(\tau) &= y_1 - u(t_0 + \tau) = u_0 + \sum_{i=1}^m p_i k_i(\tau) - u(t_0 + \tau) = \\
 &= R(0) + R'(0) + \dots + \frac{\tau^p}{p!} \cdot R^{(p)}(0) + O(\tau^{p+1}).
 \end{aligned}
 \tag{53}$$

Для того, щоб  $R(\tau) = O(\tau^{p+1})$ , тобто метод мав  $p$ -й степінь точності необхідно і достатньо, щоб  $R(0) = R'(0) = \dots = R^{(p)}(0) = 0$ . Загального розв'язку цієї нелінійної системи немає, тому розглянемо частинні випадки.

1. Методи першого порядку  $m = 1$ . Невідоме  $p_1$ :

$$k_1(\tau) = \tau f(t_0, y_0), \tag{54}$$

$$R(\tau) = u_0 + p_1 k_1(\tau) - u(t_0 + \tau), \tag{55}$$

$$R(0) = u_0 - u(t_0) = 0, \tag{56}$$

$$R'(\tau) = p_1 k'_1(\tau) - \dot{u}(t_0 + \tau) = p_1 f(t_0, y_0) - \dot{u}(t_0 + \tau), \tag{57}$$

$$R'(0) = (p_1 - 1)f_0 = 0 \tag{58}$$

Звідси  $p_1 = 1$ .

Ясно, що  $R''(0) \neq 0$ . Тому  $R(\tau) = O(\tau^2)$ . Отримуємо *явний метод Ейлера*

$$y_{n+1} = y_n + \tau f(t_n, y_n). \tag{59}$$

2. Методи другого порядку  $m = 2$ . Тут отримуємо сімейство методів.

Невідомі  $p_1, p_2, \alpha_2, \beta_{2,1}$ . Вони вибираються з умови:  $R(\tau) = O(\tau^3)$ ,  $p = 2$ . Маємо

$$p_1 + p_2 = 1, \quad 2\alpha_2 p_2 = 1, \quad 2\beta_{2,1} p_2 = 1, \tag{60}$$

і один параметр залишається вільним. Далі один параметр  $p_1$  фіксуємо і отримуємо конкретний метод:

1.  $p_1 = 0, p_2 = 1, \alpha_2 = \beta_{2,1} = 1/2$  — *модифікований метод Ейлера*:

$$y_{n+1} = y_n + \tau f\left(t_{n+1/2}, y_n + \frac{\tau f_n}{2}\right) \tag{61}$$

2.  $p_1 = 1/2, p_2 = 1/2, \alpha_2 = \beta_{2,1} = 1$  — *метод Ейлера-Коші*:

$$y_{n+1} = y_n + \frac{\tau}{2} (f(t_n, y_n) + f(t_{n+1}, y_n + \tau f_n)). \tag{62}$$

3.  $p_1 = 1/3, p_2 = 2/3, \alpha_2 = \beta_{2,1} = 3/4$  — ще один метод другого порядку точності:

$$y_{n+1} = y_n + \frac{\tau}{3} \left( f(t_n, y_n) + 2f\left(t_{n+3/4}, y_n + \frac{3f_n}{4}\right) \right). \tag{63}$$

3. Методи третього порядку  $m = 3$ .  $R(\tau) = O(\tau^4)$ . Тому  $p = 3$ . Запишемо результат вибору параметрів у вигляді таблиці Батчера (один з частинних випадків):

0			
1/3	1/3		
2/3	0	2/3	

	1/4	0	3/4
--	-----	---	-----

**Задача 38:** Довести, що метод, якому відповідає ця таблиця Батчера, має третій степінь точності.

4. Методи четвертого порядку  $m = 4$ ,  $R(\tau) = O(\tau^5)$ ,  $p = 4$ . Найбільш поширені методи:

0				
1/2	1/2			
1/2	0	1/2		
1	1/2	0	1	
	1/6	1/3	1/3	1/6

i

0				
1/4	1/4			
1/2	0	1/2		
1	1	-2	2	
	1/6	0	2/3	1/6

Розглянемо як пов'язаний порядок методу  $p$  і кількість стадій  $m$ :

$m$	1	2	3	4	5	6
$p$	1	2	3	4	4	5

Ясно чому методи Рунге-Кутта при  $m > 4$  рідко застосовуються: збільшення кількості стадій (а значить збільшення об'єму обчислювальної роботи) не дає виграшу в точності.

**Теорема:** Нехай  $m$ -стадійний метод Рунге-Кутта має  $p$ -й степінь точності на кроці, а  $f(t, u)$  задовольняє умову Ліпшиця. Тоді метод має  $p$ -й степінь точності і для похибки методів Рунге-Кутта має місце така оцінка:

$$|z_n| = |y_n - u(t_n)| \leq T \cdot e^{\alpha_0 T} \cdot \max_j |\psi_j| = O(\tau^p), \quad (64)$$

де  $\psi_j = R_j(\tau)/\tau$  — похибка апроксимації метода,  $\alpha_0 = pLm(1 + BL\tau_0)^{m-1}$ ,  $L$  — стала Ліпшиця,  $B = \max_{i,j} |f_{ij}|$ ,  $T = \max_n |t_n|$ .

## 10.4. Методи з контролем точності на кроці

Часто в ході розрахунків необхідно змінювати крок інтегрування, контролюючи величину похибки методу на кроці. При практичній оцінці цієї величини можна, наприклад, поступати так.



Перший підхід використовує принцип Рунге. Головний член похибки на одному кроці  $(t_n, t_{n+1})$  інтегрування має вигляд

$$\frac{\varphi^{(p+1)}(0)\tau^{p+1}}{(p+1)!}. \quad (65)$$

В результаті двох кроків інтегрування однокроковим методом, наприклад, методом Рунге-Кутта  $p$ -го степеня точності, буде отримано наближення  $y^\tau(t+2\tau)$  до значення  $u(t_n+2\tau)$  таке, що

$$y^\tau(t_n+2\tau) - u(t_n+2\tau) \approx 2 \frac{\varphi^{(p+1)}(0)\tau^{p+1}}{(p+1)!} = 2\overset{\circ}{R}(\tau). \quad (66)$$

Якщо тепер застосувати метод Рунге-Кутта  $p$ -го степеня з одним кроком довжини  $2\tau$  на інтервалі  $(t_n, t_{n+2})$ , то отримаємо наближене значення  $y^{2\tau}(t_n+2\tau)$ , для якого

$$y^{2\tau}(t_n+2\tau) - u(t_n+2\tau) \approx \frac{\varphi^{(p+1)}(0)(2\tau)^{p+1}}{(p+1)!} = \overset{\circ}{R}(2\tau) = 2^{p+1}\overset{\circ}{R}(\tau). \quad (67)$$

З цих співвідношень випливає представлення головного члена похибки на кроці

$$y^\tau(t_n+2\tau) - y(t_n+2\tau) \approx 2\overset{\circ}{R}(\tau) = \frac{y_{n+1}^{2\tau} - y_{n+1}^\tau}{2^p - 1}. \quad (68)$$

При необхідності можна уточнити отримане наближене значення, додавши до нього величину головного члена похибки, тобто покласти

$$u(t_n+2\tau) \approx y^\tau + \frac{y^\tau - y^{2\tau}}{2^p - 1}. \quad (69)$$

Позначимо

$$g(\tau) = \frac{y_{n+1}^{2\tau} - y_{n+1}^\tau}{2^p - 1} \quad (70)$$

Якщо  $|g(\tau)| \leq \varepsilon$ , де  $\varepsilon$  — деяка задана мала величина (похибка на одному кроці), то  $\tau$  — успішний крок і  $y_{n+1}^\tau \approx u(t_n+2\tau)$ . Якщо  $|g(\tau)| > \varepsilon$ , то зменшуємо крок  $t := \tau/2$ . Треба ще передбачити умову збільшення кроку. Задаємо деяке  $\delta \ll \varepsilon$  і якщо  $|g(\tau)| \leq \delta$ , то  $y_{n+1}^\tau \approx u(t_n+2\tau)$  і далі беремо  $\tau := 2\tau$ . Параметр  $\delta$  вибирають, наприклад, так:  $\delta = \alpha \cdot \varepsilon / 2^p$ ,  $0 < \alpha < 1$ , де  $p$  — порядок точності методу.

Інший підхід вибору кроку інтегрування заключається в використанні методів різного степеня точності. Отже, якщо в нас є два методи степеня точності на кроці  $p$  та  $p+1$ :

$$y_{n+1}^{(1)} - u(t_n) = O(\tau^{p+1}), \quad (71)$$

$$y_{n+1}^{(2)} - u(t_n) = O(\tau^{p+2}), \quad (72)$$

то головний член похибки першого методу

$$g(\tau) = y_{n+1}^{(1)} - y_{n+1}^{(2)} = O(\tau^{p+1}). \quad (73)$$

Далі з головним членом похибки  $g(\tau)$  методу степені  $p$  оперуємо так як і в першому підході.

Бажано мати можливість здійснювати крок інтегрування і оцінювати похибку при меншій кількості обчислення значень правих частин. Виграш досягається, якщо використовують методи, які називаються *вкладеними*. Таблиця Батчера для них має вигляд:

$\alpha_2$	$\beta_{2,1}$			
$\vdots$	$\vdots$	$\ddots$		
$\alpha_m$	$\beta_{m,1}$	$\dots$	$\beta_{m,m-1}$	
	$p_1$	$\dots$	$p_{m-1}$	$p_m$
	$\bar{p}_1$	$\dots$	$\bar{p}_{m-1}$	$\bar{p}_m$

Метод з параметрами  $p_1, \dots, p_{m-1}, p_m$  має порядок точності  $p$ , а з параметрами  $\bar{p}_1, \dots, \bar{p}_{m-1}, \bar{p}_m$  —  $p + 1$ . Коефіцієнти  $\alpha_i, \beta_{i,j}$  у обох методів однакові.

Найпростіший приклад вкладених методів для  $m = 2$  має таблицю Батчера:

0		
1	1	
	1	0
	0	1

Перший метод, якому відповідають коефіцієнти  $p_1 = 1, p_2 = 0$ , це метод Ейлера,  $p = 1$ . Другий —  $p_1 = 0, p_2 = 1$  — метод Ейлера-Коші,  $p = 2$ .

Іншим прикладом може служити сукупність формул шестістадійного  $m = 6$  методу Рунге-Кутта-Фельберга

$$k_1 = \tau f(t_n, y_n), \quad (74)$$

$$k_2 = \tau f\left(t_n + \frac{\tau}{2}, y_n + \frac{k_1}{2}\right), \quad (75)$$

$$k_3 = \tau f\left(t_n + \frac{\tau}{2}, y_n + \frac{k_1 + k_2}{4}\right), \quad (76)$$

$$k_4 = \tau f\left(t_n + \tau, y_n - k_2 + 2k_3\right), \quad (77)$$

$$k_5 = \tau f\left(t_n + \frac{\tau}{3}, y_n + \frac{7k_1 + 10k_2 + k_4}{27}\right), \quad (78)$$

$$k_6 = \tau f\left(t_n + \frac{\tau}{5}, y_n + \frac{28k_1 - 125k_2 + 546k_3 + 54k_4 - 378k_5}{625}\right), \quad (79)$$

тоді

$$\Delta y_n = \frac{k_1 + 4k_3 + k_4}{6} \quad (80)$$

з головним членом похибки

$$y_n(t_n + \tau) - u(t_n + \tau) = g(\tau) + O(\tau^5), \quad (81)$$

з

$$g(\tau) = -\frac{42k_1 + 224k_3 + 21k_4 - 162k_5 - 125k_6}{366}. \quad (82)$$

Методи Рунге-Кутта-Фельберга мають четвертий та п'ятий степінь точності. Порівняємо кількість обчислень правих частин в методах Рунге-Кутта та Рунге-Кутта-Фельберга ( $p = 4$  для обох). Згідно схеми кроків по змінній  $t$  та по стадіях  $i = \overline{1, 4}$  для методу Рунге-Кутта необхідно обчислити для оцінки похибки 11 значень функції, а для методу Рунге-Кутта-Фельберга — 6 значень функції.

- Основний **недолік** методів Рунге-Кутта: щоб отримати досить високий степінь точності потрібно багато раз обчислювати значення функції.
- **Достойністю** методів Рунге-Кутта є можливість зміни кроку інтегрування  $\tau = \tau u_n$  і за рахунок цього задовольняти умову точності на кроці.

## 10.5. Багатокрокові методи розв'язання задачі Коші. Методи Адамса

Недолік методів Рунге-Кутта: велика кількість обчислень значень функцій на одному кроці (особливо, чутливо це для систем). Висока точність в цих методах досягається за рахунок обчислень для  $m$  стадій коефіцієнтів  $k_i(\tau) = \tau f(\xi_i, \eta_i)$  в проміжних точках між  $t_n$  та  $t_{n+\tau}$ . А чи не можна для цього використати попередні значення  $f(t_n, y_n), f(t_{n-1}, y_{n-1}), \dots, f(t_{n-m}, y_{n-m})$ ?

Для розв'язання задачі Коші

$$\frac{du}{dt} = f(t, u), \quad t > 0, \quad u(0) = u_0 \quad (83)$$

введемо сітку

$$\omega_\tau = t_n = n\tau, n = 0, 1, \dots \quad (84)$$

з постійним кроком  $\tau > 0$ . Позначимо через  $y_n = y(t_n)$ ,  $f_n = f(t_n)$ , у функції, визначені на сітці  $\omega_\tau$ .

**Означення:** Лінійним  $m$ -кроковим різницеvim методом називається рівняння:

$$\frac{a_0 y_n + a_1 y_{n-1} + \dots + a_m y_{n-m}}{\tau} = b_0 f_n + b_1 f_{n-1} + \dots + b_m f_{n-m}, \quad (85)$$

для  $n = m, m+1, \dots$ , де  $a_k, b_k$  — числові коефіцієнти, які не залежать від  $n$ , причому  $a_0 \neq 0$ .

Рівняння (85) слід розглядати як рекурентне співвідношення, яке виражає нове значення  $y_n = y(t_n)$  через знайдені раніше значення  $y_{n-1}, y_{n-2}, \dots, y_{n-m}$ .

Розрахунок починається з  $n = m$ , тобто з рівняння

$$\frac{a_0 y_m + a_1 y_{m-1} + \dots + a_m y_0}{\tau} = b_0 f_m + b_1 f_{m-1} + \dots + b_m f_0, \quad (86)$$

Як бачимо з рівняння, для початку розрахунків необхідно задати  $m$  початкових значень. Значення  $y_0$  визначається початковою умовою задачі (83), а саме покладають  $y_0 = u_0$ . Величини  $y_1, \dots, y_{m-1}$  можна обчислити, наприклад, за методом Рунге-Кутта, або за методом рядів Тейлора. В подальшому будемо вважати, що початкові значення  $y_0, y_1, \dots, y_{m-1}$  задані.

З рівняння (85) видно, що на відміну від методу Рунге-Кутта багатокрокові різницеві методи допускають обчислення правих частин тільки в точках основної сітки  $\omega_t$ .

**Означення:** Метод (85) називається **явним**, якщо  $b_0 = 0$ , і відповідно, шукане значення  $y_n$  виражається явно через  $y_{n-1}, \dots, y_{n-m}$ . Якщо  $b_0 \neq 0$ , то метод (85) називається **неявним**.

В неявному методі для пошуку  $y_n$  потрібно розв'язувати нелінійне рівняння:

$$\frac{a_0}{\tau} \cdot y_n - b_0 f(t_n, y_n) = F(y_{n-1}, y_{n-2}, \dots, y_{n-m}), \quad (87)$$

де

$$F(y_{n-1}, y_{n-2}, \dots, y_{n-m}) = \sum_{k=1}^m \left( b_k f_{n-k} - \frac{a_0}{\tau} \cdot y_{n-k} \right) \quad (88)$$

Як правило це рівняння розв'язують методом Ньютона, обираючи, наприклад, початкове наближення  $y_n^{(0)}$  рівним  $y_{n-1}$ .

Коефіцієнти рівняння (85) визначені з точністю до множника. Для уникнення цієї неоднозначності, будемо вважати, що виконується умова:

$$\sum_{k=0}^m b_k = 1, \quad (89)$$

яка означає, що права частина різницевого рівняння (85) апроксимує праву частину диференціального рівняння (83).

В практичному використанні найбільш поширені методи Адамса, які являють собою частинний випадок багатокрокових методів (85), коли похідна  $u'(t)$  апроксимується тільки по двох точках,  $t_n$  і  $t_{n-1}$ , тобто

$$a_0 = -a_1 = 1, \quad a_k = 0, \quad k = \overline{2, m}. \quad (90)$$

Таким чином, методи Адамса мають вигляд:

$$\frac{y_n - y_{n-1}}{\tau} = \sum_{k=0}^m b_k f_{n-k}. \quad (91)$$

**Означення:** У випадку, коли  $b_0 = 0$ , методи Адамса називають *явними*, А при  $b_0 \neq 0$ , методи Адамса називають *неявними*.

Розглянемо детальніше процедуру побудови методів Адамса. Інтегруємо рівняння (83) по  $t \in [t_{n-1}, t_n]$ :

$$\frac{u(t_n) - u(t_{n-1})}{\tau} = \frac{1}{\tau} \int_{t_{n-1}}^{t_n} f(t, u(t)) dt. \quad (92)$$

Замінімо  $f(t, u(t))$  на інтерполяційний поліном. Виберемо вузлами інтерполявання точки  $t_n, t_{n-1}, \dots, t_{n-m}$ . Використаємо багаточлен степеня  $m$  за формулою Ньютона по рівновіддалених вузлах, тобто

$$f(t, u(t)) \approx L_m(t) = f_n + s \Delta f_{n-1} + \dots + \frac{s(s+1) \dots (s+m-1)}{m!} \Delta^m f_{n-m}. \quad (93)$$

де  $s = \frac{t-t_n}{\tau}$ , крок  $\tau$  сталий. В результаті підстановки в (92)  $f(t, u(t)) \approx L_m(t)$  отримуємо метод

$$\frac{y_n - y_{n-1}}{\tau} = f_n + \beta_1 \Delta f_{n-1} + \dots + \beta_m \Delta^m f_{n-m}, \quad (94)$$

де

$$\beta_k = \int_{-1}^0 \frac{s(s+1) \dots (s+k-1)}{k!} ds, \quad (95)$$

для  $k = \overline{1, m}$ , де  $\Delta f_{n-1} = f_n - f_{n-1}$ ,  $\Delta^2 f_{n-2} = f_n - 2f_{n-1} + f_{n-2}$ .

Похибка методу на одному кроці

$$R(\tau) = \int_{t_{n-1}}^{t_n} r_m(t) dt = O(\tau^{m+2}), \quad (96)$$

де  $r_m(t) = f(t) - L_m(r)$ , звідки степінь точності на одному кроці  $p = m + 1$ .

Формула (95) називається *формулою Адамса-Мултона* (неявний метод Адамса, інтерполяційний метод Адамса).

Виберемо вузли інтерполявання точки  $t_{n-1}, \dots, t_{n-m}$ . Отримаємо багаточлен степеня  $m - 1$ :

$$f(t, u(t)) \approx L_{m-1}(t) = f_{n-1} + v\Delta f_{n-2} + \dots + \frac{v(v+1) \dots (v+m-2)}{(m-1)!} \Delta^{m-1} f_{n-m}, \quad (97)$$

де  $v = \frac{t-t_{n-1}}{\tau}$ . Підставляючи в (95), маємо метод

$$\frac{y_n - y_{n-1}}{\tau} = f_{n-1} + \gamma_1 \Delta f_{n-1} + \dots + \gamma_{m-1} \Delta^{m-1} f_{n-m}, \quad (98)$$

де

$$\gamma_k = \int_0^1 \frac{v(v+1) \dots (v+k-1)}{k!} dv, \quad (99)$$

для  $k = \overline{1, m-1}$ .

Формула (98) називається *формулою Адамса-Башфорта* (явний метод Адамса, екстраполяційний метод Адамса). Похибка цього методу на одному кроці  $R(\tau) = O(\tau^{m+1})$ . Степінь точності на одному кроці  $p = m$ .

**Задача 39:** Побудувати явний та неявний двокрокові методи Адамса. Який степінь точності вони мають?

## 10.6. Метод невизначених коефіцієнтів побудови багатокрокових методів для розв'язання задачі Коші

Розглянемо задачу Коші

$$\frac{du}{dt} = f(t, u), \quad t > 0, \quad u(0) = u_0 \quad (100)$$

і багатокроковий метод

$$\sum_{k=0}^m \frac{a_k}{\tau} \cdot y_{n-k} = \sum_{k=0}^m b_k f_{n-k} \quad (101)$$

Підберемо коефіцієнти  $a_k, b_k$  так, щоб досягти найвищої точності методу (101).

Введемо похибку апроксимації  $\psi(\tau)$  для формули (101). Похибкою апроксимації на розв'язку або нев'язкою різницевого методу (101) називається функція

$$\psi(\tau) = - \sum_{k=0}^m \frac{a_k}{\tau} \cdot u_{n-k} + \sum_{k=0}^m b_k f(t_{n-k}, u_{n-k}), \quad (102)$$

де  $u_{n-k} = u(t_n - k\tau)$ . Її отримують підстановкою точного розв'язку  $u(t)$  задачі (100) в різницеве рівняння (101). Розглянемо питання про порядок похибки апроксимації при  $\tau \rightarrow 0$  в залежності від вибору коефіцієнтів  $a_k, b_k, k = \overline{0, m}$ .

Розкладемо функції  $u_{n-k} = u(t_n - k\tau)$  в точці  $t = t_n$  по формулі Тейлора:

$$u(t_n - \tau k) = \sum_{j=0}^p u^{(j)}(t_n) \cdot \frac{(-k\tau)^j}{j!} + O(\tau^{p+1}). \quad (103)$$

Тут  $u^{(j)}$  —  $j$ -та похідна. Далі

$$f(t_n - k\tau, u(t_n - k\tau)) = \frac{du}{dt}(t_n - k\tau) = \sum_{j=0}^{p-1} u^{(j+1)}(t_n) \cdot \frac{(-k\tau)^j}{j!} + O(\tau^p), \quad (104)$$

для  $k = \overline{1, m}$ . Підставляючи ці вирази в формулу (102), отримаємо:

$$\begin{aligned} \psi(\tau) &= - \sum_{k=0}^m \frac{a_k}{\tau} \left( \sum_{j=0}^p \frac{(-k\tau)^j u^{(j)}(t_n)}{j!} \right) + \sum_{k=0}^m b_k \left( \sum_{j=0}^{p-1} \frac{(-k\tau)^j u^{(j+1)}(t_n)}{j!} \right) + O(\tau^p) = \\ &= - \sum_{j=0}^p \sum_{k=0}^m \left( \frac{a_k}{\tau} \cdot \frac{(-k\tau)^j u^{(j)}(t_n)}{j!} \right) + \sum_{j=0}^{p-1} \sum_{k=0}^m \left( b_k \cdot \frac{(-k\tau)^j u^{(j+1)}(t_n)}{j!} \right) + O(\tau^p). \end{aligned} \quad (105)$$

Після перетворень, приходимо до розкладу:

$$\psi(\tau) = - \left( \sum_{k=0}^m \frac{a_k}{\tau} \right) u(t_n) + \sum_{j=1}^p \left( \sum_{k=0}^m (-k\tau)^{j-1} \left( a_k \cdot \frac{k}{j} + b_k \right) \right) \frac{u^{(j)}(t_n)}{(j-1)!} + O(\tau^p). \quad (106)$$

Звідки видно, що похибка апроксимації має порядок  $p$ , якщо виконуються умови:

$$E_0 = \sum_{k=0}^m a_k = 0, \quad E_l = \sum_{k=0}^m \frac{k^{l-1}}{l!} \cdot (ka_k + lb_k) = 0, \quad l = \overline{1, p}. \quad (107)$$

Разом з умовою нормування

$$\sum_{k=0}^m b_k = 1, \quad (108)$$

рівняння (107) утворюють систему з  $p + 2$  лінійних алгебраїчних рівнянь відносно  $2(m + 1)$  невідомих  $a_0, a_1, \dots, a_m, b_0, b_1, \dots, b_m$ .

Умови нормування запишуться у вигляді:

$$\lim_{\tau \rightarrow 0} \sum_{k=0}^m \frac{a_k}{\tau} \cdot u_{n-k} = \frac{du}{dt}(t_n), \quad (109)$$

$$\lim_{\tau \rightarrow 0} \sum_{k=0}^m b_k f(t_{n-k}, u_{n-k}) = f(t_n, u_n). \quad (110)$$

Зурахуванням (103), (105) маємо

$$\sum_{k=0}^m b_k = 1, \quad \sum_{k=0}^m k a_k = 1. \quad (111)$$

З рівняння (107) при  $l = 0$  отримуємо

$$\sum_{k=0}^m k a_k - \sum_{k=0}^m b_k = 0. \quad (112)$$

Тому умови нормування додатково до (107) дають тільки рівняння

$$\sum_{k=0}^m b_k = 1 \quad (113)$$

Для того, щоб система (105) не була перевизначеною, необхідно вимагати, щоб  $p \leq 2m$ . Ця вимога означає, що порядок апроксимації лінійних  $m$ -крокових різницевих методів не може перевищувати  $2m$  (неявні методи). Найвищий порядок апроксимації явних методів —  $2m - 1$ .

Для методів Адамса умови  $p$ -го порядку апроксимації мають вигляд:

$$\sum_{k=1}^m k^{l-1} b_k = 1, \quad l = \overline{2, p}, \quad n_0 = 1 - \sum_{j=1}^m b_j \quad (114)$$

Звідки бачимо, що найвищий порядок апроксимації неявного  $m$ -крокового методу Адамса дорівнює  $m + 1$ , а найвищий порядок апроксимації явного методу Адамса ( $b_0 = 0$ ) дорівнює  $m$ .

**Задача 40:** Побудувати явний та неявний двокрокові методи найвищого степеня точності.

## 10.7. Питання реалізації багатокрокових методів

1. Перша проблема, яка виникає при застосування багатокрокових методів це вибір додаткових початкових умов. Треба знайти  $(m - 1)$  додаткове початкове значення  $y_1, \dots, y_{m-1}$ . Шляхи вирішення проблеми такі.
  1. Можна використовувати методи Рунге-Кутта для пошуку цих початкових значень. Але треба, щоб ці методи мали або точність апроксимації  $p$ , або точність похибки на кроці  $p$ . Недолік такого способу: сама процедура обчислення за методами Рунге-Кутта займає великий об'єм, та й багатокрокові методи об'ємні.
  2. Можна використовувати метод рядів Тейлора. Знову ж таки потрібно узгоджувати точність  $p$ . Функцію розкладають в ряд:

$$u^{(p-1)}(t) = \sum_{k=0}^{p-1} \frac{u^{(k)}(t_0)}{k!} \cdot (t - t_0)^k, \quad (115)$$

де індекс  $(p - 1)$  над  $u$  вказує кількість членів ряду. Похибка

$$u - u^{(p-1)} = O(\tau^p). \quad (116)$$

2. Друга проблема: реалізація неявних методів

$$y_n = \tau \cdot \frac{b_0}{a_0} \cdot f(t_n, y_n) + \Phi(y_{n-1}, \dots, y_{n-m}) = \varphi(y_n). \quad (117)$$

1. Можна використовувати для знаходження розв'язку нелінійного рівняння метод простої ітерації:

$$y_n^{(k+1)} = \varphi(y_n^{(k)}) \quad (118)$$

де індекс  $k$  означає  $k$ -ту ітерацію. Умова збіжності методу простої ітерації  $\varphi'(u) \leq q < 1$ . Тоді маємо:

$$|\varphi'(u)| = \left| \tau \cdot \frac{b_0}{a_0} \cdot f'_u(t, u) \right| \leq q < 1, \quad (119)$$

звідки

$$\tau < \frac{a_0}{b_0 \cdot L}, \quad (120)$$

де  $|f'_u(t, u)| \leq L$ ,  $L$  — максимум похідної чи стала Ліпшиця. Тобто для збіжності методу необхідно, щоб крок  $\tau$  був досить малим.

2. Можна використовувати метод Ньютона. Як відомо, умови збіжності методу Ньютона залежать від вдалого вибору початкових умов. Хороше початкове наближення таке: або  $y_n^{(0)} = y_{n-1}$ , або  $y_n^{(0)}$  обчислюється за явною  $m$ -кроковою формулою.

В результаті умови збіжності методу Ньютона менш жорсткі, ніж у методу простої ітерації.

3. можна використовувати метод предиктор-коректор. Запишемо формули неявного методу ( $C$  — коректор,  $P$  — предиктор):

$$C: y_n = \tau \cdot \frac{b_0}{a_0} \cdot f(t_n, y_n) + \Phi(y_{n-1}, \dots, y_{n-m}) \text{ — неявний } m\text{-кроковий метод, } p = 2m.$$

$$P: y_n = \bar{\Phi}(y_{n-1}, \dots, f_{n-1}, \dots) \text{ — явний } m\text{-кроковий метод, } p = 2m - 1.$$

Далі виконується така процедура:

$$P: \bar{y}_n = \bar{\Phi}(y_{n-1}, \dots, f_{n-1}, \dots);$$

$$E: \bar{f}_n = f(t_n, \bar{y}_n);$$

$$C: y_n = \tau \cdot \frac{b_0}{a_0} \cdot f(t_n, \bar{y}_n) + \Phi(y_{n-1}, \dots);$$

$$E: f_n = f(t_n, y_n),$$

де  $E$  — підрахування правої частини рівняння (Equation). Схема  $PECE$  — це метод *предиктор-коректор*.

Іноді, щоб підвищити точність, використовують схему  $PE(CE)^j$ . Ця схема аналогічна методу простої ітерації, де  $j$  — це максимальна кількість ітерацій.