

УДК 519.612

В. С. Абрамчук*, канд. фіз.-мат. наук,

І. В. Абрамчук**, старший викладач

*Вінницький державний педагогічний університет
імені М. Коцюбинського, м. Вінниця,

**Вінницький національний технічний університет, м. Вінниця

ОПТИМІЗАЦІЙНІ МЕТОДИ РОЗВ'ЯЗУВАННЯ СИСТЕМ

$A\vec{x} = \vec{b}$ З ПОГАНО ЗУМОВЛЕНИМИ МАТРИЦЯМИ

Досліджені проблемні задачі, пов'язані з оцінкою похибок у розв'язках зашумлених систем, мінімізацією похибок обчислень при гауссових перетвореннях, прискоренням швидкості збіжності ітераційних методів.

Ключові слова: оцінка похибок розв'язку, багатошаровий метод гауссового виключення, локальний базис методу напрямленого пошуку, відношення Релея.

Вступ. СЛАР з погано зумовленими матрицями необхідно розв'язувати у багатьох задачах практики: в задачах математичної обробки і прогнозування експериментів, в задачах оптимального управління і оптимального прогнозування технічних систем, при дослідженні природних явищ і катаклізмів [1–10, 13]. Одним з підходів до розв'язування таких задач є різнищеві рівняння і метод скінченних елементів. При цьому матриці таких систем є розрідженими, великих порядків і погано зумовленими. Для створення ефективних методів розв'язування систем з погано зумовленими матрицями необхідно вяснити причини, що перешкоджають збіжності ітераційних методів та розробити нові принципи прискорення їх збіжності та мінімізації похибок обчислень при гауссових перетвореннях матриці. Основою таких досліджень є роботи [11, 12].

Постановка проблеми. Нехай необхідно розв'язати систему лінійних алгебричних рівнянь (СЛАР), яка має вид $A\vec{x} = \vec{b}$, $A \in M_{n \times n}(\mathbb{R})$, $\vec{b} \in \mathbb{R}^n$, з матрицею A довільної структури, погано зумовленою, великого порядку. Нерозв'язаними або недостатньо дослідженими є проблемні задачі: 1) вплив похибок на розв'язок СЛАР $\vec{\Delta x}(A, \vec{b}, E, \vec{\delta b})$: $(A + E)(\vec{x} + \vec{\Delta x}) = \vec{b} + \vec{\delta b}$, E , $\vec{\delta b}$ — похибки даних і обчислень; 2) причини, що не дозволяють існуючим методам бути ефективними при розв'язуванні таких СЛАР; 3) стратегія розробки стійких алгоритмів розв'язування СЛАР в умовах похибок обчислень.

Аналіз актуальних досліджень. СЛАР з погано зумовленими матрицями необхідно розв'язувати у багатьох задачах практики: в задачах

математичної обробки і прогнозування експериментів, в задачах оптимального управління і оптимального прогнозування технічних систем, при дослідженні природних явищ і катаклізмів [1–10]. При розв'язуванні СЛАР з погано зумовленими матрицями, в літературі не існує єдиного підходу: з однієї сторони, обґрунтовується, що необхідно відмовитися від розв'язання систем з матрицями типу матриці Гільберта, оскільки число зумовленості таких матриць наближається до $e^{3.5n}$ [3]; з другої сторони матриці Гільберта широко застосовуються у регресійному аналізі, тощо, тому необхідно розробляти методи розв'язування таких систем [2; 8]; по-третє, існує цілий клас методів спряжених напрямів, заснованих на погано зумовленому базисі Крилова, які ефективно використовуються для розв'язування різницевих рівнянь [2; 4]. Причиною таких протиріч є недостатній аналіз факторів, які впливають на ефективність методів розв'язування систем з погано зумовленими матрицями.

Мета статті. 1. Виконати аналіз впливу похибок обчислень на похибку розв'язку СЛАР. 2. Розробити багатоступовий метод гауссових перетворень, що мінімізує похибку обчислень. 3. Розробити стратегію побудови ітераційних методів розв'язування СЛАР на основі системи повних базисів Крилова.

Основна частина. 1. Аналіз похибок у розв'язках зашумлених систем. Однією з основних проблем теорії лінійних алгебраїчних систем є аналіз впливу похибок заокруглення і похибок початкових даних (похибок невизначеності) на розв'язок, обчислений на ЕОМ у скінченорозрядній арифметиці [2–4; 10]. На основі оберненого аналізу, на місце розв'язку системи $A\vec{x} = \vec{b}$, $A \in M_{n \times n}(\mathbb{R})$, $\text{rank } A = n$, $\vec{b} \in \text{im } A$, досліджується точний розв'язок зашумленої системи $(A + E)(\vec{x} + \vec{\delta x}) = \vec{b} + \vec{\delta b}$, де $E, \vec{\delta b}$ — похибки. Якщо $\rho(A^{-1}E) < 1$, то матриця $A + E$ має обернену і можна оцінити відносну похибку розв'язку зашумленої системи через число зумовленості матриці [2; 3; 10]. Ця оцінка є апіорною, оскільки в неї не входить ні обчислений розв'язок $\vec{\bar{x}} = \vec{x} + \vec{\delta x}$, ні розв'язок системи $A\vec{x} = \vec{b}$ [10]. Вона не дає способу, як керувати обчислювальним процесом, щоб зменшити похибку розв'язку.

Допустимо, що похибки коефіцієнтів матриці є неперервними параметрами, тоді квадратну матрицю $A + E$ можна вважати лінійним оператором деякої системи лінійних диференціальних рівнянь. Задача дослідження поведінки матриці і розв'язків системи рівнянь зведеться до структурної стійкості матриці при збуренні параметрів (теорія канонічних форм Жордана-Арнольді [9]). Проте цей підхід також не дає способу мінімізації похибки розв'язку.

У роботі [12] пропонується розв'язок «зашумленої» системи досліджувати як дробово-раціональну функцію похибок на основі формул Крамера. Skorистаємось тим, що визначник матриці є лінійною функцією вектор-стовпців (рядків) матриці [4].

Запишемо розв'язки систем $A\bar{x} = \bar{b}$, $(A + E)(\bar{x} + \delta\bar{x}) = \bar{b} + \delta\bar{b}$ через визначники: $x_i = D_i/D_0$, $x_i + \delta x_i = (D_i + \delta D_i)/(D_0 + \delta D_0)$. Тоді похибки компонент розв'язку $\delta x_i = (\delta D_i - x_i \cdot \delta D_0)/(D_0 + \delta D_0)$.

Не втрачаючи загальності, запишемо зашумлену систему другого порядку у матрично-векторній формі $[A_1 + E_1 \quad A_2 + E_2](\bar{x} + \delta\bar{x}) = \bar{b} + \delta\bar{b}$. Тоді в силу лінійності визначників відносно вектор-стовпців матриці, матимемо

$$\begin{aligned} D_0 + \delta D_0 &= \det[A_1 + E_1 \quad A_2 + E_2] = \det[A_1 \quad A_2] + (\det[A_1 \quad E_2] + \\ &+ \det[E_1 \quad A_2]) + \det[E_1 \quad E_2] = D_0 + L_0(E) + K_0^2(E), \\ D_1 + \delta D_1 &= \det[\bar{b} + \delta\bar{b} \quad A_2 + E_2] = \det[\bar{b} \quad A_2] + (\det[\delta\bar{b} \quad A_2] + \\ &+ \det[\bar{b} \quad E_2]) + \det[\delta\bar{b} \quad E_2] = D_1 + L_1(\delta\bar{b}, E_2) + K_1^2(\delta\bar{b}, E_2), \end{aligned} \quad (1)$$

аналогічно $D_2 + \delta D_2 = D_2 + L_2(\delta\bar{b}, E_1) + K_2^2(\delta\bar{b}, E_1)$.

Тут позначено: $\det[\bar{c}_1 \quad \bar{c}_2]$ — визначник матриці, складеної з двох вектор-стовпців \bar{c}_1 , \bar{c}_2 ; L — лінійна форма похибок; K^m — похибки вищих порядків.

Компоненти похибок розв'язку зашумленої системи другого порядку

$$\delta x_i = \left[(L_i - x_i L_0) + (K_i^2 - x_i K_0^2) \right] / \left[D_0 + (L_0 + K_0^2) \right], \quad i = 1, 2.$$

Для довільної зашумленої системи n -го порядку похибки компонент розв'язку $\forall i = [1 : n]$ матимуть вид:

$$\delta x_i = \left[L_i - x_i L_0 + \sum_{j=2}^{n-1} (K_i^j - x_i K_0^j) \right] / \left[D_0 + L_0 + \sum_{j=2}^{n-1} K_0^j \right], \quad (2)$$

де L_0 — лінійна форма похибок матриці, L_0 — сума n визначників матриць $(A \setminus A_i) \cup E_i$, де кожен i -й стовпчик A_i замінюється вектором похибок E_i ; L_i — лінійні форми похибок матриць $[\delta\bar{b}, E \setminus E_i]$, в яких на місце вектора E_i ставиться вектор $\delta\bar{b}$ (є сумами відповідних визначників).

Теорема 1. Для того, щоб зашумлена система $(A + E)(\bar{x} + \delta\bar{x}) = \bar{b} + \delta\bar{b}$ була стійкою до похибок $(E, \delta\bar{x})$, необхідно, щоб допущені похибки задовольняли нерівність

$$\operatorname{sgn} D_0 \cdot \left| L_0 + \sum_{j=2}^{n-1} K_0^j \right| / D_0 < 1, \quad (3)$$

необхідно і достатньо, щоб виконувалась нерівність (3) і $\forall k = [1:n]$ нерівності $|\delta x_k| < \varepsilon$, що еквівалентно нерівностям

$$\left| L_k - x_k L_0 + \sum_{j=2}^{n-1} (K_k^j - x_k K_0^j) \right| < \varepsilon \cdot \left| D_0 + L_0 + \sum_{j=2}^{n-1} K_0^j \right|, \quad (4)$$

($|\cdot|$ — абсолютна величина).

Доведення теореми 1 впливає з вищенаведених формул (1), (2).

Висновок. Зростом розмірності матриці похибки у розв'язках зростають, оскільки похибки матриці і правої частини належать простору \mathbb{R}^{n^2+n} (у той час, як розв'язок належить простору \mathbb{R}^n).

На основі проведеного аналізу побудуємо метод гауссового перетворення матриці, що мінімізує похибку обчислень.

2. Багатошаровий алгоритм гауссового перетворення (виключення). При розв'язуванні систем $A\bar{x} = \bar{b}$, $A \in M_{n \times n}(\mathbb{R})$, $\operatorname{rank} A = n$, $\bar{b} \in \operatorname{im} A$, з допомогою машинної арифметики методом гауссового виключення (або іншим прямим методом) необхідно досліджувати дві проблемні задачі, пов'язані з запобіганням: 1) катастрофічної втрати точності (втрати старших розрядів [2]) у компонентах перетвореної матриці або, навпаки, у непомірному зростанні значень компонент [3]; 2) зростання числа зумовленості при перетвореннях матриці (наближення матриці до виродженої). Ці проблеми тісно пов'язані між собою і вимагають оцінки наближених розв'язків [2–4; 10].

Із загальної теорії збурень впливає, що повністю усунути похибки обчислень при лінійних перетвореннях неможливо [2–10], їх можна лише зменшити за рахунок оптимізації методу перетворень системи. У роботі [11] доведено, що метод гауссового виключення на кожному i -му кроці, $i \in [1:n-1]$, є процесом ортогоналізації вектор-рядків A_k $\forall k \in [i+1:n]$ до одиничного вектора \bar{e}_i . Тому для розв'язання першої проблемної задачі, необхідно на кожному i -му кроці знаходити розв'язуючий вектор-рядок A_{i_0} , що утворює найменший кут з вектором \bar{e}_i [11].

Похибка обчислень при перетвореннях $(\tilde{A}_k, \bar{e}_i) = 0 \Rightarrow \Rightarrow (A_k + \alpha_k A_{i_0}, \bar{e}_i) = 0 \Rightarrow A_{k,i} + \alpha_k A_{i_0,i} = 0$ буде найменшою для всіх

$k \in [i+1 : n]$, якщо абсолютне значення $|A_{i,i}|$ найбільше (такий елемент $A_{i,i}$ називають головним [2; 3]). Оскільки головний елемент повинен належати розв'язуючому вектор-рядку, то задача його визначення еквівалентна задачі [11]: обчислити індекси i_0, j_0 для яких досягається

$$\max_{j \in [i:n]} \max_{k \in [1:n]} |\cos \{A_{k,j}, \vec{e}_i\}| = \max_{j \in [i:n]} \max_{k \in [1:n]} |A_{k,j}| / \|A_{k,j}\|_2.$$

Отже, задача зменшення похибки обчислень на i -му кроці в методі гауссового виключення (який назвемо одношаровим) полягає у знаходженні як розв'язуючого вектор-рядка A_{i_0} так і розв'язуючого вектор-стовпця A_{j_0} , на перетині яких знаходиться головний елемент A_{i_0, j_0} (який обмінюється з елементом $A_{i,i}$ шляхом перестановки вектор-рядків, вектор-стовпців). Щоб одночасно розв'язати обидві проблемні задачі, побудуємо багатошаровий метод гауссового перетворення шляхом одночасного обнулення піддіагональних елементів у стовпцях $A_{i_0}, \dots, A_{(i+m-1)}$, $m \leq n$.

Виконаємо обернений аналіз для $m = 2$, тоді для $\forall k \in [i+2 : n]$:

$$(\tilde{A}_{k,i}, \vec{e}_i) = 0, (\tilde{A}_{k,i}, \vec{e}_{i+1}) = 0, \quad (5)$$

$$\tilde{A}_{k,i} := A_{k,i} + \alpha_k A_{i,i} + \beta_k A_{(i+1)} \Rightarrow \begin{cases} A_{k,i} + \alpha_k A_{i,i} + \beta_k A_{i,i+1} = 0, \\ A_{k,i+1} + \alpha_k A_{i+1,i} + \beta_k A_{i+1,i+1} = 0. \end{cases}$$

Мінімізація похибки обчислень при розв'язуванні сукупності систем (5) $\forall k \in [i+2 : n]$ досягається тоді і лише тоді, якщо головний визначник $D_0^{(2)} = A_{i,i} A_{i+1,i+1} - A_{i+1,i} A_{i,i+1}$ приймає найбільше абсолютне значення. Отже необхідно знайти нову пару розв'язуючих векторів $A_{i,i}, A_{j,i}$ з умови $\max_{k \in [i+1:n]} \max_{j \in [i+1:n]} |A_{i,i} A_{k,j} - A_{k,i} A_{i,j}|$, що реалізується простою алгоритмічною процедурою.

Якщо уже виконане упорядкування матриці $A[i+1:n, i+1:n]$ шляхом перестановки вектор-рядків $A_{(i+1)}, A_{i,i}$, вектор-стовпців $A_{(i+1)}, A_{j,i}$, то наступна тришарова процедура гауссового виключення: $\forall k \in [i+3 : n]$, $(\tilde{A}_{k,i}, \vec{e}_i) = 0$, $(\tilde{A}_{k,i}, \vec{e}_{i+1}) = 0$, $(\tilde{A}_{k,i}, \vec{e}_{i+2}) = 0$,

$$\tilde{A}_{k,i} := A_{k,i} + \alpha_k A_{i,i} + \beta_k A_{(i+1)} + \gamma_k A_{(i+2)}, \quad (6)$$

зведеться до пошуку нової пари розв'язуючих векторів A_{i_2}, A_{j_2} . шляхом максимізації абсолютного значення визначника третього порядку системи (6). Оскільки процедура введення нових розв'язуючих векторів $A_{i_t}, A_{j_t}, t \in [1:m-1]$ є вкладеною (рекурентною) (без перетворень матриці), то алгоритм m -шарового, $m \geq 2$, гауссового перетворення легко здійснюється для довільних щільно заповнених матриць. Параметри m -шарових гауссових перетворень після визначення розв'язуючих векторів можна знайти методом Крамера, оскільки головні визначники обчислені. З використанням, наприклад, тришарового процесу (6) одночасно обнуляються піддіагональні елементи в трьох вектор-стовпцях $A_{i_t}, A_{(i+1)_t}, A_{(i+2)_t}$ для $\forall k \in [i+3:n]$ без попереднього використання одношарових і двошарових процедур. Після завершення тришарового процесу для $k \in [i+3:n]$ перетворюється рядок $A_{(i+2)_k}$ двошаровою процедурою, і нарешті $A_{(i+1)_k}$ — одношаровою.

Теорема 2. m -шаровий процес гауссового виключення коректний: мінімізує похибку обчислень, однозначно визначає параметри перетворень матриці, не приводить до виродження системи.

Доведення теореми 2. Оскільки за умовою матриця A не вироджена, то для неї існують відмінні від нуля визначники всіх порядків $2 \leq m < n$. Їх найбільше абсолютне значення гарантується скінченим числом переборів по $k, j \in [i+t:n], t \in [1:m-1]$ (тобто простою двоциклічною програмною реалізацією). Процес обчислень економний, оскільки використовує лінійні векторні перетворення типу (6). Мінімізація похибки обчислень гарантується вибором розв'язуючих векторів.

3. Модифікований алгоритм гауссового виключення для щільно заповнених погано зумовлених матриць. Для того, щоб процес гауссового виключення був стійким при машинній реалізації, пропонується масштабування матриці (A, \vec{b}) [2; 3].

Запропонований в п. 2 алгоритм гауссового виключення дозволяє модифікувати обчислювальний процес в залежності від типу матриці. Нехай розв'язується система з матрицею типу матриці Гільберта $H = (1/(i+j-1))_{i,j=1}^n$, яка виникає у багатьох практичних дослідженнях (поліноміальна регресія [2, 3, 10]) і є погано зумовленою (число зумовленості асимптотично наближається до e^{cn} , $c \approx 3.5$ [10], наприклад, похибки в даних поліноміальної регресії 10-го порядку зростають на множник, більший, ніж $3 \cdot 10^{12}$ [3].

Виникає питання, у чому причина сильного зростання похибок у розв'язках систем з погано зумовленими матрицями. Очевидна відповідь, що матриця типу матриці Гільберта є погано зумовленою — одностороння. Більш глибока причина полягає у тому, що не дивлячись на частковий або повний вибір головного елемента, гауссове перетворення матриці Гільберта приводить до зростання числа зумовленості перетворюваних векторів і не зростання значень визначників у перетворених матрицях $A[i:n, i:n]$, $i \in [1:n-1]$ (наближення матриці $A[i:n, i:n]$ до виродженої — до структурної нестійкості).

Побудуємо алгоритм гауссового виключення для матриць типу матриці Гільберта або матриці Коші [2], що мінімізує похибку обчислень. Нехай на кожному i -му кроці перетворень матрицю можна масштабувати, наприклад, за умовою $A_{k,i} = 1$, $k \in [i:n]$. Тоді, на основі п.2, розв'язуючим вектор-рядком є той, для якого досягається $\max |A_{k,i+1}|$, $k \in [i:n]$. Для матриці Гільберта порядку $n \geq 2$, розв'язуючим вектором для $i=1$ є не A_1 (що визначається стандартним методом гауссового виключення [2], а вектор $\tilde{A}_n := A_{..n}/A_{n,1}$. Всі наступні перетворення з умовою масштабування $A_{k,i} = 1$, $k \in [i:n]$, для матриці Гільберта є стійкими.

Нехай матриця A не може бути задовільно масштабована умовою $A_{k,i} = 1$, $k \in [i:n]$ (частковий вибір), або $A_{k,j} = 1$, $k, j \in [i:n]$ (повний вибір) [3]. На початковому кроці розширимо матрицю (A, \bar{b}) , $A \in M_{n \times n}$ до матриці $(A1, \bar{b}1)$, $A1 \in M_{(n+1) \times (n+1)}$ введенням першого стовпця з одиничних елементів $A1_{k,1} = 1$, $k \in [1:n+1]$, першого рядка з елементів одиничного вектора $A1_{..1} = \bar{e}_1$, $\forall i, j \in [2:n+1]$, $A1_{i,j} = A_{i-1,j-1}$; з елементами правих частин: $b1 = 1$, $\forall i \in [2:n+1]$, $b1_i = b_{i-1} + 1$.

Алгоритм. Сформувати матрицю $(A1, \bar{b}1)$.

Для $i \in [1:n+1]$ вибрати пару розв'язуючих векторів \bar{a} , \bar{c} з вектор-рядків $A_{..k}$, $k \in [i:n+1]$ таких, що для всіх векторів $\bar{d} \in \{A_{(i+1)}, \dots, A_{(n+1)}\}$, $\bar{d} \neq \alpha \bar{a} + \beta \bar{c}$, виконуються умови: $\bar{d} := \bar{d} + \alpha \bar{a} + \beta \bar{c}$, $(\bar{d}, \bar{e}_i) = 0$, $(\bar{d}, \bar{e}_{i+1}) = 1$, де $\alpha = -(1 + \beta)$, $\beta = (1 - d_{i+1} + a_{i+1}) / (c_{i+1} - a_{i+1})$.

Розв'язуючі вектори \vec{a} , \vec{c} вибираються з умови $\max_{j \in [i:n]} |c_{i+1} - a_{i+1}| =$
 $= \max_{k, j \in [i:n+1]} |A_{k,i+1} - A_{j,i+1}|$ (частковий вибір), або $\max_{k, j \in [i:n+1]} |A_{k,t} - A_{j,t}|$
 (повний вибір для всіх стовпчиків $t \in [i+1:n+1]$), $\max_{t \in [i+1:n+1]} |c_t - a_t| > 0$
 існує (у протилежному випадку матриця A вироджена).

4. Конус K_{\min} . При розв'язуванні систем $A\vec{x} = \vec{b}$ ітераційними методами основна увага приділена асимптотичній швидкості збіжності, з розрахунку що буде досягнута на деякому кроці задана точність $\|\vec{r}\| < \varepsilon$. Ця теоретична оцінка лежить, як правило, в основі порівняння методів на ефективність. Але в дійсності оцінка не виконується, особливо при розв'язуванні систем з погано зумовленими матрицями. На перешкоді стає конус K_{\min} . Простір \mathbb{R}^n розбивається гіперплощинами $A_i \vec{x} = \vec{b}_i$, $i = \{1, \dots, n\}$, на 2^n конусів з вершиною у точці $\vec{x}^* = A^{-1}\vec{b}$. Пару конусів, що містять сингулярну пряму $\vec{x} = \vec{x}^* + t \cdot \vec{e}_{\min}$, \vec{e}_{\min} — власний вектор матриці $A^T A$, що відповідає $\lambda_{\min}(A^T A)$, назвемо конусом K_{\min} (конусом мінімальних нев'язок).

Щоб пояснити деталі, розглянемо приклад у просторі \mathbb{R}^2 .

Приклад. Проаналізувати систему

$$A\vec{x} = \vec{b}, \quad A = \begin{bmatrix} 1 & 1 \\ 1 & 0.9999 \end{bmatrix}, \quad \vec{b} = \begin{bmatrix} 0 \\ 0.01 \end{bmatrix}, \quad \vec{x}^* = \begin{bmatrix} 100 \\ -100 \end{bmatrix}.$$

Обчислимо основні характеристики:

$$\mathbf{cond} A = \sqrt{\lambda_{\max}(A^T A) / \lambda_{\min}(A^T A)} = 12648.79453, \quad \det A = -0.0001,$$

$$\lambda_{\max}(A^T A) = 3.99980075, \quad \lambda_{\min}(A^T A) = 0.000000025,$$

$$S_{\min} : \begin{cases} x_1 = 100 - 0.999950012t, \\ x_2 = -100 + t; \end{cases} \quad S_{\max} : \begin{cases} x_1 = 100 + 1.000967t, \\ x_2 = -100 + t; \end{cases}$$

$t \in \mathbb{R}$, S_{\min} , S_{\max} , — сингулярні прямі, що визначаються власними векторами \vec{e}_{\min} , \vec{e}_{\max} матриці $A^T A$.

Малість нев'язки в області K_{\min} не дозволяє оцінити реальний розв'язок, оскільки норма нев'язки може стати малою величиною, а норма похибки залишатись великою величиною. Будь-який ітераційний процес, що заснований на мінімізації норми вектора нев'язки або

норми вектора похибки, швидко приводить наближення в область K_{\min} . Дійсно, виберемо дві початкові точки $\bar{x}^{(1)} = [150; 50]^T \notin K_{\min}$, $\|\bar{x}^{(1)} - \bar{x}^*\|_2 \approx 158$, $\bar{x}^{(2)} = [0; 0]^T \in K_{\min}$, $\|\bar{x}^{(2)} - \bar{x}^*\|_2 \approx 14$.

Застосуємо один і той же ітераційний процес $\bar{x}^{(k)} = \bar{x}^{(k)} + \alpha_k A^T \bar{r}^{(k)}$, $\alpha_k = -\|\bar{r}^{(k)}\|_2^2 / \|A^T \bar{r}^{(k)}\|_2^2$, що мінімізує норму вектора похибки, відстанню $\bar{x}^{(1)} \approx [49.999 - 49.996]^T$; $\bar{x}^{(2)} \approx [0.0050005 \ 0.005]^T$; $\bar{x}^{(1)}, \bar{x}^{(2)} \in K_{\min}$, (обчислення проведені на 10-розрядному калькуляторі). Незалежно від того, що $\bar{x}^{(1)} \notin K_{\min}$, $\bar{x}^{(2)} \in K_{\min}$, ітераційний процес за один крок привів нові наближення в область K_{\min} (точніше, точка $\bar{x}^{(1)}$, яка не належить K_{\min} і знаходилась на великій відстані від розв'язку, перейшла в ближчу точку $\bar{x}^{(1)} \in K_{\min}$, ніж точка $\bar{x}^{(2)}$, що знаходилась в K_{\min} на меншій відстані від розв'язку, яка майже не змінила своє положення). Обчислені нев'язки $\bar{r}^{(1)}$, $\bar{r}^{(2)}$ не вносять ніякого роз'яснення в дану ситуацію.

5. Ітераційні методи розв'язування систем $A\bar{x} = \bar{b}$ великих порядків. Щоб обґрунтувати можливу ефективність довільного методу розв'язування системи лінійних алгебричних рівнянь з невідірженими погано зумовленими матрицями великих порядків, необхідно виконати аналіз впливу похибок заокруглення на результат обчислювального процесу реалізації алгоритму методу у машинній арифметиці (в умовах скінчено-розрядної сітки). Оскільки прямі методи пов'язані з накопиченням похибок при перетвореннях матриці, то вони відносяться до неефективних методів розв'язування систем з погано зумовленими матрицями великих порядків.

Нерозв'язаною проблемою є: «чи існує стратегія побудови ітераційних методів розв'язування таких систем, яка забезпечує як теоретичну збіжність, так і її практичну алгоритмічну реалізацію у скінчено-розрядній арифметиці при економній кількості обчислень». Існують дві принципово різні стратегії побудови ітераційних методів.

Перша — ітераційний процес задається у формі $\forall k = 0, 1, \dots$ $\bar{x}^{(k+1)} = \bar{x}^{(k)} + \alpha_k \bar{c}^{(k)}$ (або у більш загальній формі $\bar{x}^{(k+1)} = B^{(k)} \bar{x}^{(k)} + \beta_k \bar{w}^{(k)}$), де $B^{(k)}$, $\bar{c}^{(k)}$, $\bar{w}^{(k)}$, задаються в залежності від класу методів, як функції матриці A , правої частини \bar{b} , наближення до A^{-1} , можуть враховувати структуру матриці та попередні наближення

за t кроків [2; 4]. Збіжність цих процесів заснована на мінімізації вектора похибки $\vec{\varepsilon}^{(k+1)} = \vec{x}^{(k+1)} - \vec{x}^*$ або вектора нев'язки $\vec{r}^{(k+1)} = A\vec{x}^{(k+1)} - \vec{b}$, або на іншій процедурі, наприклад, на побудові спряжених, біспряжених напрямів, тощо. Проте всі ці процедури, за умови збіжності, приводять в область K_{\min} , де процес збіжності сповільнюється і може стати, в силу похибок обчислень, неконтрольованим відносно монотонності процесу. Більш глибокі результати у розв'язуванні систем з погано зумовленими матрицями засновані на принципі регуляризації — побудові псевдорозв'язку з мінімальною нормою $\|\cdot\|_2$ [8]. Але, якщо шукати нормальний розв'язок відносно $\vec{x}^* = A^{-1}\vec{b}$, то знову виникне проблемна задача впливу похибок заокруглення на збіжність методу (оскільки абсолютний нуль у машинній арифметиці відсутній).

Запропонуємо двоциклічну стратегію розв'язування цієї проблеми на основі методу напрямленого пошуку [12]:

$$\vec{x}^{(k,j)} = \vec{x}^{(j)} + \alpha_{k,j} \vec{c}^{(k,j)}, \quad k = 0, 1, \dots, m, \quad j = 1, \dots, m_2, \quad (7)$$

де у внутрішньому циклі по $k \in \{0, 1, \dots, m\}$ наближення $\vec{x}^{(j)}$, нев'язки $\vec{r}^{(j)} = A\vec{x}^{(j)} - \vec{b}$ не змінюються і а) наближення $\vec{x}^{(j+1)}$, $j \in \{0, 1, \dots, m_2\}$ (зовнішній цикл) оптимізується у внутрішньому циклі за рахунок напрямного вектора $\vec{c}^{(k,j)} \in K_m$, $(\vec{c}^{(k,j)}, \vec{r}^{(j)}) \neq 0$, шляхом мінімізації норми $\|\cdot\|_2$ вектора похибки $\vec{\varepsilon}^{(k,j)} = \vec{x}^{(k,j)} - \vec{x}^*$, але вектор $\vec{x}^{(k,j)}$ не обчислюється); б) нев'язка $\vec{r}^{(k+1)}$ оптимізується у внутрішньому циклі за рахунок вектора-поправки $\vec{c}^{(k,j)} \in K_m$, $(A\vec{c}^{(k,j)}, \vec{r}^{(j)}) \neq 0$ шляхом мінімізації норми $\|\cdot\|_2$ вектора нев'язки $\vec{r}^{(k,j)} = A\vec{x}^{(k,j)} - \vec{b}$, але вектори $\vec{x}^{(k,j)}$, $\vec{r}^{(k,j)}$ не обчислюються. Ітераційні процеси а), б) формально ідентичні з точки зору організації обчислень, проте мають різний геометричний зміст.

Побудуємо алгоритм процедури б) (для п. а) будується аналогічно), який складається з двох етапів: 1) обчислення початкового вектора $\vec{x}^{(0)}$; 2) уточнення розкладу вектора \vec{b} .

Алгоритм обчислення початкового вектора $\vec{x}^{(0)}$. Вектор \vec{x} , розв'язок системи $A\vec{x} = \vec{b}$, визначає розклад вектора \vec{b} по неортогональному базису $\{A_i = A\vec{e}_i\}_{i=1}^n$. На місце фактичного розкладу, виконаємо ортогональні проєкції \vec{b} на вектори A_i , дістанемо наближені

значення компонент $x_i^{(0)} = (\vec{b}, A\vec{e}_i) / \|A\vec{e}_i\|_2^2$ розкладу. Вектор $\vec{b}' = \sum_{i=1}^n A_i x_i^{(0)}$ за напрямом і евклідовою довжиною не збігається з вектором \vec{b} . Масштабуємо вектор \vec{b}' за умови $\|\alpha \vec{b}'\|_2 = \|\vec{b}\|_2 \Rightarrow \alpha = \|\vec{b}\|_2 / \|\vec{b}'\|_2$. Масштабований вектор $\vec{y}^{(0)} = \alpha \vec{x}^{(0)}$ приймемо за початковий розв'язок системи $A\vec{x} = \vec{b}$.

Побудуємо алгоритм уточнення розв'язку системи $A\vec{x} = \vec{b}$ (з матрицею – чорний ящик) на основі базисних векторів $Kr_m = \{\vec{r}^{(0)}, A\vec{r}^{(0)}, \dots, A^{m-1}\vec{r}^{(0)}\} = \{\vec{u}_0, \vec{u}_1, \dots, \vec{u}_{m-1}\}$ підпростору Крилова і системи повних базисів $\{Ke_i\}_{i=1}^n$. Базисні вектори підпростору Крилова є сильнозумовленими, тому їх ортогоналізація методом Грама-Шміда (Арнольдї або Ланцоша) приводить до швидкої розортогоналізації в силу похибок обчислень [2; 4], отже $m \ll n$. Систему повних базисів найпростіше формувати на основі одиничного базису $\{\vec{e}_i\}_{i=1}^n$ у формі криловських базисів $Ke_i = \{\vec{e}_i, A\vec{e}_i, \dots, A^{m-1}\vec{e}_i\}$, $i \in [1:n]$. Розмірність m криловських базисів $K_m = \{\vec{p}, A\vec{p}, \dots, A^{m-1}\vec{p}\}$ вибиратимемо в залежності від зумовленості матриці A , за правилом: $\forall j = 0, 1, \dots, t \leq m-2$ виконується $|\cos\{A^j \vec{p}, A^{j+1} \vec{p}\}| < 1 - 10^{-s}$ і $|\cos\{A^{m-1} \vec{p}, A^m \vec{p}\}| \geq 1 - 10^{-s}$ ($s > 1$ — ціле додатне число, що вибирається в залежності від машинної точності).

Базис $Kr_m = \{\vec{r}^{(0)}, A\vec{r}^{(0)}, \dots, A^{m-1}\vec{r}^{(0)}\} = \{\vec{u}_0, \vec{u}_1, \dots, \vec{u}_{m-1}\}$, $m \ll n$ як не повний базис, використаємо для побудови початкового вектора $\vec{c}^{(0)}$. Методом Грама-Шміда з векторів системи $\{\vec{u}_i\}_{i=0}^{m-1}$ побудуємо A -ортогональну систему $\{\vec{u}_i\}_{i=0}^{m-1}$, $(A\vec{u}_i, A\vec{u}_j) = 0$, $i \neq j$. Вектор $\vec{c}^{(0)}$, що оптимізує норму $\|\circ\|_2$ вектора нев'язки $\vec{r} = \vec{r}_0 + \sum_{i=0}^{m-1} \alpha_i A\vec{u}_i$ запишемо у формі

$$\vec{c}_0 = \sum_{i=0}^{m-1} \alpha_i \vec{u}_i, \alpha_i = -(\vec{r}^{(0)}, A\vec{u}_i) / \|A\vec{u}_i\|_2^2. \quad (8)$$

Результатом є вектор $\vec{c}^{(0)}$, що забезпечує монотонність процесу: $\vec{x} = \vec{x}^{(0)} + \gamma \vec{c}^{(0)}$, $\gamma = -(\vec{r}^{(0)}, \vec{c}^{(0)}) / \|A\vec{c}^{(0)}\|_2^2$. Дійсно $\|\vec{r}\|_2^2 = \|\vec{r}^{(0)}\|_2^2 -$

$-\left(\bar{r}^{(0)}, A\bar{c}^{(0)}\right)^2 / \left\|A\bar{c}^{(0)}\right\|_2^2 < \left\|\bar{r}^{(0)}\right\|_2^2$ Оптимальні параметри α_i розраховані теоретично в абсолютній арифметиці. Чим менше похибка розортогоналізації, тим ближче реально обчислений вектор $\bar{c}^{(0)}$ до оптимального.

Уточнимо вектор $\bar{c}^{(0)}$, якщо $\varphi = \cos^2 \left\{\bar{r}^{(0)}, A\bar{c}^{(0)}\right\} < 1 - 10^{-s}$ на основі системи повних базисів Ke_i у формі $\bar{c}^{(i)} := \alpha \bar{c}^{(i)} + \beta_{k,i} \bar{s}^{(k,i)}$, $i \in I\left(\bar{r}^{(0)}\right)$, $k \in \{0, 1, \dots, m-1\}$, де $I\left(\bar{r}^{(0)}\right)$ — множина індексів $i \in \{1, \dots, n\}$ упорядкована за неспаданням послідовності $\left|\bar{r}^{(0)}, A\bar{e}_i\right| / \left\|A\bar{e}_i\right\|_2^2$. Вектори

$\bar{s}^{(k,i)}$ формуються з векторів $Ke_i = \left\{\bar{w}_0^{(i)}, \dots, \bar{w}_{m-1}^{(i)}\right\}$ послідовно за правилом: $\bar{s}^{(0,i)} = \gamma \bar{c}^{(i)} + \bar{w}_0^{(i)}$, $\left(A\bar{s}^{(0,i)}, A\bar{c}^{(i)}\right) = 0 \Rightarrow \gamma = -\left(A\bar{w}_0^{(i)}, A\bar{c}^{(i)}\right) / \left\|A\bar{c}^{(i)}\right\|_2^2$;

$$\forall k, j \in [1 : m-1], \bar{s}^{(k,i)} = \gamma \bar{c}^{(i)} + \sum_{j=0}^{k-1} \eta_j \bar{s}^{(j,i)} + \bar{w}^{(k)}, \quad \left(A\bar{s}^{(k,i)}, A\bar{c}^{(i)}\right) = 0,$$

$$\forall j \in [1 : m-1] \quad \left(A\bar{s}^{(k,i)}, A\bar{s}^{(j,i)}\right) = 0, \quad k \neq j, \quad (9)$$

$$\Rightarrow \gamma = -\left(A\bar{w}^{(k)}, A\bar{c}^{(i)}\right) / \left\|A\bar{c}^{(i)}\right\|_2^2, \quad \eta_j = -\left(A\bar{w}^{(k)}, A\bar{s}^{(j,i)}\right) / \left\|A\bar{s}^{(j,i)}\right\|_2^2.$$

Оскільки $\forall k \in [0, 1, \dots, m-1]$ система векторів $\left\{A\bar{s}^{(k,i)}\right\}$ ортогональна і $\left(A\bar{s}^{(k,i)}, A\bar{c}^{(i)}\right) = 0$, то $\bar{s}^{(i)} = \sum_{k=0}^{m-1} \delta_k \bar{s}^{(k,i)}$ при $\delta_k = -\left(\bar{r}^{(0)}, A\bar{s}^{(k,i)}\right) / \left\|A\bar{s}^{(k,i)}\right\|_2^2$ мінімізує норму $\|\cdot\|_2$ вектора нев'язки.

Внутрішній цикл по k завершено. Оскільки вектор $\bar{s}^{(i)}$ A -ортогональний до вектора $\bar{c}^{(i)}$, тому вектор $\bar{c}^{(i)} := \alpha \bar{c}^{(i)} + \beta \bar{s}^{(i)}$ на i -му кроці зовнішнього циклу мінімізує норму вектора нев'язки: $\bar{r}^{(i)} = \bar{r}^{(0)} + \alpha_i A\bar{c}^{(i)} + \beta_i A\bar{s}^{(i)}$,

$$\alpha_i = -\left(\bar{r}^{(0)}, A\bar{c}^{(i)}\right) / \left\|A\bar{c}^{(i)}\right\|_2^2, \quad \beta_i = -\left(\bar{r}^{(0)}, A\bar{s}^{(i)}\right) / \left\|A\bar{s}^{(i)}\right\|_2^2 \quad (10)$$

(це забезпечує строгу монотонність процесу обчислень згладжує похибку обчислень).

Якщо $\cos^2 \left\{\bar{r}^{(0)}, A\bar{c}^{(i)}\right\} < 1 - 10^{-s}$, то продовжити процес для $i \in [1 : n]$. Після завершення циклу по i вектор $\bar{x}^{(1)}$ задаємо у формі

$$\bar{x}^{(1)} = \bar{x}^{(0)} + \xi \bar{c}^{(t)}, \quad \bar{r}^{(1)} = \bar{r}^{(0)} + \xi A\bar{c}^{(t)}, \quad (11)$$

де t — значення з множини $[1: n]$ після завершення процесу. Параметр ξ знайдемо з умови мінімізації вектора нев'язки: $\xi = -(\vec{r}^{(0)}, A\vec{c}^{(t)}) / \|A\vec{c}^{(t)}\|_2^2$. Якщо $\cos^2(\vec{b}^{(1)}, \vec{b}) < 1 - 10^{-s}$, $\vec{b}^{(1)} = \sum_{i=1}^n A_i \cdot \vec{x}_i^{(1)}$, то замінивши $\vec{x}^{(0)} := \vec{x}^{(1)}$ процес повторити для нового наближення $\vec{x}^{(1)}$, нев'язки $\vec{r}^{(1)} = A\vec{x}^{(1)} - \vec{b}$.

Алгоритм побудований для матриць типу «чорний ящик». Якщо змінні системи упорядковані так, що матриця розбивається на укрупнені підсистеми ортогональних векторів-рядків (векторів-стовпців), то метод напрямленого пошуку стає економним і прискорюється швидкість збіжності та зменшується вплив похибок заокруглення на швидкість збіжності. Такими системами є різниці еліптичні рівняння [5; 6; 12].

Теорема 3. Алгоритм оптимізації напрямного вектора (методом мінімізації норми вектора похибок або нев'язок) коректний: збігається, визначає параметри однозначно, мінімізує похибку обчислень (похибку розортогоналізації базисних векторів).

Доведення теореми впливає з формул (7)–(11).

6. Прискорення збіжності ітераційних процесів розв'язування систем $A\vec{x} = \vec{b}$ на основі мінімізації відношення Релея $\rho(A^T A, \vec{x} - \vec{x}^*)$. Нехай $\vec{x} \in K_{\min}$ — наближений розв'язок системи $A\vec{x} = \vec{b}$. Якщо вектор похибки $\vec{\varepsilon} = \vec{y} - \vec{x}^*$ є власним вектором матриці $A^T A$, то обчислений вектор $\vec{p} = A^T \vec{r} = A^T A(\vec{y} - \vec{x}^*) = \lambda(\vec{y} - \vec{x}^*)$ є напрямним вектором і розв'язок системи здійсниться за один крок.

Побудуємо алгоритм мінімізації відношення Релея $\rho(A^T A, \vec{x} - \vec{x}^*) = \|\vec{r}\|_2^2 / \|\vec{x} - \vec{x}^*\|_2^2$ на сфері $\|\vec{x} - \vec{x}^*\|_2^2 = \|\vec{\tilde{x}} - \vec{x}^*\|_2^2$ на основі підпростору Крилова $Kr_4 = \text{span}\{\vec{r}, A\vec{r}, A^2\vec{r}, A^3\vec{r}\} = \text{span}\{\vec{u}_1, \vec{u}_2, \vec{u}_3, \vec{u}_4\}$, $\vec{r} = A\vec{\tilde{x}} - \vec{b}$. Вектор $\vec{\tilde{x}}$, що мінімізує відношення Релея, шукатимемо у формі $\vec{x} = \vec{\tilde{x}} + \alpha A^T \vec{p} + \beta A^T \vec{s} \Rightarrow \vec{r} = \vec{\tilde{r}} + \alpha AA^T \vec{p} + \beta AA^T \vec{s}$, де $\vec{s} = \gamma \vec{u}_1 + \delta \vec{u}_2 + t \vec{u}_3 + \vec{u}_4$, $\vec{p} = \vec{u}_1$, параметри γ, δ знайдемо з умов

$$(A^T \vec{s}, A^T \vec{p}) = 0, (AA^T \vec{s}, AA^T \vec{p}) = 0. \quad (12)$$

Розкриємо вирази $\|\vec{x} - \vec{x}^*\|_2^2$, $\|\vec{r}\|_2^2$, матимемо

$$\varphi = \|\vec{r}\|_2^2 = \|\vec{\tilde{r}}\|_2^2 + 2\alpha(\vec{\tilde{r}}, AA^T \vec{p}) + 2\beta(\vec{\tilde{r}}, AA^T \vec{s}) + \alpha^2 \|AA^T \vec{p}\|_2^2 + \beta^2 \|AA^T \vec{s}\|_2^2.$$

Рівняння сфери $\|\bar{x} - \bar{x}^*\|_2^2 = \|\bar{x} - \bar{x}^*\|_2^2$ еквівалентне заданню еліпса відносно параметрів α, β :

$$L_1: 2\alpha(\bar{r}, \bar{p}) + 2\beta(\bar{r}, \bar{s}) + \alpha^2 \|A^T \bar{p}\|_2^2 + \beta^2 \|A^T \bar{s}\|_2^2 = 0. \quad (13)$$

Оскільки визначник системи (12) відносно параметрів γ, δ відмінний від нуля: $D_0 = \|A^T \bar{u}_1\|_2^2 (AA^T \bar{u}_1, AA^T \bar{u}_2) - \|AA^T \bar{u}_1\|_2^2 (A^T \bar{u}_1, A^T \bar{u}_2) \neq 0$ (для векторів підпростору Крилова), то вектор визначиться з системи (12) у формі $\bar{s} = \bar{s}_0 t + \bar{s}_1$.

Дві неперервні замкнуті опуклі лінії — лінія рівня L_2 функції $\varphi = \|\bar{r}\|_2^2$

$$L_2: 2\alpha(\bar{r}, AA^T \bar{p}) + 2\beta(\bar{r}, AA^T \bar{s}) + \alpha^2 \|AA^T \bar{p}\|_2^2 + \beta^2 \|AA^T \bar{s}\|_2^2 = 0$$

і еліпс L_1 , що мають спільну точку $(\alpha_0, \beta_0) = (0, 0)$, де $\bar{s} = \bar{s}_0 t + \bar{s}_1$, $t \in \mathbb{R}$, або перетинаються, або дотикаються. Якщо лінії перетинаються, то існує найближча точка $(\alpha_1, \beta_1) \in L_1$, в якій опукла функція φ приймає найменше значення: $\varphi(\alpha_1, \beta_1) < \varphi(\alpha_0, \beta_0)$. Щоб довести перетин ліній L_1, L_2 , допустимо протилежне. Нехай $(0, 0)$ є точкою дотику, тоді вектори \bar{u}, \bar{v} з спільним початком у точці $(0, 0)$ і кінцями у центрах ліній L_1, L_2 повинні бути ортогональними: $\cos^2(\bar{u}, \bar{v}) = 0$.

Дослідимо функцію $w(t) = \cos^2\{\bar{u}(t), \bar{v}(t)\}$ по параметру $t \in \mathbb{R}$,

$$\text{де } \bar{u} = \left[-(\bar{r}, \bar{p}) / \|A^T \bar{p}\|_2^2 - (\bar{r}, \bar{s}(t)) / \|A^T \bar{s}(t)\|_2^2 \right]^T,$$

$$\bar{v} = \left[-(\bar{r}, AA^T \bar{p}) / \|AA^T \bar{p}\|_2^2 - (\bar{r}, AA^T \bar{s}(t)) / \|AA^T \bar{s}(t)\|_2^2 \right]^T.$$

Оскільки функція $w(t)$, визначена, неперервна і не є сталою, то існує принаймні одна точка t_0 , для якої $w(t_0) \neq 0$ (отже лінії L_1, L_2 перетинаються).

Пошук точок на еліпсі з найменшим (найбільшим) значенням опуклої функції легко здійснити методами умовного екстремуму.

Існування мінімального значення відношення Релея в області, що є перетином конуса K_{\min} і сфери $\|\bar{x} - \bar{x}^*\|_2 = \|\bar{x} - \bar{x}^*\|_2$ впливає з умови, що область K_{\min} містить сингулярну пряму S_{\min} .

7. Критерій зупинки ітерацій і прийняття наближеного рішення системи $A\vec{x} = \vec{b}$. При розв'язуванні систем $A\vec{x} = \vec{b}$ з погано зумовленими матрицями в умовах машинної арифметики критерії зупинки ітерацій: $\|\vec{x}^{(k+1)} - \vec{x}^{(k)}\| < \varepsilon$ або $\|\vec{r}^{(k+1)}\| < \varepsilon$, або на основі теоретичного розрахунку числа ітерацій, що забезпечують асимптотичну збіжність з похибкою $\|\vec{r}^{(k)}\| < \varepsilon$ втрачають зміст. Область K_{\min} і похибки обчислень стають на перешкоді застосуванню критерія прийняття нормального рішення [8]: $\min_{\vec{x} \in \mathbb{R}^n} \left(\|A\vec{x} - \vec{b}\|_2^2 + \alpha \|\vec{x} - \vec{x}^*\|_2^2 \right)$.

За критерій зупинки ітерацій і прийняття рішення приймемо умову: якщо для деякого наближення $\vec{x}^{(k)} \in \mathbb{R}^n$ виконується нерівність $\cos^2 \left\{ \vec{c} \left(\vec{x}^{(k)}, \vec{b} \right) \right\} \geq 1 - 10^{-t}$, $t > 1$ — пов'язане з машинною точністю, $\vec{c} \left(\vec{x}^{(k)} \right) = \sum_{i=1}^n A_i \cdot \vec{x}_i^{(k)}$, то ітераційний процес зупинити і $\vec{x}^{(k)}$ — наближене рішення системи.

Теорема 4. Функціонал $\cos^2 \left\{ \vec{c} \left(\vec{x}, \vec{b} \right) \right\}$ стійкий до похибок збурень.

Доведення теореми. Нехай

$$\begin{aligned} \vec{x} &= \vec{x}^* + \Delta \vec{x} \Rightarrow \vec{c} \left(\vec{x} \right) = \sum_{i=1}^n A_i \left(x_i^* + \Delta x_i \right), \\ \vec{c} \left(\vec{x} \right) &= \vec{b} + \vec{w} \Rightarrow \cos^2 \left\{ \vec{c}, \vec{b} \right\} = \left(\vec{b} + \vec{w}, \vec{b} \right)^2 / \left(\|\vec{b}\|_2^2 \cdot \|\vec{b} + \vec{w}\|_2^2 \right) = \\ &= (1 + \alpha)^2 / \left(1 + 2\alpha + \beta^2 \right)^2 = \varphi, \text{ де } \alpha = \left(\vec{w}, \vec{b} \right) / \|\vec{b}\|_2^2, \beta = \|\vec{w}\|_2 / \|\vec{b}\|_2. \end{aligned}$$

Максимальне значення функція φ приймає тоді і лише тоді, якщо $\alpha = \beta = 0$. Отже, похибка $\|\vec{w}\|_2 = 0$, що можливо лише за умови $\forall i \in [1:n] \Delta x_i = 0$.

8. Бі-спряжені базиси. Якщо матриця A системи $A\vec{x} = \vec{b}$ близька до виродженої ($|\det A|$ — мале додатне число, $\mathit{cond} A$ — велике число), або A , \vec{b} взяті з експерименту (задані з наближенням), то наближений розв'язок (квазірозв'язок) такої системи можна отримати шляхом мінімізації функціоналу $\varphi = \|A\vec{x} - \vec{b}\|_2^2 + \delta \|\vec{x} - \vec{y}^{(0)}\|_2^2$ (δ — параметр регуляризації [8]). Побудуємо у внутрішньому циклі процесу (7) із базису Ke_i , $i \in [1:n]$ бі-спряжену систему векторів з матрицями A^T та AA^T [13]:

$$\left(AA^T \bar{c}^{(t,j)}, AA^T \bar{c}^{(p,j)} \right) = 0, \quad \left(A^T \bar{c}^{(t,j)}, A^T \bar{c}^{(p,j)} \right) = 0, \quad \forall t \neq p \in \{0, 1, \dots, m-1\}.$$

Тоді мінімум φ існує і єдиний, оптимальні параметри $\alpha_{k,j}$ визначаються за формулами

$$\alpha_{k,j} = -\left(\bar{r}^{(j)}, AA^T \bar{c}^{(k,j)} \right) / \left\| AA^T \bar{c}^{(k,j)} \right\|_2^2 - \delta \left(\bar{\rho}^{(j)}, \bar{c}^{(k,j)} \right) / \left\| A^T \bar{c}^{(k,j)} \right\|_2^2,$$

де $\bar{\rho}^{(j)} = A\bar{x}^{(j)} - A\bar{y}^{(0)}$, $\bar{y}^{(0)} \in \mathbb{R}^n$ — вектор відносно якого шукаємо нормальний розв'язок. Якщо $\bar{y}^{(0)} = \bar{x}^*$, то $\bar{\rho}^{(j)} = \bar{r}^{(j)}$.

Висновки. У роботі розв'язані проблемні задачі теорії лінійних систем: 1) встановлена залежність похибок зашумленої системи від похибки матриці, похибки правої частини, від розв'язку, порядку та величини визначника матриці; 2) побудований багатосаровий алгоритм гауссового виключення, що одночасно обнулює піддіагональні елементи в m стовпцях і мінімізує похибку обчислень та модифікований метод для погано зумовлених матриць типу матриці Гільберта та матриць, що не можуть задовільно масштабуватись; 3) показано, що перепоною для відшукування з високою точністю наближених розв'язків у системах з погано зумовленими матрицями стають області K_{\min} , що містять сингулярні прямі з найменшими власними значеннями; 4) запропонований двоциклічний ітераційний метод напрямленого пошуку для розв'язування систем з матрицями великих порядків на основі системи повних базисів криловського типу, що мінімізує похибку обчислень за рахунок вибору оптимальних параметрів у підпросторах Крилова малих розмірностей і проведення обчислень поза областю K_{\min} ; 5) запропонований алгоритм прискорення збіжності на основі мінімізації відношення Релея на сфері.

Список використаних джерел:

1. Хейгеман Л. Прикладные итерационные методы / Л. Хейгеман, Д. Янг ; пер. с англ. — М. : Мир, 1986. — 448 с.
2. Деммель Дж. Вычислительная линейная алгебра / Дж. Деммель. — М. : Мир, 2001. — 429 с.
3. Райс Дж. Матричные вычисления и математическое обеспечение / Дж. Райс ; пер. с англ. — М. : Мир, 1984. — 264 с.
4. Воеводин В. В. Матрицы и вычисления / В. В. Воеводин, Ю. А. Кузнецов. — М. : Наука, 1984. — 320 с.
5. Патанкар С. Численные методы решения задач теплообмена и динамики жидкости / С. Патанкар. — М. : Энергоатомиздат, 1984. — 125 с.
6. Зверев В. Г. Модифицированный полилинейный метод решения разностных эллиптических уравнений / В. Г. Зверев // ЖВМ и МФ. — 1998. — Т. 38. — № 9. — С. 1553–1562.
7. Зенкевич О. Конечные элементы и аппроксимация / О. Зенкевич, К. Морган. — М. : Мир, 1986. — 318 с.

8. Тихонов А. Н. Методы решения некорректных задач / А. Н. Тихонов, В. Я. Арсенин. — М. : Наука, 1974. — 223 с.
9. Гилмор Р. Прикладная теория катастроф / Р. Гилмор. — М. : Мир, 1984. — 350 с. (кн.1), 282 с. (кн.2).
10. Хорн Р. Матричный анализ / Р. Хорн, Ч. Джонсон ; пер. с англ. — М. : Мир, 1989. — 655 с.
11. Абрамчук В. С. Итерационные методы направленного поиска решения систем $Ax = fc$ сингулярно-естественным упорядочением переменных / В. С. Абрамчук // Доп. НАН України. — 1996. — № 8. — С. 4–8.
12. Абрамчук В. С. Проблеми, методи, алгоритми розв'язування систем лінійних рівнянь з погано зумовленими матрицями / В. С. Абрамчук, І. В. Абрамчук // Математичне та комп'ютерне моделювання. Серія: Фізико-математичні науки : зб. наук. праць. — Кам'янець-Подільський : Кам'янець-Подільський національний університет імені Івана Огієнка, 2016. — Вип. 13. — С. 5–16.
13. Ильин В. П. Методы би-сопряженных направлений в подпространствах Крылова / В. П. Ильин // СибЖИМ. — 2008. — Т. 11, № 4 (36). — С. 47–60.

The problems, that investigates, was related with estimation of errors in the solutions of noised systems and minimizations errors of calculations that arise in Gaussian transformations and acceleration of convergence of the iterative methods.

Key words: *the estimation of solution's error, multilayer method of Gaussian elimination, the local basis of the direct search method, Rayleigh–Ritz quotient.*

Отримано: 23.03.2017

УДК 517.97

С. М. Бак, канд. фіз.-мат. наук, доцент

Вінницький державний педагогічний університет
імені Михайла Коцюбинського, м. Вінниця

ІСНУВАННЯ СТОЯЧИХ ХВИЛЬ В ДИСКРЕТНОМУ НЕЛІНІЙНОМУ РІВНЯННІ ШРЕДІНГЕРА З КУБІЧНОЮ НЕЛІНІЙНІСТЮ НА ДВОВИМІРНІЙ ГРАТЦІ

Стаття присвячена вивченню дискретного нелінійного рівняння типу Шредінгера з кубічною нелінійністю на двовимірній ґратці. Одержано результат про існування стоячих хвиль для таких рівнянь.

Ключові слова: *дискретне нелінійне рівняння Шредінгера, двовимірна ґратка, стоячі хвилі, критичні точки, теорема про зачеплення.*

Вступ. Останнім часом значну увагу приділяють моделям, дискретним за просторовою змінною. Серед рівнянь, які описують такі моделі, найбільш відомими є рівняння ланцюгів осциляторів, дискре-