

WYKŁAD 2 : Klasyczne metody iteracyjne dla układów równań liniowych

1. Wstęp

Wiele zagadnień numerycznych prowadzi do układów równań liniowych o wielkiej liczbie niewiadomych. Typowe zagadnienia z dziedziny mechaniki ośrodków ciągłych (mechanika płynów, dynamika gazów, mechanika odkształcalnych ciał stałych, zagadnienia transportu ciepła i masy, itp.) wymagają rozwiązania (często - wielokrotnego) układów liniowych o rozmiarach od rzędu kilkudziesięciu / kilkuset tysięcy (zagadnienia 2D) do kilkudziesięciu / kilkuset milionów (a nawet więcej!) w złożonych geometrycznie zagadnieniach 3D.

Cechą tych układów liniowych jest to, że ich macierze współczynników zawierają relatywnie mało elementów ni zerowych. O takich macierzach mówimy, że są RZADKIE (ang. sparse matrices). W typowych zastosowaniach współczynnik „nadalności” (sparseness factor) - tj. stosunek liczby elementów ni zerowych do wszystkich elementów - może być naprawdę bardzo mały, szczególnie dla dużych rozmiarów macierzy. Np. przy zastosowaniu metody różnic skończonych do rozwiązania z operatorem Laplace'a w 2D będzie on wynosił $\sim 5/N$ gdzie N to wymiar macierzy (powiemy $N \sim 10^4 \div 10^5$).

2

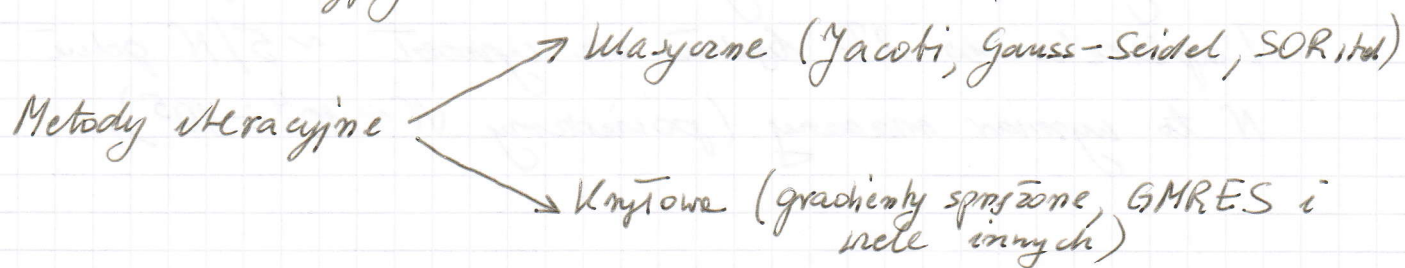
Do rozwiązywania układów liniowych o takich rozmiarach i strukturze nie stosuje się metod „dokładnych” (czyli nieiteracyjnych) typu eliminacji Gaussa (lub jej pochodnych)

Istnieją zasadniczo dwa podstawowe powody:

- 1) koszt numeryczny metod dokładnych - nawet wstawię zaimplementowanych pod kątem efektywnego użycia dla macierzy nadek - jest gigantyczny, tj. proporcjonalny do kwadratu rozmiaru układu (wsp. proporcjonalności może być różnie znany, a jego dokładna wartość zależy od tzw. rozkładu „pasma” macierzy, czyli maksymalnej odległości elementów niezerowych od głównej diagonalnej)
- 2) podczas procesu rozwiązywania potrzebne zwykle znaczące więcej pamięci operacyjnej niż do przewidzenia zeroch niezerowych elementów macierzy. Powodem tego jest efekt pojawiania się niezerowych wartości w miejscach gdzie wartości elementów macierzy były zerowe (tzw. fill-in). Efekt ten może zwiększyć zapotrzebowanie na pamięć nawet o rząd wielkości!

Reasumując: układów liniowe z macierzami rzadkimi rozwiązyjemy metodami iteracyjnymi!

Podział metod iteracyjnych:



③

2. Wybieramy informacje nt. stacjonarnych liniowych procesów iteracyjnych.

Stacjonarnym liniowym procesem iteracyjnym nazywamy proces zadany wzorem:

$$x^{n+1} = Gx^n + W \quad (1)$$

Diagram explaining the components of equation (1):

- x^{n+1} : wektor x w iteracji $n+1$ -szej (indicated by an upward arrow from the text below)
- G : stała, zadana macierz (indicated by an upward arrow from the text below)
- x^n : wektor x w iteracji n -tej (indicated by a curved arrow from the text below)
- W : zadany wektor (indicated by a downward arrow from the text below)

pono tym zadany jest wektor startowy x^0 .

Punktem (wektorem) stałym procesu (1) nazywamy wektor x_* taki, że

$$x_* = Gx_* + W \quad (2)$$

Zauważmy, że $(I - G)x_* = W$. Wyznaczenie punktu stałego iteracji (1) oznacza zatem rozwiązanie układu z macierzą $I - G$ i wektorem prawych stron W (albo zadania równowrotnego).

Zobaczymy od tego zależy właściwość procesu (1).

Wprowadźmy wektor błędu:

$$e^n := x^n - x_* \quad ; \quad n = 0, 1, 2, \dots$$

Zauważmy, że odejmując stronami (1) i (2) otrzymamy

$$x^{n+1} - x_* = G(x^n - x_*) \Rightarrow e^{n+1} = Ge^n \quad (3)$$

(4)

Wektoru x spełniają zatem naszą formułę iteracyjną
z $W=0$.

Udowodnimy następujące.

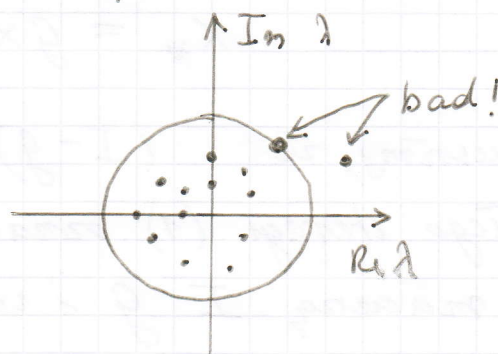
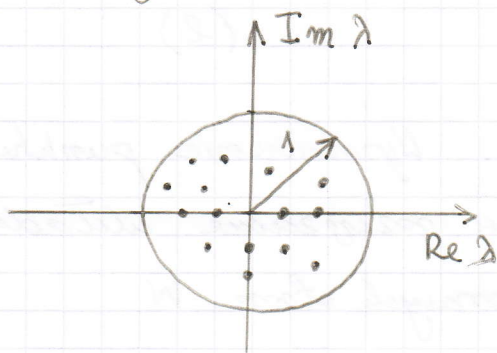
Twierdzenie: W warunkach koniecznym i wystarczającym
zbieżności procesu iteracyjnego

$$x^{n+1} = Gx^n + W$$

do punktu stałego x_* (tj. $x_* = Gx_* + W$) jest, aby
wszystkie wartości własne macierzy G spełniały
warunek

$$|\lambda| < 1. \quad (\text{nierówność ostra!})$$

Komentarz: innymi słowy wymaga się, aby wszystkie
wartości własne G (są to - ma ogół - liczby zespolone)
znajdowały się we wnętrzu koła o promieniu 1



Wielkość $\rho(G) = \max |\lambda(G)|$ nazywamy **PROMIENIEM**
SPEKTRALNYM macierzy G . Twierdzenie mówi, że dla zbieżności
potrzeba i wystarcza, aby $\rho(G) < 1$!

Dowód:

Dowodnimy tylko konieczności warunku. Założymy - wbrew
założeniu - że istnieje wartość własna λ macierzy G
taka, że $|\lambda| \geq 1$. Z teorii macierzy wynika, że
istnieje wówczas wektor (własny) \hat{v} taki, że:

(5)

$$g\hat{v} = \hat{\lambda}\hat{v} \quad (4)$$

Przyjmijmy wektor startowy $x^0 = x_* + \hat{v}$. Wówczas $e^0 = \hat{v}$ i

$$e^1 = ge^0 = g\hat{v} = \hat{\lambda}\hat{v}$$

Dalej $e^2 = ge^1 = g(\hat{\lambda}\hat{v}) = \hat{\lambda}g\hat{v} = \hat{\lambda}^2\hat{v}$ itd.

Ogólnie
$$e^n = \hat{\lambda}^n \hat{v} = \hat{\lambda}^n e^0$$

a stąd, przy dowolnej dobranej normie wektorowej mamy

$$\|e^n\| = |\hat{\lambda}|^n \cdot \|e^0\|.$$

Ponieważ $|\hat{\lambda}| \geq 1$ więc zachodzi dla dowolnego $n > 0$ nierówność $\|e^n\| \geq \|e^0\|$ co oczywiście zaprzecza możliwości, że $\lim_{n \rightarrow \infty} \|e^n\| = 0$. KONIEC

Wyjaśnimy znaczenie wielkości promienia spektralnego $\rho(g)$.

Dla uproszczenia założymy, że macierz g ma pełny układ wektorów własnych $\{v_1, v_2, \dots, v_n\}$, $n = \dim g$.

Niech ponadto $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$ będzie układem odpowiadających tym wektorom wartości własnych, tj.

$$gv_j = \lambda_j v_j, \quad j = 1, 2, \dots, n$$

Ponieważ $\{v_1, \dots, v_n\}$ to układ liniowo niezależny to powyższy wektor będzie można zapisać jako ich kombinację liniową

$$e^0 = x^0 - x_* = \alpha_1 v_1 + \alpha_2 v_2 + \dots + \alpha_n v_n$$

$\alpha_1, \alpha_2, \dots, \alpha_n$ — jednoznacznie określone liczby.

Wówczas

$$e^1 = \sum_{j=1}^n \alpha_j g v_j = \sum_{j=1}^n \alpha_j \lambda_j v_j$$

$$e^2 = \sum_{j=1}^n \alpha_j \lambda_j g v_j = \sum_{j=1}^n \alpha_j \lambda_j^2 v_j \quad \text{itd}$$

Ogólnie:
$$e^n = \sum_{j=1}^n \alpha_j \lambda_j^n v_j$$

Dla dużej wartości n otrzymamy sumę zdominowaną przez składnik odpowiadający największej co do modułu wartości własnej (założamy, że stojący przy tym składnik współczynnik α nie jest „przypadkiem” zerowy). Oznacza to, że dla dużej n

$$e^n \approx \alpha_{\max} \lambda_{\max}^n v_{\max}$$

$$e^{n+1} \approx \alpha_{\max} \lambda_{\max}^{n+1} v_{\max}$$

Stąd:
$$\mu = \frac{\|e^{n+1}\|}{\|e^n\|} \approx \frac{|\alpha_{\max}| \lambda_{\max}^{n+1}(g) \|v_{\max}\|}{|\alpha_{\max}| \lambda_{\max}^n(g) \|v_{\max}\|} = \lambda(g)$$

gdzie μ nazywamy asymptotycznym tempem zbieżności. Jak widać im mniejsze $\lambda(g)$ tym lepiej!

Zadanie pomocnicze zmniejsza do uzyskania jak największej redukcji promienia spektralnego $\lambda(g)$ nazywamy poprawianiem warunkowania procesu iteracyjnego (ang. preconditioning)

(7)

3. Klasyczne metody iteracyjne

a) metoda Jacobiego

$AX = b$ - układ do rozwiązania.

Załóżmy, że dla dowolnego $j = 1, 2, \dots, n = \dim A$ elementy diagonalne $a_{jj} \neq 0$. Jeśli było macierz A jest nieosobliwa da się to osiągnąć metodą zmiany numeracji równań i/lub niewiadomych.

Definiujemy macierze:

$$U: u_{ij} = \begin{cases} a_{ij}, & i < j \\ 0, & i \geq j \end{cases}$$

$$L: l_{ij} = \begin{cases} 0, & i \leq j \\ a_{ij}, & i > j \end{cases}$$

$$D: d_{ij} = \begin{cases} a_{ii}, & i = j \\ 0, & i \neq j \end{cases}$$

Obrazowo ...

$$A = U + L + D$$

Skoro $A = U + L + D$ to układ przyjmuje postać

$$(U + L + D)x = b$$

$$Dx = -(U + L)x + b$$

$$x = -D^{-1}(U + L)x + D^{-1}b \quad (5)$$

⑧

Otrzymaliśmy układ równań dla wektora (punktu) stałego dla następującego procesu iteracyjnego:

$$x^{k+1} = -D^{-1}(U+L)x^k + D^{-1}b \quad (6)$$

Jest to właśnie metoda Jacobiego. Ułamy dla niej

$$G_J = -D^{-1}(U+L), \quad w_J = D^{-1}b. \quad (7)$$

Rozpisując (6) na składowe, otrzymamy pętlę do wykonania w każdej iteracji metody Jacobiego

```
for i ∈ {1, 2, ..., n} do:
    
$$x_i^{k+1} = \frac{1}{a_{ii}} \left( b_i - \sum_{j=1}^{i-1} a_{ij} x_j^k - \sum_{j=i+1}^n a_{ij} x_j^k \right) \quad (8)$$

end;
```

Zauważmy, że sposób zapisania pętli w (8) ma sugerować, że kolejność indeksu i ma znaczenie – aktualizacja każdego elementu wektora x^{k+1} jest kompletnie niezależna od pozostałych elementów. Wynika stąd wniosek, że pętlę można wykonać RÓWNOLEGLE!

b) metoda Gaussa - Seidela

Zauważmy, że w pętli (8) indeks i zmieniany jest „zwykcyjnie” tzn. rośnie od 1 do n . Wówczas elementy aktualizowanego wektora x^{k+1} o numerach od 1 do $i-1$ mogą być od razu wykorzystane do obliczenia x_i^{k+1} (bo zostały policzone wcześniej w tej właśnie iteracji)

9

Otrzymaliśmy proces iteracyjny opisany następująco

$$\begin{aligned} &\text{for } i = 1:n \text{ do} \\ &\quad x_i^{k+1} = \frac{1}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{k+1} - \sum_{j=i+1}^n a_{ij} x_j^k \right) \quad (9) \\ &\text{end;} \end{aligned}$$

Jest to metoda iteracyjna Gaussa-Seidela.

Wyprowadzimy postać maciernowo-wektorową tej metody.

Z formuły (9) wynika, że

$$x^{k+1} = -D^{-1} [L x^{k+1} + U x^k] + D^{-1} b$$

czyli

$$D x^{k+1} = -L x^{k+1} - U x^k + b$$

$$(D+L) x^{k+1} = -U x^k + b$$

$$x^{k+1} = -\underbrace{(D+L)^{-1}}_{||} U x^k + (D+L)^{-1} b$$

Mamy zatem:

$$x^{k+1} = G_{gs} x^k + W_{gs} \quad (10)$$

gdzie

$$\begin{cases} G_{gs} = -(D+L)^{-1} U \\ W_{gs} = (D+L)^{-1} b \end{cases} \quad (11)$$

Zauważmy, że zmieniając kolejność pętli w (9) na odwrotną (tj. for $i = n:1$ step -1 do) otrzymamy alternatywny wariant metody G-S o następującej postaci maciernowo-wektorowej.

$$\begin{cases} x^{k+1} = \tilde{G}_{gs} x^k + \tilde{W}_{gs} \\ \tilde{G}_{gs} = -(D+U)^{-1} L, \quad \tilde{W}_{gs} = (D+U)^{-1} b \end{cases} \quad (12)$$

c) Metoda nadrelaksacji (SOR - Successive Over Relaxation)

Najlepszym ulepszeniem jest wprowadzenie nadrelaksacji.

Punktem wyjścia jest metoda Gaussa-Seidela w formie rekurencyjnym wzorem (9). Metoda SOR przedstawia się następująco:

$$\begin{aligned} &\text{for } i = 1:n \text{ do} \\ &\quad t_i = \frac{1}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{k+1} - \sum_{j=i+1}^n a_{ij} x_j^k \right) \\ &\quad x_i^{k+1} = \omega t_i + (1-\omega) x_i^k \\ &\text{end;} \end{aligned} \quad (13)$$

W metodzie nadrelaksacji, parametr ω przyjmuje wartość większą od 1. Zauważmy, że dla $\omega = 1$ metoda SOR sprowadza się do metody G-S.

Celem stosowania nadrelaksacji jest przyspieszenie zbieżności procesu iteracyjnego. Wartość parametru ω wpływa na wielkość promienia spektralnego macierzy G_{SOR} odpowiadającej metodzie SOR - wybierając odpowiednio ω można ten promień zmniejszyć (w porównaniu z metodą G-S) i znacząco przyspieszyć zbieżność. Dla pewnych klas macierzy (opisywanych w zastosowaniach) istnieją twierdzenia o optymalnym wyborze parametru ω .

Zadanie: wyproceduj macierzowo-wektorową postać metody SOR

$$x^{k+1} = G_{SOR} x^k + W_{SOR} \quad (14)$$

Odp. $G_{SOR} = (D - \omega L)^{-1} [(1-\omega)D + \omega U]$, $W_{SOR} = \omega (D - \omega L)^{-1} b$.

4. Zbieżność klasycznych metod iteracyjnych

Ogólne kryterium zbieżności stacjonarnego procesu iteracyjnego $x^{k+1} = Gx^k + W$ wymaga, aby promień spektralny macierzy G był mniejszy od jedności. W praktyce, potrzebujemy kryteriów dotyczących bezpośrednio macierzy układu $AX=b$, a nie macierzy procesu iteracyjnego. Innymi słowy, interesują nas warunki, które powinna spełniać macierz A takie, aby macierz G tej czy innej metody iteracyjnej miała promień spektralny mniejszy od 1.

Zacznijmy od zdefiniowania kilku własności.

DEF: Ułóżmy, że macierz kwadratowa A o wymiarze n jest:

1) diagonalnie dominująca jeżeli

$$|a_{ii}| \geq \sum_{j \neq i} |a_{ij}|, \text{ dla } i = 1, 2, \dots, n \quad (15)$$

Jeżeli powyższa nierówność jest OSTRA dla każdej wartości indeksu i to macierz A jest SILNIE diagonalnie dominująca

2) redukowalna jeśli istnieje taka macierz permutacyjna P (jest to macierz jednostkowa o poprzeszaniach kolumnach lub - co równoważne - wierszów), że:

$$PAP^T = \begin{bmatrix} B_{11} & B_{12} \\ 0 & B_{22} \end{bmatrix} \rightarrow \text{blok kwadratowy}$$

i nieredukowalna jeśli takiej macierzy P nie ma.

- 3) nieredukowalna i diagonalnie dominująca jeśli jest diagonalnie dominująca, nieredukowalna i przynajmniej dla jednej wartości „i” nierówność (15) jest ostra.
- 4) nieujemna (co piszemy $A \geq 0$) jeśli wszystkie elementy macicy są nieujemne, i dodatnia (co piszemy $A > 0$) jeśli wszystkie elementy są dodatnie.
- 5) M-macierz jeśli wszystkie elementy pozadiagonalne ($a_{ij}, i \neq j$) są niedodatnie i jednocześnie macierz odwrotna do A jest nieujemna ($A^{-1} \geq 0$).

Podamy bez dowodu kluczowe twierdzenia:

Tw1: Niech macierz A będzie macierzą nieredukowalną i diagonalnie dominującą lub silnie diagonalnie dominującą. Wówczas metody Jacobi'ego i Gaussa-Seidela są zbieżne do rozwiązania $AX=b$ przy dowolnym wektorze startowym x^0 .

Tw2: Załóżmy, że macierz A jest M-macierzą. Wówczas metody Jacobi'ego i Gaussa-Seidela są zbieżne do rozwiązania układu $AX=b$ przy dowolnym wektorze startowym x^0 .

Dalej, mamy ciekawe twierdzenie o metodzie SOR.

Tw3 (Kahan): Załóżmy, że macierz A ma rzeczywiste elementy diagonalne ($a_{ii} \neq 0, i = 1, 2, \dots, n$). Wówczas promień spektralny macierzy G_{SOR} spełnia nierówność $\rho(G_{SOR}) \geq |\omega - 1|$.

Ponieważ warunkiem zbieżności jest, aby $\rho(y_{SOR}) < 1$ to musi zachodzić nierówność

$$|\omega - 1| < 1 \Rightarrow \omega \in (0, 2) \quad (16)$$

Wniosek: Warunkiem koniecznym zbieżności metody SOR jest aby parametr ω należał do przedziału $(0, 2)$ (w praktyce $\omega \in [1, 2)$).

Ważnym przypadkiem układu liniowego jest przypadek, w którym macierz A jest (jednocześnie):

- symetryczna ($A = A^T$, czyli $a_{ij} = a_{ji}$, $i, j = 1, 2, \dots, n$)
- dodatnio określona tj. dla dowolnego wektora $x \neq 0$ ma miejsce nierówność $x^T A x = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j > 0$.

Mamy dwa ważne twierdzenia:

Tw4 (Ostrowskiego-Reicha): Jeżeli macierz A jest macierzą symetryczną i nieosobliwą (tj. $\det A \neq 0$), posiadającą dodatnie elementy diagonalne (tj. $a_{ii} > 0$, $i = 1, \dots, n$) oraz $\omega \in (0, 2)$ to metoda iteracyjna SOR zbiega do rozwiązania układu liniowego $Ax = b$ przy dowolnym wektorze startowym x^0 wtedy i tylko wtedy, gdy macierz A jest dodatnio określona.

Wniosek: Dla macierzy symetrycznej i dodatnio określonej metoda Gaussa-Seidela jest zbieżna do rozwiązania przy dowolnym wektorze startowym x^0 .

Ciekawe, że symetria i dodatnia określoność NIE WYSTARCZY do zapewnienia zbieżności metody Jacobiego. Dla tej metody mamy nieco „bardziej wymagające” twierdzenie

TW 5: Załóżmy, że macierz $A = D + B$ jest symetryczna (tj. $B = L + U = L + L^T$) i macierz diagonalna D

zawiera wyłącznie liczby dodatnie ($a_{ii} > 0, i = 1, 2, \dots, n$).

Wówczas metoda iteracyjna Jacobiego zbiega do rozwiązania układu $AX = b$ przy dowolnym wektorze startowym x^0 wtedy i tylko wtedy, gdy zarówno A jak i macierz $\tilde{A} = D - B$ są macierzami dodatnio określonymi.

5. UWAGI nt. kryterium zatrzymywania iteracji

Każdy proces iteracyjny musi mieć określone kryterium zakończenia (stopu). Naturalnym kryterium stopu wydaje się być wielkość normy wektora błędu (np. norma maximum czyli $\|e\| := \max_{i=1, \dots, n} |e_i|$)

$$\|e^k\| = \|x^k - x_*\|$$

Problem w tym, że oczywiście x_* nie jest a priori znane!

W praktyce badamy zmianę rozwiązania w kolejnych iteracjach i stosujemy kryterium

$$\|x^{k+1} - x^k\| \leq \varepsilon \quad (\varepsilon - \text{mała liczba, np. } \varepsilon = 10^{-6})$$

ewentualnie

$$\|x^{k+1} - x^k\| \leq \tilde{\varepsilon} \cdot \|x^k\|$$

Oba powyższe kryteria mogą być zawodne - w przypadku powolnej zbieżności, zmiany rozwiązania z reguły są niewielkie. Z tego powodu, po stwierdzeniu spełnienia jednej z powyższych nierówności, warto jest sprawdzić ile wynosi szacowana wartość residuum układu, tj.

$$\frac{1}{\|b\|} \|b - Ax^k\|.$$