

# Wstęp do metod numerycznych

## SVD, metody iteracyjne i metoda gradientów sprzężonych

P. F. Góra

<http://th-www.if.uj.edu.pl/zfs/gora/>

2011

## Współczynnik uwarunkowania macierzy symetrycznej

**Twierdzenie 1.** Niech  $\mathbf{A} \in \mathbb{R}^{N \times N}$  będzie macierzą symetryczną,  $\mathbf{A} = \mathbf{A}^T$ , i niech liczby  $\{\lambda_i\}_{i=1}^N$  będą jej wartościami własnymi. Jeżeli  $\det \mathbf{A} \neq 0$ , współczynnik uwarunkowania tej macierzy spełnia

$$\kappa = \frac{\max_i |\lambda_i|}{\min_i |\lambda_i|}. \quad (1)$$

*Dowód.* W celu udowodnienia tego twierdzenia obliczmy normę macierzy  $\mathbf{A}$ . Ponieważ jest to macierz symetryczna i rzeczywista, jej wartości własne są rzeczywiste, natomiast jej unormowane wektory własne tworzą bazę w  $\mathbb{R}^N$ . Oznaczmy przez  $\mathbf{y}_i$  jej  $i$ -ty wektor własny,  $\mathbf{A}\mathbf{y}_i = \lambda_i\mathbf{y}_i$ . Każdy

wektor  $\mathbf{x} \in \mathbb{R}^N$ ,  $\|\mathbf{x}\| = 1$ , można przedstawić jako kombinację liniową

$$x = \sum_{i=1}^N \alpha_i y_i, \quad (2)$$

przy czym warunek unormowania prowadzi do następującej więzi na współczynniki tej kombinacji:

$$\sum_{i=1}^N \alpha_i^2 = 1. \quad (3)$$

Obliczmy teraz

$$\begin{aligned}\|\mathbf{Ax}\|^2 &= \left\| \mathbf{A} \sum_{i=1}^N \alpha_i \mathbf{y}_i \right\|^2 = \left\| \sum_{i=1}^N \alpha_i \mathbf{A} \mathbf{y}_i \right\|^2 = \left\| \sum_{i=1}^N \alpha_i \lambda_i \mathbf{y}_i \right\|^2 \\ &\leq \sum_{i=1}^N \|\alpha_i \lambda_i \mathbf{y}_i\|^2 = \sum_{i=1}^N \alpha_i^2 \lambda_i^2 \leq \sum_{i=1}^N \alpha_i^2 \max_i \lambda_i^2 \\ &= \max_i \lambda_i^2 \sum_{i=1}^N \alpha_i^2 = \left( \max_i |\lambda_i| \right)^2\end{aligned}\tag{4}$$

Widzimy zatem, iż  $\forall \mathbf{x} \in \mathbb{R}^N$ ,  $\|\mathbf{x}\|^2 = 1$ , zachodzi  $\|\mathbf{A}\mathbf{x}\| \leq \max_i |\lambda_i|$ , a zatem na mocy definicji normy indukowanej,  $\|\mathbf{A}\| = \max_i |\lambda_i|$ .

Ponieważ  $\det \mathbf{A} \neq 0$ , macierz  $\mathbf{A}^{-1}$  istnieje, jest symetryczna i rzeczywista, a jej wartościami własnymi są liczby  $\{1/\lambda_i\}_{i=1}^N$ . Zupełnie analogicznie dowodzimy, iż  $\|\mathbf{A}^{-1}\| = \max_i (1/|\lambda_i|) = 1 / \left( \min_i |\lambda_i| \right)$ , skąd natychmiast wynika teza (1).  $\square$

## Singular Value Decomposition

**Twierdzenie 2.** Dla każdej macierzy  $\mathbf{A} \in \mathbb{R}^{M \times N}$ ,  $M \geq N$ , istnieje rozkład

$$\mathbf{A} = \mathbf{U} [\text{diag}(w_i)] \mathbf{V}^T, \quad (5)$$

gdzie  $\mathbf{U} \in \mathbb{R}^{M \times N}$  jest macierzą kolumnowo ortogonalną,  $\mathbf{V} \in \mathbb{R}^{N \times N}$  jest macierzą ortogonalną oraz  $w_i \in \mathbb{R}$ ,  $i = 1, \dots, N$ . Rozkład ten nazywamy rozkładem względem wartości osobliwych (*Singular Value Decomposition, SVD*). Jeżeli  $M = N$ , macierz  $\mathbf{U}$  jest macierzą ortogonalną.

## Jądro i zasięg operatora

Niech  $\mathbf{A} \in \mathbb{R}^{M \times N}$ . *Jądrem operatora  $\mathbf{A}$*  nazywam

$$\text{Ker } \mathbf{A} = \{\mathbf{x} \in \mathbb{R}^N : \mathbf{A}\mathbf{x} = \mathbf{0}\}. \quad (6)$$

*Zasięgiem operatora  $\mathbf{A}$*  nazywam

$$\text{Range } \mathbf{A} = \{\mathbf{y} \in \mathbb{R}^M : \exists \mathbf{x} \in \mathbb{R}^N : \mathbf{A}\mathbf{x} = \mathbf{y}\}. \quad (7)$$

Jądro i zasięg operatora są przestrzeniami liniowymi. Jeśli  $M = N < \infty$ ,  
 $\dim(\text{Ker } \mathbf{A}) + \dim(\text{Range } \mathbf{A}) = N$ .

## Sens SVD

Sens SVD najlepiej widać w przypadku, w którym co najmniej jedna z wartości  $w_i = 0$ . Dla ustalenia uwagi przyjmijmy  $w_1 = 0$ ,  $w_{i \neq 1} \neq 0$ .

Po pierwsze, co to jest  $\mathbf{z} = [z_1, z_2, \dots, z_n]^T = \mathbf{V}^T \mathbf{x}$ ? Ponieważ  $\mathbf{V}$  jest macierzą ortogonalną,  $\mathbf{z}$  jest rozkładem wektora  $\mathbf{x}$  w bazie kolumn macierzy  $\mathbf{V}$ . Korzystając z (5), dostajemy

$$\mathbf{Ax} = \mathbf{U} [\text{diag}(w_i)] \mathbf{V}^T \mathbf{x} = \mathbf{U} [\text{diag}(0, w_2, \dots, w_N)] \mathbf{z} = \mathbf{U} \begin{bmatrix} 0 \\ w_2 z_2 \\ \vdots \\ w_N z_N \end{bmatrix}. \quad (8)$$

Wynikiem ostatniego mnożenia będzie pewien wektor z przestrzeni  $\mathbb{R}^M$ . Ponieważ pierwszym elementem wektora  $[0, w_2 z_2, \dots, w_N z_N]^T$  jest zero, **wynik ten nie zależy od pierwszej kolumny macierzy  $\mathbf{U}$** . Widzimy zatem, że **kolumny macierzy  $\mathbf{U}$ , odpowiadające niezerowym współczynnikom  $w_i$ , stanowią bazę w zasięgu operatora  $\mathbf{A}$** .

Co by zaś się stało, gdyby  $\mathbf{x}$  był równoległy do wektora stanowiącego pierwszą kolumnę  $\mathbf{V}$ ? Wówczas  $\mathbf{z} = 0$ , a wobec tego  $\mathbf{Ax} = 0$ . Ostatecznie więc widzimy, że **kolumny macierzy  $\mathbf{V}$ , odpowiadające zerowym współczynnikom  $w_i$ , stanowią bazę w jądrze operatora  $\mathbf{A}$** .

## SVD i odwrotność macierzy

Niech  $\mathbf{A} \in \mathbb{R}^{N \times N}$ . Zauważmy, że  $|\det \mathbf{A}| = \prod_{i=1}^N w_i$ , a zatem  $\det \mathbf{A} = 0$  wtedy i tylko wtedy, gdy co najmniej jeden  $w_i = 0$ . Niech  $\det \mathbf{A} \neq 0$ . Wówczas równanie  $\mathbf{A}\mathbf{x} = \mathbf{b}$  ma rozwiązanie postaci

$$\mathbf{x} = \mathbf{A}^{-1}\mathbf{b} = \mathbf{V} [\text{diag}(w_i^{-1})] \mathbf{U}^T \mathbf{b}. \quad (9)$$

Niech teraz  $\det \mathbf{A} = 0$ . Równanie  $\mathbf{A}\mathbf{x} = \mathbf{b}$  *także* ma rozwiązanie, o ile tylko  $\mathbf{b} \in \text{Range } \mathbf{A}$ . Rozwiązanie to dane jest wzorem

$$\mathbf{x} = \tilde{\mathbf{A}}^{-1}\mathbf{b} = \mathbf{V} [\text{diag}(\tilde{w}_i^{-1})] \mathbf{U}^T \mathbf{b}. \quad (10a)$$

gdzie

$$\tilde{w}_i^{-1} = \begin{cases} w_i^{-1} & \text{gdy } w_i \neq 0, \\ 0 & \text{gdy } w_i = 0. \end{cases} \quad (10b)$$

## SVD i współczynnik uwarunkowania

**Twierdzenie 3.** Jeżeli macierz  $\mathbf{A} \in \mathbb{R}^{N \times N}$  posiada rozkład (5) oraz  $\det \mathbf{A} \neq 0$ , jej współczynnik uwarunkowania spełnia

$$\kappa = \frac{\max_i |w_i|}{\min_i |w_i|}. \quad (11)$$

Jeśli macierz jest źle uwarunkowana, ale *formalnie* odwracalna, numeryczne rozwiązanie równania  $\mathbf{A}\mathbf{x} = \mathbf{b}$  może być zdominowane przez wzmocniony błąd zaokrąglenia. Aby tego uniknąć, często zamiast (bezużytecznego!) rozwiązania dokładnego (9), używa się *przybliżonego* (i użytecznego!) rozwiązania w postaci (10) z następującą modyfikacją

$$\tilde{w}_i^{-1} = \begin{cases} w_i^{-1} & \text{gdy } |w_i| > \tau, \\ 0 & \text{gdy } |w_i| \leq \tau, \end{cases} \quad (12)$$

gdzie  $\tau$  jest pewną zadaną tolerancją.

## Metody iteracyjne

W metodach dokładnych otrzymane rozwiązanie jest dokładne z dokładnością do błędów zaokrąglenia, które, dodajmy, dla układów źle uwarunkowanych mogą być *znaczne*.

W metodach iteracyjnych rozwiązanie dokładne otrzymuje się, teoretycznie, w granicy nieskończenie wielu kroków — w praktyce liczymy na to, że po skończonej (i niewielkiej) ilości kroków zbliżymy się do wyniku ścisłego w granicach błędu zaokrąglenia.

Rozpatrzmy układ równań:

$$a_{11}x_1 + a_{12}x_2 + a_{13}x_3 = b_1 \quad (13a)$$

$$a_{21}x_1 + a_{22}x_2 + a_{23}x_3 = b_2 \quad (13b)$$

$$a_{31}x_1 + a_{32}x_2 + a_{33}x_3 = b_3 \quad (13c)$$

Przepiszmy ten układ w postaci

$$x_1 = (b_1 - a_{12}x_2 - a_{13}x_3)/a_{11} \quad (14a)$$

$$x_2 = (b_2 - a_{21}x_1 - a_{23}x_3)/a_{22} \quad (14b)$$

$$x_3 = (b_3 - a_{31}x_1 - a_{32}x_2)/a_{33} \quad (14c)$$

Gdyby po prawej stronie (14) były “stare” elementy  $x_j$ , a po lewej “nowe”, dostalibyśmy metodę iteracyjną

$$x_i^{(k+1)} = \left( b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k)} - \sum_{j=i+1}^N a_{ij}x_j^{(k)} \right) / a_{ii} \quad (15)$$

Górny indeks  $x^{(k)}$  oznacza, że jest to przybliżenie w  $k$ -tym kroku. Jest to tak zwana **metoda Jacobiego**.

Zauważmy, że w metodzie (15) nie wykorzystuje się najnowszych przybliżeń: Powiedzmy, obliczając  $x_2^{(k+1)}$  korzystamy z  $x_1^{(k)}$ , mimo iż znane jest już wówczas  $x_1^{(k+1)}$ . (Za to metodę tę łatwo można zrównoleglić.) Sugeruje to następujące ulepszenie:

$$x_i^{(k+1)} = \left( b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k+1)} - \sum_{j=i+1}^N a_{ij}x_j^{(k)} \right) / a_{ii} \quad (16)$$

Jest to tak zwana **metoda Gaussa-Seidela**.

Jeżeli macierz  $A = \{a_{ij}\}$  jest rzadka, obie te metody iteracyjne będą efektywne *tylko i wyłącznie* wówczas, gdy we wzorach (15), (16) uwzględnimy ich strukturę, to jest uniknie redundantnych mnożeń przez zera.

## Trochę teorii

Metody Jacobiego i Gaussa-Seidela należą do ogólnej kategorii

$$\mathbf{M}\mathbf{x}^{(k+1)} = \mathbf{N}\mathbf{x}^{(k)} + \mathbf{b} \quad (17)$$

gdzie  $\mathbf{A} = \mathbf{M} - \mathbf{N}$  jest *podziałem (splitting)* macierzy. Dla metody Jacobiego  $\mathbf{M} = \mathbf{D}$  (część diagonalna),  $\mathbf{N} = -(\mathbf{L} + \mathbf{U})$  (części pod- i ponad-diagonalne, bez przekątnej). Dla metody Gaussa-Seidela  $\mathbf{M} = \mathbf{D} + \mathbf{L}$ ,  $\mathbf{N} = -\mathbf{U}$ . Rozwiązanie równania  $\mathbf{A}\mathbf{x} = \mathbf{b}$  jest punktem stałym iteracji (17).

**Definicja** *Promieniem spektralnym* (diagonalizowalnej) macierzy  $G$  nazywam

$$\rho(G) = \max\{|\lambda| : \exists y \neq 0 : Gy = \lambda y\} \quad (18)$$

**Twierdzenie 4.** *Iteracja (17) jest zbieżna jeśli  $\det M \neq 0$  oraz  $\rho(M^{-1}N) < 1$ .*

*Dowód.* Przy tych założeniach iteracja (17) jest odwzorowaniem zwężającym. □

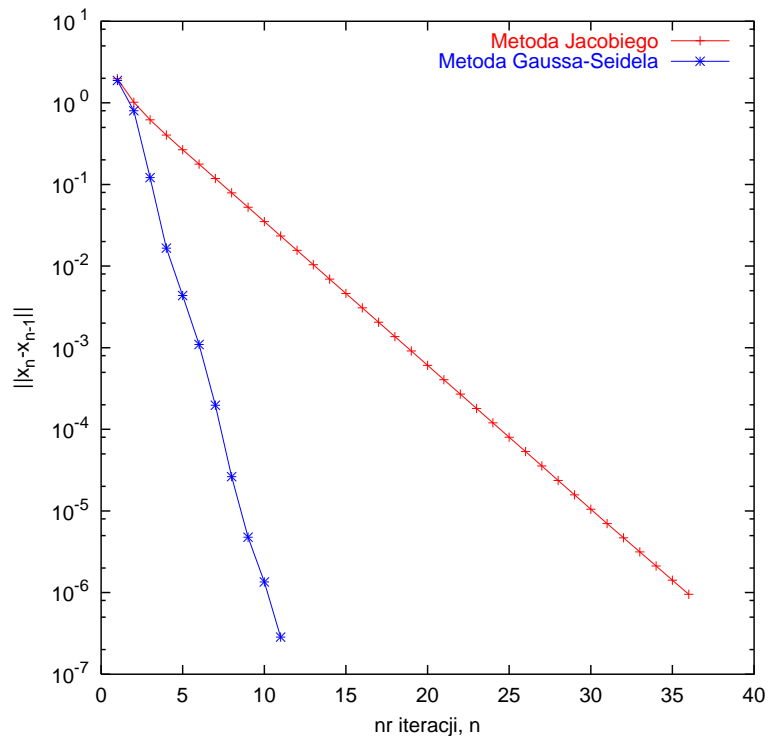
**Twierdzenie 5.** *Metoda Jacobiego jest zbieżna jeśli macierz  $A$  jest silnie diagonalnie dominująca.*

**Twierdzenie 6.** *Metoda Gaussa-Seidela jest zbieżna jeśli macierz  $A$  jest symetryczna i dodatnio określona.*

## Przykład

Rozwiązujemy układ równań:

$$\begin{array}{rcccccc} 3x & + & y & + & z & = & 1 \\ x & + & 3y & + & z & = & 1 \\ x & + & y & + & 3z & = & 1 \end{array}$$



## SOR

Jeśli  $\rho(\mathbf{M}^{-1}\mathbf{N})$  w metodzie Gaussa-Seidela jest bliskie jedności, zbieżność

metody jest bardzo wolna. Można próbować ją poprawić:

$$x_i^{(k+1)} = w \left( b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k+1)} - \sum_{j=i+1}^N a_{ij}x_j^{(k)} \right) / a_{ii} + (1-w)x_i^{(k)}, \quad (19)$$

gdzie  $w \in \mathbb{R}$  jest *parametrem relaksacji*. Metoda ta zwana jest *successive over-relaxation*, SOR. W postaci macierzowej

$$\mathbf{M}_w \mathbf{x}^{(k+1)} = \mathbf{N}_w \mathbf{x}^{(k)} + w \mathbf{b} \quad (20)$$

$\mathbf{M}_w = \mathbf{D} + w\mathbf{L}$ ,  $\mathbf{N}_w = (1-w)\mathbf{D} - w\mathbf{U}$ . *Teoretycznie* należy dobrać takie  $w$ , aby zminimalizować  $\rho(\mathbf{M}_w^{-1}\mathbf{N}_w)$ .

## Metoda gradientów sprzężonych — motywacja

Rozważmy funkcję  $f : \mathbb{R}^N \rightarrow \mathbb{R}$

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x} + c, \quad (21)$$

gdzie  $\mathbf{x}, \mathbf{b} \in \mathbb{R}^N$ ,  $c \in \mathbb{R}$ ,  $\mathbf{A} = \mathbf{A}^T \in \mathbb{R}^{N \times N}$  jest *symetryczna i dodatnio określona*. Przy tych założeniach, funkcja (21) ma dokładnie jedno minimum, będące zarazem minimum globalnym. Szukanie minimów dodatnio określonych form kwadratowych jest (względnie) łatwe i z praktycznego punktu widzenia ważne. Minimum to leży w punkcie spełniającym

$$\nabla f = 0. \quad (22)$$

Obliczmy

$$\begin{aligned}\frac{\partial f}{\partial x_i} &= \frac{1}{2} \frac{\partial}{\partial x_i} \sum_{j,k} A_{jk} x_j x_k - \frac{\partial}{\partial x_i} \sum_j b_j x_j + \underbrace{\frac{\partial c}{\partial x_i}}_0 \\ &= \frac{1}{2} \sum_{j,k} A_{jk} \left( \underbrace{\frac{\partial x_j}{\partial x_i}}_{\delta_{ij}} x_k + x_j \underbrace{\frac{\partial x_k}{\partial x_i}}_{\delta_{ik}} \right) - \sum_j b_j \underbrace{\frac{\partial x_j}{\partial x_i}}_{\delta_{ij}} \\ &= \frac{1}{2} \sum_k A_{ik} x_k + \frac{1}{2} \sum_j A_{ji} x_j - b_i = \frac{1}{2} \sum_k A_{ik} x_k + \frac{1}{2} \sum_j A_{ij} x_j - b_i \\ &= (\mathbf{Ax} - \mathbf{b})_i .\end{aligned}\tag{23}$$

Widzimy zatem, że funkcja (21) osiąga minimum w punkcie, w którym zachodzi

$$\mathbf{Ax} - \mathbf{b} = 0 \Leftrightarrow \mathbf{Ax} = \mathbf{b}. \quad (24)$$

Rozwiązywanie układu równań liniowych (24) z macierzą symetryczną, dodatnio określoną jest równoważne poszukiwaniu minimum dodatnio określonej formy kwadratowej.

Przypuśćmy, że macierz  $\mathbf{A}$  jest przy tym *rzadka* i duża (lub co najmniej średnio-duża). Wówczas metoda gradientów sprzężonych jest godną uwagi metodą rozwiązywania (24)

## Metoda gradientów sprzężonych, *Conjugate Gradients*, CG

$\mathbf{A} \in \mathbb{R}^{N \times N}$  symetryczna, dodatnio określona,  $x_1$  — początkowe przybliżenie rozwiązania równania (24),  $0 < \varepsilon \ll 1$ .

$$\begin{aligned} & \mathbf{r}_1 = \mathbf{b} - \mathbf{A}\mathbf{x}_1, \mathbf{p}_1 = \mathbf{r}_1 \\ & \mathbf{while} \quad \|\mathbf{r}_k\| > \varepsilon \\ & \quad \alpha_k = \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A}\mathbf{p}_k} \\ & \quad \mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{A}\mathbf{p}_k \\ & \quad \beta_k = \frac{\mathbf{r}_{k+1}^T \mathbf{r}_{k+1}}{\mathbf{r}_k^T \mathbf{r}_k} \\ & \quad \mathbf{p}_{k+1} = \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k \\ & \quad \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k \\ & \mathbf{end} \end{aligned} \tag{25}$$

Wówczas zachodzą twierdzenia:

**Twierdzenie 7.** Ciągi wektorów  $\{\mathbf{r}_k\}$ ,  $\{\mathbf{p}_k\}$  spełniają następujące zależności:

$$\mathbf{r}_i^T \mathbf{r}_j = 0, \quad i > j, \quad (26a)$$

$$\mathbf{r}_i^T \mathbf{p}_j = 0, \quad i > j, \quad (26b)$$

$$\mathbf{p}_i^T \mathbf{A} \mathbf{p}_j = 0, \quad i > j. \quad (26c)$$

**Twierdzenie 8.** Jeżeli  $\mathbf{r}_M = 0$ , to  $\mathbf{x}_M$  jest ścisłym rozwiązaniem równania (24).

*Dowód.* Oba (sic!) dowody przebiegają indukcyjnie. □

Ciąg  $\{x_k\}$  jest w gruncie rzeczy “pomocniczy”, nie bierze udziału w iteracjach, służy tylko do konstruowania kolejnych przybliżeń rozwiązania.

Istotą algorytmu jest konstruowanie dwu ciągów wektorów spełniających zależności (26). Wektory  $\{r_k\}$  są wzajemnie prostopadłe, a zatem *w arytmetyce dokładnej*  $r_{N+1} = 0$ , wobec czego  $x_{N+1}$  jest poszukiwanym ścisłym rozwiązaniem.

Zauważmy, że ponieważ  $A$  jest symetryczna, dodatnio określona, warunek (26c) oznacza, że wektory  $\{p_k\}$  są wzajemnie prostopadłe w metryce zadanej przez  $A$ . Ten właśnie warunek nazywa się warunkiem *sprzężenia względem  $A$* , co daje nazwę całej metodzie.

## Koszt metody

W arytmetyce dokładnej metoda zbiega się po  $N$  krokach, zatem jej koszt wynosi  $O(N \cdot \text{koszt\_jednego\_kroku})$ . Koszt jednego kroku zdominowany jest przez obliczanie iloczynu  $A p_k$ . Jeśli macierz  $A$  jest pełna, jest to  $O(N^2)$ , a zatem całkowity koszt wynosi  $O(N^3)$ , czyli tyle, ile dla metod dokładnych. Jeżeli jednak  $A$  jest rzadka, koszt obliczania iloczynu jest mniejszy (o ile obliczenie to jest odpowiednio zaprogramowane). Jeśli  $A$  jest pasmowa o szerokości pasma  $M \ll N$ , całkowity koszt wynosi  $O(M \cdot N^2)$ .

## Problem!

W arytmetyce o skończonej dokładności kolejne generowane wektory nie są *ściśle* ortogonalne do swoich poprzedników — na skutek akumulującego się błędu zaokrąglenia rzut na poprzednie wektory może stać się z czasem znaczny. Powoduje to istotne spowolnienie metody.

**Twierdzenie 9.** *Jeżeli  $\mathbf{x}$  jest ścisłym rozwiązaniem równania (24),  $\mathbf{x}_k$  są generowane w metodzie gradientów sprzężonych, zachodzi*

$$\|\mathbf{x} - \mathbf{x}_k\| \leq 2\|\mathbf{x} - \mathbf{x}_1\| \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{k-1}, \quad (27)$$

*gdzie  $\kappa$  jest współczynnikiem uwarunkowania macierzy  $\mathbf{A}$ .*

*Jeżeli  $\kappa \gg 1$ , zbieżność może być bardzo wolna.*

## “Prewarunkowana” (*preconditioned*) metoda gradientów sprzężonych

Spróbujmy przyspieszyć zbieżność odpowiednio modyfikując równanie (24) i algorytm (25), jednak tak, aby

- nie zmienić rozwiązania,
- macierz zmodyfikowanego układu pozostała symetryczna i dodatnio określona, aby można było zastosować metodę gradientów sprzężonych,
- macierz zmodyfikowanego układu pozostała rzadka, aby jeden krok iteracji był numerycznie tani,
- macierz zmodyfikowanego układu miała niski współczynnik uwarunkowania.

Czy to się w ogóle da zrobić? **Okazuje się, że tak!**

Postępujemy następująco: Niech  $C \in \mathbb{R}^{N \times N}$  będzie odwracalną macierzą symetryczną, rzeczywistą, dodatnio określoną. Wówczas  $\tilde{A} = C^{-1}AC^{-1}$  też jest symetryczna, rzeczywista, dodatnio określona.

$$C^{-1}A \underbrace{C^{-1}C}_I x = C^{-1}b, \quad (28a)$$

$$\tilde{A}\tilde{x} = \tilde{b}, \quad (28b)$$

gdzie  $\tilde{x} = Cx$ ,  $\tilde{b} = C^{-1}b$ . Do równania (28b) stosujemy teraz metodę gradientów sprzężonych.

W każdym kroku iteracji musimy obliczyć (tylko, bo odnosi się to do “tyldowanego” układu (28b))

$$\alpha_k = \frac{\tilde{\mathbf{r}}_k^T \tilde{\mathbf{r}}_k}{\tilde{\mathbf{p}}_k^T \tilde{\mathbf{A}} \tilde{\mathbf{p}}_k} = \frac{\tilde{\mathbf{r}}_k^T \tilde{\mathbf{r}}_k}{\tilde{\mathbf{p}}_k^T \mathbf{C}^{-1} \mathbf{A} \mathbf{C}^{-1} \tilde{\mathbf{p}}_k}, \quad (29a)$$

$$\tilde{\mathbf{r}}_{k+1} = \tilde{\mathbf{r}}_k - \alpha_k \tilde{\mathbf{A}} \tilde{\mathbf{p}}_k = \tilde{\mathbf{r}}_k - \alpha_k \mathbf{C}^{-1} \mathbf{A} \mathbf{C}^{-1} \tilde{\mathbf{p}}_k, \quad (29b)$$

$$\beta_k = \frac{\tilde{\mathbf{r}}_{k+1}^T \tilde{\mathbf{r}}_{k+1}}{\tilde{\mathbf{r}}_k^T \tilde{\mathbf{r}}_k}, \quad (29c)$$

$$\tilde{\mathbf{p}}_{k+1} = \tilde{\mathbf{r}}_{k+1} + \beta_k \tilde{\mathbf{p}}_k, \quad (29d)$$

$$\tilde{\mathbf{x}}_{k+1} = \tilde{\mathbf{x}}_k + \alpha_k \tilde{\mathbf{p}}_k. \quad (29e)$$

Równania (29) zawierają jawne odniesienia do macierzy  $C^{-1}$ , co nie jest zbyt wygodne. Łatwo się przekonać, iż za pomocą prostych przekształceń macierz tę można „usunąć”, tak, iż pozostaje tylko jedno jej nietrywialne wystąpienie. Zdefiniujmy mianowicie

$$\tilde{\mathbf{r}}_k = C^{-1}\mathbf{r}_k, \quad \tilde{\mathbf{p}}_k = C\mathbf{p}_k, \quad \tilde{\mathbf{x}}_k = C\mathbf{x}_k. \quad (30)$$

W tej sytuacji  $\tilde{\mathbf{r}}_k^T \tilde{\mathbf{r}}_k = (C^{-1}\mathbf{r}_k)^T C^{-1}\mathbf{r}_k = \mathbf{r}_k^T (C^{-1})^T C^{-1}\mathbf{r}_k = \mathbf{r}_k^T C^{-1}C^{-1}\mathbf{r}_k = \mathbf{r}_k^T (C^{-1})^2 \mathbf{r}_k$  etc.

Wówczas równania (29) przechodzą w

$$\alpha_k = \frac{\mathbf{r}_k^T (\mathbf{C}^{-1})^2 \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}, \quad (31a)$$

$$\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k, \quad (31b)$$

$$\beta_k = \frac{\mathbf{r}_{k+1}^T (\mathbf{C}^{-1})^2 \mathbf{r}_{k+1}}{\mathbf{r}_k^T (\mathbf{C}^{-1})^2 \mathbf{r}_k}, \quad (31c)$$

$$\mathbf{p}_{k+1} = (\mathbf{C}^{-1})^2 \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k, \quad (31d)$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k. \quad (31e)$$

W powyższych równaniach rola macierzy  $C$  sprowadza się do obliczenia — *jeden raz w każdym kroku iteracji* — wyrażenia  $(C^{-1})^2 r_k$ , co, jak wiadomo, robi się rozwiązując odpowiedni układ równań. Zdefiniujmy

$$M = C^2. \quad (32)$$

Macierz  $M$  należy rzecz jasna dobrać tak, aby równanie  $Mz = r$  można było szybko rozwiązać.

Ostatecznie otrzymujemy następujący algorytm:

$$\begin{aligned} & \mathbf{r}_1 = \mathbf{b} - \mathbf{A}\mathbf{x}_1 \\ & \text{rozwiąż } \mathbf{M}\mathbf{z}_1 = \mathbf{r}_1 \\ & \mathbf{p}_1 = \mathbf{z}_1 \\ & \text{while } \|\mathbf{r}_k\| > \varepsilon \\ & \quad \alpha_k = \frac{\mathbf{r}_k^T \mathbf{z}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k} \\ & \quad \mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k \\ & \quad \text{rozwiąż } \mathbf{M}\mathbf{z}_{k+1} = \mathbf{r}_{k+1} \\ & \quad \beta_k = \frac{\mathbf{r}_{k+1}^T \mathbf{z}_{k+1}}{\mathbf{r}_k^T \mathbf{z}_k} \\ & \quad \mathbf{p}_{k+1} = \mathbf{z}_{k+1} + \beta_k \mathbf{p}_k \\ & \quad \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k \\ & \text{end} \end{aligned} \tag{33}$$

## Incomplete Cholesky preconditioner

Niech rozkład  $QR$  macierzy  $C$  ma postać  $C = QH^T$ , gdzie  $Q$  jest macierzą ortogonalną,  $H^T$  jest macierzą trójkątną górną. Zauważmy, że

$$M = C^2 = C^T C = (QH^T)^T QH^T = HQ^T QH^T = HH^T, \quad (34)$$

a więc macierz  $H$  jest czynnikiem Cholesky'ego macierzy  $M$ . Niech rozkład Cholesky'ego macierzy  $A$  ma postać  $A = GG^T$ . *Przypuśćmy, iż  $H \simeq G$ .*

Wówczas

$$\begin{aligned} \tilde{\mathbf{A}} &= \mathbf{C}^{-1} \mathbf{A} \mathbf{C}^{-1} = (\mathbf{C}^T)^{-1} \mathbf{A} \mathbf{C}^{-1} = \left( (\mathbf{Q} \mathbf{H}^T)^T \right)^{-1} \mathbf{A} (\mathbf{Q} \mathbf{H}^T)^{-1} = \\ &= (\mathbf{H} \mathbf{Q}^T)^{-1} \mathbf{A} (\mathbf{H}^T)^{-1} \mathbf{Q}^T = \mathbf{Q} \underbrace{\mathbf{H}^{-1} \mathbf{G}}_{\simeq \mathbb{I}} \underbrace{\mathbf{G}^T (\mathbf{H}^T)^{-1}}_{\simeq \mathbb{I}} \mathbf{Q}^T \simeq \mathbf{Q} \mathbf{Q}^T = \mathbb{I}. \end{aligned} \tag{35}$$

Ponieważ  $\tilde{\mathbf{A}} \simeq \mathbb{I}$ , współczynnik uwarunkowania tej macierzy powinien być bliski jedności.

## Niepełny rozkład Cholesky'ego — algorytm w wersji GAXPY

```
for    $k = 1:N$   
       $H_{kk} = A_{kk}$   
      for    $j = 1:k-1$   
           $H_{kk} = H_{kk} - H_{kj}^2$   
      end  
       $H_{kk} = \sqrt{H_{kk}}$   
      for    $l = k+1:N$   
           $H_{lk} = A_{lk}$   
          if    $A_{lk} \neq 0$   
              for    $j = 1:k-1$   
                   $H_{lk} = H_{lk} - H_{lj}H_{kj}$   
              end  
               $H_{lk} = H_{lk}/H_{kk}$   
          endif  
      end  
end
```

## Uwagi

- Ponieważ  $\mathbf{A}$  jest rzadka, powyższy algorytm na obliczanie przybliżonego czynnika Cholesky'ego wykonuje się szybko. Wykonuje się go *tylko raz*.
- Równanie  $\mathbf{Mz} = \mathbf{r}$  rozwiązuje się szybko, gdyż znamy czynnik Cholesky'ego  $\mathbf{M} = \mathbf{H}\mathbf{H}^T$ .
- Obliczone  $\mathbf{H}$  jest rzadkie, a zatem równanie  $\mathbf{Mz} = \mathbf{r}$  rozwiązuje się szczególnie szybko.
- Mamy nadzieję, że macierz  $\tilde{\mathbf{A}}$  ma współczynnik uwarunkowania bliski jedności, a zatem nie potrzeba wielu iteracji (33).

## Przykład — macierz pasmowa z pustymi diagonalami

Rozważmy macierz o następującej strukturze:

$$\mathbf{A} = \begin{bmatrix} a_1 & 0 & b_3 & 0 & c_5 & 0 & 0 & 0 & 0 & 0 & \dots \\ 0 & a_2 & 0 & b_4 & 0 & c_6 & 0 & 0 & 0 & 0 & \dots \\ b_3 & 0 & a_3 & 0 & b_5 & 0 & c_7 & 0 & 0 & 0 & \dots \\ 0 & b_4 & 0 & a_4 & 0 & b_6 & 0 & c_8 & 0 & 0 & \dots \\ c_5 & 0 & b_5 & 0 & a_5 & 0 & b_7 & 0 & c_9 & 0 & \dots \\ 0 & c_6 & 0 & b_6 & 0 & a_6 & 0 & b_8 & 0 & c_{10} & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \end{bmatrix} \quad (36)$$

Macierz ta jest symetryczna, zakładamy też, że jest dodatnio określona.

*Niepełny* czynnik Cholesky'ego macierzy (36) ma postać

$$\mathbf{H} = \begin{bmatrix} p_1 & & & & & & \\ 0 & p_2 & & & & & \\ q_3 & 0 & p_3 & & & & \\ 0 & q_4 & 0 & p_4 & & & \\ r_5 & 0 & q_5 & 0 & p_5 & & \\ 0 & r_6 & 0 & q_6 & 0 & p_6 & \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \end{bmatrix} \quad (37)$$

(W pełnym czynniku Cholesky'ego macierzy (36) zera leżące w (37) *po-  
między* diagonalą “*p*” a diagonalną “*r*” znikłyby — w ogólności mogłyby  
tam znajdować się jakieś niezerowe liczby.)

Zgodnie z podanym algorytmem, elementy ciągów  $\{p_k\}$ ,  $\{q_k\}$ ,  $\{r_k\}$  wyliczamy z następujących wzorów:

$$\begin{array}{ll}
 p_1 = \sqrt{a_1}, & p_2 = \sqrt{a_2}, \\
 q_3 = b_3/p_1, & q_4 = b_4/p_2, \\
 r_5 = c_5/p_1, & r_6 = c_6/p_2, \\
 p_3 = \sqrt{a_3 - q_3^2}, & p_4 = \sqrt{a_4 - q_4^2}, \\
 q_5 = (b_5 - r_5q_3)/p_3, & q_6 = (b_6 - r_6q_4)/p_4, \\
 r_7 = c_7/p_3, & r_8 = c_7/p_4, \\
 p_5 = \sqrt{a_5 - q_5^2 - r_5^2}, & p_6 = \sqrt{a_6 - q_5^2 - r_6^2}, \\
 q_7 = (b_7 - r_7q_5)/p_5, & q_8 = (b_8 - r_8q_6)/p_6, \\
 r_9 = c_9/p_5, & r_{10} = c_{10}/p_6, \\
 p_7 = \sqrt{a_7 - q_7^2 - r_7^2}, & p_8 = \sqrt{a_8 - q_8^2 - r_8^2}, \\
 q_9 = (b_9 - r_9q_7)/p_7, & q_{10} = (b_{10} - r_{10}q_8)/p_8, \\
 r_{11} = c_{11}/p_7, & r_{12} = c_{12}/p_8, \\
 \dots & \dots
 \end{array}$$

## Macierze niesymetryczne

Jeżeli w równaniu

$$\mathbf{Ax} = \mathbf{b} \quad (38)$$

macierz  $\mathbf{A}$  nie jest symetryczna i dodatnio określona, sytuacja się komplikuje. Zakładając, że  $\det \mathbf{A} \neq 0$ , równanie (38) możemy “zsymetryzować” na dwa sposoby.

CGNR:

$$\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}, \quad (39)$$

lub CGNE:

$$\mathbf{AA}^T \mathbf{y} = \mathbf{b}, \quad (40a)$$

$$\mathbf{x} = \mathbf{A}^T \mathbf{y}. \quad (40b)$$

Do dwu powyższych równań formalnie rzecz biorąc *można* używać metody gradientów sprzężonych. Trzeba jednak pamiętać, że nawet jeśli macierz  $\mathbf{A}$  jest rzadka, macierze  $\mathbf{A}^T \mathbf{A}$ ,  $\mathbf{A} \mathbf{A}^T$  nie muszą być rzadkie, a co gorsza, ich współczynnik uwarunkowania jest kwadratem współczynnika uwarunkowania macierzy wyjściowej.

Alternatywnie, zamiast “symetryzować” macierz, można zmodyfikować algorytm, tak aby zamiast dwu, generował on *cztery* ciągi wektorów. Należy jednak pamiętać, że dla wielu typów macierzy taki algorytm bywa bardzo wolno zbieżny, a niekiedy nawet dochodzi do kompletnej stagnacji przed uzyskaniem rozwiązania:

## Metoda gradientów bi-sprzężonych (*Bi-Conjugate Gradients, Bi-CG*)

$$\begin{aligned} & \mathbf{r}_1 = \mathbf{b} - \mathbf{A}\mathbf{x}_1, \mathbf{p}_1 = \mathbf{r}_1, \bar{\mathbf{r}}_1 \neq 0 \text{ dowolny, } \bar{\mathbf{p}}_1 = \bar{\mathbf{r}}_1 \\ & \mathbf{while} \quad \|\mathbf{r}_k\| > \varepsilon \\ & \quad \alpha_k = \frac{\bar{\mathbf{r}}_k^T \mathbf{r}_k}{\bar{\mathbf{p}}_k^T \mathbf{A}\mathbf{p}_k} \\ & \quad \mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{A}\mathbf{p}_k \\ & \quad \bar{\mathbf{r}}_{k+1} = \bar{\mathbf{r}}_k - \alpha_k \mathbf{A}^T \bar{\mathbf{p}}_k \\ & \quad \beta_k = \frac{\bar{\mathbf{r}}_{k+1}^T \mathbf{r}_{k+1}}{\bar{\mathbf{r}}_k^T \mathbf{r}_k} \\ & \quad \mathbf{p}_{k+1} = \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k \\ & \quad \bar{\mathbf{p}}_{k+1} = \bar{\mathbf{r}}_{k+1} + \beta_k \bar{\mathbf{p}}_k \\ & \quad \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k \\ & \mathbf{end} \end{aligned} \tag{41}$$

Wektory wygenerowane w algorytmie (41) spełniają następujące relacje:

$$\bar{\mathbf{r}}_i^T \mathbf{r}_j = \mathbf{r}_i^T \bar{\mathbf{r}}_j = 0, \quad i > j, \quad (42a)$$

$$\bar{\mathbf{r}}_i^T \mathbf{p}_j = \mathbf{r}_i^T \bar{\mathbf{p}}_j = 0, \quad i > j, \quad (42b)$$

$$\bar{\mathbf{p}}_i^T \mathbf{A} \mathbf{p}_j = \mathbf{p}_i^T \mathbf{A}^T \bar{\mathbf{p}}_j = 0, \quad i > j. \quad (42c)$$

Jeżeli w algorytmie (41) weźmiemy  $\bar{\mathbf{r}}_1 = \mathbf{A} \mathbf{r}_1$ , we wszystkich krokach zachodzić będzie  $\bar{\mathbf{r}}_k = \mathbf{A} \mathbf{r}_k$  oraz  $\bar{\mathbf{p}}_k = \mathbf{A} \mathbf{p}_k$ . Jest to wersja przydatna dla rozwiązywania układów równań z macierzami symetrycznymi, ale nieokreślonymi dodatnio. Jest to przy okazji szczególny wariant algorytmu GMRES (*generalised minimum residual*), formalnie odpowiadającego minimalizacji funkcjonau

$$\Phi(\mathbf{x}) = \frac{1}{2} \|\mathbf{A} \mathbf{x} - \mathbf{b}\|^2. \quad (43)$$