

## 1 Układy równań liniowych

### 1.1 Podstawowe pojęcia

•  $Ax = b$  ma rozwiązanie gdy:  $\det(A) \neq 0$ ,  $b = 0$  to  $x = 0$ ,  $\ker(A) = \{0\}$ ,  $A$  nie ma zerowej wartości własnej,  $A$  jest odwracalna.

• Macierz z diagonalą rozwiązuje się trywialnie:  $x_i = \frac{b_i}{a_{ii}}$ . Koszt  $O(n)$ .

• FLOP - Floating point Operation - dodawanie, mnożenie, dzielenie, odejmowanie, sqrt to jeden flop.

• Macierz trójkątna dolna, układ równań rozwiązany w przód:  $x_i = \frac{b_i - \sum_{j=1}^{i-1} a_{ji}x_j}{a_{ii}}$ . Koszt  $O(n^2)$ .

• Macierz trójkątna górna rozwiązanie analogiczne (w tył), koszt  $O(n^2)$ .

• Macierz ortogonalna  $Q$  spełnia  $Q^T Q = I$ ,  $Q^{-1} = Q^T$ .

• Gdy  $A$  ortogonalna, to  $Ax = b$  ma rozwiązanie  $x = A^T b$ . Koszt  $O(n^2)$ .

### 1.2 Rozkład LU

•  $PA = LU$  (z wyborem) lub  $A = LU$  (bez wyboru),  $P$  - permutacja,  $L$  - trójkątna dolna,  $U$  - trójkątna górna.

• Standardowy algorytm rozkładu LU z bez wyboru elementu głównego:  $A = LU$ .  $PA = \left[ \begin{array}{c|c} a_{1,1} & a_{1,2}^T \\ \hline a_{2,1} & A_{2,2} \end{array} \right] = \left[ \begin{array}{c|c} 1 & 0^T \\ \hline l_{2,1} & L_{2,2} \end{array} \right]$ .

$$\left[ \begin{array}{c|c} u_{1,1} & u_{1,2}^T \\ \hline 0 & U_{2,2} \end{array} \right]$$

$u_{1,1} = a_{1,1}$ ,  $u_{1,2} = a_{1,2}$ ,  $l_{2,1} = \frac{a_{2,1}}{u_{1,1}}$ , aktualizuj  $A_{2,2} = A_{2,2} - l_{2,1}u_{1,2}^T$ , wyznacz  $A_{2,2} = L_{2,2}U_{2,2}$ .

Działa gdy wszystkie minory główne macierzy  $A$  są niezerowe.

• GEPP - Gauss Elimination with Partial Pivoting - GE z częściowym wyborem elementu głównego.

$PA = LU$ ,  $P$  - permutacja,  $L$  - trójkątna dolna,  $U$  - trójkątna górna.

$PAx = LUx = Pb$ ,  $\tilde{b} = Pb$ ,  $Lx = \tilde{b}$ ,  $Uy = y$ . Koszt  $O(n^3)$ .

### 2 Arytmetyka zmiennopozycyjna

•  $6.63 \cdot 10^{-34}$  - 6.63 - mantysa, -34 - cecha, 10 - podstawa

•  $x = (-1)^s \cdot m \cdot \beta^e$  -  $s$  - znak,  $m = (f_0.f_1f_2\dots f_{p-1})_2$  - mantysa,  $b$  - podstawa,  $e$  - cecha

• W liczbach maszynowych  $\beta = 2$

•  $(1 - e_{max} \leq e \leq e_{max})$

• Liczby maszynowe są normalizowane  $f_0 = 1$  i nie jest zapisywane.

• W 6-bitowej arytmetyce zmiennopozycyjnej liczby subnormalne zachodzą dla  $e = -2$  i  $f_0 = 0$ .

• Metody aproksymacji liczb maszynowych:

RN - do najbliższej (domyślnie)

RD - w dół, tzn w stronę  $-\infty$

RU - w górę, tzn w stronę  $\infty$

RZ - do zera  $RZ(x) = RD(x)$  gdy  $x \geq 0$ ,  $RU(x)$  gdy  $x \leq 0$

• Jeśli  $|x| \in [realmin, realmax]$  to  $\frac{|x - RN(x)|}{|x|} \leq \frac{1}{2^p} = v$

•  $fl(x) = x(1 \cdot \epsilon)$ , gdy  $|\epsilon| \leq v$ .

• dla float32  $v = 6.0 \cdot 10^{-8}$  oraz dla float64  $v = 1.1 \cdot 10^{-16}$

### 3 Uwarunkowanie zad. i numeryczna poprawność

•  $\|\tilde{x} - x\|$  - błąd bezwzględny,  $\|\tilde{x} - x\|/\|x\|$  - błąd względny ( $x \neq 0$ )

• Wskaźnik uwarunkowania (bezwzględny) na poziomie  $\epsilon$

$$cond_{abs}(P, x, \epsilon) = \sup_{\|\Delta\| \leq \epsilon} \frac{\|P(x+\Delta) - P(x)\|}{\|\Delta\|}$$

$\|P(\tilde{x}) - P(x)\| \leq cond_{abs}(P, x, \epsilon) \cdot \|\tilde{x} - x\|$ , dla  $\|\tilde{x} - x\| \leq \epsilon$

• Idealizacja punktowy wskaźnik uwarunkowania

$$cond_{abs}(P, x) = \lim_{\epsilon \rightarrow 0} cond_{abs}(P, x, \epsilon)$$

• Gdy  $P$  różniczkowalna to  $cond_{abs}(P, x) = \|P'(x)\|$

• Wrażliwość dla błędu względnego

$$cond_{rel}(P, x, \epsilon) = \sup_{\|\Delta\| \leq \epsilon} \frac{\|P(x+\Delta) - P(x)\|}{\|P(x)\|} \cdot \frac{\|\Delta\|}{\|x\|} = cond_{abs}(P, x, \epsilon \|x\|) \cdot \frac{\|x\|}{\|P(x)\|}$$

Punktowo:

$$cond_{rel}(P, x) = \lim_{\epsilon \rightarrow 0} cond_{rel}(P, x, \epsilon)$$

Gdy  $P$  jest różniczkowalna to  $cond_{rel}(P, x) = \frac{\|P'(x)\|}{\|P(x)\|}$

• Zadanie dobrze uwarunkowanie:  $cond(P, x)$  nieduże

Zadanie źle uwarunkowane:  $cond(P, x)$  bardzo duże

• Normy wektorowe

$$\|x\|_1 = \sum_{i=1}^n |x_i|, \|x\|_2 = \sqrt{\sum_{i=1}^n |x_i|^2}, \|x\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}, \|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$$

$$\|x\|_\infty \leq \|x\|_1 \leq N\|x\|_\infty, \|x\|_\infty \leq \|x\|_2 \leq \sqrt{n}\|x\|_\infty, \|x\|_2 \leq \|x\|_1 \leq \sqrt{n}\|x\|_2$$

$$\|A\| = \max_{\|x\| \neq 0} \|Ax\|/\|x\| = \max_{\|x\|=1} \|Ax\| = \max_{\|x\| \leq 1} \|Ax\|$$

$$\|Ax\| \leq \|A\|\|x\|, \|AB\| \leq \|A\|\|B\|, \|I\| = 1, \|A\|_1 = \max_j \sum_i |a_{ij}|, \|A\|_\infty = \max_i \sum_j |a_{ij}|, \|A\|_2 = \max\{\sqrt{\mu} : \mu \text{ jest w. wł. } A^T A\}$$

$$\|cond(A) = \|A\|\|A^{-1}\|$$

$$\|Ay = b. \text{ Jeśli } \epsilon cond(A) \leq 1/2 \text{ to } \|\tilde{y} - y\|/\|y\| \leq 4cond(A) \cdot \epsilon$$

$$\| \text{Jeśli } \|\Delta\| < 1 \text{ to } I + \Delta \text{ nieosobliwa i } 1/(1 + \|\Delta\|) \leq \|(I + \Delta)^{-1}\| \leq 1/(1 - \|\Delta\|)$$

• Algorytm poprawnie numeryczny - dla każdego  $x \in X$  wynik algorytmu  $A$  zrealizowanego w fl  $fl(A(fl(x)))$  jest dokładnym rozwiązaniem zadania dla danych  $x$  zaburzonych na poziomie błędu reprezentacji.

• Algorytm NP daje wynik, którego błąd można oszacować na podstawie własności zadania obliczeniowego:

$$\|\tilde{y} - y\|/\|y\| = \|P(\tilde{x}) - P(x)\|/\|P(x)\| \lesssim cond_{rel}(P, x)\|\tilde{x} - x\|/\|x\| \leq K \cdot cond_{rel}(P, x) \cdot v$$

### 4 LZNK

$A \in \mathbb{R}^{m \times n}$ ,  $m \geq n$ ,  $\text{rank}(A) = n$ ,  $\vec{b} \in \mathbb{R}^m$ , znaleźć  $\vec{x} \in \mathbb{R}^n$  taki, że  $\|\vec{b} - A\vec{x}\|_2 \rightarrow \min$ .

$$\|\vec{b} - A\vec{x}\| \leq \|\vec{b} - A\vec{y}\| \forall y.$$

$$\bullet QR - A = QR = \begin{bmatrix} \hat{Q} & \tilde{Q} \end{bmatrix} \begin{bmatrix} \hat{R} \\ 0 \end{bmatrix} \Rightarrow$$

$$\|\vec{b} - A\vec{x}\|_2^2 = \|\hat{Q}^T \vec{b} - \hat{R}\vec{x}\|_2^2 + \|\tilde{Q}^T \vec{b}\|_2^2 = \min \iff \hat{R}\vec{x} = \hat{Q}^T \vec{b}$$

•  $A^T A$  - mała, symetryczna, dodatnio określona  $A^T A \vec{x} = A^T \vec{b}$

• Algorytm: oblicz  $B = A^T A \rightarrow$  wyznacz Cholesky'ego  $B = LL^T \rightarrow$  rozwiąż  $LL^T \vec{x} = A^T \vec{b}$

•  $SVD - A \in \mathbb{R}^{m \times n}$ ,  $m \geq n$ , istnieje  $A = U\Sigma V^T$ , że  $U \in \mathbb{R}^{m \times n} : U^T U = I$ ,  $V \in \mathbb{R}^{n \times n} : V^T V = I$ ,  $\Sigma \in \mathbb{R}^{n \times n} : \Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$ ,  $\sigma_1 \geq \dots \geq \sigma_n \geq 0$

$$\lambda(A^T A) = (\sigma_1^2, \dots, \sigma_n^2), \sigma_1 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_n = 0 \implies \text{rank} A = r$$

$A = \sum_{i=1}^n \sigma_i \vec{u}_i \vec{v}_i^T$  gdzie  $\vec{u}_i, \vec{v}_i$  - kolumny  $U, V$

$$\|\vec{b} - A\vec{x}\|_2^2 = \sum_{i=1}^r \left( \sigma_i y_i - \vec{u}_i^T \vec{b} \right)^2 + \sum_{i=r+1}^n \left( \vec{u}_i^T \vec{b} \right)^2 \rightarrow \min, \text{ gdy } y_i = \frac{\vec{u}_i^T \vec{b}}{\sigma_i},$$

To oznacza, że rozwiązanie LZNK jest jednoznacznie określone  $\vec{x}^* = \sum_{i=1}^r \frac{\vec{u}_i^T \vec{b}}{\sigma_i} \vec{v}_i$

### 5 Metody iteracyjne

• Ciąg  $x_{k+1} = Bx_k + c$  jest zbieżny do  $x^*$  dla każdego  $x_0$  wtedy i tylko wtedy, gdy  $\rho(B) < 1$ . gdzie  $\rho(B)$  to promień spektralny macierzy  $B$  ( $\rho(B) = \max\{|\lambda| : \lambda \text{ jest w. wł. } B\}$ )

• Jeśli  $\|B\| < 1$  to ciąg  $x_{k+1} = Bx_k + c$  jest zbieżny do  $x^*$  dla każdego  $x_0$ . oraz  $\|x^* - x_k\| \leq \|B\| \cdot \|x_k - x^*\|$

• Niech  $A = M - Z$  oraz  $A, M$  nieosobliwe.  $Ax^* = b$ . Jeśli  $\rho(M^{-1}Z) < 1$  to metoda iteracyjna  $Mx_{k+1} = Zx_k + b$  jest zbież-

na do  $x^*$  dla każdego  $x_0$ .  
 Jeśli dodatkowo  $\gamma = \|M^{-1}Z\|_\infty < 1$  to  $\|x^* - x_k\| \leq \gamma \cdot \|x_k - x^*\|$   
 • Niech  $A = L + D + U$ ,  $D = \text{diag}(A)$ ,  $L$  (odp. U) - dolna (odp. górna) trójkątna z zerową diagonalą.  
 $x_{k+1} = x_k + M^{-1}(b - Ax_k)$   
 Metoda Jacobiego:  $M = D$ , Metoda Gaussa-Seidela:  $M = D + L$ ,  
 Metoda SOR:  $M = 1/\omega D + L$ .

• Jeśli  $A$  jest diagonalnie dominująca to m. Jacobiego jest zbieżna do  $x^*$  dla każdego  $x_0$ .

•  $x_{k+1} = x_k + \delta_k$ . Jak wybrać  $\delta_k$ ?. Idealna poprawka  $\delta_k^*$ :  
 $A\delta_k^* = r_k \rightarrow x_{k+1} = x^*$   
 Wyznaczamy idealną poprawkę  $\delta_k^* = V_k a_k$ , gdzie  $V_k, U_k \in \mathbb{R}^{N \times r}$   
 max rzędu t. że  $a_k \in \mathbb{R}^r$  spełnia:  
 $U_k^T A V_k a_k = U_k^T r_k$

### 6 Zagadnienie własne

• Znaleźć  $\vec{v} \in \mathbb{C}^n$  oraz  $\lambda \in \mathbb{C}$  t. że  $A\vec{v} = \lambda\vec{v} \wedge \|\vec{v}\| = 1$ .  
 • Każda symetryczna macierz  $A \in \mathbb{R}^{n \times n}$  ma rozkład

$$A = Q\Lambda Q^T, \Lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix},$$

$\lambda_i \in \mathbb{R}$  — wartości własne  $A$ , kolumny  $Q$  — wektory własne  $A$ .

• Metoda potęgowa — zakładamy ze  $|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$ ,  
 losuj  $\vec{x}_0$ ,  $\vec{x}_{k+1} = A\vec{x}_k / \|A\vec{x}_k\|$ , znajduje  $\lambda_{\max}$ .  
 • Wyznaczanie wartości własnej na podstawie wektora własnego.  
 $\|A\vec{v} - \lambda\vec{v}\| \rightarrow \min, \lambda = \frac{\vec{x}^H A \vec{v}}{\vec{x}^H \vec{x}}$   
 • Odwrotna metoda potęgowa — zakładamy ze  $|\lambda_1| \geq \dots \geq |\lambda_n| > 0$ , wyznacz  $PA = LU$ , losujemy  $\vec{x}_0$ , rozwiąż  $LU\vec{x}_{k+1} = P\vec{x}_k$ , przeskaluj  $\vec{x}_{k+1} = \vec{x}_{k+1} / \|\vec{x}_{k+1}\|$ , znajduje  $1/\lambda_{\min}$ .  
 • Odwrotna metoda potęgowa z parametrem — to samo, co odwrotna ale dla  $(A - \mu)^{-1}$ , znajduje  $\lambda$  najbliższe  $\mu$ .

•  $\forall_i \lambda_i(A) \leq \|A\|$   
 • Koła Greszgorina —  $\forall_\lambda \exists_i |\lambda - a_{ii}| \leq \sum_{j \neq i} |a_{ij}|$

	Macierz	wart. wł.	wekt. wł.	zastrz.
	$A$	$\lambda$	$\vec{v}$	—
•	$A - \mu I$	$\lambda - \mu$	$\vec{v}$	—
	$A^{-1}$	$1/\lambda$	$\vec{v}$	$A$ nieos.
	$(A - \mu I)^{-1}$	$1/(\lambda - \mu)$	$\vec{v}$	$A - \mu$ nieos.

### 7 Dopiski

#### 7.1 Rozkład QR

•  $A = QR$ ,  $Q$  — ortogonalna,  $R$  — trójkątna górna.  
 • Housholder — dany  $\vec{x} \in \mathbb{R}^n$ ,  $\alpha = \|\vec{x}\| \cdot -\text{sgn}(x_1)$

$$\vec{u} = \vec{x} - \alpha \vec{e}_1, \quad \vec{v} = \frac{\vec{u}}{\|\vec{u}\|}, \quad Q = I - 2\vec{v}\vec{v}^T, \quad Qx = \begin{bmatrix} \alpha \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\bullet Q_1 A = \begin{bmatrix} \alpha_1 & * & \dots & * \\ 0 & & & \\ \vdots & & & \\ 0 & & & \end{bmatrix} A'$$

$$\bullet Q_k = \begin{bmatrix} I_{k-1} & 0 \\ 0 & Q'_k \end{bmatrix}$$

•  $R = Q_t \dots Q_2 Q_1 A$   
 •  $Q^T = Q_t \dots Q_1$   
 •  $Q = Q_1^T \dots Q_t^T = Q_1 \dots Q_t$   
 • Zamiast tego można użyć ortogonalizacji Grama-Schmidta, która nie jest tak numerycznie stabilna, ale można ją zastosować kilkakrotnie (reortogonalizacja) żeby poprawić wynik.