

1 Iterační metody pro řešení soustav lineárních rovnic

Přímé metody řešení soustav lineárních rovnic (LU, LDLT, Choleského faktORIZACE atd.) vyžadují $\mathcal{O}(n^3)$ operací a velikost soustavy, kterou jimi dokážeme vyřešit je značně omezená. V případě husté matice $A \in \mathbb{R}^{n \times n}$ se přibližná velikost řešitelné soustavy historicky vyvíjela takto:

- 1950: $n = 20$ (Wilkinson)
- 1965: $n = 200$ (Forsythe a Moler)
- 1980: $n = 2000$ (LINPACK)
- 1995: $n = 20000$ (LAPACK)
- 2010: $n = 200000$ (HDSS)

Pracujeme-li s řídkými maticemi, jsme schopni řešit i systémy s mnohem větší dimenzí (milióny, desítky miliónů neznámých) – zejména, pokud je řešič schopen pracovat paralelně. V případě použití přímého řešiče na řídkou matici však může dojít k jejímu zaplnění. Např. využitím prvního řádku následující matice k vynulování prvního sloupce dojde k zaplnění všech ostatních prvků v matici:

$$\begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & 0 & 0 & 0 \\ \times & 0 & \times & 0 & 0 \\ \times & 0 & 0 & \times & 0 \\ \times & 0 & 0 & 0 & \times \end{bmatrix} \rightarrow \begin{bmatrix} \times & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \end{bmatrix}$$

Tento problém lze částečně řešit vhodnou pivotizací.

Mezi další nevýhody přímých řešičů patří:

- Známe-li již přibližné řešení soustavy, nedokážeme tuto znalost využít ke snížení celkového počtu operací a zkrácení doby výpočtu.
- Naopak, pokud nám postačuje znalost pouze přibližného řešení, nemůžeme výpočet pomocí přímého řešiče ukončit předčasně.

Alternativou k přímým řešičům jsou iterační řešiče, které generují posloupnost přibližných řešení $\{\mathbf{x}^k\}$ a pracují téměř výhradně s násobením matice-vektor, které má náročnost $\mathcal{O}(n^2)$. Důležitou vlastností každé iterační metody je rychlost konvergence posloupnosti $\{\mathbf{x}^k\}$ k řešení. Může se totiž stát, že pro některé matice A iterační metoda konverguje velmi pomalu nebo vůbec.

1.1 Lineární iterační metody

Prvním typem iteračních metod, kterým se budeme zabývat, jsou tzv. lineární iterační metody. Ty hledají posloupnost řešení soustavy $A\mathbf{x} = \mathbf{b}$ ve tvaru

$$\mathbf{x}^{k+1} := M\mathbf{x}^k + N\mathbf{b}, \quad (1.1)$$

kde M a N jsou nějaké matice odpovídajících rozměrů¹.

Definice Lineární iterační metodu nazveme konzistentní, řeší-li rovnici $M\mathbf{x} + N\mathbf{b} = \mathbf{x}$ právě jeden vektor $\mathbf{x} = A^{-1}\mathbf{b}$. Je možné ukázat, že metoda je konzistentní právě tehdy, je-li splněno $M = I - NA$.

Definice Iterační metodu nazveme konvergentní platí-li $\mathbf{x}^k \rightarrow \mathbf{x} = A^{-1}\mathbf{b}$ pro $k \rightarrow \infty$. Je možné ukázat, že metoda je konvergentní, právě tehdy, je-li splněno $\|M\| < 1$, kde $\|M\| = \max_{\mathbf{v} \in \mathbb{R}^n} \frac{\|M\mathbf{v}\|_v}{\|\mathbf{v}\|_v}$ je maticová norma indukovaná vektorovou normou.

Při odvozování následujících iteračních metod budeme využívat rozkladu matice A na součet dolní trojúhelníkové, diagonální a horní trojúhelníkové matice, tedy $A = L + D + U$ (pozor, nepleťte si matice L, D, U se stejně nazvanými maticemi, které se vyskytovaly u přímých řešičů). Např. matici

$$A = \begin{bmatrix} 2 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 2 \end{bmatrix}$$

rozložíme na

$$L = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, D = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix}, U = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}.$$

1.1.1 Jacobiho metoda

Vyděme z rovnice $A\mathbf{x} = \mathbf{b}$. Dosazením $A = L + D + U$ dostaneme

$$(L + D + U)\mathbf{x} = \mathbf{b}.$$

Roznásobme a přezávorkujme výraz na levé straně

$$D\mathbf{x} + (L + U)\mathbf{x} = \mathbf{b}$$

Jacobiho metodu odvodíme tak, že přidáme indexy $k + 1$ a k k příslušným vektorům \mathbf{x}

$$D\mathbf{x}^{k+1} + (L + U)\mathbf{x}^k = \mathbf{b}.$$

Osamostatněním \mathbf{x}^{k+1} dostaneme předpis pro $k + 1$ aproximaci vektoru \mathbf{x}

$$\mathbf{x}^{k+1} := D^{-1}(\mathbf{b} - (L + U)\mathbf{x}^k) = \underbrace{-D^{-1}(L + U)}_{=M} \mathbf{x}^k + \underbrace{D^{-1}}_{=N} \mathbf{x}^k.$$

¹Index k označuje číslo aktuální iterace.

Jednotlivé složky vektoru \mathbf{x}^{k+1} můžeme vyjádřit jako

$$(\mathbf{x}^{k+1})_i := \frac{1}{a_{i,i}} \left(b_i - \sum_{j=1, j \neq i}^n a_{i,j} (\mathbf{x}^k)_j \right) = \quad (1.2)$$

$$= \frac{1}{a_{i,i}} \left(b_i - \sum_{j=1}^{i-1} a_{i,j} (\mathbf{x}^k)_j - \sum_{j=i+1}^n a_{i,j} (\mathbf{x}^k)_j \right) \quad (1.3)$$

Pro snadnější a přehlednější zápis budeme i -tý prvek vektoru \mathbf{x}^{k+1} také značit jako x_i^{k+1} (horní index tedy označuje číslo iterace, dolní index značí pořadí prvku ve vektoru).

Konzistence metody vyplývá z jejího odvození, můžeme však ještě ověřit, že $M = I - NA$. V případě Jacobiho metody je $M = -D^{-1}(L + U)$ a $N = D^{-1}$ (viz výše). Platí tedy

$$\begin{aligned} I - NA &= I - D^{-1}A = I - D^{-1}(L + D + U) = \\ &= I - D^{-1}(L + U) - D^{-1}D = -D^{-1}(L + U) = M. \end{aligned}$$

Metoda je tedy konzistentní.

Lze dokázat, že metoda je konvergentní právě tehdy, když $\|M\| = \|D^{-1}(L + U)\| < 1$. To splňují např. striktně diagonálně dominantní matice (tedy matice, pro které platí $\forall i : |a_{i,i}| > \sum_{j=1, j \neq i}^n |a_{i,j}|$). Konvergenci metody pro diagonálně dominantní matice se dá poměrně snadno dokázat. Vyjděme ze vztahu pro i -tý prvek aproximovaného vektoru v iteraci $k + 1$:

$$x_i^{k+1} = \frac{1}{a_{i,i}} \left(b_i - \sum_{j=1, j \neq i}^n a_{i,j} x_j^k \right).$$

Jelikož je metoda konzistentní, musí tuto rovnost splňovat i prvky vektoru přesného řešení \mathbf{x} :

$$x_i = \frac{1}{a_{i,i}} \left(b_i - \sum_{j=1, j \neq i}^n a_{i,j} x_j \right).$$

Odečteme-li od první rovnosti druhou, dostaneme

$$\underbrace{x_i^{k+1} - x_i}_{=e_i^{k+1}} = -\frac{1}{a_{i,i}} \sum_{j=1, j \neq i}^n a_{i,j} \underbrace{(x_j^k - x_j)}_{=e_j^k},$$

kde e^k je vektor chyby v k -tém kroce. Chybu v kroce $k + 1$ můžeme tedy odhadnout pomocí vlastností striktně diagonálně dominantní matice:

$$\begin{aligned} |e_i^{k+1}| &\leq \frac{1}{|a_{i,i}|} \sum_{j=1, j \neq i}^n |a_{i,j}| |e_j^k| \leq \frac{1}{|a_{i,i}|} \sum_{j=1, j \neq i}^n |a_{i,j}| \max_{j=1, \dots, n, j \neq i} |e_j^k| = \\ &= \max_{j=1, \dots, n, j \neq i} |e_j^k| \underbrace{\frac{1}{|a_{i,i}|} \sum_{j=1, j \neq i}^n |a_{i,j}|}_{<1} < \max_{j=1, \dots, n, j \neq i} |e_j^k| \end{aligned}$$

Každý prvek vektoru chyby v kroce $k + 1$ je tedy v absolutní hodnotě menší než maximální prvek vektoru chyby v předchozím kroce. Vektor chyby tedy konverguje k nulovému vektoru.

1.1.2 Gaussova-Seidelova metoda

Všimněme si, že při výpočtu x_i^{k+1} využíváme v sumě $\sum_{j=1}^{i-1} a_{i,j} x_j^k$ ve výrazu (1.3) pouze prvky x_1^k, \dots, x_{i-1}^k . Tyto prvky tedy můžeme nahradit již vypočtenými prvky aktuálními iterace $x_1^{k+1}, \dots, x_{i-1}^{k+1}$. Dostaneme tak předpis Gaussovy-Seidelovy metody:

$$x_i^{k+1} := \frac{1}{a_{i,i}} \left(b_i - \sum_{j=1}^{i-1} a_{i,j} x_j^{k+1} - \sum_{j=i+1}^n a_{i,j} x_j^k \right) \quad (1.4)$$

Podobně jako v předchozím případě můžeme metodu odvodit, nahradíme-li v soustavě $A\mathbf{x} = \mathbf{b}$ matici A součtem $L + D + U$. Tentokrát ovšem vektor s indexem $k + 1$ ponecháme u součtu $L + D$

$$(L + D)\mathbf{x}^{k+1} = \mathbf{b} - U\mathbf{x}^k, \quad (1.5)$$

tedy

$$\mathbf{x}^{k+1} = \underbrace{-(L + D)^{-1}U}_{=M} \mathbf{x}^k + \underbrace{(L + D)^{-1}\mathbf{b}}_{=N}.$$

Vztah mezi maticovým zápisem a zápisem po prvcích (1.4) je nejlépe vidět na rovnosti (1.5). Jedná se o soustavu rovnic s dolní trojúhelníkovou maticí $L + D$, vektorem pravé strany $\mathbf{b} - U\mathbf{x}^k$ a neznámým vektorem \mathbf{x}^{k+1} . Všimněte si, že výraz (1.4) pak přesně odpovídá algoritmu pro dopřednou substituci pro řešení takovéto soustavy.

Podobně jako v případě Jacobiho metody můžeme ověřit, zda je metoda konzistentní porovnáním $I - NA$ a M .

$$\begin{aligned} I - NA &= I - (L + D)^{-1}A = I - (L + D)^{-1}(L + D + U) = \\ &= I - (L + D)^{-1}(L + D) - (L + D)^{-1}U = -(L + D)^{-1}U = M. \end{aligned}$$

Metoda je tedy konzistentní.

Metoda je konvergentní, právě když $\|(L + D)^{-1}U\| < 1$.

1.1.3 Richardsonova metoda

Iterace Richardsonovy metody je dána předpisem

$$\mathbf{x}^{k+1} := \mathbf{x}^k + \omega \mathbf{r}^k,$$

kde $\omega \in \mathbb{R}_+$ a $\mathbf{r}^k = \mathbf{b} - \mathbf{A}\mathbf{x}^k$ je reziduum, které určuje, jak dobře je splněna původní rovnice. Vztah mezi reziduem a chybou $\mathbf{e}^k = \mathbf{x}^k - \mathbf{x}$ lze odvodit přenásobením definice chyby maticí \mathbf{A}

$$\mathbf{A}\mathbf{e}^k = \mathbf{A}\mathbf{x}^k - \mathbf{A}\mathbf{x} = \mathbf{A}\mathbf{x}^k - \mathbf{b} = -\mathbf{r}^k.$$

Studujme konvergenci metody pro symetrickou pozitivně definitní matici \mathbf{A} . V takovém případě vlastní čísla λ_i a vlastní vektory \mathbf{v}_i (tedy skalární vektory, pro které platí $\mathbf{A}\mathbf{v}_i = \lambda_i\mathbf{v}_i$, $\|\mathbf{v}_i\| = 1$) splňují

$$0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$$

a

$$\text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\} = \mathbb{R}^n.$$

Zde span značí lineární obal. Jinými slovy, vlastní vektory tvoří bázi \mathbb{R}^n .

Díky předchozímu poznatku můžeme reziduum vyjádřit jako lineární kombinaci prvků báze tvořené vlastními vektory, tedy $\mathbf{r}^{k+1} = \sum_{i=1}^n \alpha_i^{k+1} \mathbf{v}_i$. Studujme nyní, jak se chová reziduum (a tedy i chyba) v jednotlivých iteracích. Na základě toho se později pokusíme odvodit optimální hodnotu ω pro co nejrychlejší konvergenci.

$$\begin{aligned} \sum_{i=1}^n \alpha_i^{k+1} \mathbf{v}_i &= \mathbf{r}^{k+1} = \mathbf{b} - \mathbf{A}\mathbf{x}^{k+1} = \underbrace{\mathbf{b} - \mathbf{A}(\mathbf{x}^k + \omega \mathbf{r}^k)}_{=\mathbf{r}^k} = \\ &= \mathbf{r}^k - \mathbf{A}\omega \mathbf{r}^k = (1 - \omega \mathbf{A}) \underbrace{\mathbf{r}^k}_{=\sum_{i=1}^n \alpha_i^k \mathbf{v}_i} = (1 - \omega \mathbf{A}) \sum_{i=1}^n \alpha_i^k \mathbf{v}_i = \\ &= \sum_{i=1}^n \alpha_i^k \mathbf{v}_i - \sum_{i=1}^n \alpha_i^k \omega \underbrace{\mathbf{A}\mathbf{v}_i}_{=\lambda_i \mathbf{v}_i} = \sum_{i=1}^n \alpha_i^k \mathbf{v}_i - \sum_{i=1}^n \alpha_i^k \omega \lambda_i \mathbf{v}_i = \\ &= \sum_{i=1}^n (1 - \omega \lambda_i) \alpha_i^k \mathbf{v}_i. \end{aligned}$$

Všimněme si, že jsme vyjádřili koeficienty rozvoje rezidua \mathbf{r}^{k+1} v bázi $\{\mathbf{v}_i\}_{i=1}^n$ pomocí násobků koeficientů v předchozím kroce (viz podtržené části předchozího výrazu). Koeficienty se tedy budou zmenšovat (a jednotlivé složky vektoru rezidua budou konvergovat k nule) právě tehdy, když $|1 - \omega \lambda_i| < 1$ pro všechna $i = 1, 2, \dots, n$. Rychlost konvergence bude záviset na největší hodnotě $|1 - \omega \lambda_i|$. Shrňme tento poznatek do následující věty.

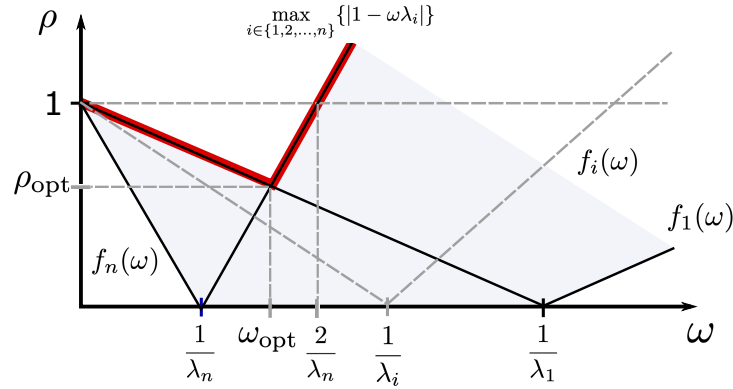


Figure 1.1: Konvergenční faktor Richardsonovy metody v závislosti na ω .

Věta Richardsonova metoda konverguje, právě když $\forall i \in \{1, 2, \dots, n\} : |1 - \omega\lambda_i| < 1$. Konvergenční faktor $\rho = \max_{i \in \{1, 2, \dots, n\}} \{|1 - \omega\lambda_i|\}$ určuje rychlost konvergence: $\|\mathbf{r}^{k+1}\| \leq \rho \|\mathbf{r}^k\|$.

Čím menší bude konvergenční faktor $\rho = \max_{i \in \{1, 2, \dots, n\}} \{|1 - \omega\lambda_i|\}$, tím rychleji bude metoda konvergovat. Vzhledem k tomu, že vlastní čísla matice A jsou daná, můžeme konvergenční faktor ovlivnit pouze vhodnou volbou parametru ω . Odvození ideální hodnoty ω ilustrujeme na Obrázku 1.1. Jsou na něm znázorněny funkce $f_1(\omega) = |1 - \omega\lambda_1|$ a $f_n(\omega) = |1 - \omega\lambda_n|$. Protože směrnice funkcí $f_i(\omega) = |1 - \omega\lambda_i|$ jsou určeny vlastními čísly matice A a ta jsou seřazena od nejmenšího po největší, budou grafy všech funkcí $f_i, i = 2, \dots, n-1$, ležet "mezi" grafy f_1 a f_n (v šedě vyznačené oblasti). Funkci

$$\max_{i \in \{1, 2, \dots, n\}} \{|1 - \omega\lambda_i|\}$$

tedy můžeme vykreslit jako červeně zvýrazněnou lomenou čáru tvořenou částí funkce f_1 a částí funkce f_n . Z grafu této funkce tedy rovnou můžeme odvodit:

1. Interval, ve kterém musí ω ležet. Aby metoda konvergovala, musí platit $\rho = \max_{i \in \{1, 2, \dots, n\}} \{|1 - \omega\lambda_i|\} < 1$. Červená funkce tedy musí ležet pod zakreslenou konstantní funkcí $\rho = 1$. Levý krajní bod intervalu je 0, pravý určíme jako průsečík příslušné části funkce f_n s konstantní funkcí 1:

$$-(1 - \omega\lambda_n) = 1 \quad \Rightarrow \quad \omega = \frac{2}{\lambda_n}.$$

Metoda tedy konverguje pro $\omega \in (0, 2/\lambda_n)$.

2. Optimální ω je bod, ve kterém červeně vyznačená funkce dosahuje minima. Tento bod dostaneme jako průsečík funkcí f_1 a f_n :

$$1 - \omega_{\text{opt}}\lambda_1 = -(1 - \omega_{\text{opt}}\lambda_n) \quad \Rightarrow \quad \omega_{\text{opt}} = \frac{2}{\lambda_1 + \lambda_n}.$$

Zjistili jsme tedy, že nejlepší konvergence dosáhneme, zvolíme-li $\omega_{\text{opt}} = \frac{2}{\lambda_1 + \lambda_n}$. Konvergenční faktor bude v tomto případě

$$\begin{aligned}\rho_{\text{opt}} &= 1 - \omega_{\text{opt}} \lambda_1 = 1 - \frac{2\lambda_1}{\lambda_1 + \lambda_n} = \frac{\lambda_1 + \lambda_n - 2\lambda_1}{\lambda_1 + \lambda_n} = \\ &= \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \frac{\frac{1}{\lambda_1}}{\frac{1}{\lambda_1}} = \frac{\frac{\lambda_n}{\lambda_1} - 1}{\frac{\lambda_n}{\lambda_1} + 1} = \frac{\kappa(\mathbf{A}) - 1}{\kappa(\mathbf{A}) + 1},\end{aligned}$$

kde $\kappa(\mathbf{A}) = \lambda_n/\lambda_1$ je číslo podmíněnosti matice \mathbf{A} .

Můžeme také odvodit, kolik iterací je třeba, abychom dosáhli požadované relativní změny normy rezidua. Hledáme tedy k , pro které platí

$$\frac{\|\mathbf{r}^k\|}{\|\mathbf{r}^0\|} \leq \varepsilon, \quad \text{tedy} \quad \|\mathbf{r}^k\| \leq \varepsilon \|\mathbf{r}^0\|$$

Využijme toho, že $\|\mathbf{r}^k\| \leq \rho_{\text{opt}}^k \|\mathbf{r}^0\|$ a přepíšme nerovnici na

$$\rho_{\text{opt}}^k \|\mathbf{r}^0\| \leq \varepsilon \|\mathbf{r}^0\|.$$

Vykrácením normy a zlogaritmováním obou stran nerovnice dostaneme řešení $k \geq \frac{\log \varepsilon}{\log \rho_{\text{opt}}}$ (nezapomeňte, že protože $\rho_{\text{opt}} \in (0, 1)$, je třeba otočit znaménko nerovnosti).

1.2 Gradientní iterační metody

Věta Řešení soustavy $\mathbf{Ax} = \mathbf{b}$ se symetrickou pozitivně definitní maticí \mathbf{A} je ekvivalentní s minimalizací kvadratické formy

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}.$$

Důkaz Dokažme nejdříve implikaci $\mathbf{Ax} = \mathbf{b} \Rightarrow \mathbf{x} = \arg \min_{\mathbf{v} \in \mathbb{R}^n} f(\mathbf{v})$. Podíváme se, jak se změní funkční hodnota f , posuneme-li se z bodu \mathbf{x} o nějaký nenulový vektor \mathbf{c} :

$$\begin{aligned}f(\mathbf{p}) &= f(\mathbf{x} + \mathbf{c}) = \frac{1}{2} (\mathbf{x} + \mathbf{c})^T \mathbf{A} (\mathbf{x} + \mathbf{c}) - \mathbf{b}^T (\mathbf{x} + \mathbf{c}) = \\ &= \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{c}^T \underbrace{\mathbf{A} \mathbf{x}}_{=\mathbf{b}} + \frac{1}{2} \mathbf{c}^T \mathbf{A} \mathbf{c} - \mathbf{b}^T \mathbf{x} - \mathbf{b}^T \mathbf{c} = \\ &= \underbrace{\frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}}_{=f(\mathbf{x})} + \underbrace{\mathbf{c}^T \mathbf{b} - \mathbf{b}^T \mathbf{c}}_{=0} + \frac{1}{2} \mathbf{c}^T \mathbf{A} \mathbf{c} = f(\mathbf{x}) + \underbrace{\frac{1}{2} \mathbf{c}^T \mathbf{A} \mathbf{c}}_{>0}.\end{aligned}$$

Díky pozitivní definitnosti \mathbf{A} je výraz $\mathbf{c}^T \mathbf{A} \mathbf{c}$ kladný. Posuneme-li se tedy z bodu \mathbf{x} v libovolném směru, hodnota funkce f se zvětší. V bodě \mathbf{x} tedy nastává minimum.

K důkazu opačné implikace $\mathbf{x} = \arg \min_{\mathbf{v} \in \mathbb{R}^n} f(\mathbf{v}) \Rightarrow \mathbf{Ax} = \mathbf{b}$ je třeba si uvědomit nutnou podmínku minima funkce $f : \mathbb{R}^n \rightarrow \mathbb{R}$, tedy nulovost gradientu:

$$\mathbf{x} = \arg \min_{\mathbf{v} \in \mathbb{R}^n} f(\mathbf{v}) \Rightarrow \nabla f(\mathbf{x}) = \mathbf{o}. \quad (1.6)$$

Lze ukázat, že pro gradient funkce f platí

$$\nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1}(\mathbf{x}), \dots, \frac{\partial f}{\partial x_n}(\mathbf{x}) \right]^T = \frac{1}{2} \mathbf{A}^T \mathbf{x} + \frac{1}{2} \mathbf{Ax} - \mathbf{b} = \mathbf{Ax} - \mathbf{b}.$$

Z podmínky (1.6) tedy vyplývá $\mathbf{Ax} - \mathbf{b} = \mathbf{o}$. \square

2 LDMT, LDLT a Choleského dekompozice

2.1 LDMT (LDU) dekompozice

Uvažujme rozklad obecné čtvercové matice $\mathbf{A} \in \mathbb{R}^{n \times n}$

$$\mathbf{A} = \mathbf{LDM}^T,$$

kde $\mathbf{L} \in \mathbb{R}^{n \times n}$ a $\mathbf{M} \in \mathbb{R}^{n \times n}$ jsou dolní trojúhelníkové matice s jedničkami na diagonálách a $\mathbf{D} \in \mathbb{R}^{n \times n}$ je diagonální matice. Matici tedy rozkládáme do následujícího tvaru (\times značí nenulové prvky matice):

$$\mathbf{A} = \mathbf{LDM}^T = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ \times & 1 & 0 & 0 & 0 \\ \times & \times & 1 & 0 & 0 \\ \times & \times & \times & 1 & 0 \\ \times & \times & \times & \times & 1 \end{bmatrix} \begin{bmatrix} \times & 0 & 0 & 0 & 0 \\ 0 & \times & 0 & 0 & 0 \\ 0 & 0 & \times & 0 & 0 \\ 0 & 0 & 0 & \times & 0 \\ 0 & 0 & 0 & 0 & \times \end{bmatrix} \begin{bmatrix} 1 & \times & \times & \times & \times \\ 0 & 1 & \times & \times & \times \\ 0 & 0 & 1 & \times & \times \\ 0 & 0 & 0 & 1 & \times \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Při odvozování vhodného algoritmu budeme vycházet z předpokladu, že již známe $j - 1$ sloupců matice \mathbf{L} , řádků matice \mathbf{M} (tedy sloupců matice \mathbf{M}^T) a diagonálních prvků d_1, d_2, \dots, d_{j-1} matice \mathbf{D} , $1 \leq j \leq n$. Chceme odvodit vzorec pro j -tý sloupec matice \mathbf{L} (tedy prvky $\mathbf{L}(j+1:n, j)$), j -tý řádek matice \mathbf{M} (tedy prvky $\mathbf{M}(j, j+1:n)$ neboli $\mathbf{M}^T(j+1:n, j)$) a diagonální prvek d_j .

Sloupec j matice \mathbf{A} můžeme vyjádřit pomocí našeho rozkladu jako

$$\mathbf{A}(1:n, j) = (\mathbf{LDM}^T)(1:n, j) = \mathbf{LDM}^T \mathbf{e}_j = \mathbf{L}(\mathbf{DM}^T \mathbf{e}_j) = \mathbf{Lv},$$

kde \mathbf{e}_j je j -tý vektor standardní báze prostoru \mathbb{R}^n (tedy vektor, který má 1 na j -té pozici a jinde nuly). Součin $\mathbf{DM}^T \mathbf{e}_j$ jsme si označili \mathbf{v} a z tvarů matic \mathbf{D}, \mathbf{M}^T není těžké ověřit, že tento vektor má nenulové prvky pouze na prvních j pozicích. Uvědomme si, že \mathbf{v} reprezentuje j -tý sloupec součinu matic \mathbf{DM}^T a vynásobíme-li jím matici \mathbf{L} , získáme j -tý sloupec matice \mathbf{A} .

Rozdělme nyní sloupec $\mathbf{A}(1:n, j)$ na dvě části - $\mathbf{A}(1:j, j)$ a $\mathbf{A}(j+1:n, j)$. Pro první část platí

$$A(1:j, j) = L(1:j, 1:j)\mathbf{v}(1:j).$$

Výraz na levé straně známe, submatici $L(1:j, 1:j)$ známe také (z předpokladu známe $j-1$ sloupců L a víme, že $L(j, j) = 1$). Neznámý vektor $\mathbf{v}(1:j)$ tedy získáme řešením této soustavy j rovnic o j neznámých s dolní trojúhelníkovou maticí (k řešení můžeme použít efektivní dopřednou substituci).

Po nalezení vektoru $\mathbf{v}(1:j)$ určíme příslušné prvky matice M . Vektor \mathbf{v} je j -tým sloupcem součinu DM^T . Z tvarů matic D, M^T tedy vyplývá, že pro prvky $\mathbf{v}(1:j)$ platí

$$\begin{aligned}\mathbf{v}(1:j) &= D(1:j, 1:j)M^T(1:j, j) \\ D(1:j, 1:j)^{-1}\mathbf{v}(1:j) &= M^T(1:j, j)\end{aligned}$$

Protože matice D je diagonální, k jejímu invertování stačí pouze invertovat diagonální prvky. Po přenásobení inverzní diagonální matice vektorem $\mathbf{v}(1:j)$ dostaneme pro $i = 1, \dots, j-1$ rovnost

$$\frac{1}{d_i}\mathbf{v}(i) = M^T(i, j),$$

tedy $M(j, i) = \frac{1}{d_i}\mathbf{v}(i)$. Zbývá určit neznámý prvek d_j . Protože $M(j, j) = 1$, musí platit

$$1 = \frac{1}{d_j}\mathbf{v}(j),$$

tedy $d_j = \mathbf{v}(j)$.

Určili jsme tedy j -tý řádek matice M a prvek d_j , zbývá určit j -tý sloupec matice L . K tomu využijeme druhou část sloupce $A(1:n, j)$. Tu si můžeme vyjádřit jako

$$A(j+1:n, j) = L(j+1:n, 1:j)\mathbf{v}(1:j)$$

Výraz na levé straně opět známe, známe již také vektor $\mathbf{v}(1:j)$ a všechny sloupce submatice $L(j+1:n, 1:j)$ kromě posledního, tedy kromě $L(j+1:n, j)$. Součin na pravé straně si můžeme rozepsat na

$$A(j+1:n, j) = L(j+1:n, 1:j-1)\mathbf{v}(1:j-1) + L(j+1:n, j)\mathbf{v}(j).$$

Osamostatněním neznámého sloupce získáme

$$L(j+1:n, j) = \frac{A(j+1:n, j) - L(j+1:n, 1:j-1)\mathbf{v}(1:j-1)}{\mathbf{v}(j)}.$$

Výsledný algoritmus tedy můžeme zapsat takto:

```
function LDMT(A)
    n = size(A)
    L = eye(n, n), M = eye(n, n), D = zeros(n, n)
    v(1) = A(1, 1)
    D(1, 1) = v(1)
    L(2:n, 1) = A(2:n, 1)/v(1)
```

```

for  $j = 2, \dots, n$  do
  Solve  $L(1:j, 1:j)v(1:j) = A(1:j, j)$  for  $v(1:j)$ 
  for  $i = 1, \dots, j-1$  do
     $M(j, i) = v(i)/D(i, i)$ 
  end for
   $D(j, j) = v(j)$ 
   $L(j+1:n, j) = \frac{A(j+1:n, j) - L(j+1:n, 1:j-1)v(1:j-1)}{v(j)}$ 
end for
end function

```

2.2 LDLT dekompozice

Dá se ukázat, že je-li matice A na vstupu algoritmu pro LDMT dekompozici symetrická, platí $L = M$. Matici jsme tedy schopni rozložit na součin $A = LDL^T$. Vhodnou úpravou předchozího algoritmu jsme schopni snížit jeho náročnost pro symetrické matice na polovinu.

Pro odvození tedy opět vyjděme z předpokladu, že známe $j-1$ sloupců matice L , řádků M a diagonálních prvků matice D . Navíc, protože $M = L$, známe i prvky v příslušných *sloupcích* matice M , včetně prvků $M(j, 1:j-1)$, tedy $M^T(1:j-1, j)$. Vzpomeňme si, co platilo pro vektor $v(1:j)$:

$$v(1:j) = D(1:j, 1:j)M^T(1:j, j) = D(1:j, 1:j)L^T(1:j, j)$$

Díky tomu, že D je diagonální matice, můžeme si součin na pravé straně snadno rozepsat na

$$v(1:j) = \begin{bmatrix} D(1,1)L^T(1,j) \\ D(2,2)L^T(2,j) \\ \vdots \\ D(j-1,j-1)L^T(j-1,j) \\ D(j,j) \end{bmatrix}$$

Díky předpokladům známe všechny prvky vektoru napravo kromě posledního. Získali jsme tedy předpis, který nám říká, jak snadno získat $v(1:j-1)$. Na posledním řádku jsme navíc využili toho, že $L(j, j) = 1$. Zbývá nalézt $v(j)$. K tomu si vyjádříme, čemu se rovná $A(j, j)$:

$$A(j, j) = L(j, 1:j)v(1:j)$$

Rozepišme součin napravo opět na známou a neznámou část

$$A(j, j) = L(j, 1:j-1)v(1:j-1) + L(j, j)v(j) = L(j, 1:j-1)v(1:j-1) + 1v(j)$$

a vyjádříme neznámý prvek $v(j)$

$$v(j) = A(j, j) - L(j, 1:j-1)v(1:j-1).$$

Tímto jsme získali $v(1:j)$ aniž bychom museli řešit soustavu rovnic, čímž jsme původní algoritmus výrazně urychlili. Zbývající část odvození je identická. Výsledný algoritmus pro LDLT rozklad symetrické matice tedy můžeme zapsat takto:

```

function LDLT(A)
    n = size(A)
    L = eye(n,n), D = zeros(n,n)
    v(1) = A(1,1)
    D(1,1) = v(1)
    L(2:n,1) = A(2:n,1)/v(1)
    for j = 2, ..., n do
        for i = 1, ..., j-1 do
            v(i) = L(j,i)D(i,i)
        end for
        v(j) = A(j,j) - L(j,1:j-1)v(1:j-1)
        D(j,j) = v(j)
        L(j+1:n,j) =  $\frac{A(j+1:n,j) - L(j+1:n,1:j-1)v(1:j-1)}{v(j)}$ 
    end for
end function

```

2.3 Choleského rozklad

Matici A , pro kterou platí $a_{ij} = a_{ji}$ a

$$\forall \mathbf{x} \neq 0 : \mathbf{x}^T \mathbf{A} \mathbf{x} > 0$$

nazýváme symetrickou pozitivně definitní. Tyto matice často vznikají diskretizací fyzikálních systémů. Zopakujme si některé jejich vlastnosti:

1. Každá pozitivně definitní matice má kladné prvky na diagonále ($a_{ii} = \mathbf{e}_i^T \mathbf{A} \mathbf{e}_i > 0$, nerovnost vyplývá z definice pozitivní definitnosti)
2. Matice je pozitivně definitní, je-li symetrická a má-li všechna vlastní čísla kladná.
3. Každá pozitivně definitní matice může být rozložena na

$$\mathbf{A} = \mathbf{R}^T \mathbf{R},$$

kde \mathbf{R} je horní trojúhelníková matice s kladnými prvky na diagonále.

Dekompozici z bodu 3 nazýváme Choleského rozklad. Pokusme se odvodit jeho základní tvar. Rozdělme matice \mathbf{A} a \mathbf{R} na čtyři bloky:

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}(1,1) & \mathbf{A}(1,2:n) \\ \mathbf{A}(2:n,1) & \mathbf{A}(2:n,2:n) \end{bmatrix}, \mathbf{R} = \begin{bmatrix} \mathbf{R}(1,1) & \mathbf{R}(1,2:n) \\ 0 & \mathbf{R}(2:n,2:n) \end{bmatrix}$$

Hledaný součin má tvar

$$\begin{aligned} \begin{bmatrix} \mathbf{A}(1,1) & \mathbf{A}(1,2:n) \\ \mathbf{A}(2:n,1) & \mathbf{A}(2:n,2:n) \end{bmatrix} &= \begin{bmatrix} \mathbf{R}(1,1) & 0 \\ \mathbf{R}^T(1,2:n) & \mathbf{R}^T(2:n,2:n) \end{bmatrix} \begin{bmatrix} \mathbf{R}(1,1) & \mathbf{R}(1,2:n) \\ 0 & \mathbf{R}(2:n,2:n) \end{bmatrix} = \\ &= \begin{bmatrix} \mathbf{R}^2(1,1) & \mathbf{R}(1,1)\mathbf{R}(1,2:n) \\ \mathbf{R}(1,1)\mathbf{R}^T(1,2:n) & \mathbf{R}^T(1,2:n)\mathbf{R}(1,2:n) + \mathbf{R}^T(2:n,2:n)\mathbf{R}(2:n,2:n) \end{bmatrix} \end{aligned}$$

Porovnáním bloků na odpovídajících pozicích na levé a pravé straně nejdříve snadno vypočítáme první řádek matice R . Musí platit $R^2(1,1) = A(1,1)$, takže $R(1,1) = \sqrt{A(1,1)}$ (uvědomme si, že díky pozitivní definitnosti A víme, že diagonální prvek je kladný). Odsud a z $A(1,2:n) = R(1,1)R(1,2:n)$ snadno vyjádříme $R(1,2:n) = A(1,2:n)/R(1,1)$. Zbývá určit blok $R(2:n,2:n)$. Porovnáním bloků na pozicích $(2,2)$ získáme rovnost

$$A(2:n,2:n) - R^T(1,2:n)R(1,2:n) = R^T(2:n,2:n)R(2:n,2:n).$$

Dá se ukázat, že výraz na levé straně je opět symetrická pozitivně definitní matice. Výraz na pravé straně je její Choleského rozklad. K nalezení matice $R(2:n,2:n)$ tedy rekurzivně aplikujeme tento postup na levou stranu předchozí rovnosti. V dalším kroce tedy nalezneme řádek $R(2,2:n)$ a rekurzivně budeme pokračovat, dokud nenalezneme zbývající řádky matice R . Promyslete si, že se algoritmus dá zapsat takto:

```

function CHOLESKY(A)
    m = size(A)
    R = A
    for k = 1, ..., m do
        for j = k + 1, ..., m do
            R(j, j : m) = R(j, j : m) - R(k, j : m)R(k, j) / R(k, k)
        end for
        R(k, k : m) = R(k, k : m) / sqrt(R(k, k))
    end for
end function

```

3 Výpočetní náročnost

3.1 Gaussova eliminace (LU rozklad)

Zopakujme si algoritmus LU faktorizace:

```

function LU(A)
    m = size(A)
    U = A
    for k = 1, ..., m - 1 do
        for j = k + 1, ..., m do
            L(j, k) = U(j, k) / U(k, k)
            U(j, k : m) = U(j, k : m) - L(j, k)U(k, k : m)
        end for
    end for
end function

```

Výpočtu dominuje vnitřní smyčka, konkrétně řádek

$$U(j, k : m) = U(j, k : m) - L(j, k)U(k, k : m) \quad (3.1)$$

Vektor $U(k, k : m)$ má v k -tém kroce délku $l = m - k + 1$. Na řádku tento vektor přenásobíme skalárem (l násobení) a odečteme od vektoru $U(j, k : m)$ (l rozdílů). Celkem tedy na tomto řádku provedeme $2l$ operací. Zajímá nás, kolikrát se (3.1) provede v průběhu celého běhu algoritmu a s jak dlouhými vektory při tom pracuje. Rozepišme si postupně počet volání (3.1) a délku l pro jednotlivé kroky k :

$$\begin{aligned} k = 1 : \text{počet: } m - 1, \quad l = m \\ k = 2 : \text{počet: } m - 2, \quad l = m - 1 \\ \vdots \\ k = m - 1 : \text{počet: } 1, \quad l = m - (m - 1) + 1 = 2. \end{aligned}$$

Celkový počet operací q je tedy dán výrazem

$$q = 2((m - 1)m + (m - 2)(m - 1) + \dots + 2 \cdot 3 + 1 \cdot 2),$$

tedy

$$q = \sum_{i=1}^m 2(m-i)(m-i+1) = 2 \sum_{j=0}^{m-1} j(j+1) = 2 \sum_{j=0}^{m-1} (j^2 + j) = 2 \left(\sum_{j=0}^{m-1} j^2 + \sum_{j=0}^{m-1} j \right).$$

K vyčíslení těchto sum můžeme použít vzorce $\sum_{j=0}^m j^2 = \frac{1}{6}(2m+1)(m+1)m$ a $\sum_{j=0}^m j = \frac{m(m+1)}{2}$. Dostaneme tedy (m^2 a m odčítáme, protože naše sumy sčítají pouze po $j = m - 1$):

$$q = 2 \left(\frac{1}{6}(2m+1)(m+1)m - m^2 + \frac{m(m+1)}{2} - m \right) = 2 \left(\frac{1}{3}m^3 - \frac{1}{3}m \right)$$

Výpočtu dominuje třetí mocnina, proto můžeme říct, že náročnost Gaussovy eliminace (LU rozkladu), je $\approx \frac{2}{3}m^3$.

3.2 LU rozklad s částečnou a úplnou pivotizací

Určeme nyní náročnost částečné pivotizace při výpočtu LU rozkladu. V každém kroce $k = 1, \dots, k = m - 1$ hledáme největší pivot ve sloupci délky $m - k + 1$. Celkový počet porovnání je tedy

$$m + (m - 1) + (m - 2) + \dots + 2 = \frac{(m + 2)(m + 1)}{2} = \frac{1}{2}(m^2 + m - 2)$$

Výpočtu dominuje druhá mocnina, můžeme tedy říct, že celkový počet porovnání je $\approx \frac{1}{2}m^2$, tedy náročnost částečné pivotizace je $O(m^2)$. Vzhledem k tomu, že náročnost samotného LU rozkladu je kubická, nepřináší částečná pivotizace významný overhead.

Podobně bychom mohli odvodit, že náročnost úplné pivotizace je $O(m^3)$, což výrazně přispěje k původní náročnosti metody. Z tohoto důvodu se úplná pivotizace často nepoužívá.

3.3 Další algoritmy

Stejným způsobem lze odvodit, že náročnost LDMT rozkladu je $2/3m^3$, LDLT a Choleského rozklad symetrické matice vyžadují poloviční počet operací, tedy $1/3m^3$. Řešení soustavy s dolní nebo horní trojúhelníkovou maticí pomocí dopředné nebo zpětné substituce vyžaduje $1/2m^2$ operací. Pro srovnání výpočet inverzní matice vyžaduje m^3 operací a řešení systému pomocí inverzní matice (násobení inverzní maticí) vyžaduje m^2 operací.

4 Iterační metody

Jak jsme si ukázali v předchozí kapitole, přímé řešiče vyžadují $O(m^3)$ operací, což je pro praktické použití s velkými maticemi příliš mnoho. Proto velikost soustav s plnými maticemi, které jsme schopni řešit pomocí přímých řešičů, nerostla v historii příliš rychle:

- 1950: $m = 20$ (Wilkinson)
- 1965: $m = 200$ (Forsythe, Moler)
- 1980: $m = 2000$ (LINPACK)
- 1995: $m = 20000$ (LAPACK)
- 2010: $m = 200000$ (HDSS)

V případě řídkých systémů jsme sice schopni řešit soustavy s ještě větším počtem neznámých, nicméně použití přímých řešičů může vést k zaplnění systémů. Uvědomme si např., co se stane pokud eliminujeme první sloupec následující matice pomocí prvního řádku:

$$\begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & 0 & 0 & 0 \\ \times & 0 & \times & 0 & 0 \\ \times & 0 & 0 & \times & 0 \\ \times & 0 & 0 & 0 & \times \end{bmatrix} \Rightarrow \begin{bmatrix} \times & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \\ 0 & \times & \times & \times & \times \end{bmatrix}$$

Alternativou je použít iterační řešiče, ty generují posloupnost přibližných řešení $\{x^{(k)}\}$ a pracují téměř výhradně s násobením matice a vektoru (což vyžaduje $O(m^2)$ operací). Podstatnou vlastností každého iteračního řešiče je tedy rychlost konvergence $\{x^{(k)}\}$ ke správnému řešení. Pokud bychom totiž k řešení dokonvergovali až po m a více iteracích, potřebovali bychom opět minimálně m^3 operací.

4.1 Lineární iterační řešiče

Hledejme posloupnost přibližných řešení rovnice $A\mathbf{x} = \mathbf{b}$, kde $A \in \mathbb{R}^{n \times n}$ je regulární matice, ve tvaru

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

(někdy se v literatuře setkáte také s tvarem $R\mathbf{x}^{(k+1)} := W\mathbf{x}^{(k)} + \mathbf{b}$). Horní index (k) zde znamená k -tou iteraci.

Řekneme, že metoda je *konzistentní*, pokud rovnici $\mathbf{x} := M\mathbf{x} + N\mathbf{b}$ řeší právě jeden vektor $\mathbf{x} = A^{-1}\mathbf{b}$. Lze ukázat, že toto platí právě tehdy, když $M = I - NA$.

Řekneme, že metoda je *konvergentní*, platí-li $\mathbf{x}^{(k)} \rightarrow \mathbf{x}$ pro $k \rightarrow \infty$. Toto platí právě tehdy, když $\|M\| < 1$.

Rozložme nyní matici A na součet $A = L + D + U$. Např. pro matici

$$A = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix}$$

by rozklad vypadal takto

$$L = \begin{bmatrix} 0 & 0 & 0 \\ -1 & 0 & 0 \\ 0 & -1 & 0 \end{bmatrix}, D = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix}, U = \begin{bmatrix} 0 & -1 & 0 \\ 0 & 0 & -1 \\ 0 & 0 & 0 \end{bmatrix}$$

4.1.1 Jacobiho metoda

V případě Jacobiho metody vyjdeme ze zápisu soustavy $A\mathbf{x} = \mathbf{b}$ ve tvaru

$$(L + D + U)\mathbf{x} = \mathbf{b}.$$

Roznásobme závorku do následujícího tvaru

$$D\mathbf{x} + (L + U)\mathbf{x} = \mathbf{b}$$

a převedme člen s $(L + U)$ na pravou stranu

$$D\mathbf{x} = -(L + U)\mathbf{x} + \mathbf{b}.$$

Jacobiho iterace pak definujeme takto

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

Přemnožením inverzí k D získáme konečný tvar

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

Zde $M = -D^{-1}(L + U)$ a $N = D^{-1}$. Vzhledem k tomu, že D je diagonální matice a že $(L + U)\mathbf{x}^{(k)}$ je násobením vektoru $\mathbf{x}^{(k)}$ původní maticí A bez její diagonály, můžeme napsat předpis pro i -tý prvek vektoru $\mathbf{x}^{(k+1)}$ takto:

$$\begin{aligned} (\mathbf{x}^{(k+1)})_i &:= \frac{1}{a_{ii}} \left(b_i - \sum_{j=1, j \neq i}^n a_{ij} (\mathbf{x}^{(k)})_j \right) = \\ &= \frac{1}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij} (\mathbf{x}^{(k)})_j - \sum_{j=i+1}^n a_{ij} (\mathbf{x}^{(k)})_j \right) \end{aligned} \quad (4.1)$$

Konzistence metody vyplývá ze způsobu, jakým jsme ji odvodili, tedy

$$\mathbf{Ax} = (\mathbf{L} + \mathbf{D} + \mathbf{U})\mathbf{x} = \mathbf{b},$$

tedy

$$\mathbf{x} = -\mathbf{D}^{-1}(\mathbf{L} + \mathbf{U})\mathbf{x} + \mathbf{D}^{-1}\mathbf{b}. \quad (4.2)$$

Řeší-li $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ původní rovnici $\mathbf{Ax} = \mathbf{b}$, řeší také (4.2). Metoda je tedy konzistentní. Konzistenci můžeme také ověřit zjištěním, zda platí rovnost $\mathbf{M} = \mathbf{I} - \mathbf{NA}$:

$$\begin{aligned} \mathbf{I} - \mathbf{NA} &= \mathbf{I} - \mathbf{D}^{-1}\mathbf{A} = \mathbf{I} - \mathbf{D}^{-1}(\mathbf{L} + \mathbf{D} + \mathbf{U}) \\ &= \mathbf{I} - \mathbf{D}^{-1}\mathbf{D} - \mathbf{D}^{-1}(\mathbf{L} + \mathbf{U}) = -\mathbf{D}^{-1}(\mathbf{L} + \mathbf{U}) = \mathbf{M}. \end{aligned}$$

Jacobiho metoda je konvergentní, platí-li $\|\mathbf{M}\| = \|\mathbf{D}^{-1}(\mathbf{L} + \mathbf{U})\| < 1$, což splňují např. diagonálně dominantní matice, tedy takové matice, pro které platí: $\forall i : |a_{ii}| > \sum_{j \neq i} a_{ij}$.

4.1.2 Gaussova-Seidelova metoda

Všimneme si, že počítáme-li Jacobiho metodou i -tý prvek vektoru $\mathbf{x}^{(k+1)}$, tedy $(\mathbf{x}^{(k+1)})_i$, známe již všech $i - 1$ předchozích prvků tohoto vektoru. V předpisu (4.1) můžeme tedy člen $\sum_{j=1}^{i-1} a_{ij}(\mathbf{x}^{(k)})_j$ nahradit členem $\sum_{j=1}^{i-1} a_{ij}(\mathbf{x}^{(k+1)})_j$. Touto změnou dostaneme předpis pro Gaussovou-Seidelovu metodu:

$$(\mathbf{x}^{(k+1)})_i := \frac{1}{a_{ii}}(b_i - \sum_{j=1}^{i-1} a_{ij}(\mathbf{x}^{(k+1)})_j - \sum_{j=i+1}^n a_{ij}(\mathbf{x}^{(k)})_j), \quad (4.3)$$

což lze zapsat maticově jako

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

Převedením $\mathbf{x}^{(k+1)}$ na levou stranu dostaneme

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

odtud

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

Pro Gaussovou-Seidelovu metodu tedy platí $\mathbf{M} = -(\mathbf{D} + \mathbf{L})^{-1}\mathbf{U}$, $\mathbf{N} = (\mathbf{D} + \mathbf{L})^{-1}$. Ověřme její konzistenci:

$$\begin{aligned} \mathbf{I} - \mathbf{NA} &= \mathbf{I} - (\mathbf{D} + \mathbf{L})^{-1}\mathbf{A} = \mathbf{I} - (\mathbf{D} + \mathbf{L})^{-1}(\mathbf{L} + \mathbf{D} + \mathbf{U}) \\ &= \mathbf{I} - (\mathbf{D} + \mathbf{L})^{-1}(\mathbf{D} + \mathbf{L}) - (\mathbf{D} + \mathbf{L})^{-1}\mathbf{U} = -(\mathbf{D} + \mathbf{L})^{-1}\mathbf{U} = \mathbf{M}. \end{aligned}$$

Metoda je tedy konzistentní. Gaussova-Seidelova metoda je konvergentní, právě když $\|\mathbf{M}\| = \|-(\mathbf{D} + \mathbf{L})^{-1}\mathbf{U}\| < 1$.

4.1.3 Richardsonova metoda

Předpis pro $k + 1$. krok Richardsonovy metody je dán

$$\mathbf{x}^{k+1} := \mathbf{x}^{(k)} + \omega \mathbf{r}^{(k)},$$

kde $\mathbf{r}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}$ je tzv. reziduum a $\omega > 0$.

Metoda je konzistentní pro libovolné ω ($\mathbf{x} = \mathbf{x} + \omega(\mathbf{b} - \mathbf{A}\mathbf{x})$), $\mathbf{x} = \mathbf{x} + \omega\mathbf{0}$).

Studujme nyní konvergenci metody pro \mathbf{A} symetrickou pozitivně definitní matici. V takovém případě vlastní čísla λ_i a vlastní vektory \mathbf{v}_i (tedy λ_i a \mathbf{v}_i , pro které platí $\mathbf{A}\mathbf{v}_i = \lambda_i\mathbf{v}_i$) splňují

$$0 < \lambda_1 \leq \lambda_2 \leq \lambda_3 \leq \dots \leq \lambda_n$$

a $\{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \dots, \mathbf{v}_n\}$ tvoří ortonormální bázi prostoru \mathbb{R}^n . Každý prvek toho prostoru tedy můžeme vyjádřit jako lineární kombinaci vektorů \mathbf{v}_i .

Studujme nyní, jak se šíří reziduum $\mathbf{r}^{(k+1)}$. Nejdříve si vyjádříme $\mathbf{r}^{(k+1)}$ a $\mathbf{r}^{(k)}$ v bázi dané vektory \mathbf{v}_i .

$$\mathbf{r}^{(k+1)} = \sum_{i=1}^n \alpha_i^{(k+1)} \mathbf{v}_i, \quad \mathbf{r}^{(k)} = \sum_{i=1}^n \alpha_i^{(k)} \mathbf{v}_i$$

a dále postupujme takto:

$$\begin{aligned} \mathbf{r}^{(k+1)} &= \mathbf{b} - \mathbf{A}\mathbf{x}^{(k+1)} = \mathbf{b} - \mathbf{A}(\mathbf{x}^{(k)} + \omega \mathbf{r}^{(k)}) = (\mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}) - \omega \mathbf{A}\mathbf{r}^{(k)} \\ &= \mathbf{r}^{(k)} - \omega \mathbf{A}\mathbf{r}^{(k)} = (I - \omega \mathbf{A})\mathbf{r}^{(k)} = (I - \omega \mathbf{A}) \sum_{i=1}^n \alpha_i^{(k)} \mathbf{v}_i \\ &= \sum_{i=1}^n \alpha_i^{(k)} \mathbf{v}_i - \sum_{i=1}^n \omega \lambda_i \alpha_i^{(k)} \mathbf{v}_i = \sum_{i=1}^n (\alpha_i^{(k)} \mathbf{v}_i - \omega \lambda_i \alpha_i^{(k)} \mathbf{v}_i) \\ &= \sum_{i=1}^n (1 - \omega \lambda_i) \alpha_i^{(k)} \mathbf{v}_i. \end{aligned}$$

Všimneme si tedy, že jsme dostali rovnost

$$\mathbf{r}^{(k+1)} = \sum_{i=1}^n \alpha_i^{(k+1)} \mathbf{v}_i = \sum_{i=1}^n (1 - \omega \lambda_i) \alpha_i^{(k)} \mathbf{v}_i.$$

Vyjádřili jsme tedy koeficienty $\alpha_i^{(k+1)}$ pomocí koeficientů v předchozím kroce, tedy $\alpha_i^{(k)}$, přenásobených konstantami $(1 - \omega \lambda_i)$. Konvergence rezidua k nulovému vektoru tedy bude zajištěna, právě když $\forall i : |1 - \omega \lambda_i| < 1$. Platí $\|\mathbf{r}^{(k+1)}\| \leq \varrho \|\mathbf{r}^{(k)}\|$, kde ϱ je konvergenční faktor $\varrho = \max_{i \in \{1, \dots, n\}} |1 - \omega \lambda_i|$.

Jednoduchou úvahou můžeme odvodit vhodné hodnoty parametru ω .