

Poznámky k předmětu Numerická lineární algebra I.

Michal Merta*

*Katedra aplikované matematiky, VŠB-Technická univerzita Ostrava, e-mail: michal.merta@vsb.cz

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 $A\mathbf{x} = \mathbf{b}$ s regulární maticí 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 $Mx + Nx = b$ právě jeden vektor $x = A^{-1}b$. Je možné ukázat, že metoda je konzistentní právě tehdy, je-li splněno $M = I - NA$.

Definice Iterační metodu nazveme konvergentní platí-li $\forall b, \forall x^0 : x^k \rightarrow x = A^{-1}b$ pro $k \rightarrow \infty$. Je možné ukázat, že metoda je konvergentní, právě tehdy, je-li splněno $\|M\| < 1$, kde $\|M\| = \max_{v \in \mathbb{R}^n, v \neq 0} \frac{\|Mv\|_v}{\|v\|_v}$ 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

Vyjděme z rovnice $Ax = b$. Dosazením $A = L + D + U$ dostaneme

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

Roznásobme a přezávorkujme výraz na levé straně

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

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

$$Dx^{k+1} + (L + U)x^k = b.$$

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

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

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 $\mathbf{M} = \mathbf{I} - \mathbf{N}\mathbf{A}$. V případě Jacobiho metody je $\mathbf{M} = -\mathbf{D}^{-1}(\mathbf{L} + \mathbf{U})$ a $\mathbf{N} = \mathbf{D}^{-1}$ (viz výše). Platí tedy

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

Metoda je tedy konzistentní.

Lze dokázat, že metoda je konvergentní právě tehdy, když $\|\mathbf{M}\| = \|\mathbf{D}^{-1}(\mathbf{L} + \mathbf{U})\| < 1$. 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 kroku. Chybu v kroku $k + 1$ 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 kroku $k + 1$ je tedy v absolutní hodnotě menší než maximální prvek vektoru chyby v předchozím kroku. 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 $\|(L + D)^{-1}U\| < 1$, což opět platí pro diagonálně dominantní matice.

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}

$$\mathbf{A}\mathbf{e}^k = \mathbf{A}\mathbf{x}^k - \mathbf{A}\mathbf{x} = \mathbf{A}\mathbf{x}^k - \mathbf{b} = -\mathbf{r}^k. \quad (1.6)$$

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

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

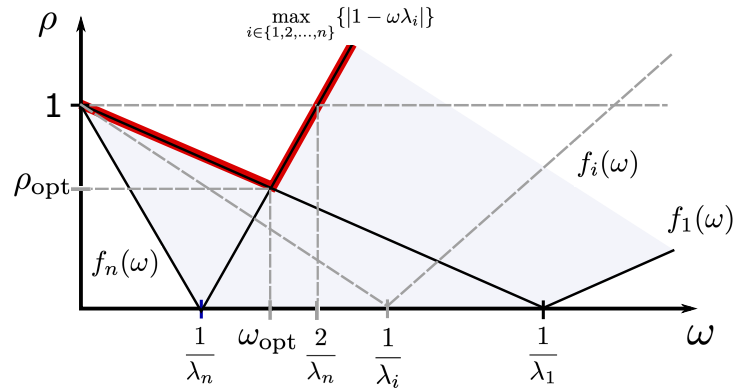
a z vlastních vektorů lze utvořit ortonormální bázi \mathbb{R}^n .

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.

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

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 kroku (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.

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: $\|\mathbf{r}^{k+1}\| \leq \rho \|\mathbf{r}^k\|$.



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

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).

1.1.4 Ukončovací podmínky

Při použití iteračního řešiče většinou nemáme předem zadaný počet iterací, které mají proběhnout. Chceme výpočet ukončit ve chvíli, kdy se s odhadem řešení dostaneme dostatečně blízko přesnému řešení. Vzhledem k tomu, že přesné řešení (tedy ani přesnou chybu v dané iteraci) neznáme, musíme si pomoci jinak.

Jednou z možností je ukončit cyklus ve chvíli, kdy se s novým odhadem řešení příliš nepohneme od předchozího odhadu (tzn. $\|\mathbf{x}^{k+1} - \mathbf{x}^k\| < \varepsilon$). Tato podmínka ale nijak nebere v potaz velikost prvků v matici soustavy a vektoru pravé strany (jiná situace nastane, pokud jsou prvky matice a vektoru v řádech tisíců, jiná pokud jsou v řádech tisícín). Proto je vhodné tuto podmínku zvolit relativně např. vzhledem k normě vektoru pravé strany (tzn. $\|\mathbf{x}^{k+1} - \mathbf{x}^k\| < \|\mathbf{b}\| \varepsilon$, tedy $\|\mathbf{x}^{k+1} - \mathbf{x}^k\|/\|\mathbf{b}\| < \varepsilon$).

Nejčastěji se ovšem k výpočtu ukončovací podmínky používá normy vektoru rezidua $\mathbf{r}^{k+1} = \mathbf{b} - \mathbf{A}\mathbf{x}^{k+1}$. To nám poskytuje přirozený odhad toho, jak dobře je splněna původní rovnice. Ukončovací podmínku lze tedy volit ve tvaru $\|\mathbf{b} - \mathbf{A}\mathbf{x}^{k+1}\| < \varepsilon$. Podobně jako v předchozím případě je i zde vhodnější použít relativní změnu rezidua oproti vektoru pravé strany ($\|\mathbf{b} - \mathbf{A}\mathbf{x}^{k+1}\|/\|\mathbf{b}\| < \varepsilon$) nebo počátečnímu reziduu ($\|\mathbf{b} - \mathbf{A}\mathbf{x}^{k+1}\|/\|\mathbf{b} - \mathbf{A}\mathbf{x}^0\| < \varepsilon$).

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 $\mathbf{Ax} = \mathbf{b} \Rightarrow \mathbf{x} = \arg \min_{\mathbf{v} \in \mathbb{R}^n} f(\mathbf{v})$. Podívejme se, jak se změni 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:

$$\mathbf{x} = \arg \min_{\mathbf{v} \in \mathbb{R}^n} f(\mathbf{v}) \Rightarrow \nabla f(\mathbf{x}) = \mathbf{o}. \quad (1.7)$$

Lze ukázat, že pro gradient funkce f platí

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

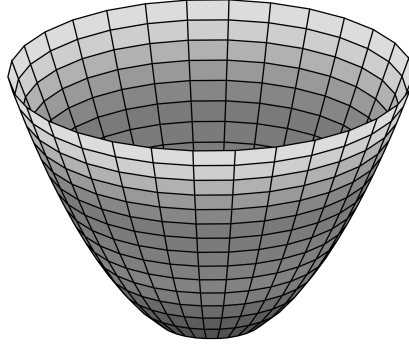
Z podmínky (1.7) tedy vyplývá $\mathbf{Ax} - \mathbf{b} = \mathbf{o}$. \square

V případě symetrické pozitivně definitní matice $\mathbf{A} \in \mathbb{R}^{n \times n}$ máme tedy dvě možnosti, jak geometricky nahlížet na řešení soustavy lineárních rovnic. První přístup je chápat každou rovnici jako předpis nadroviny v n -rozměrném prostoru. Řešení soustavy pak odpovídá hledání průsečíku těchto rovin. Druhý přístup, který využijeme při odvozování následujících algoritmů, odpovídá minimalizaci příslušné pozitivně definitní kvadratické formy. Grafem pozitivně definitní kvadratické formy $f : \mathbb{R}^n \rightarrow \mathbb{R}$ je n -dimenzionální paraboloid, který má minimum (viz Obrázek 1.2 pro $n = 2$).

1.2.1 Metoda největšího spádu

Metoda největšího spádu je iterační metoda s předpisem

$$\mathbf{x}^{k+1} := \mathbf{x}^k + \alpha^k \mathbf{v}^k, \quad (1.8)$$



Obrázek 1.2: Graf kvadratické formy s pozitivně definitní maticí \mathbf{A} (zdroj Wikipedia)

kde \mathbf{v}^k volíme jako směr největšího poklesu funkce f . Všimněme si, že pro gradient platí $\nabla f(\mathbf{x}^k) = \mathbf{A}\mathbf{x}^k - \mathbf{b}$ a pro reziduum k -tém kroku $\mathbf{r}^k = \mathbf{b} - \mathbf{A}\mathbf{x}^k$. Tedy $\mathbf{r}^k = -\nabla f(\mathbf{x}^k)$. Protože gradient odpovídá směru největšího růstu funkce v daném bodě, reziduum je směr největšího spádu. Logicky, protože chceme dosáhnout minima dané funkce, vydáváme se v každém kroku ve směru rezidua, tedy $\mathbf{v}^k = \mathbf{r}^k$.

Otázkou je, jak daleko se v každém kroku v tomto směru vydat, tedy jak zvolit koeficient α^k . Metoda největšího spádu volí tento koeficient tak, aby v každém kroku dosáhla minima funkce f ve směru rezidua. Definujme si tedy pomocnou funkci $F : \mathbb{R} \rightarrow \mathbb{R}$:

$$\begin{aligned} F(\alpha) &= f(\mathbf{x}^k + \alpha \mathbf{r}^k) = \frac{1}{2}(\mathbf{x}^k + \alpha \mathbf{r}^k)^T \mathbf{A}(\mathbf{x}^k + \alpha \mathbf{r}^k) - \mathbf{b}^T(\mathbf{x}^k + \alpha \mathbf{r}^k) = \\ &= \frac{1}{2}(\mathbf{x}^k)^T \mathbf{A}\mathbf{x}^k + \alpha(\mathbf{x}^k)^T \mathbf{A}\mathbf{r}^k + \frac{1}{2}\alpha^2(\mathbf{r}^k)^T \mathbf{A}\mathbf{r}^k - \mathbf{b}^T \mathbf{x}^k - \alpha \mathbf{b}^T \mathbf{r}^k \end{aligned}$$

Hledáme α , ve kterém tato funkce dosahuje minima, její derivace se tedy musí rovnat nule:

$$F'(\alpha) = \alpha(\mathbf{r}^k)^T \mathbf{A}\mathbf{r}^k + (\mathbf{r}^k)^T \underbrace{\mathbf{A}\mathbf{x}^k}_{=(\mathbf{b}-\mathbf{r}^k)} - \mathbf{b}^T \mathbf{r}^k = \alpha(\mathbf{r}^k)^T \mathbf{A}\mathbf{r}^k - (\mathbf{r}^k)^T \mathbf{r} = 0$$

Odtud

$$\alpha^k = \frac{(\mathbf{r}^k)^T \mathbf{r}^k}{(\mathbf{r}^k)^T \mathbf{A}\mathbf{r}^k}. \quad (1.9)$$

Stejný předpis můžeme odvodit, použijeme-li místo pomocné funkce funkci f a položíme její derivaci ve směru \mathbf{r}^k rovnu nule (vzpomeňme si, že platí $\frac{df(\mathbf{x})}{d\mathbf{h}} =$

$(\nabla f(\mathbf{x}))^T \mathbf{h}$:

$$\frac{df(\mathbf{x}^{k+1})}{d\mathbf{r}^k} = 0 \quad (1.10)$$

$$(\nabla f(\mathbf{x}^{k+1}))^T \mathbf{r}^k = 0 \quad (1.11)$$

$$(-\mathbf{r}^{k+1})^T \mathbf{r}^k = 0 \quad (1.12)$$

Dosazením $\mathbf{r}^{k+1} = \mathbf{r}^k - \alpha^k \mathbf{A} \mathbf{r}^k$ do předchozí rovnice a jednoduchou úpravou dostaneme stejný předpis pro α^k jako v předchozím případě (1.16). Předchozí odvození nám také prozradilo důležitou vlastnost metody největšího spádu – každý směr $\mathbf{v}^k = \mathbf{r}^k$ je kolmý na předchozí směr. Jak brzy uvidíme, není to vždy žádaná vlastnost.

Algoritmus tedy počítá jednotlivé aproximace pomocí následujících předpisů:

$$\begin{aligned} \mathbf{r}^k &:= \mathbf{b} - \mathbf{A} \mathbf{x}^k = \mathbf{b} - \mathbf{A}(\mathbf{x}^{k-1} + \alpha^{k-1} \mathbf{r}^{k-1}) = \underbrace{\mathbf{b} - \mathbf{A} \mathbf{x}^{k-1}}_{\mathbf{r}^{k-1}} - \alpha^{k-1} \mathbf{A} \mathbf{r}^{k-1} = \\ &= \mathbf{r}^{k-1} - \alpha^{k-1} \mathbf{A} \mathbf{r}^{k-1} \\ \alpha^k &:= \frac{(\mathbf{r}^k)^T \mathbf{r}^k}{(\mathbf{r}^k)^T \mathbf{A} \mathbf{r}^k} \\ \mathbf{x}^{k+1} &:= \mathbf{x}^k + \alpha^k \mathbf{r}^k \end{aligned}$$

Díky úpravě předpisu pro výpočet \mathbf{r}^k jsme ušetřili jedno násobení matice-vektor ($\mathbf{A} \mathbf{x}^k$) – výsledek $\mathbf{A} \mathbf{r}^{k-1}$ si totiž můžeme zapamatovat z předchozí iterace.

Ukažme si nyní, že metoda konverguje. Konvergenci budeme dokazovat v tzv. energetické normě $\|\cdot\|_A$, tedy normě indukované skalárním součinem $(\mathbf{A} \mathbf{x}, \mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x} : \|\mathbf{x}\|_A = \sqrt{\mathbf{x}^T \mathbf{A} \mathbf{x}}$.

$$\begin{aligned} \|e^{k+1}\|_A^2 &= (e^{k+1})^T \mathbf{A} e^{k+1} = (e^k + \alpha^k \mathbf{r}^k)^T \mathbf{A} (e^k + \alpha^k \mathbf{r}^k) = \\ &= (e^k)^T \mathbf{A} e^k + 2\alpha^k (\mathbf{r}^k)^T \underbrace{\mathbf{A} e^k}_{=-\mathbf{r}^k} + (\alpha^k)^2 (\mathbf{r}^k)^T \mathbf{A} \mathbf{r}^k = \\ &= \|e^k\|_A^2 - 2 \underbrace{\frac{(\mathbf{r}^k)^T \mathbf{r}^k}{(\mathbf{r}^k)^T \mathbf{A} \mathbf{r}^k}}_{=\alpha^k} (\mathbf{r}^k)^T \mathbf{r}^k + \left(\frac{(\mathbf{r}^k)^T \mathbf{r}^k}{(\mathbf{r}^k)^T \mathbf{A} \mathbf{r}^k} \right)^2 (\mathbf{r}^k)^T \mathbf{A} \mathbf{r}^k = \\ &= \|e^k\|_A^2 - \frac{((\mathbf{r}^k)^T \mathbf{r}^k)^2}{(\mathbf{r}^k)^T \mathbf{A} \mathbf{r}^k} = \\ &= \|e^k\|_A^2 \left(1 - \frac{((\mathbf{r}^k)^T \mathbf{r}^k)^2}{((\mathbf{r}^k)^T \mathbf{A} \mathbf{r}^k) \underbrace{((e^k)^T \mathbf{A} e^k)}_{=(\mathbf{A}^{-1} \mathbf{r}^k)^T \mathbf{A} (\mathbf{A}^{-1} \mathbf{r}^k) = (\mathbf{r}^k)^T \mathbf{A}^{-1} \mathbf{r}^k}} \right). \end{aligned}$$

Rovnost v poslední závorce vyplývá ze vztahu (1.6) mezi reziduem a chybou. Zjistili jsme tedy, že chyba v kroku $k+1$ je nějakým násobkem chyby v před-

chozím kroku. Podívejme se na tento násobek podrobněji:

$$\begin{aligned} 1 - \frac{((\mathbf{r}^k)^T \mathbf{r}^k)^2}{((\mathbf{r}^k)^T \mathbf{A} \mathbf{r}^k)((\mathbf{r}^k)^T \mathbf{A}^{-1} \mathbf{r}^k)} &\leq 1 - \frac{((\mathbf{r}^k)^T \mathbf{r}^k)^2}{\|\mathbf{A}\| \|\mathbf{r}^k\|^2 \|\mathbf{A}^{-1}\| \|\mathbf{r}^k\|^2} = \\ &= 1 - \frac{1}{\|\mathbf{A}\| \|\mathbf{A}^{-1}\|} = q, \end{aligned}$$

kde první nerovnost vyplývá z Cauchyho-Schwarzovy nerovnosti a definice maticové normy indukované vektorovou normou. Výraz $\|\mathbf{A}\| \|\mathbf{A}^{-1}\|$ nazýváme číslem podmíněnosti (viz také Richardsonova metoda) a v případě symetrické pozitivně definitní matice jej vypočítáme jako podíl největšího a nejmenšího vlastního čísla $\kappa(\mathbf{A}) = \frac{\lambda_{\max}}{\lambda_{\min}}$.

Protože $\kappa(\mathbf{A}) = \lambda_{\max}/\lambda_{\min} \geq 1$, platí $q = 1 - 1/\kappa(\mathbf{A}) < 1$. Chyba se tedy v každém kroku zmenšuje a metoda konverguje k řešení. Navíc

$$\|\mathbf{e}^{k+1}\|_A^2 \leq \left(1 - \frac{1}{\kappa(\mathbf{A})}\right) \|\mathbf{e}^k\|_A^2$$

a je možné dokázat vztah mezi počáteční chybou a chybou v k -té iteraci:

$$\|\mathbf{e}^{k+1}\|_A \leq \left(\frac{\kappa(\mathbf{A}) - 1}{\kappa(\mathbf{A}) + 1}\right)^k \|\mathbf{e}^0\|_A. \quad (1.13)$$

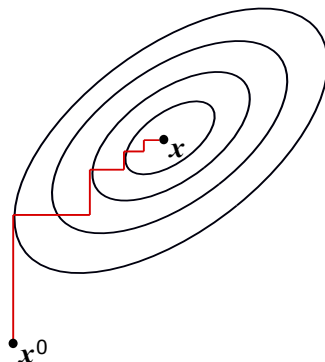
Metoda tedy konverguje k řešení pomalu. Zejména je-li $\kappa(\mathbf{A})$ velké, je zlomek ve výrazu (1.13) blízký 1 a dostaneme velmi pomalou konvergenci.

Celý algoritmus můžeme napsat na několika řádcích:

```
function STEEPESTD_DESCENT( $A$ ,  $\mathbf{b}$ ,  $\mathbf{x}^0$ )
   $\mathbf{r}^0 = \mathbf{b} - \mathbf{A}\mathbf{x}^0$ 
   $k = 0$ 
  while  $\|\mathbf{r}^k\|/\|\mathbf{r}^0\| > \varepsilon$  do
     $\alpha_k = ((\mathbf{r}^k)^T \mathbf{r}^k)/((\mathbf{r}^k)^T \mathbf{A} \mathbf{r}^k)$ 
     $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{r}^k$ 
     $\mathbf{r}^{k+1} = \mathbf{r}^k - \alpha_k \mathbf{A} \mathbf{r}^k$ 
     $k = k + 1$ 
  end while
end function
```

Poznámka Z (1.12) vyplývá, že směr \mathbf{r}^{k+1} , ve kterém se v každém kroku vydáme, je kolmý na předchozí směr. To může vést k situacím jako na Obrázku 1.3, kdy procházíme „cik-cak“ údolí tvořené minimalizovanou kvadratickou formou.

Vydáváme se tedy několikrát v tom samém směru a postupně zkracujeme krok. Ideální by bylo najít metodu, která by minimalizovala celkovou chybu v daném směru vždy pouze jednou (v každé iteraci vynulovala chybu v daném směru). Taková metoda by byla schopná vyřešit soustavu o n neznámých během n iterací. Zkusme si tedy pro $\mathbf{A} \in \mathbb{R}^{n \times n}$ zvolit n navzájem ortogonálních směrů $\mathbf{d}^0, \mathbf{d}^1, \dots, \mathbf{d}^{n-1}$ (např. rovnoběžných s osami souřadnic) a použít předpis: $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{d}^k$.



Obrázek 1.3: Pomalá konvergence metody největšího spádu (vrstevnice minimalizované kvadratické formy a jednotlivé směry \mathbf{r}^k).

Jak ale získat α_k ? Ilustrujme to ve dvou dimenzích na Obrázku 1.4. Na začátku (Obrázek 1.4 vlevo) si zvolíme nějaký počáteční odhad řešení \mathbf{x}^0 ; počáteční chybu \mathbf{e}^0 pak lze rozložit na dvě kolmé složky rovnoběžné s osami souřadnic ($\mathbf{e}^0 = \mathbf{e}_x^0 + \mathbf{e}_y^0$). V první iteraci si zvolme vektor \mathbf{d}^0 rovnoběžný s osou x (Obrázek 1.4 vpravo). Z obrázku je patrné, že abychom vynulovali chybu ve směru osy x , musíme bod \mathbf{x}^1 zvolit tak, aby vektor nové chyby \mathbf{e}^1 byl kolmý k \mathbf{d}^0 .

Obecně tedy chceme, aby nová chyba byla kolmá k aktuálnímu směru:

$$(\mathbf{d}^k)^T \mathbf{e}^{k+1} = 0.$$

Z definice chyby dostaneme

$$0 = (\mathbf{d}^k)^T (\mathbf{x}^{k+1} - \mathbf{x}) = (\mathbf{d}^k)^T (\mathbf{x}^k + \alpha_k \mathbf{d}^k - \mathbf{x}) = (\mathbf{d}^k)^T (\mathbf{e}^k + \alpha_k \mathbf{d}^k)$$

a odtud

$$\alpha_k = -\frac{(\mathbf{d}^k)^T \mathbf{e}^k}{(\mathbf{d}^k)^T \mathbf{d}^k}.$$

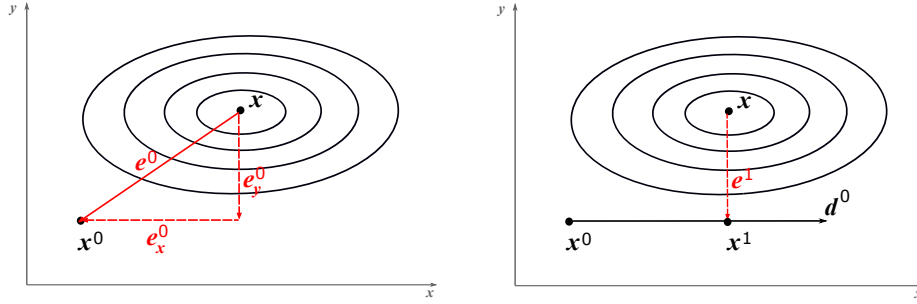
Je jasné, že jsme si příliš nepomohli. Abychom vypočítali α_k , potřebovali bychom znát chybu \mathbf{e}^k . Tu ale bez znalosti přesného řešení nespočítáme. Jak si ale ukážeme v příští lekci, pokud požadavek na ortogonalitu směrů \mathbf{d}^k nahradíme za tzv. A-ortogonalitu, budeme skutečně schopni odvodit metodu, která (alespoň teoreticky) dokonverguje k řešení během n iterací.

1.2.2 Metoda sdružených gradientů

Definice Buď A symetrická pozitivně definitní matice. Nenulové vektory $\{\mathbf{p}_i\}_{i=0}^{n-1}$ nazveme sdružené (A-ortogonální), platí-li

$$\forall i, j \in \{0, 2, \dots, n-1\} : i \neq j \Rightarrow \mathbf{p}_i^T A \mathbf{p}_j = 0.$$

Dva navzájem A-ortogonální vektory někdy označujeme $\mathbf{a} \perp_A \mathbf{b}$.



Obrázek 1.4: Ortogonální složky počáteční chyby (vlevo). Chceme-li vynulovat složku chyby ve směru \mathbf{d}^0 , musí být \mathbf{e}^1 kolmý k \mathbf{d}^0 (vpravo).

Lemma Prvky množiny sdružených vektorů jsou lineárně nezávislé.

Důkaz Chceme ukázat, že platí

$$\alpha_0 \mathbf{p}_0 + \alpha_1 \mathbf{p}_1 + \dots + \alpha_{n-1} \mathbf{p}_{n-1} = \mathbf{0} \Rightarrow \forall i : \alpha_i = 0.$$

Přenasobme rovnici výrazem $(\mathbf{A}\mathbf{p}_k)^T, k \in \{0, 1, \dots, n-1\}$ zleva:

$$\alpha_0 \mathbf{p}_k^T \mathbf{A}\mathbf{p}_0 + \alpha_1 \mathbf{p}_k^T \mathbf{A}\mathbf{p}_1 + \dots + \alpha_{n-1} \mathbf{p}_k^T \mathbf{A}\mathbf{p}_{n-1} = 0.$$

Protože jsou vektory sdružené, všechny členy v této sumě až na jeden jsou nulové:

$$\alpha_k \mathbf{p}_k^T \mathbf{A}\mathbf{p}_k = 0.$$

Současně víme, že matice \mathbf{A} je symetrická pozitivně definitní a součin $\mathbf{p}_k^T \mathbf{A}\mathbf{p}_k$ musí být kladný. Takže $\alpha_k = 0$. \square

Vratme se k poznámce na konci předchozí kapitoly a ortogonální směry nahradme za \mathbf{A} -ortogonální. Předpokládejme tedy, že známe n \mathbf{A} -ortogonálních směrů $\{\mathbf{d}^0, \mathbf{d}^1, \dots, \mathbf{d}^{n-1}\}$ a podobně jako na konci předchozí kapitoly se pokusme odvodit iterační metodu ve tvaru $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{d}^k$. Požadavek na ortogonalitu nové chyby k aktuálnímu směru nahradme za požadavek na jejich vzájemnou \mathbf{A} -ortogonalitu, tzn. $\mathbf{e}^{k+1} \perp_{\mathbf{A}} \mathbf{d}^k$. Ukažme si, že toto nastane v minimu naší kvadratické formy $f(\mathbf{x}) = 1/2 \mathbf{x}^T \mathbf{A}\mathbf{x} - \mathbf{b}^T \mathbf{x}$ ve směru \mathbf{d}^k . Hledejme tedy nový bod \mathbf{x}^{k+1} tak, aby se v něm derivace f ve směru \mathbf{d}^k rovnala nule (nutná podmínka minima):

$$\frac{df(\mathbf{x}^{k+1})}{d\mathbf{d}^k} = \text{grad} f(\mathbf{x}^{k+1})^T \mathbf{d}^k = -(\underbrace{\mathbf{r}^{k+1}}_{= -\mathbf{A}\mathbf{e}^{k+1}})^T \mathbf{d}^k = (\mathbf{e}^{k+1})^T \mathbf{A}\mathbf{d}^k = 0.$$

Chyba \mathbf{e}^{k+1} je tedy v bodu minima skutečně \mathbf{A} -ortogonální ke směru \mathbf{d}^k . Nyní snadno můžeme získat předpis pro koeficient α_k (použijeme již několikrát zmíněný vztah $\mathbf{e}^{k+1} = \mathbf{e}^k + \alpha_k \mathbf{d}^k$):

$$(\mathbf{e}^{k+1})^T \mathbf{A}\mathbf{d}^k = (\mathbf{e}^k + \alpha_k \mathbf{d}^k)^T \mathbf{A}\mathbf{d}^k = 0 \Rightarrow \alpha_k = -\frac{(\mathbf{e}^k)^T \mathbf{A}\mathbf{d}^k}{(\mathbf{d}^k)^T \mathbf{A}\mathbf{d}^k} = \frac{(\mathbf{r}^k)^T \mathbf{d}^k}{(\mathbf{d}^k)^T \mathbf{A}\mathbf{d}^k}$$

Věta Algoritmus spočte přesné řešení \mathbf{x} v nejvýše n krocích.

Důkaz Vyjádříme si počáteční chybu \mathbf{e}^0 v bázi tvořené \mathbf{A} -ortogonálními (tedy nezávislými) vektory $\{\mathbf{d}^i\}_{i=0}^{n-1}$:

$$\mathbf{e}^0 = \sum_{j=0}^{n-1} \delta_j \mathbf{d}^j.$$

Souřadnice δ_j najdeme tak, že rovnici přenásobíme zleva výrazem $(\mathbf{d}^k)^T \mathbf{A}$, $k \in \{0, 1, \dots, n-1\}$ a využijeme toho, že většina členů výsledné sumy bude díky \mathbf{A} -ortogonalitě nulová.

$$(\mathbf{d}^k)^T \mathbf{A} \mathbf{e}^1 = \sum_{j=0}^{n-1} \delta_j (\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^j = \delta_k (\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^k.$$

Odtud

$$\delta_k = \frac{(\mathbf{d}^k)^T \mathbf{A} \mathbf{e}^1}{(\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^k} = \frac{(\mathbf{d}^k)^T \mathbf{A} (\mathbf{e}^1 + \sum_{i=0}^{k-1} \alpha_i \mathbf{d}^i)}{(\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^k} = \frac{(\mathbf{d}^k)^T \mathbf{A} \mathbf{e}^k}{(\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^k} = -\frac{(\mathbf{d}^k)^T \mathbf{r}^k}{(\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^k}.$$

Zde si uvědomme, že přičtením sumy $\sum_{i=0}^{k-1} \alpha_i \mathbf{d}^i$ jsme výraz nijak nezměnili (všechna \mathbf{d}^i v sumě jsou \mathbf{A} -ortogonální k \mathbf{d}^k před závorkou, součin je tedy nulový). Dále rovnost $\mathbf{e}^0 + \sum_{i=0}^{k-1} \alpha_i \mathbf{d}^i = \mathbf{e}^k$ snadno vyplývá z již známého vztahu $\mathbf{e}^{k+1} = \mathbf{e}^k + \alpha_k \mathbf{d}^k$.

Všimněme si, že $\delta_k = -\alpha_k$. Vyjádříme si tedy chybu v i -té iteraci jako

$$\mathbf{e}^i = \mathbf{e}^0 + \sum_{j=0}^{i-1} \alpha_j \mathbf{d}^j = \sum_{j=0}^{n-1} \delta_j \mathbf{d}^j + \sum_{j=0}^{i-1} \alpha_j \mathbf{d}^j = \sum_{j=0}^{n-1} \delta_j \mathbf{d}^j - \sum_{j=0}^{i-1} \delta_j \mathbf{d}^j = \sum_{j=i}^{n-1} \delta_j \mathbf{d}^j.$$

V každé iteraci se tedy odstraní jedna složka počáteční chyby \mathbf{e}^0 . Po n iteracích je každá složka chyby v bázi $\{\mathbf{d}_i\}_{i=0}^{n-1}$ vynulovaná $\Rightarrow \mathbf{e}^n = \mathbf{o}$. \square

Zbývá ukázat, jak najít jednotlivé \mathbf{A} -ortogonální směry $\{\mathbf{d}^i\}_{i=0}^{n-1}$. Předpokládejme, že známe $\{\mathbf{d}^0, \mathbf{d}^1, \dots, \mathbf{d}^k\}$ a hledíme \mathbf{d}^{k+1} . Vyjdeme z vektoru rezidua \mathbf{r}^{k+1} a pomocí Gramova-Schmidtova procesu (připomeňte si) jej \mathbf{A} -ortogonalizujeme vůči předchozím směrům. Hledíme nový vektor ve tvaru

$$\mathbf{d}^{k+1} = \mathbf{r}^{k+1} - \sum_{j=0}^k \beta_{k,j} \mathbf{d}^j. \quad (1.14)$$

Určíme koeficienty $\beta_{k,j}$ tak, aby výsledný vektor byl \mathbf{A} -ortogonální k předchozím vektorům $\{\mathbf{d}^0, \mathbf{d}^1, \dots, \mathbf{d}^k\}$. Přenásobíme proto rovnost zleva výrazem $(\mathbf{d}^i)^T \mathbf{A}$, $i \in \{0, 1, \dots, k\}$

$$(\mathbf{d}^i)^T \mathbf{A} \mathbf{d}^{k+1} = (\mathbf{d}^i)^T \mathbf{A} \mathbf{r}^{k+1} - \sum_{j=0}^k \beta_{k,j} (\mathbf{d}^i)^T \mathbf{A} \mathbf{d}^j$$

a využijme toho, že díky \mathbf{A} -ortogonalitě je výraz na levé straně rovnosti stejně jako většina členů sumy napravo roven nule. Tedy

$$0 = (\mathbf{d}^i)^T \mathbf{A} \mathbf{r}^{k+1} - \beta_{k,i} (\mathbf{d}^i)^T \mathbf{A} \mathbf{d}^i \Rightarrow \beta_{k,i} = \frac{(\mathbf{d}^i)^T \mathbf{A} \mathbf{r}^{k+1}}{(\mathbf{d}^i)^T \mathbf{A} \mathbf{d}^i}.$$

Dosazením vypočtených koeficientů do (1.14) dostaneme

$$\mathbf{d}^{k+1} = \mathbf{r}^{k+1} - \sum_{j=0}^k \frac{(\mathbf{d}^j)^T \mathbf{A} \mathbf{r}^{k+1}}{(\mathbf{d}^j)^T \mathbf{A} \mathbf{d}^j} \mathbf{d}^j. \quad (1.15)$$

Všimněme si, že abychom pomocí tohoto vzorce vypočítali nový směr, museli bychom si pamatovat všechny předchozí směry, což by z paměťového hlediska bylo značně nevýhodné. Ukážeme si ale nyní, že všechny koeficienty $\beta_{k,j}$ kromě $\beta_{k,k}$ jsou nulové. Dokažme tedy následující tři lemmata.

Lemma Reziduum \mathbf{r}^{k+1} je ortogonální k $\mathbf{d}^0, \mathbf{d}^1, \dots, \mathbf{d}^k$.

Důkaz Přenásobme \mathbf{r}^{k+1} vektory $\mathbf{d}^i, i \in \{0, 1, \dots, k\}$.

$$(\mathbf{d}^i)^T \underbrace{\mathbf{r}^{k+1}}_{= -\mathbf{A} \mathbf{e}^{k+1}} = -(\mathbf{d}^i)^T \mathbf{A} \mathbf{e}^{k+1} = -(\mathbf{d}^i)^T \mathbf{A} \sum_{j=k+1}^{n-1} \delta_j \mathbf{d}^j = 0,$$

kde poslední rovnost vyplývá z \mathbf{A} -ortogonalitě všech vektorů \mathbf{d}^j v sumě k vektoru \mathbf{d}^i . \square

Lemma Reziduum \mathbf{r}^{k+1} je ortogonální k $\mathbf{r}^0, \mathbf{r}^1, \dots, \mathbf{r}^k$.

Důkaz Využijme předpis (1.14) a vyjádřeme \mathbf{r}^i :

$$\mathbf{r}^i = \mathbf{d}^i + \sum_{j=0}^{i-1} \beta_{i,j} \mathbf{d}^j.$$

Pak díky předchozímu lemmatu dostaneme pro všechna $i \leq k$:

$$(\mathbf{r}^{k+1})^T \mathbf{r}^i = (\mathbf{r}^{k+1})^T \left(\mathbf{d}^i + \sum_{j=0}^{i-1} \beta_{i,j} \mathbf{d}^j \right) = (\mathbf{r}^{k+1})^T \mathbf{d}^i + \sum_{j=0}^{i-1} \beta_{i,j} (\mathbf{r}^{k+1})^T \mathbf{d}^j = 0.$$

\square

Lemma Necht je dána množina reziduí $\{\mathbf{r}^i\}_{i=0}^{n-1}$ a sdružených směrů $\{\mathbf{d}^i\}_{i=0}^{n-1}$. Potom platí:

$$\begin{aligned} j < k : (\mathbf{r}^{k+1})^T \mathbf{A} \mathbf{d}^j &= 0 \\ j = k : (\mathbf{r}^{k+1})^T \mathbf{A} \mathbf{d}^j &\neq 0 \end{aligned}$$

Důkaz Pro reziduum platí

$$\mathbf{r}^{j+1} = \mathbf{b} - \mathbf{A}\mathbf{x}^{j+1} = \mathbf{b} - \mathbf{A}(\mathbf{x}^j + \alpha_j \mathbf{d}^j) = \mathbf{r}^j - \alpha_j \mathbf{A}\mathbf{d}^j$$

Odtud

$$\mathbf{A}\mathbf{d}^j = -\frac{1}{\alpha_j}(\mathbf{r}^{j+1} - \mathbf{r}^j)$$

a

$$(\mathbf{r}^{k+1})^T \mathbf{A}\mathbf{d}^j = (\mathbf{r}^{k+1})^T \left(\frac{1}{\alpha_j}(\mathbf{r}^j - \mathbf{r}^{j+1}) \right) = \frac{1}{\alpha_j}(\mathbf{r}^{k+1})^T \mathbf{r}^j - \frac{1}{\alpha_j}(\mathbf{r}^{k+1})^T \mathbf{r}^{j+1}$$

Pro $j < k$ je tento výraz díky ortogonalitě reziduí nulový. Pro $j = k$ dostaneme

$$(\mathbf{r}^{k+1})^T \mathbf{A}\mathbf{d}^k = \frac{1}{\alpha_k}(\mathbf{r}^{k+1})^T \mathbf{r}^k - \frac{1}{\alpha_k}(\mathbf{r}^{k+1})^T \mathbf{r}^{k+1} = -\frac{1}{\alpha_k}(\mathbf{r}^{k+1})^T \mathbf{r}^{k+1}. \quad (1.16)$$

□

Vratme se nyní zpět k předpisu pro výpočet koeficientů $\beta_{k,i}$ v Gramově-Schmidtově procesu. Díky předchozímu lemmatu jsme zjistili, že pro $j = 0, \dots, k-1$ jsou čitatele v sumě v (1.15) nulové. Předpis se tedy zjednoduší na

$$\mathbf{d}^{k+1} = \mathbf{r}^{k+1} - \frac{(\mathbf{d}^k)^T \mathbf{A}\mathbf{r}^{k+1}}{(\mathbf{d}^k)^T \mathbf{A}\mathbf{d}^k} \mathbf{d}^k.$$

K výpočtu nového směru nám stačí znát nové reziduum a předchozí směr. Protože všechny ostatní koeficienty v dané sumě jsou nulové, budeme pro jednoduchost psát místo $\beta_{k,k}$ jen β_k .

Využijeme-li rovnosti $\mathbf{r}^{k+1} = \mathbf{r}^k - \alpha_k \mathbf{A}\mathbf{d}^k$, tedy $\mathbf{A}\mathbf{d}^k = \frac{1}{\alpha_k}(\mathbf{r}^k - \mathbf{r}^{k+1})$, můžeme přepsat čitatele na

$$(\mathbf{r}^{k+1})^T \mathbf{A}\mathbf{d}^k = \frac{1}{\alpha_k}(\mathbf{r}^{k+1})^T (\mathbf{r}^k - \mathbf{r}^{k+1}) = -\frac{1}{\alpha_k}(\mathbf{r}^{k+1})^T \mathbf{r}^{k+1}.$$

Podobně lze upravit jmenovatele

$$(\mathbf{d}^k)^T \mathbf{A}\mathbf{d}^k = (\mathbf{r}^k - \beta_{k-1} \mathbf{d}^{k-1})^T \mathbf{A}\mathbf{d}^k = \frac{1}{\alpha_k}(\mathbf{r}^k)^T (\mathbf{r}^k - \mathbf{r}^{k+1}) = \frac{1}{\alpha_k}(\mathbf{r}^k)^T \mathbf{r}^k.$$

Po dosazení do předpisu pro β_k dostaneme

$$\beta_k = -\frac{(\mathbf{r}^{k+1})^T \mathbf{r}^{k+1}}{(\mathbf{r}^k)^T \mathbf{r}^k}$$

V podobném duchu můžeme upravit i čitatele v předpisu α_k

$$(\mathbf{r}^k)^T \mathbf{d}^k = (\mathbf{r}^k)^T \mathbf{d}^k = \mathbf{r}^k - \beta_{k-1} \mathbf{d}^{k-1} = (\mathbf{r}^k)^T \mathbf{r}^k - \beta_{k,k-1} \underbrace{(\mathbf{r}^k)^T \mathbf{d}^{k-1}}_{=0},$$

čímž získáme

$$\alpha_k = \frac{(\mathbf{d}^k)^T \mathbf{r}^k}{(\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^k} = \frac{(\mathbf{r}^k)^T \mathbf{r}^k}{(\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^k}.$$

Pokud použijme tyto předpisy pro α_k, β_k , ušetříme násobení matice-vektor a skalární součin.

Zapišme nyní celý algoritmus:

```
function CONJUGATE_GRADIENT( $A, \mathbf{b}, \mathbf{x}^0$ )
   $\mathbf{d}^0 = \mathbf{r}^0 = \mathbf{b} - \mathbf{A}\mathbf{x}^0$ 
   $k = 0$ 
  while  $\|\mathbf{r}^k\|/\|\mathbf{r}^0\| > \varepsilon$  do
     $\alpha_k = \frac{(\mathbf{r}^k)^T \mathbf{r}^k}{(\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^k}$ 
     $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{d}^k$ 
     $\mathbf{r}^{k+1} = \mathbf{r}^k - \alpha_k \mathbf{A} \mathbf{r}^k$ 
     $\beta_k = \frac{(\mathbf{r}^{k+1})^T \mathbf{r}^{k+1}}{(\mathbf{r}^k)^T \mathbf{r}^k}$ 
     $\mathbf{d}^{k+1} = \mathbf{r}^{k+1} + \beta_k \mathbf{d}^k$ 
     $k = k + 1$ 
  end while
end function
```

Doplňme ještě informace o konvergenci. Chybu můžeme odhadnout pomocí vztahů

$$\|\mathbf{e}^{k+1}\|_A \leq \frac{\kappa(A) - 1}{\kappa(A) + 1} \|\mathbf{e}^k\|_A$$

a

$$\|\mathbf{e}^k\|_A \leq 2 \left(\frac{\sqrt{\kappa(A)} - 1}{\sqrt{\kappa(A)} + 1} \right)^k \|\mathbf{e}^0\|_A.$$

Můžeme také odhadnout maximální počet iterací nutných k dosažení relativní přesnosti ε (tedy $\|\mathbf{e}^k\| \leq \varepsilon \|\mathbf{e}^0\|$):

$$k \leq \left\lceil \frac{1}{2} \sqrt{\kappa(A)} \ln \frac{2}{\varepsilon} \right\rceil.$$

Porovnejme to s metodou největšího spádu, pro kterou platí

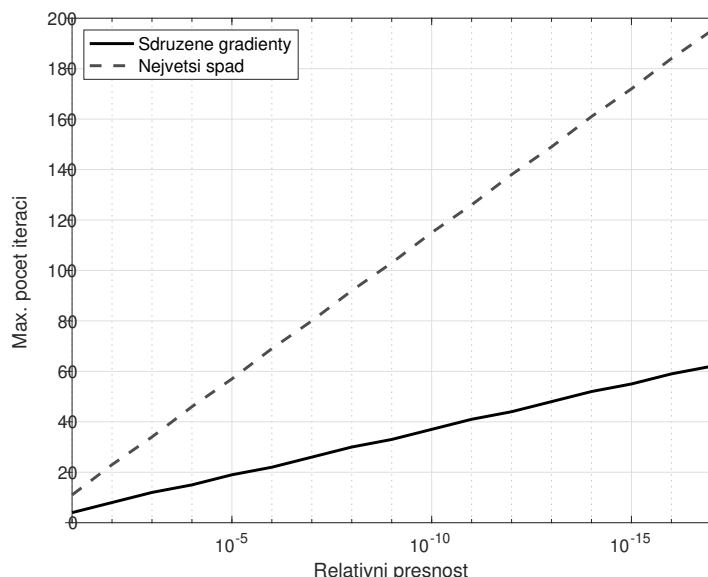
$$k \leq \left\lceil \frac{1}{2} \kappa(A) \ln \frac{1}{\varepsilon} \right\rceil.$$

Rozdíl mezi metodami je dobře patrný na Obrázku 1.5.

Na závěr poznamenejme, že metoda sdružených gradientů patří mezi tzv. Krylovovské metody. Tyto metody generují posloupnost Krylovových podprostorů ve tvaru $\mathcal{K}^k(A, \mathbf{b}) = \text{span}\{\mathbf{b}, \mathbf{A}\mathbf{b}, \mathbf{A}^2\mathbf{b}, \dots, \mathbf{A}^k\mathbf{b}\}$ a v každé iteraci minimalizují chybu v energetické normě na daném podprostoru, tedy

$$\mathbf{x}^{k+1} = \arg \min_{\bar{\mathbf{x}} \in \mathbf{x}^0 + \mathcal{K}^k(A, \mathbf{r}^0)} \|\mathbf{x} - \bar{\mathbf{x}}\|_A,$$

kde \mathbf{x} je přesné řešení.



Obrázek 1.5: Porovnání maximálního počtu iterací nutných k dosažení dané relativní přesnosti ($\kappa(A) = 10$).

1.3 Předpodmínění

V předchozích kapitolách jsme viděli, že rychlost konvergence iteračních metod k řešení závisí na čísle podmíněnosti $\kappa(A)$ matice soustavy. Předpodmínění je nenáročná ekvivalentní úprava soustavy, jejímž cílem je snížení čísla podmíněnosti, tedy snížení počtu potřebných iterací.

Předpokládejme, že známe matici M takovou, že platí $\kappa(M^{-1}A) \ll \kappa(A)$. Pak místo původní soustavy $Ax = b$ můžeme řešit ekvivalentní systém

$$M^{-1}Ax = M^{-1}b. \quad (1.17)$$

Vzhledem k menšímu číslu podmíněnosti matice $M^{-1}A$ by měl oproti původní soustavě klesnout počet iterací nutných k jejímu vyřešení. Matici M nazýváme *levým předpodmiňovačem*.

Je-li matice soustavy A symetrická pozitivně definitní, je výhodné původní soustavu řešit metodou největšího spádu nebo metodou sdružených gradientů. Po předpodmiňovači M pak budeme chtít, aby byl také symetrický pozitivně definitní. Problémem však je, že součin dvou symetrických pozitivně definitních matic obecně nemusí být symetrická pozitivně definitní matice. Systém (1.17) tedy můžeme vyřešit pomocí iteračního řešiče, který nevyžaduje tyto vlastnosti matice, ale na rozdíl od původní soustavy $Ax = b$ nemůžeme použít metodu největšího spádu ani metodu sdružených gradientů.

Pokusme se tedy odvodit takový tvar předpokládání, abychom zachovali symetrii a pozitivní definitnost matice systému. Využijeme toho, že každou symetrickou pozitivně definitní matici můžeme rozložit (např. pomocí Choleského rozkladu) na součin $M = LL^T$, kde L je dolní trojúhelníková matice. Dále využijeme toho, že matice $M^{-1}A$ a $L^{-1}AL^{-T}$ mají stejná vlastní čísla (a tedy stejné číslo podmíněnosti). Je-li totiž \mathbf{v} vlastní vektor matice $M^{-1}A$ s vlastním číslem λ , pak $L^T\mathbf{v}$ je vlastní vektor $L^{-1}AL^{-T}$ se stejným vlastním číslem λ :

$$\begin{aligned} (L^{-1}AL^{-T})(L^T\mathbf{v}) &= L^{-1}A \underbrace{L^{-T}L^T}_{=I} \mathbf{v} = L^{-1}A\mathbf{v} = \underbrace{(L^TL^{-T})}_{=I} L^{-1}A\mathbf{v} \\ &= L^T \underbrace{L^{-T}L^{-1}}_{=M^{-1}} A\mathbf{v} = L^T \underbrace{M^{-1}A\mathbf{v}}_{=\lambda\mathbf{v}} = \lambda L^T\mathbf{v} \end{aligned}$$

Přenásobme nyní původní soustavu $A\mathbf{x} = \mathbf{b}$ zleva maticí L^{-1} :

$$L^{-1}A\mathbf{x} = L^{-1}\mathbf{b}$$

Abychom zachovali symetrii a pozitivní definitnost, musíme matici A přenásobit L^{-T} zprava. “Vložme” tedy mezi A a \mathbf{x} jednotkovou matici ve tvaru $I = L^{-T}L^T$, tím soustavu nijak nezměníme:

$$L^{-1}AL^{-T}L^T\mathbf{x} = L^{-1}\mathbf{b}.$$

Místo původní soustavy $A\mathbf{x} = \mathbf{b}$ a místo soustavy (1.17) řešíme ekvivalentní systém

$$(L^{-1}AL^{-T})(L^T\mathbf{x}) = L^{-1}\mathbf{b}.$$

Zavedme značení $\hat{A} = L^{-1}AL^{-T}$, $\hat{\mathbf{x}} = L^T\mathbf{x}$, $\hat{\mathbf{b}} = L^{-1}\mathbf{b}$. Soustavu pak můžeme zapsat jako

$$\hat{A}\hat{\mathbf{x}} = \hat{\mathbf{b}}. \quad (1.18)$$

V tomto případě je matice soustavy $\hat{A} = L^{-1}AL^{-T}$ symetrická a pozitivně definitní, k řešení (1.18) tedy můžeme použít metodu největšího spádu či metodu sdružených gradientů. Po nalezení $\hat{\mathbf{x}}$ získáme řešení původní soustavy jako $\mathbf{x} = L^{-T}\hat{\mathbf{x}}$.²

Poznamenejme, že nevýhodou tohoto přístupu je nutnost znalosti rozkladu $M = LL^T$. Vhodnou manipulací s předpisy jednotlivých iteračních metod se však této potřeby zbavíme. Demonstrujeme to pro jednoduchost na Richardsonově metodě, která má při aplikaci na (1.18) předpis

$$\hat{\mathbf{x}}^{k+1} = \hat{\mathbf{x}}^k + \omega(\hat{\mathbf{b}} - \hat{A}\hat{\mathbf{x}}^k),$$

²Inverzní matice L^{-1} a L^{-T} samozřejmě nemusíme explicitně sestavovat. Např. násobením $\mathbf{w} = L^{-1}\mathbf{u}$ je ekvivalentní s řešením soustavy $L\mathbf{w} = \mathbf{u}$. Matice L je dolní trojúhelníková, řešení tedy snadno získáme pomocí dopředné substituce. Obdobně postupujeme i v případě násobení maticí L^{-T} . Není také třeba a není výhodné, abychom explicitně sestavovali matici $L^{-1}AL^{-T}$ jakou maticový součin. Součin $L^{-1}AL^{-T}\mathbf{u}$ vypočítáme jako posloupnost tří po sobě jdoucích součinů matice-vektor.

tedy

$$\mathbf{L}^T \mathbf{x}^{k+1} = \mathbf{L}^T \mathbf{x}^k + \omega(\mathbf{L}^{-1} \mathbf{b} - \mathbf{L}^{-1} \mathbf{A} \mathbf{L}^{-T} \mathbf{L}^T \mathbf{x}^k).$$

Přenásobme tento předpis zleva maticí \mathbf{L}^{-T} :

$$\mathbf{L}^{-T} \mathbf{L}^T \mathbf{x}^{k+1} = \mathbf{L}^{-T} \mathbf{L}^T \mathbf{x}^k + \omega(\mathbf{L}^{-T} \mathbf{L}^{-1} \mathbf{b} - \mathbf{L}^{-T} \mathbf{L}^{-1} \mathbf{A} \mathbf{L}^{-T} \mathbf{L}^T \mathbf{x}^k).$$

Odtud snadno získáme předpis předpodmíněné Richardsonovy metody, který již nevyžaduje znalost rozkladu \mathbf{M} :

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \omega \mathbf{M}^{-1}(\mathbf{b} - \mathbf{A} \mathbf{x}^k).$$

Podobným způsobem bychom mohli upravit předpisy metody největšího spádu či metody sdružených gradientů. Všimněme si, že při použití tohoto předpisu potřebujeme aplikovat \mathbf{M}^{-1} (tedy řešit soustavu s maticí \mathbf{M}). Předpodmiňovač tedy musíme volit tak, abychom toto řešení byli schopni najít rychle a efektivně.

Jak tedy zvolit \mathbf{M} ? Požadavky na předpodmiňovač se dají shrnout takto: inverze \mathbf{M} by měla aproximovat inverzi \mathbf{A} ($\mathbf{M}^{-1} \approx \mathbf{A}^{-1}$), řešení soustavy s \mathbf{M} by mělo být snadné a efektivní a pokud je matice \mathbf{A} symetrická a pozitivně definitní, měl by i předpodmiňovač být symetrický pozitivně definitní. Dva extrémní případy jsou:

- $\mathbf{M} = \mathbf{A}$ – v takovém případě inverze \mathbf{M} aproximuje inverzi \mathbf{A} přesně, ale aplikace samotného předpodmínění je stejně náročná jako řešení původní úlohy;
- $\mathbf{M} = \mathbf{I}$ – v tomto případě je sice aplikace předpodmiňovače triviální, nedosáhneme ale žádného zlepšení čísla podmíněnosti (získáme původní nepředpodmíněný algoritmus).

Mezi těmito nepoužitelnými extrémy existuje celé spektrum sofistikovaných metod předpodmínění, pro jednoduchost však zmiňme dva dobře známé přístupy:

- diagonální předpodmiňovač – volíme $\mathbf{M} = \text{diag } \mathbf{A}$. V tomto případě je výpočet inverze \mathbf{M} velmi snadný. Tento předpodmiňovač je efektivní např. pro soustavy s diagonálně dominantní maticí.
- neúplný LU/Choleského rozklad – je vhodný pro soustavy s řídkou maticí. Problémem klasického LU/Choleského rozkladu je, že přestože je matice \mathbf{A} řídká, mohou být její faktory \mathbf{L}, \mathbf{U} plné nebo více zaplněné. Princip neúplného LU/Choleského rozkladu je jednoduchý – použijeme stejný algoritmus jako u klasického LU/Choleského rozkladu, ale nenulový prvek do matice \mathbf{L} nebo \mathbf{U} uložíme pouze tehdy, existuje-li nenulový prvek na stejné pozici v rozkládané matici \mathbf{A} . Dostaneme neúplné faktory $\tilde{\mathbf{L}}, \tilde{\mathbf{U}}$ a příslušné předpodmiňovače definujeme jako $\mathbf{M} = \tilde{\mathbf{L}} \tilde{\mathbf{U}}$ nebo $\mathbf{M} = \tilde{\mathbf{L}} \tilde{\mathbf{L}}^T$. Řešit soustavu s \mathbf{M} je pak snadné, protože známe její rozklad na součin trojúhelníkových matic a můžeme použít dopřednou a zpětnou substituci.

References

- [1] Trefethen, L. N, Bau, D. Numerical Linear Algebra. SIAM. 1997.
- [2] Schewchuk, J. R. An Introduction to the Conjugate Gradient Method Without the Agonizing Pain. 1994. Dostupné z <https://www.cs.cmu.edu/~quake-papers/painless-conjugate-gradient.pdf>
- [3] Lukáš, D. Zápisky z přednášek. Dostupné z <https://home1.vsb.cz/~luk76>