



Hierarchické matice: Moderní přístup k práci s velkými hustými maticemi

Bakalářská práce

Studijní program: B1101 – Matematika
Studijní obory: 7504R015 – Matematika se zaměřením na vzdělávání
7507R036 – Anglický jazyk se zaměřením na vzdělávání

Autor práce: **Barbora Košková**
Vedoucí práce: Martin Plešinger





Hierarchical matrices: A contemporary approach for large-scale dense matrices

Bachelor thesis

Study programme: B1101 – Mathematics
Study branches: 7504R015 – Mathematics for Education
7507R036 – English for Education

Author: **Barbora Košková**
Supervisor: Martin Plešinger



Tento list nahrad'te
originálem zadání.

Prohlášení

Byla jsem seznámena s tím, že na mou bakalářskou práci se plně vztahuje zákon č. 121/2000 Sb., o právu autorském, zejména § 60 – školní dílo.

Beru na vědomí, že Technická univerzita v Liberci (TUL) nezasahuje do mých autorských práv užitím mé bakalářské práce pro vnitřní potřebu TUL.

Užiji-li bakalářskou práci nebo poskytnu-li licenci k jejímu využití, jsem si vědoma povinnosti informovat o této skutečnosti TUL; v tomto případě má TUL právo ode mne požadovat úhradu nákladů, které vynaložila na vytvoření díla, až do jejich skutečné výše.

Bakalářskou práci jsem vypracovala samostatně s použitím uvedené literatury a na základě konzultací s vedoucím mé bakalářské práce a konzultantem.

Současně čestně prohlašuji, že tištěná verze práce se shoduje s elektronickou verzí, vloženou do IS STAG.

Datum:

Podpis:

Anotace

Tato bakalářská práce se zabývá problematikou práce s hustými maticemi pomocí konceptu tzv. hierarchických matic. Takové husté matice často vznikají např. jako inverze matic řídkých. Na příkladu třídiagonálních matic a s využitím jejich spektrálních vlastností ukážeme, že jejich hustá inverze má veškeré mimodiagonální bloky hodnoti nejvýše jedna. Toto pozorování lze v jistém smyslu zobecnit na řadu dalších řídkých matic, které můžeme nalézt v mnoha úlohách z reálného světa, fyzikálních, inženýrských, atd.

V práci zavedeme koncept hierarchických matic, jejichž základním kamenem je stromová struktura (typicky např. binární strom, jak jej známe z teorie grafů), která popisuje rekurzivní členění matice na bloky. Dále se v práci soustředíme na základní operace s takovými maticemi. Ukážeme, v jakém smyslu je lze zejména sčítat a násobit. Hlavní nástroj, který při popisu operací používáme, je tzv. low-rank aritmetika matic. Ta využívá maticových rozkladů (zejména QR rozklad a singulární rozklad (SVD)) k chytré manipulaci s bloky nízké hodnoti a ke kompresi výsledku operací.

Klíčová slova:

třídiagonální matice; vlastní čísla; husté matice; husté inverze řídkých matic; hierarchické matice; low-rank aritmetika matic; stromy (teorie grafů)

Abstract

The bachelor thesis focuses on a work and manipulation with dense matrices using the concept of so-called hierarchical matrices. Such dense matrices often appear, e.g., as inversions of sparse matrices. On the example of tridiagonal matrices and by employing their spectral properties, we demonstrate that their dense inverses have off-diagonal blocks of a low rank (not more than one). This observation can be generalized to a lot of other cases of sparse matrices, which can be found in many real-world problems in physics, engineering, etc.

In the thesis we introduce the concept of hierarchical matrices, where the key idea is the tree structure (e.g., a binary tree, which we know from graph theory) that describes a recursive partitioning of the matrix into blocks. We also focus on basic operations with these matrices. We show in which way it is possible to do the matrix addition and multiplication. The main tool that we use for describing these operations is so-called low-rank arithmetic of matrices. It employs matrix decompositions (especially the QR decomposition and singular value decomposition (SVD)) for smart manipulation with low-rank blocks and for compression of the result of operations.

Key words:

tridiagonal matrices; eigenvalues; dense matrices; dense inverses of sparse matrices; hierarchical matrices; low-rank arithmetic of matrices; trees (graph theory)

Poděkování

Ráda bych poděkovala, svému vedoucímu bakalářské práce, panu Martinu Plešingerovi, za cenné rady, věcné připomínky a vstřícnost při konzultacích a vypracování bakalářské práce.

Obsah

Anotace	5
Abstract	6
Seznam obrázků	10
Seznam tabulek	11
Použité značení a zkratky	12
Úvod	14
1 Vybrané základní pojmy	16
1.1 Spektrum a vlastní čísla matic	16
1.2 (Neorientované) grafy, zejména stromy	16
2 Symetrické třídiagonální matice a jejich spektrální vlastnosti	19
2.1 Třídiagonální matice	19
2.2 Násobnosti vlastních čísel symetrické třídiagonální matice	20
2.3 Vlastní čísla hlavních podmatic	21
2.4 Prokládání vlastních čísel	23
3 Hodnosti mimodiagonálních bloků inverzní matice	24
3.1 Obecná regulární matice a její inverze	24
3.1.1 Vztah hodnotí bloků B a F	24
3.1.2 Regularita bloků A a H	25
3.1.3 Vztah hodnotí bloků C a G	26
3.1.4 Regularita bloků D a E	26
3.2 Symetrická třídiagonální regulární matice	27
3.3 Symetrická třídiagonální pozitivně definitní matice	28
3.4 Upřesnění na závěr	29
4 Hierarchické matice	32
4.1 Hierarchické dělení matice do stromové struktury	32
4.2 Paměťové náklady na uložení hierarchické struktury a její omezení	33
4.2.1 Poznámka k velikosti nejmenších bloků	34
4.3 Hierarchický přístup pro obecnou matici	35

5	Základní operace	38
5.1	Součet dvou hierarchických matic	38
5.1.1	Součet dvou bloků nízké hodnoti	38
5.1.2	Komprese	39
5.2	Součet hierarchické matice, low-rank matice a husté matice	41
5.2.1	Součet hierarchické matice a low-rank matice	41
5.2.2	Součet hierarchické a husté, resp. low-rank a husté matice	41
5.3	Součin hierarchické matice a vektoru	42
5.3.1	Součin hierarchické matice a husté matice	42
5.4	Součin dvou hierarchických matic	43
5.4.1	Součiny různých bloků	43
5.4.2	Rekapitulace součinů bloků a součty součinů bloků	44
5.5	Poznámka k součtu a součinu dvou hierarchických matic	45
	Závěr	46
	Reference	48

Seznam obrázků

1	Příklady řídkých matic z Harwell–Boeing Collection	15
2	Náhodná symetrická třídiagonální matice a její inverze.	15
1.1	Příklady grafů.	18
3.1	Hodnota mimodiagonálního bloku.	28
4.1	Hierarchická reprezentace třídiagonální matice	33
4.2	Hierarchická matice a její inverze	36
4.3	Vliv velikosti největšího hustě ukládaného bloku	37
5.1	Operace na dvou hierarchických maticích s rozdílným stromem	45
5.2	Schéma výpočtu inverze řídké matice v hierarchickém formátu	46

Seznam tabulek

5.1	Výpočetní náklady součtu dvou matic nízké hodnoti	40
5.2	Výpočetní náklady součinu matice s vektorem	42
5.3	Tabulka rekapituluje součiny bloků a jejich výsledky.	45

Použité značení a zkratky

V textu značíme

vektory	pomocí malých písmen $u_1, u_2, u_r, v_1, v_2, v_r, x$, atd.,
matice a jejich bloky	pomocí velkých písmen (latinských i řeckých) $A, B, C, D, E, F, U, V, \Sigma$, atd.,

Pomocí malých písmen (latinských i řeckých) také značíme prvky matic a také skaláry. Speciální význam pak mají písmena i, j, ℓ , jimiž zpravidla indexujeme prvky matic, a k, m, n, r , která používáme k označení dimenze matice, resp. hodnosti (ranku) matice.

Matice a vektory

Značení	Význam
$A \in \mathbb{R}^{n \times m}$	reálná matice s rozměry n krát m , s prvky $a_{i,j}$
A^T	transpozice matice A
$\text{rank}(A)$	hodnost matice definovaná jako počet lineárně nazávislých řádků, resp. sloupců matice A
A^{-1}	inverze čtvercové regulární matice A
$\text{sp}(A)$	spektrum čtvercové matice A
λ	vlastní číslo čtvercové matice A
$\det(A)$	determinant čtvercové matice A
I, I_n	jednotková matice, resp. jednotková matice řádu n
e_j	j -tý sloupec jednotkové matice I vhodného řádu
T_n	symetrická třídiagonální matice
$\alpha_1, \dots, \alpha_n$	prvky na diagonále T_n
β_2, \dots, β_n	prvky na první nad- a poddiagonále T_n

Použité zkratky a akronymy

Zkratka	Význam
QR	QR rozklad matice, $A = QR$
SVD	singulární rozklad matice (singular value decomposition), $W = U_w \Sigma_w V_w^T$
Ds	hustá (dense) matice
Sp	řídká (sparse) matice
LR	low-rank matice
Hi	hierarchická matice

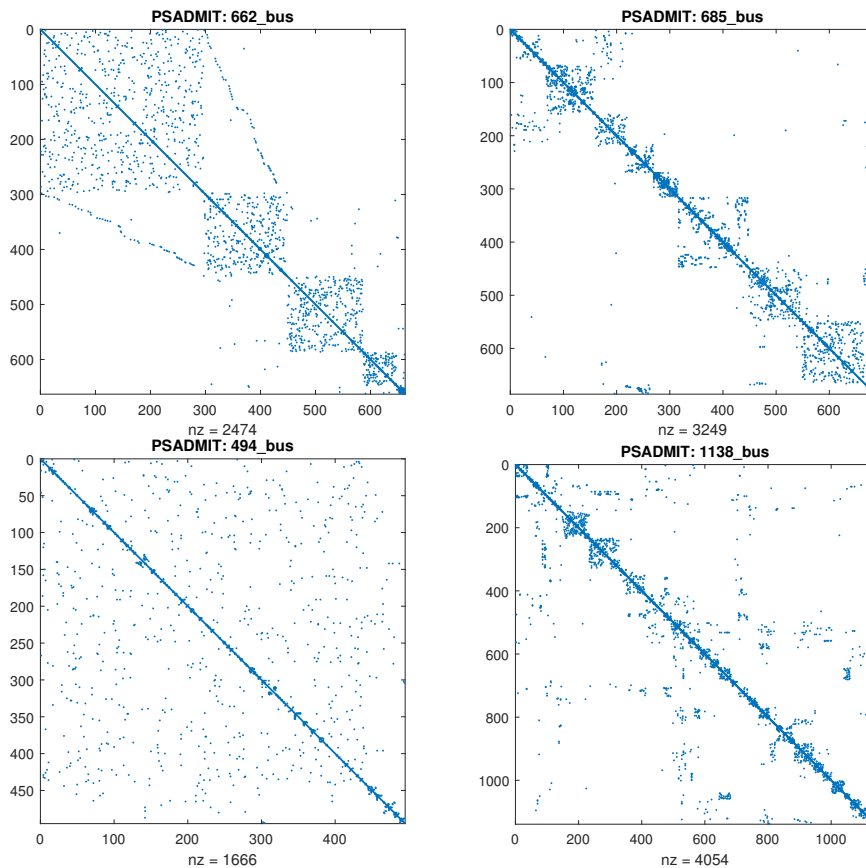
Úvod

Ruku v ruce s rozmachem výpočetní techniky vznikají i možnosti a potřeby řešit stále rozsáhlejší matematické úlohy, které jsou navíc často zformulovány jazykem lineární algebry, tedy pomocí matic; viz např. úvod v [5].

Taková matematická úloha často vzniká diskretizací nějakého reálného fyzikálního problému a typicky zahrnuje práci s rozsáhlými (protože se často snažíme o co nejjemnější diskretizaci, abychom co nejpřesněji podchytili nuance původního problému) a *řídkými* (protože diskretizace typicky pracuje s lokální komunikací) maticemi. Tak nazýváme matice, které mají zpravidla velké množství nulových prvků. S rozsáhlými řídkými maticemi se pracuje (v počítači) snadno, zejména proto, že u takových matic ukládáme do paměti jen nenulové prvky a jejich polohu v matici (tj. index řádku a sloupce), a tedy můžeme uložit i matici poměrně velkých rozměrů. Příklady řídkých matic z Harwell–Boeing Collection (viz [18] nebo [19]) jsou vidět na obrázku 1.

U některých úloh se však můžeme setkat s tím, že takovou matici (předpokládáme, že je regulární) potřebujeme explicitně invertovat; viz např. [10, Sekce 2.3.2, Algoritmus 3]. Není-li však řídká matice (blokově) diagonální, je její inverze obecně hustá, to lze vidět na příkladu náhodné třídiagonální matice; viz obrázek 2. Je tedy zřejmé, že pokud pracujeme s řídkou maticí příliš velkého řádu, kterou navíc potřebujeme explicitně invertovat, snadno se může stát, že při této operaci vyčerpáme zdroje dostupné na daném výpočetním systému. (Stačí si představit hustou matici řádu např. milion. Vzhledem k tomu, že k uložení jednoho reálného čísla v počítači (ve formátu `double`) je zapotřebí 8 bytů, pak by k uložení celé matice bylo potřeba $8 \cdot (10^6)^2$ bytů, což je přibližně 7 451 gigabytů.) Z tohoto důvodu vznikl koncept tzv. *hierarchických*, nebo také hierarchicky uložených matic, viz např. monografie Wolfganga Hackbushe [7], Maria Bebendorfa [1] a Steffena Börma [2], případně práce Stefana Pauliho [13]; k dispozici je též software `HLib` napsaný v jazyce `C` kolektivem kolem prof. Hackbushe, který je volně (resp. na požádání) dostupný, viz [21].

V kapitole 1 shrneme vybrané základní pojmy, jako je spektrum, vlastní čísla matic a grafy, zejména stromy. V kapitole 2 vyložíme vlastnosti symetrických třídiagonálních matic. Primárně zaměřené na vlastní čísla těchto matic. V kapitole 3 rozebereme podrobněji mimodiagonální bloky a jejich hodnoty. Kapitola 4 je věnována konceptu hierarchických matic. Zejména pak vhodnému dělení a ukládání těchto matic do paměti počítače. V poslední kapitole 5 se zaměříme na základní operace, zejména součet a součin dvou hierarchických matic.



Obrázek 1: Příklady řídkých matic mající původ ve výpočtech rozvodů silového elektrického vedení z Harwell-Boeing Collection.

```

MATLAB Command Window
>> N = 10; D0 = rand(N,1); D1 = rand(N-1,1); T = diag(D0) + diag(D1,1) + diag(D1,-1)
T =
    0.1966    0.7572         0         0         0         0         0         0         0         0
    0.7572    0.2511    0.7537         0         0         0         0         0         0         0
         0    0.7537    0.6160    0.3804         0         0         0         0         0         0
         0         0    0.3804    0.4733    0.5678         0         0         0         0         0
         0         0         0    0.5678    0.3517         0.0759         0         0         0         0
         0         0         0         0    0.0759    0.8308    0.0540         0         0         0
         0         0         0         0         0    0.0540    0.5853    0.5308         0         0
         0         0         0         0         0         0    0.5308    0.5497    0.7792         0
         0         0         0         0         0         0         0    0.7792    0.9172    0.9340
         0         0         0         0         0         0         0         0    0.9340    0.2858

>> inv(T)
ans =
    0.5594    1.1754   -0.9535   -0.7847    1.2930   -0.1197    0.0261   -0.0166   -0.0061    0.0198
    1.1754   -0.3052    0.2476    0.2037   -0.3357    0.0311   -0.0068    0.0043    0.0016   -0.0051
   -0.9535    0.2476    0.8754    0.7205   -1.1871    0.1099   -0.0240    0.0153    0.0056   -0.0182
   -0.7847    0.2037    0.7205   -1.5703   -2.5873   -0.2396    0.0522   -0.0332   -0.0121    0.0396
    1.2930   -0.3357   -1.1871    2.5873   -1.3612    0.1261   -0.0275    0.0175    0.0064   -0.0209
   -0.1197    0.0311    0.1099   -0.2396    0.1261    1.2092   -0.2636    0.1677    0.0612   -0.2001
    0.0261   -0.0068   -0.0240    0.0522   -0.0275   -0.2636    4.0981   -2.6079   -0.9519    3.1103
   -0.0166    0.0043    0.0153   -0.0332    0.0175    0.1677   -2.6079    2.8585    1.0433   -3.4091
   -0.0061    0.0016    0.0056   -0.0121    0.0064    0.0612   -0.9519    1.0433   -0.0876    0.2864
    0.0198   -0.0051   -0.0182    0.0396   -0.0209   -0.2001    3.1103   -3.4091    0.2864    2.5627

```

Obrázek 2: Náhodná symetrická třídiagonální matice a její inverze.

1 Vybrané základní pojmy

V této úvodní kapitole shrneme a zopakujeme několik vybraných základních pojmů a konceptů z lineární algebry a teorie grafů, které budou užitečné v následujícím textu.

1.1 Spektrum a vlastní čísla matic

Nejdůležitější pojem známý již z úvodních kurzů lineární algebry, se kterým budeme pracovat, bude pojem vlastního čísla a pojmy s ním související. Proto začneme následující definicí.

Definice 1. *Nechť $A \in \mathbb{R}^{n \times n}$ je reálná čtvercová matice, poté se rovnice $\det(A - \lambda I) = 0$ nazývá charakteristickou rovnicí matice A a $\lambda \in \mathbb{C}$ splňující tuto rovnici jejím vlastním číslem.*

Poznamenejme, že $\det(A - \lambda I)$ je polynom v proměnné λ stupně n . Protože u tohoto polynomu se mění znaménko u nejvyšší mocniny λ v závislosti na n , bývá obvyklejší pracovat s polynomem $(-1)^n \det(A - \lambda I) = \det(\lambda I - A)$, který se nazývá *charakteristický polynom*, viz [5, str. 8]. My budeme pro jednoduchost pracovat s polynomem $\det(A - \lambda I)$.

Množinu všech vlastních čísel matice A budeme značit $\text{sp}(A)$, nazývá se *spektrém* matice.

Alternativně můžeme spektrum matice $A \in \mathbb{R}^{n \times n}$ popsat například pomocí hodnoty. Konkrétně $\text{rank}(A - \lambda I)$ může být buď roven n (tedy $(A - \lambda I)$ je regulární matice), nebo ostře menší než n (tedy $(A - \lambda I)$ je singulární matice). Druhý případ nastane tehdy a jen tehdy, je-li $\lambda \in \text{sp}(A)$.

1.2 (Neorientované) grafy, zejména stromy

Zejména v druhé části této práce bude potřeba pracovat s (neorientovanými) grafy, konkrétně se stromy. Proto bude vhodné tyto pojmy zavést. Naší ambicí není hovořit o teorii grafů jako takové, budeme tedy pojmy zavádět nepatrně zjednodušeně.

Definice 2. *Graf je uspořádaná dvojice (V, H) , kde V je množina vrcholů a H množina hran. Přičemž hrana je neuspořádaná dvojice vrcholů, tedy platí $H \subseteq \binom{V}{2}$, kde pomocí symbolu $\binom{V}{2}$ značíme množinu všech dvouprvkových podmnožin V .*

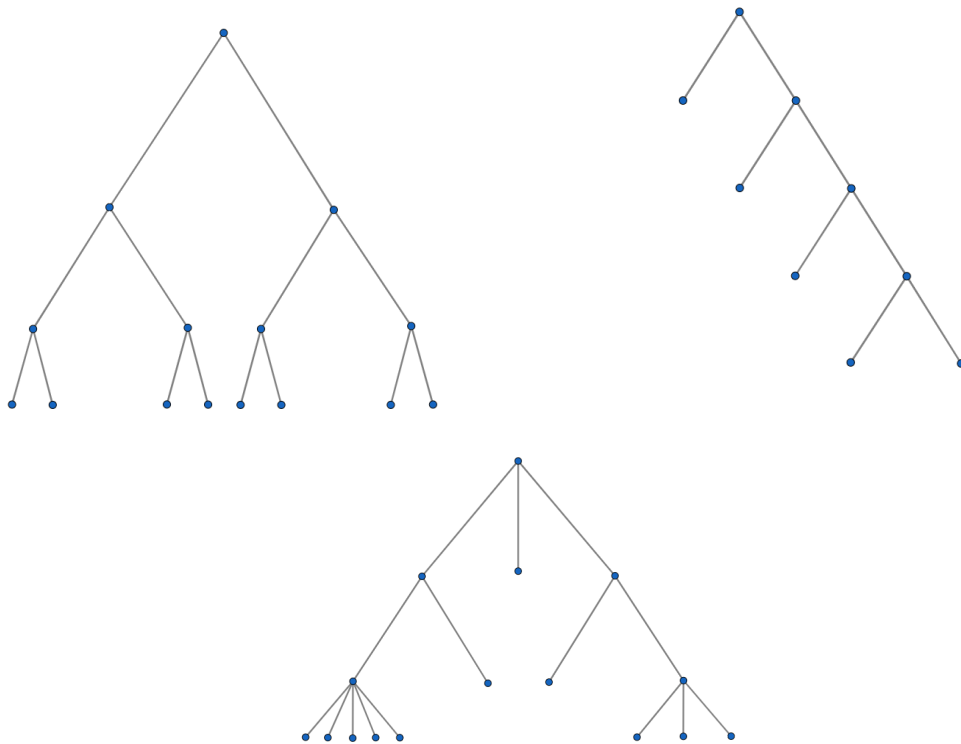
My speciálně budeme potřebovat graf, který se nazývá *strom*. Nejprve se však podíváme na pojem *cesta*. Konkrétně cesta délky l z vrcholu v_1 do v_2 je uspořádaná l -tice hran tak, že dvě sousední hrany mají společný vrchol. Dále se hodí pojem *kružnice*. Kružnicí budeme nazývat každou cestu z vrcholu v zpět do vrcholu v délky alespoň 3 obsahující různé hrany.

Definice 3. *Stromem rozumíme graf, u kterého vede cesta mezi libovolnými dvěma vrcholy a neobsahuje kružnici. Navíc obsahuje jeden významný vrchol, který se nazývá kořen.*

V dalším textu budeme pracovat pouze se stromy a z definice stromu vyplývá, že mezi libovolnými dvěma vrcholy vede cesta (takový graf je tzv. *souvislý*). Dále budeme pod pojmem *cesta* rozumět nejkratší cestu mezi danými vrcholy. Protože stromy navíc neobsahují kružnice, bude tato nejkratší cesta vždy dána jednoznačně.

U stromů kromě *kořene* zavádíme ještě pojem *list*, který má do jisté míry opačnou funkci. Listem je takový vrchol v , že cestu vedoucí mezi kořenem a v nelze prodloužit na straně vrcholu v , aniž by došlo k opakovanému užití hrany. Tedy, list je vrchol, do kterého vede jen jedna hrana a zároveň není kořenem.

Vrcholy na cestách mezi kořenem a listy pojmenováváme následujícím způsobem. O libovolných dvou vrcholech spojených hranou říkáme, že jsou *předek* (*rodič*) a *potomek*, přičemž předek je ten, který je blíž kořenu, a potomek ten, který je blíž listu. Alternativně bychom kořen a listy mohli definovat tak, že kořenem je vrchol, který nemá předky, a listem vrchol, který nemá potomky, viz obrázek 1.1. Pro bližší seznámení se s pojmy teorie grafů doporučujeme např. [9] nebo české učebnice [3], [4], [12], [14], [15], případně rozsáhlejší text [11].



Obrázek 1.1: Příklady stromů. První dva stromy nazýváme *binární*, neboť každý vrchol, který není listem má právě dva potomky. První z nich je navíc tzv. *vyvážený binární strom*. Třetí strom je obecný.

2 Symetrické třídiagonální matice a jejich spektrální vlastnosti

Třídiagonální matice vznikají přirozeně v řadě úloh spojených s obory, jako je například aplikovaná matematika, moderní fyzika a inženýrství. Nejprve zavedeme pojmy rozlišující jednotlivé diagonály matice. Nechť $A \in \mathbb{R}^{n \times n}$ je čtvercová matice s prvky $a_{i,j}$. Prvky jejichž indexy splňují $i - j = 0$, nazýváme prvky (hlavní) *diagonály*. Pokud indexy prvků splňují $i - j = -1$, hovoříme o prvcích (první) *naddiagonály* a pokud indexy splňují $i - j = 1$, hovoříme o prvcích (první) *poddiagonály*. (Obecně pokud indexy prvků $a_{i,j}$ splňují $i - j = \pm k, k > 0$, pak hovoříme o prvcích k -té nad- resp. poddiagonály).

2.1 Třídiagonální matice

Důležitým pojmem pro další výklad je tzv. třídiagonální matice, kterou definujeme nyní.

Definice 4. Nechť $A \in \mathbb{R}^{n \times n}$ je čtvercová matice s prvky $a_{i,j}$. Když platí $a_{i,j} = 0$ pro všechna $|i - j| \geq 2$, pak matici A nazýváme *třídiagonální*.

Definice 4 tedy říká, že třídiagonální matice může obsahovat nenulové prvky pouze na hlavní diagonále, na první naddiagonále a první poddiagonále

$$A = \begin{bmatrix} a_{1,1} & a_{1,2} & 0 & 0 & \cdots & 0 \\ a_{2,1} & a_{2,2} & a_{2,3} & 0 & \cdots & 0 \\ 0 & a_{3,2} & a_{3,3} & a_{3,4} & & \vdots \\ 0 & 0 & a_{4,3} & a_{4,4} & \ddots & 0 \\ \vdots & \vdots & & \ddots & \ddots & a_{n-1,n} \\ 0 & 0 & \cdots & 0 & a_{n,n-1} & a_{n,n} \end{bmatrix}. \quad (2.1)$$

Poznamenejme, že prvky na těchto diagonálách nenulové být nemusí, tedy i nulová matice je formálně třídiagonální.

Pokud $A \in \mathbb{R}^{n \times n}$ je třídiagonální matice a platí $A = A^T$, poté hovoříme o symetrické třídiagonální matici. Alternativně, je to třídiagonální matice, kde navíc platí $a_{i,j} = a_{j,i}$ pro libovolné $i, j \in \{1, 2, \dots, n\}$. Pro jednoduchost budeme její prvky

značit takto

$$T_n = \begin{bmatrix} \alpha_1 & \beta_2 & & & \\ \beta_2 & \ddots & \ddots & & \\ & \ddots & \ddots & \beta_n & \\ & & & \beta_n & \alpha_n \end{bmatrix}. \quad (2.2)$$

2.2 Násobnosti vlastních čísel symetrické třídiagonální matice

Nejprve se podíváme na několik zajímavých vlastností matice T_n souvisejících s vlastními čísly. Začneme násobností vlastních čísel. Budeme tedy studovat matici

$$(T_n - \lambda I) = \begin{bmatrix} \alpha_1 - \lambda & \beta_2 & & & \\ \beta_2 & \alpha_2 - \lambda & \beta_3 & & \\ & \beta_3 & \alpha_3 - \lambda & \ddots & \\ & & \ddots & \ddots & \beta_{n-1} \\ & & & \beta_{n-1} & \alpha_{n-1} - \lambda & \beta_n \\ & & & & \beta_n & \alpha_n - \lambda \end{bmatrix}. \quad (2.3)$$

Vybereme-li z ní šikvou podmatici, vynecháním např. prvního řádku a posledního sloupce, dostaneme

$$\begin{bmatrix} \beta_2 & \alpha_2 - \lambda & \beta_3 & & \\ & \beta_3 & \alpha_3 - \lambda & \ddots & \\ & & \ddots & \ddots & \beta_{n-1} \\ & & & \beta_{n-1} & \alpha_{n-1} - \lambda \\ & & & & \beta_n \end{bmatrix}. \quad (2.4)$$

Tato podmatice je zřejmě horní trojúhelníková, a pokud $\beta_i \neq 0, \forall i = 2, \dots, n$, pak je také regulární. Tedy má právě $n - 1$ lineárně nezávislých řádků (respektive sloupců). Proto také původní matice $T_n - \lambda I$ má alespoň $n - 1$ lineárně nezávislých řádků (respektive sloupců). Jinými slovy $\text{rank}(T_n - \lambda I)$ je roven buď n , nebo $n - 1$ v závislosti na čísle λ .

Protože je matice T_n symetrická, je i normální (viz [5, kapitola 2, obrázek 2.1]), a proto číslo

$$n - \text{rank}(T_n - \lambda I)$$

je přímo rovno násobnosti vlastního čísla λ . Vidíme tedy, že číslo λ z (2.3) je vlastním číslem matice T_n s násobností 1, nebo není vůbec kořenem jejího charakteristického polynomu. Shrňme tyto úvahy v následujícím závěru.

Důsledek 1. *Nechť $T_n \in \mathbb{R}^{n \times n}$ je reálná symetrická třídiagonální matice (2.2) splňující navíc $\beta_2 \beta_3 \cdots \beta_n \neq 0$. Pak je násobnost vlastního čísla λ rovna jedné.*

neboli

$$\det(T_n - \lambda I) = (\alpha_n - \lambda) \det(T_{n-1} - \lambda I) - \beta_n^2 \det(T_{n-2} - \lambda I).$$

Předpokládejme nyní, že dva po sobě jdoucí charakteristické polynomy $P_n(\lambda)$ matice T_n a $P_{n-1}(\lambda)$ matice T_{n-1} mají stejný kořen. Tedy, že

$$\exists \varphi \in \mathbb{C}, \quad P_n(\varphi) = 0 \quad \wedge \quad P_{n-1}(\varphi) = 0.$$

Poté z rovnice

$$P_n(\varphi) = (\alpha_n - \varphi)P_{n-1}(\varphi) - \beta_n^2 P_{n-2}(\varphi)$$

vyplývá, že

$$\beta_n^2 P_{n-2}(\varphi) = 0.$$

Za předpokladu, že $\beta_n \neq 0$, nutně platí $P_{n-2}(\varphi) = 0$, tedy φ je kořenem také charakteristického polynomu $P_{n-2}(\lambda)$ matice T_{n-2} . Tedy $\varphi \in \text{sp}(T_{n-2})$ a stejným postupem (tj. za předpokladu $\beta_2 \beta_3 \cdots \beta_n \neq 0$) zřejmě dostaneme

$$P_n(\varphi) = P_{n-1}(\varphi) = P_{n-2}(\varphi) = P_{n-3}(\varphi) = \dots = P_3(\varphi) = P_2(\varphi) = P_1(\varphi) = 0.$$

Podívejme se podrobně na poslední dva polynomy. Zřejmě

$$P_2(\lambda) = \det \left(\begin{bmatrix} \alpha_1 - \lambda & \beta_2 \\ \beta_2 & \alpha_2 - \lambda \end{bmatrix} \right) = (\alpha_1 - \lambda)(\alpha_2 - \lambda) - \beta_2^2 \quad (2.6)$$

a zároveň

$$P_2(\lambda) = (\alpha_2 - \lambda)P_1(\lambda) - \beta_2^2 P_0(\lambda),$$

kde $P_1(\lambda) = (\alpha_1 - \lambda)$. Porovnáním předchozích rovnic vidíme, že a $P_0(\lambda)$ je polynom stupně nula splňující $P_0(\lambda)\beta_2^2 = \beta_2^2$, tedy $P_0(\lambda) = 1$.

Pro $\lambda = \varphi$ by ale mělo platit

$$P_2(\varphi) = (\alpha_2 - \varphi)(\alpha_1 - \varphi) - \beta_2^2 = 0$$

a zároveň $P_1(\varphi) = (\alpha_1 - \varphi) = 0$. Z toho vyplývá, že $\beta_2^2 = 0$, což ale odporuje předpokladu, že $\beta_2 \beta_3 \cdots \beta_n \neq 0$. Zformulujme toto pozorování jako důsledek.

Důsledek 2. *Nechť*

$$T_n = \begin{bmatrix} T_{n-1} & \beta_n e_{n-1} \\ e_{n-1}^T \beta_n & \alpha_n \end{bmatrix} \in \mathbb{R}^{n \times n} \quad (2.7)$$

je reálná symetrická třídiagonální matice splňující navíc $\beta_2 \beta_3 \cdots \beta_n \neq 0$. Pak matice T_n a T_{n-1} mají různá vlastní čísla.

Analogický výsledek dostaneme také odstraněním prvního řádku a prvního sloupce z matice

$$T_n = \begin{bmatrix} \alpha_1 & e_1^T \beta_2 \\ \beta_2 e_1 & \tilde{T}_{n-1} \end{bmatrix}. \quad (2.8)$$

2.4 Prokládání vlastních čísel

Nechť $M \in \mathbb{R}^{n \times n}$ je reálná symetrická matice

$$M = \begin{bmatrix} A & b \\ b^T & \delta \end{bmatrix}. \quad (2.9)$$

Vlastní čísla takovéto symetrické matice jsou reálná (viz např. [5, věta 2.8]) a můžeme je tedy seřadit

$$\lambda_1(M) \geq \lambda_2(M) \geq \dots \geq \lambda_n(M).$$

Pak také $A \in \mathbb{R}^{(n-1) \times (n-1)}$ je reálná symetrická matice a její vlastní čísla lze seřadit obdobným způsobem. Z věty o prokládání (eigenvalue interlacing; viz např. [5], nebo [16]) vyplývá, že jejich vlastní čísla se prokládají následujícím způsobem

$$\lambda_1(M) \geq \lambda_1(A) \geq \lambda_2(M) \geq \lambda_2(A) \geq \dots \geq \lambda_{n-1}(A) \geq \lambda_n(M).$$

Poznamenejme, že v článku [16] je tato nerovnost ve skutečnosti odvozena pro singulární čísla. Tato věta nás v kombinaci z předchozími dvěma tvrzeními vede k následujícímu pozorování.

Důsledek 3. *Z důsledků 1, 2 a z věty o prokládání plyne, že vlastní čísla třídiagonálních matic*

$$T_n = \begin{bmatrix} T_{n-1} & \beta_n e_{n-1} \\ e_{n-1}^T \beta_n & \alpha_n \end{bmatrix} \in \mathbb{R}^{n \times n}$$

a T_{n-1} (viz (2.7)) se prokládají následujícím způsobem

$$\lambda_1(T_n) > \lambda_1(T_{n-1}) > \lambda_2(T_n) > \lambda_2(T_{n-1}) > \dots > \lambda_{n-1}(T_{n-1}) > \lambda_n(T_n).$$

Analogický výsledek opět dostaneme také pro dvojici matic

$$T_n = \begin{bmatrix} \alpha_1 & e_1^T \beta_2 \\ \beta_2 e_1 & \tilde{T}_{n-1} \end{bmatrix}$$

a \tilde{T}_{n-1} , viz (2.8). Tento důsledek je nejdůležitějším závěrem této kapitoly.

Na závěr poznamenejme, že determinant matice je součin jejích vlastních čísel $\det(M) = \prod_{j=1}^n \lambda_j(M)$, a zároveň je nulový právě tehdy, když je matice M singulární.

Tedy, protože dvě po sobě jdoucí matice nemohou mít stejné vlastní číslo, nejvýše jedna z matic T_{n-1} , T_n může být singulární. Jinými slovy, v posloupnosti symetrických třídiagonálních matic $T_1, T_2, T_3, \dots, T_{n-1}, T_n, \dots$ s nenulovými mimodiagonálními prvky β_j a zkonstruovanými jako ve (2.7) nemohou být dvě po sobě singulární.

3 Hodnosti mimodiagonálních bloků inverzní matice

Nyní budeme studovat regulární matice a jejich inverze zapsané pomocí bloků. Speciálně nás budou zajímat hodnosti mimodiagonálních bloků. Nejprve se na celou věc podíváme obecně.

3.1 Obecná regulární matice a její inverze

Nechť $M \in \mathbb{R}^{n \times n}$ je reálná čtvercová *regulární* matice

$$M = \begin{bmatrix} A & B \\ C & D \end{bmatrix}, \quad (3.1)$$

rozdělená na čtyři bloky a nechť je její inverze

$$M^{-1} = \begin{bmatrix} E & F \\ G & H \end{bmatrix} \quad (3.2)$$

rozdělena na bloky tak, aby obě matice šly násobit po blocích. Pokud tyto matice vynásobíme, dostaneme zřejmě

$$MM^{-1} = \begin{bmatrix} AE + BG & AF + BH \\ CE + DG & CF + DH \end{bmatrix} = \begin{bmatrix} I & 0 \\ 0 & I \end{bmatrix} = I, \quad (3.3)$$

kde jednotkové matice I a nulové matice 0 jsou vhodných rozměrů.

Nás bude zajímat v jakém vztahu jsou hodnosti dvojice matic B a F a také dvojice C a G . Z porovnání mimodiagonálních bloků součinu MM^{-1} plyne, že

$$AF = -BH \quad \text{a} \quad CE = -DG. \quad (3.4)$$

3.1.1 Vztah hodností bloků B a F

Pokud bude matice A čtvercová a regulární, můžeme první rovnici vynásobit zleva A^{-1} . Dostaneme tak rovnost

$$F = -A^{-1}BH.$$

Z tohoto vztahu triviálně plyne, že hodnost matice F se rovná hodnosti matice $A^{-1}BH$. Matice A^{-1} je čtvercová a regulární, proto hodnost součinu $A^{-1}BH$ se

od hodnoty BH neliší. Hodnost matice B se však po vynásobení maticí H může zmenšit, nebo být stejná; pokud totiž násobíme dvě matice, počet lineárně nezávislých řádků se nikdy nezvýší, pouze se může snížit či zůstat stejný. Obecně platí

$$\text{rank}(BH) \leq \min\{\text{rank}(B), \text{rank}(H)\}.$$

Celkově tak dostáváme nerovnost

$$\text{rank}(F) = \text{rank}(A^{-1}BH) = \text{rank}(BH) \leq \text{rank}(B). \quad (3.5)$$

Dále je-li matice H čtvercová a regulární, dostaneme z (3.4) vztah

$$AFH^{-1} = -B.$$

Pokud se podíváme na hodnoty těchto matic, analogicky dostaneme

$$\text{rank}(B) = \text{rank}(AFH^{-1}) = \text{rank}(AF) \leq \text{rank}(F). \quad (3.6)$$

Pokud jsou tedy regulární obě matice A i H , pak z nerovností (3.5) a (3.6) plyne

$$\text{rank}(B) = \text{rank}(F). \quad (3.7)$$

3.1.2 Regularita bloků A a H

Nakonec zbývá ukázat, že tato situace (3.7) nastane vždy, když je A nebo H regulární. Následující lemma nám totiž říká, že regularita těchto dvou bloků musí vždy nastat současně.

Lemma 1. *Nechť M je regulární matice a nechť M a M^{-1} jsou rozdělené na bloky tak, aby šly násobit po blocích, následujícím způsobem*

$$M = \begin{bmatrix} A & B \\ C & D \end{bmatrix}, \quad M^{-1} = \begin{bmatrix} E & F \\ G & H \end{bmatrix}.$$

Pak matice A je čtvercová regulární tehdy a jen tehdy, když matice H je čtvercová regulární.

Důkaz. Aby násobení bylo proveditelné po blocích z obou stran, tj. ve tvaru MM^{-1} i $M^{-1}M$, musí být velikosti bloků M a $(M^{-1})^T$ stejné. Tedy speciálně, A má stejnou velikost jako E^T a D má stejnou velikost jako H^T . Matice A je čtvercová právě tehdy, když E^T a D jsou čtvercové, a tedy právě tehdy, když H^T (a tedy i H) je čtvercová.

Regularitu dokážeme pomocí výpočtu inverze pomocí Gaußovy eliminace. Předpokládejme nejprve, že A je regulární. Zřejmě platí

$$\begin{aligned} [M|I] &= \left[\begin{array}{cc|cc} A & B & I & 0 \\ C & D & 0 & I \end{array} \right] \sim \left[\begin{array}{cc|cc} I & A^{-1}B & A^{-1} & 0 \\ C & D & 0 & I \end{array} \right] \\ &\sim \left[\begin{array}{cc|cc} I & A^{-1}B & A^{-1} & 0 \\ 0 & D - CA^{-1}B & -CA^{-1} & I \end{array} \right]. \end{aligned}$$

Protože původní matice M byla regulární musí být i (blokově) trojúhelníková matice

$$\begin{bmatrix} I & A^{-1}B \\ 0 & S \end{bmatrix}, \quad \text{kde} \quad S = D - CA^{-1}B,$$

regulární. Z regularity M a A tedy ihned plyne regularita matice S (která se nazývá Schurovým doplňkem bloku A v matici M , viz např. [5, Definice 1.4, str. 6]).

Nyní dokončíme vyjádření inverzní matice M^{-1} eliminací naddiagonálního bloku $A^{-1}B$. Zřejmě platí

$$\begin{aligned} & \left[\begin{array}{cc|cc} I & A^{-1}B & A^{-1} & 0 \\ 0 & S & -CA^{-1} & I \end{array} \right] \sim \left[\begin{array}{cc|cc} I & A^{-1}B & A^{-1} & 0 \\ 0 & I & -S^{-1}CA^{-1} & S^{-1} \end{array} \right] \\ & \sim \left[\begin{array}{cc|cc} I & 0 & A^{-1} + A^{-1}BS^{-1}CA^{-1} & -A^{-1}BS^{-1} \\ 0 & I & -S^{-1}CA^{-1} & S^{-1} \end{array} \right] = [I|M^{-1}]. \end{aligned}$$

Dostáváme tedy speciálně $H = S^{-1}$. Z regularity celé matice M a jejího bloku A tedy plyne regularita bloku $H = S^{-1} = (D - CA^{-1}B)^{-1}$ inverze M^{-1} .

Opačnou implikaci, tedy že z regularity M a H plyne regularita A bychom dokázali obdobně. Zřejmě $A = (E - FH^{-1}G)^{-1}$ je inverzí Schurova doplňku bloku H v matici M^{-1} . \square

3.1.3 Vztah hodnotí bloků C a G

Z druhé rovnosti ve vztahu (3.4), $CE = -DG$, odvodíme vztahy mezi hodnotmi druhých mimodiagonálních bloků C a G . Pokud bude matice E čtvercová a regulární můžeme rovnici vynásobit zleva E^{-1} . Dostaneme tak rovnost

$$C = -E^{-1}GD,$$

ze které vyplývá

$$\text{rank}(C) = \text{rank}(E^{-1}GD) = \text{rank}(GD) \leq \text{rank}(G).$$

Dále je-li i matice D čtvercová a regulární, dostaneme rovnost

$$CED^{-1} = -G.$$

Pokud se podíváme na hodnoti těchto matic, analogicky dostaneme

$$\text{rank}(G) = \text{rank}(CED^{-1}) = \text{rank}(CE) \leq \text{rank}(C).$$

Jsou-li tedy regulární obě matice D i E , pak platí

$$\text{rank}(C) = \text{rank}(G).$$

3.1.4 Regularita bloků D a E

Lemma 1 tvrdí, že matice A je regulární tehdy a jen tehdy, když matice H je regulární. Zcela analogicky lze odvodit, že E je regulární, tehdy a jen tehdy, když D je regulární. Tedy obě dvě matice D a E jsou regulární zároveň.

3.2 Symetrická třídiagonální regulární matice

Uvažujme nyní naši symetrickou třídiagonální matici T_n , kterou rozdělíme tak, jako v (3.1), tedy následujícím způsobem

$$T_n = \begin{bmatrix} T_k & e_k \beta_{k+1} e_1^T \\ e_1 \beta_{k+1} e_k^T & \tilde{T}_{n-k} \end{bmatrix} \in \mathbb{R}^{n \times n}, \quad (3.8)$$

kde $T_k \in \mathbb{R}^{k \times k}$, $\tilde{T}_{n-k} \in \mathbb{R}^{(n-k) \times (n-k)}$ a

$$e_1 \beta_{k+1} e_k^T = \begin{bmatrix} 0 & \cdots & 0 & \beta_{k+1} \\ 0 & \cdots & 0 & 0 \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \cdots & 0 & 0 \end{bmatrix} \in \mathbb{R}^{(n-k) \times k}. \quad (3.9)$$

Ze struktury nulových a nenulových prvků matice (3.9) je zřejmé, že hodnost této matice je rovna nejvýše jedné. Přesněji je rovna právě jedné, právě tehdy, když $\beta_{k+1} \neq 0$. Necht' dále je T_n regulární, a necht' je inverze T_n zapsána ve tvaru

$$T_n^{-1} = \begin{bmatrix} E_k & F_k \\ F_k^T & H_k \end{bmatrix}. \quad (3.10)$$

Z předchozí sekce 3.1 je zřejmé, že z regularity T_k vyplývá, že $F_k \in \mathbb{R}^{k \times (n-k)}$ má (za předpokladu $\beta_{k+1} \neq 0$) také hodnost rovnou právě jedné. Poznamenejme, že pro $\beta_{k+1} = 0$ platí

$$T_n = \begin{bmatrix} T_k & 0 \\ 0 & \tilde{T}_{n-k} \end{bmatrix} \quad \text{a} \quad T_n^{-1} = \begin{bmatrix} T_k^{-1} & 0 \\ 0 & \tilde{T}_{n-k}^{-1} \end{bmatrix} \quad (3.11)$$

a tedy $F_k = 0$ a $\text{rank}(F_k) = 0$.

Pokud však T_k regulární není, pak můžeme (3.8) nahradit jinými děleními tak, že řád první matice z prvního blokového řádku bude o jedna menší či větší. Tedy

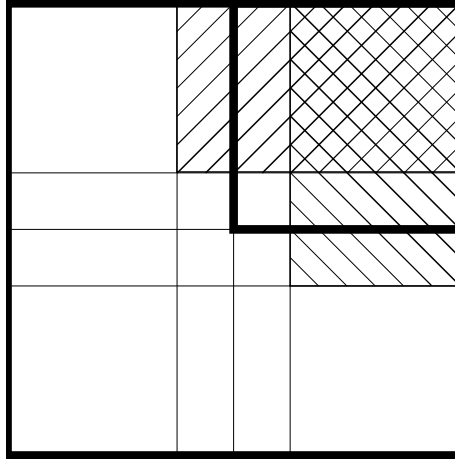
$$T_n = \begin{bmatrix} T_{k-1} & e_k \beta_k e_1^T \\ e_1 \beta_k e_k^T & \tilde{T}_{n-k-1} \end{bmatrix} = \begin{bmatrix} T_{k+1} & e_k \beta_{k+2} e_1^T \\ e_1 \beta_{k+2} e_k^T & \tilde{T}_{n-k+1} \end{bmatrix},$$

kde T_{k-1} i T_{k+1} regulární jsou, viz poznámka pod důsledkem 2.3. (Přesněji řečeno jsou regulární, pokud existují; pro $k = 1$ a $k = n - 1$ bude mít první respektive druhé blokové schéma nepatrně jinou strukturu.) Inverzi matice T_n pak zapíšeme ve tvarech

$$T_n^{-1} = \begin{bmatrix} E_{k-1} & F_{k-1} \\ F_{k-1}^T & H_{k-1} \end{bmatrix} = \begin{bmatrix} E_{k+1} & F_{k+1} \\ F_{k+1}^T & H_{k+1} \end{bmatrix}.$$

Jak již víme, hodnosti obou mimodiagonálních bloků $F_{k-1} \in \mathbb{R}^{(k-1) \times (n-k+1)}$ a $F_{k+1} \in \mathbb{R}^{(k+1) \times (n-k-1)}$ budou rovny právě jedné, za předpokladu $\beta_k \neq 0$ a $\beta_{k+2} \neq 0$. Tedy

$$\text{rank}(F_{k-1}) = \text{rank}(F_{k+1}) = 1,$$



Obrázek 3.1: Šrafované bloky F_{k-1} a F_{k+1} mají hodnotu 1. Nás zajímá hodnota tučně vyznačeného bloku.

z obrázku 3.1 je přitom vidět, že podmatice F_{k-1} a F_{k+1} vymezují odlišné části matice. Dále víme, že mají hodnotu rovnou právě jedné. Z toho dokážeme odvodit, jaké má vlastnosti mezilehlý člen. Určitě tedy platí

$$\text{rank} \left(\begin{array}{c|c} \begin{array}{c|c} \text{diagonal} & \text{cross-hatch} \\ \hline \text{diagonal} & \text{diagonal} \end{array} & \begin{array}{c|c} \text{cross-hatch} & \text{cross-hatch} \\ \hline \text{diagonal} & \text{diagonal} \end{array} \\ \hline \begin{array}{c|c} \text{diagonal} & \text{diagonal} \\ \hline \text{diagonal} & \text{diagonal} \end{array} & \begin{array}{c|c} \text{diagonal} & \text{diagonal} \\ \hline \text{diagonal} & \text{diagonal} \end{array} \end{array} \right) \leq 2; \quad (3.12)$$

později, v sekci 3.4 ukážeme, že je to pro libovolnou regulární symetrickou třídiagonální matici vždy nejvýše jedna.

3.3 Symetrická třídiagonální pozitivně definitní matice

Uvažujme nyní naši symetrickou třídiagonální matici T_n , která má všechna vlastní čísla kladná, takovou matici nazýváme *pozitivně definitní*. Odřízneme-li poslední řádek a poslední sloupec dostaneme matici T_{n-1} , její vlastní čísla jsou opět kladná. V důsledku prokládání 2.4 snadno zjistíme, že i matice T_k z (3.8) musí mít kladná vlastní čísla pro $k = 1, \dots, n$. Soubor matic T_1, \dots, T_n nazýváme hlavní rohové podmatice. V opačném případě pokud má matice T_n všechna vlastní čísla záporná, nazýváme ji *negativně definitní*.

Důsledek 4. *Nechť T_n je symetrická třídiagonální pozitivně definitní matice*

$$T_n^{-1} = \begin{bmatrix} E_k & F_k \\ F_k^T & H_k \end{bmatrix},$$

kteřá splňuje $\beta_2\beta_3\cdots\beta_n \neq 0$, pak všechny mimodiagonální bloky F_k a F_k^T , $k = 1, 2, \dots, n-1$, její inverze mají hodnotu rovnou jedné. Obecně pak mimodiagonální bloky mají hodnotu rovnou nejvýše jedné. Pokud navíc je některé $\beta_{k+1} = 0$, pak víme, že odpovídající mimodiagonální bloky budou hodnoty nula, viz (3.11).

3.4 Upřesnění na závěr

Nyní ukážeme, že každý mimodiagonální blok inverze regulární symetrické třídiagonální matice T_n^{-1} má hodnotu opravdu nejvýše jedna a zpřesníme tak výsledek (3.12). Toho ovšem musíme dosáhnout nepatrně komplikovanější cestou, konkrétně tak, že se podíváme na jednotlivé prvky daného bloku. Pro jednoduchost se budeme soustředit na naddiagonální blok. Jeho (i, s) -tý prvek je dán známým vztahem

$$(T_n^{-1})_{i,s} = (-1)^{i+s} \cdot \frac{\det([T_n]_{i,s})}{\det(T_n)},$$

kde $[T_n]_{i,s}$ vznikne z matice T_n odstraněním i -tého řádku a s -tého sloupce (připomeňme, že matice T_n je symetrická):

$$\begin{array}{c} \boxed{\begin{array}{cccc} \alpha_1 & \beta_2 & & \\ \beta_2 & \alpha_2 & \ddots & \\ & \ddots & \ddots & \beta_{i-1} \\ & & \beta_{i-1} & \alpha_{i-1} \end{array}} & \begin{array}{c} \beta_i \\ \alpha_i \\ \beta_{i+1} \end{array} & | & \\ \hline & \boxed{\begin{array}{ccc} \beta_{i+1} & \alpha_{i+1} & \beta_{i+2} \\ & \beta_{i+2} & \alpha_{i+2} & \ddots \\ & & \ddots & \ddots & \beta_{s-1} \\ & & & \beta_{s-1} & \alpha_{s-1} \\ & & & & \beta_s \end{array}} & \begin{array}{c} \beta_s \\ \alpha_s \\ \beta_{s+1} \end{array} & | & \\ & & \boxed{\begin{array}{ccc} \beta_{s+1} & \alpha_{s+1} & \beta_{s+2} \\ \beta_{s+2} & \alpha_{s+2} & \ddots \\ & \ddots & \ddots & \beta_n \\ & & \beta_n & \alpha_n \end{array}} & & \end{array}$$

Vidíme tedy, že matice $[T_n]_{i,s}$ je blokově trojúhelníková s vyznačenými čtvercovými bloky na diagonále. Nechť

$$\text{tridiag}((\alpha_j, \dots, \alpha_k), (\beta_{j+1}, \dots, \beta_k))$$

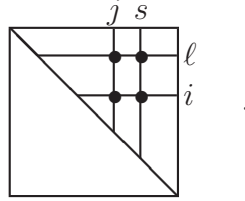
značí třídiagonální matici, která je podmaticí (resp. blokem) matice T_n obsahující dané prvky. Označme dále

$$D_j^k = \det \left(\text{tridiag}((\alpha_j, \dots, \alpha_k), (\beta_{j+1}, \dots, \beta_k)) \right).$$

Pak zřejmě

$$(T_n^{-1})_{i,s} = (-1)^{i+s} \cdot \frac{D_1^{i-1} \cdot (\beta_{i+1} \cdots \beta_s) \cdot D_{s+1}^n}{D_1^n}. \quad (3.13)$$

Zbývá si uvědomit, že na to, aby měla matice (resp. zde mimodiagonální blok) hodnost nejvýše jedna, každá její podmatice řádu 2 musí být singulární. Vyberme tedy z naddiagonálních prvků podmatice následujícím způsobem:



Pro jednotlivé řádkové a sloupcové indexy zřejmě platí

$$1 \leq \ell < i < j < s \leq n. \quad (3.14)$$

Vybraná podmatice

$$\begin{bmatrix} (T^{-1})_{\ell,j} & (T^{-1})_{\ell,s} \\ (T^{-1})_{i,j} & (T^{-1})_{i,s} \end{bmatrix}$$

má determinant roven nule právě tehdy, když platí

$$(T^{-1})_{\ell,j}(T^{-1})_{i,s} = (T^{-1})_{\ell,s}(T^{-1})_{i,j}.$$

Dosazením (3.13) za jednotlivé prvky rovnost snadno ověříme. Vlevo dostaneme

$$\begin{aligned} & (T^{-1})_{\ell,j} \cdot (T^{-1})_{i,s} \\ &= (-1)^{\ell+j+i+s} \cdot \frac{D_1^{\ell-1} \cdot D_1^{i-1} \cdot (\beta_{\ell+1} \cdots \beta_j)(\beta_{i+1} \cdots \beta_s) \cdot D_{j+1}^n \cdot D_{s+1}^n}{(D_1^n)^2}, \end{aligned}$$

vpravo

$$\begin{aligned} & (T^{-1})_{\ell,s} \cdot (T^{-1})_{i,j} \\ &= (-1)^{\ell+s+i+j} \cdot \frac{D_1^{\ell-1} \cdot D_1^{i-1} \cdot (\beta_{\ell+1} \cdots \beta_s)(\beta_{i+1} \cdots \beta_j) \cdot D_{s+1}^n \cdot D_{j+1}^n}{(D_1^n)^2}. \end{aligned}$$

Původní rovnost se tedy zredukuje na rovnost

$$(\beta_{\ell+1} \cdots \beta_j)(\beta_{i+1} \cdots \beta_s) = (\beta_{\ell+1} \cdots \beta_s)(\beta_{i+1} \cdots \beta_j).$$

Z nerovnosti mezi indexy (3.14) však ihned plyne, že se obě strany rovnají součinu

$$(\beta_{\ell+1} \cdots \beta_i)(\beta_{i+1} \cdots \beta_j)^2(\beta_{j+1} \cdots \beta_s).$$

Tím jsme tedy dokázali, že mimodiagonální blok nemůže mít hodnost vyšší než jedna. Toto pozorování můžeme zformulovat do věty, kterou již nebudeme (nemusíme) dokazovat.

Věta 1. *Nechť T je reálná čtvercová regulární symetrická třídiagonální matice. Pak každý mimodiagonální blok matice T^{-1} má hodnotu nejvýše jedna.*

Poznamenejme, že hodnota mimodiagonálního bloku ovšem může být menší než jedna. Jednak může být matice T (a tudíž i její inverze) blokově diagonální, viz (3.11), to ale souvisí s tím, že matice T má některý prvek β nulový. Druhý případ nastává například zde

$$\begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}^{-1} = \begin{bmatrix} 0 & 1 & 0 & -1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ -1 & 0 & 1 & 0 \end{bmatrix}.$$

Matice napravo obsahuje nulové mimodiagonální bloky tvaru $[0, 0]$ a $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$. Nejsou to ovšem bloky stejného tvaru jako výše používaný blok F_k .

Podobné výsledky je možné odvodit např. pro symetrickou pětidiagonální matici, jejíž inverze bude mít mimodiagonální bloky hodnoty nejvýše dva. Inverze symetrické sedmidiagonální matice bude mít hodnotu mimodiagonálních bloků omezenou třemi, atp.

4 Hierarchické matice

Nyní je dobré si uvědomit, že symetrická třídiagonální matice T_n (2.2) je jednoznačně určena právě $2n - 1$ čísly $\alpha_1, \dots, \alpha_n, \beta_2, \dots, \beta_n$. Chceme-li si takovou matici zapamatovat, nebo ji uložit v počítači, není si tedy třeba pamatovat všech n^2 prvků, ale stačí $(2n - 1) \sim n$ výše zmíněných. Naopak inverze třídiagonální matice je obecně hustá, jak je naznačeno např. na obrázku 2. Pro její zapamatování je (na první pohled) třeba uložit všech n^2 , resp. s využitím symetrie $\binom{n+1}{2} = \frac{1}{2}(n^2 + n) \sim n^2$ prvků. V následujícím textu však ukážeme, že (nejen) v případě inverze třídiagonální matice lze při chytrém způsobu ukládání prvků místo v paměti výrazně ušetřit.

4.1 Hierarchické dělení matice do stromové struktury

Následně budeme řešit, zda se dá její inverze zapsat také tímto způsobem. Pro jednoduchost si zvolíme symetrickou třídiagonální pozitivně definitní matici o velikosti $n = 2^\ell$. Díky tomu můžeme matici rozdělit přesně „na půl“ následujícím způsobem

$$T_n = \begin{bmatrix} T_{n/2}^{(1)} & e_{n/2}\beta_{n/2+1}e_1^T \\ e_1\beta_{n/2+1}e_{n/2}^T & T_{n/2}^{(2)} \end{bmatrix} \in \mathbb{R}^{n \times n}. \quad (4.1)$$

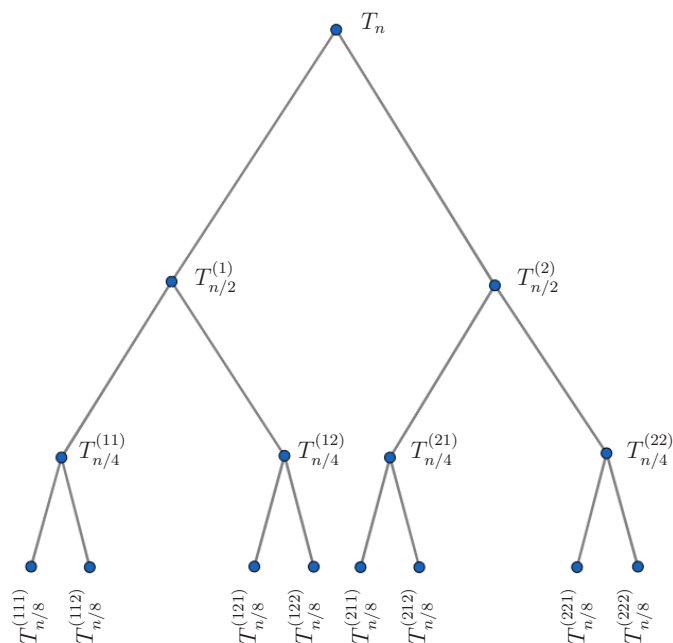
Dále se zaměříme na matice $T_{n/2}^{(1)}$ a $T_{n/2}^{(2)}$ na diagonále, a ty dále rozdělujeme opět napůl

$$T_{n/2}^{(1)} = \begin{bmatrix} T_{n/4}^{(11)} & e_{n/4}\beta_{n/4+1}e_1^T \\ e_1\beta_{n/4+1}e_{n/4}^T & T_{n/4}^{(12)} \end{bmatrix} \in \mathbb{R}^{(n/2) \times (n/2)} \quad \text{a}$$

$$T_{n/2}^{(2)} = \begin{bmatrix} T_{n/4}^{(21)} & e_{n/4}\beta_{n/4+1}e_1^T \\ e_1\beta_{n/4+1}e_{n/4}^T & T_{n/4}^{(22)} \end{bmatrix} \in \mathbb{R}^{(n/2) \times (n/2)},$$

a dále bychom stejným způsobem pokračovali při dalším rozdělování. Vytvoříme tím dělení původní matice do struktury vyváženého binárního stromu viz obrázek 1.1, resp. 4.1 (díky tomu, že $n = 2^\ell$), kde T_n je kořen a je předkem dvou potomků $T_{n/2}^{(1)}$ a $T_{n/2}^{(2)}$ a tak dále. Analogicky se dělí všechny větve, například matici $T_{n/64}^{(121121)}$ rozdělíme (je-li n dostatečně velké) tak, že vzniknou následující bloky na její diagonále

$$T_{n/64}^{(121121)} \longrightarrow \left(T_{n/128}^{(1211211)}, T_{n/128}^{(1211212)} \right).$$



Obrázek 4.1: Hierarchická reprezentace symetrické pozitivně definitní třídiagonální matice strukturovaná do vyváženého binárního stromu. Ve schématu nejsou vyznačené mimodiagonální bloky hodnoty jedna, které musíme ukládat.

Mimodiagonální bloky, které při rozdělování vznikají, mají všechny hodnotu vždy rovnu nejvýše jedné.

Inverzi matice (4.1) lze, jak již víme, zapsat ve tvaru

$$T_n^{-1} = \begin{bmatrix} E_{n/2}^{(1)} & F_{n/2} \\ F_{n/2}^T & E_{n/2}^{(2)} \end{bmatrix}, \quad (4.2)$$

přičemž hodnota mimodiagonálního bloku $F_{n/2}$ (a ze symetrie triviálně i $F_{n/2}^T$) je rovna jedné. Zcela analogicky bychom se nyní zaměřili hierarchicky na matice na matice $E_{n/2}^{(1)}$ a $E_{n/2}^{(2)}$ a ty hierarchicky „půlili“. Protože všechny mimodiagonální bloky matice T_n^{-1} mají hodnotu jedna, budou mít hodnotu jedna i mimodiagonální bloky vzniklé „půlením“ matic $E_{n/2}^{(1)}$ a $E_{n/2}^{(2)}$.

4.2 Paměťové náklady na uložení hierarchické struktury a její omezení

Pokud máme matici hodnoty jedna, tak to znamená, že má všechny řádky stejné až na násobek skalárem (lineárně závislé), to samé platí pro sloupce. Jediné co potřebujeme znát, abychom tuto matici mohli zapsat, je jediný její (nenulový) řádek, řekněme například první, a jediný její (nenulový) sloupec, řekněme opět například první. Pro jednoduchost můžeme tedy říct, že pro uložení matice $m \times k$

hodnosti jedna je potřeba si zapamatovat $(m + k)$ čísel. Poznamenejme, že analogicky, pro stejně velkou matici hodnosti dva je potřeba si zapamatovat dva její lineárně nezávislé sloupce a dva lineárně nezávislé řádky, tj. $2(m + k)$ čísel. Pro matici hodnosti r to bude $r(m + k)$ čísel.

Viděli jsme, že u hierarchického dělení symetrické třídiagonální pozitivně definitní matice T_n mají mimodiagonální bloky typu $e_1 \beta_{n/2+1} e_{n/2}^T$ (viz (4.1)) hodnost jedna. Stejně tak odpovídající mimodiagonální bloky typu $F_{n/2}$ matice T_n^{-1} (viz (4.2)) jsou také hodnosti jedna, jak jsme ukázali v kapitole 3, viz zejména důsledek 4.

Podívejme se tedy, kolik čísel je třeba pro hierarchické uložení zejména matice T_n^{-1} (u matice T_n máme přímější a jednodušší způsob jak ji uložit, tj. pouze její diagonálu a naddiagonálu). Soustředit se budeme primárně na uložení mimodiagonálních bloků typu $F_{n/2}$ (viz (4.2)).

V našem dělení jsou všechny tyto mimodiagonální bloky čtvercové a řádů $n/(2^j)$, pro různá přirozená j . První řádek z největšího mimodiagonálního bloku má $\frac{n}{2}$ čísel, u menšího bloku $\frac{n}{4}$, pak $\frac{n}{8}$ a tak dále. To samé platí pro sloupce. Pro zapamatování takových bloků (hodnosti jedna) je tedy potřeba si uložit právě $n = \frac{n}{2} + \frac{n}{2}$, resp. $\frac{n}{2} = \frac{n}{4} + \frac{n}{4}$ a $\frac{n}{4} = \frac{n}{8} + \frac{n}{8}$ atd. čísel.

To ale není vše, protože matic určitých rozměrů vzniká v dělení obecně více. Na první úrovni dělení vzniká jen jedna mimodiagonální matice (druhou z důvodu symetrie nemusíme uvažovat). Na druhé úrovni již dvě, na třetí úrovni dělení čtyři a tak dále. Pro uložení všech *mimodiagonálních bloků* (resp. naddiagonálních s využitím symetrie) inverze T_n^{-1} tedy potřebujeme

$$1 \left(\frac{n}{2} + \frac{n}{2} \right) + 2 \left(\frac{n}{4} + \frac{n}{4} \right) + 4 \left(\frac{n}{8} + \frac{n}{8} \right) + \dots$$

čísel. Protože jsme uvažovali původní matici T_n řádu $n = 2^\ell$, poslední sčítanec bude formálně

$$2^{\ell-1} \left(\frac{n}{2^\ell} + \frac{n}{2^\ell} \right).$$

Sčítanců je tedy právě ℓ . Hodnota každého sčítance je navíc rovna n . Celkem tedy potřebujeme

$$\underbrace{n + n + n + \dots + n}_\ell + n = n\ell + n = n(\log_2(n) + 1),$$

čísel, kde prvních $n\ell$ čísel má původ v mimodiagonálních blocích hodnosti jedna a posledních n čísel jsou prvky na diagonále.

4.2.1 Poznámka k velikosti nejmenších bloků

Vzhledem k tomu, že n jsme si zvolili jako 2^ℓ , tak je možné matici dělit vždy na poloviny. Tímto rozkladem můžeme postupovat až na „základní úroveň“, kde nám na blokové diagonále zůstanou právě pouze diagonální prvky. Otázkou však je, jestli je nutné, resp. výhodné, provádět dělení až do takto jemných detailů.

Uvažujeme-li (symetrickou) matici 2×2 , je zřejmé, že její (jediný možný) mimodiagonální blok nemá smysl ukládat jako dvojici (*nenulový řádek, nenulový sloupec*), tj. jako dvě čísla, když je sám o sobě jediným číslem. Fakticky tedy matici 2×2 není vhodné ukládat hierarchicky (protože by to vyžadovalo uložení čtyř čísel a navíc nějakou režii pro složitější organizaci dat v paměti, kterou zde vůbec nediskutujeme) ale uložíme ji přímo po prvcích (což s využitím symetrie vyžaduje uložení pouze tří čísel).

Podobnou úvahu můžeme provést i pro matici 4×4 , kde (když uvažujeme půlení rozměrů, jak bylo naznačeno) existuje také jen jeden mimodiagonální blok velikosti 2×2 . U takového bloku je zapamatování dvojice řádek-sloupec stejně náročné jako zapamatování bloku celého; vždy si musíme zapamatovat čtyři čísla. Teoreticky tedy nemá smysl ukládat matice 4×4 (a menší) hierarchickým způsobem.

Prakticky může být hranice ještě výš a i větší (diagonální) bloky budeme ukládat přímo, viz obrázek 4.3. Důvodem může být v první řadě to, že při uložení takové hierarchické matice musíme také nějakou paměť věnovat na uložení stromu (struktury samotné), a dále také to, že uložení matice do této struktury a další práce s ní vyžaduje určitý čas a výpočetní prostředky, tj. již výše zmíněnou režii.

4.3 Hierarchický přístup pro obecnou matici

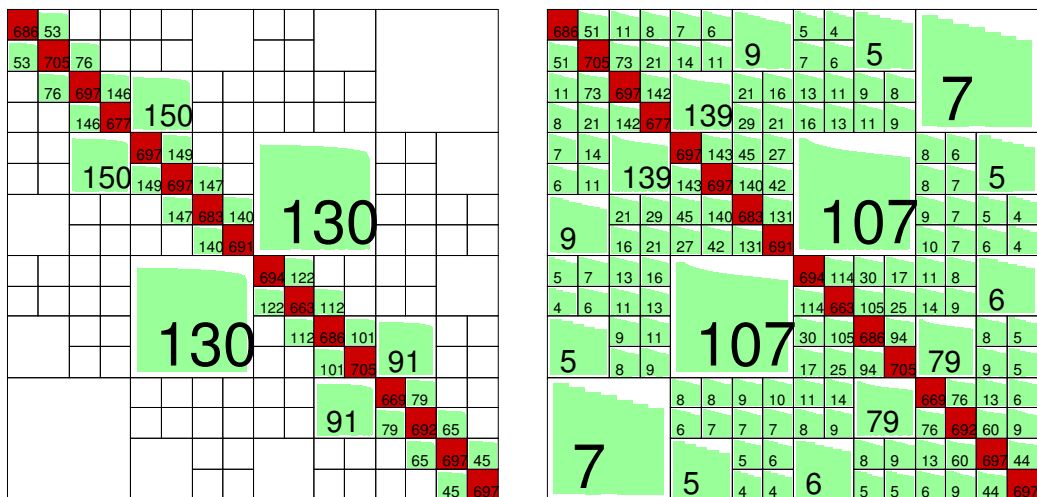
V případě čtvercových matic obecných rozměrů, nemůžeme takto postupovat. Proto musíme zvolit jiný způsob dělení. Ve skutečnosti je potřeba v každém kroku hierarchie rozdělit množinu indexů na disjunktní podmnožiny po sobě jdoucích index, jejichž sjednocením dostaneme původní množinu, a poté analogicky postupovat u takto získaných podmnožin. Je v zásadě jedno jakým způsobem toto budeme provádět. Zápisem takového rozdělení získáme strom.

Námi výše zvolená struktura dělení na poloviny (která je možná jen při $n = 2^\ell$) vede na vyvážený binární strom, viz první graf na obrázku 1.1. Obecně při dělení na dvě disjunktní podmnožiny dostaneme nějaký obecný binární strom. Mimodiagonální bloky v dělení již nebudou čtvercové, vše ale bude fungovat zcela analogicky. Jen je potřeba si uvědomit, že pokud zvolíme velmi špatné dělení (výrazně nevyvážený strom, viz druhý graf na obrázku 1.1) nemusí výsledná hierarchická struktura vést k úspoře místa při ukládání matice (pokud budou mít všechny mimodiagonální bloky např. jen jeden řádek; tj. pokud bychom matice $T^{(1)}$, $T^{(21)}$, $T^{(221)}$, atd. volili řádu jedna). Z tohoto úhlu pohledu tedy bude vždy nejvýhodnější dělit matici napůl, nebo alespoň přibližně napůl, není-li jiná možnost.

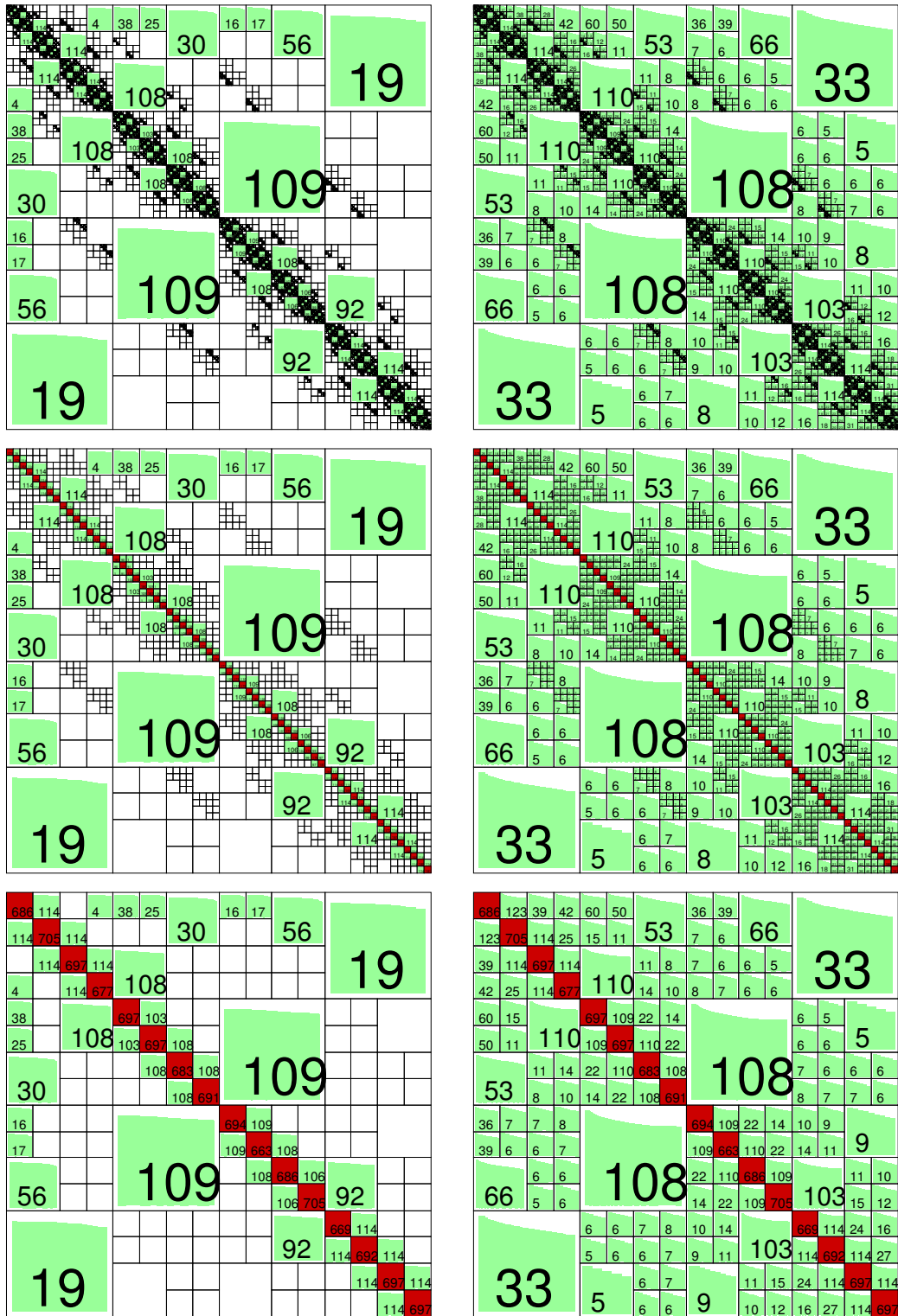
Poznamenejme, že předchozí odvození byla udělána pro symetrickou třídiagonální matici, resp. její inverzi, ale v praxi je můžeme provést vždy, když bude možné hodnoty všech mimodiagonálních bloků dané matice nějak rozumně omezit (např. když budeme mít obecnou matici řádu milion, o které víme, že všechny její mimodiagonální bloky jsou hodnoty nejvýše sto).

Poslední věcí, kterou je dobré poznamenat, je, že v praxi nejsme omezeni binárním dělením (stromem) matice a můžeme ji na každé úrovni rozdělit na obecný počet bloků (disjunktních podmnožin indexů). Význam a smysl takového dělení může mít

původ v úloze (např. ve fyzické geometrii oblasti, na které úlohu řeším), ze které matice vznikla, viz např. [7], [1] a [2]. Pro příklad obecného dělení obecné matice a její inverze viz obrázky 4.2 a 4.3. Poznamenejme, že matice na obrázku 4.2 a na obrázku 4.3 se liší pouze v permutaci řádků a sloupců.



Obrázek 4.2: Vlevo je obecná matice řádu $n = 11036$, vpravo je její inverze. Červené bloky jsou uloženy přímo (hustě; číslo v bloku označuje jeho řád). Zelenobílé bloky jsou hodnoty r (číslo v bloku), uloženy jako r lineárně nezávislých řádků a r lineárně nezávislých sloupců. Matice pocházejí z úlohy [10] a jsou vytvořené pomocí softwaru `hlib` [21] a Matlab.



Obrázek 4.3: Vlevo je obecná matice řádu $n = 11036$, vpravo je její inverze. V jednotlivých řádcích je velikost nejmenšího hierarchicky ukládaného bloku omezená na řády 16, 256 a 1024 (tj. na diagonále ležící bloky menších řádů se ukládají hustě). Matice pocházejí z úlohy [10] a jsou vytvořené pomocí softwaru hlib [21] a Matlab.

5 Základní operace

Vytvořit ze symetrické pozitivně definitní třídiagonální matice hierarchickou maticí ve výše uvedeném smyslu je snadné, ať už je jakkoliv velká. Méně snadné je získat její inverzi v tomto tvaru. Zřejmě není cílem inverzi nejprve spočítat (neboť je hustá), a pak ji rozkládat do stromové struktury. Inverzi bude potřeba provést přímo v hierarchickém tvaru. K tomu budeme potřebovat umět s takovými maticemi manipulovat. Typicky budeme potřebovat operace jako jsou *součet* dvou (hierarchicky uložených) matic, *součin matice s vektorem*, případně *součin dvou matic*, nebo i složitější operace jako je *LU rozklad matice*, atd., viz [13] nebo [21] (v obou případech je k dispozici i software, který takové operace implementuje za částečného použití Matlabu).

5.1 Součet dvou hierarchických matic

Začneme nejjednodušší operací a tou je součet dvou hierarchických matic. Uvažujme matice M_1 a M_2 stejných rozměrů, které mají stejnou hierarchickou strukturu. Tedy za prvé jejich hierarchická struktura je popsána stejným stromem a za druhé při každém větvení stromu se původní množina indexů dělí na stejné disjunktí podmnožiny. Nechť tedy

$$M_1 + M_2 = \begin{bmatrix} A_1 & B_1 \\ C_1 & D_1 \end{bmatrix} + \begin{bmatrix} A_2 & B_2 \\ C_2 & D_2 \end{bmatrix} = \begin{bmatrix} A_1 + A_2 & B_1 + B_2 \\ C_1 + C_2 & D_1 + D_2 \end{bmatrix} = M.$$

Protože jsou původní matice uloženy hierarchicky, matice A a D jsou buď opět hierarchické, nebo jsou uloženy hustě. V prvním případě součet matic $A_1 + A_2$, resp. $D_1 + D_2$ realizujeme stejně jako u původního součtu $M_1 + M_2$. Rekurzivně tak postupujeme stromovou hierachií a nakonec se dopracujeme k blokům, které jsou uloženy hustě, tj. k druhému případu. Tyto bloky pak sečteme přímo. Potřebujeme se tedy soustředit hlavně na součty mimodiagonálních bloků. Protože pracujeme s hierarchickým formátem matice, budeme předpokládat, že hodnota mimodiagonálních bloků B_1 , C_1 , B_2 a C_2 je rozumně malá.

5.1.1 Součet dvou bloků nízké hodnoti

V dalším textu se budeme soustředit pouze na součet bloků $B = B_1 + B_2$ (pro bloky C_1 a C_2 bude sčítání probíhat zcela analogicky); nechť navíc B , B_1 , $B_2 \in \mathbb{R}^{m \times k}$ pro jednoduchost. Označme tedy, že

$$r_1 = \text{rank}(B_1) \quad \text{a} \quad r_2 = \text{rank}(B_2).$$

Přičemž bloky B_j , $j = 1, 2$ nízké hodnoti zpravidla ukládáme jako sadu r_j lineárně nezávislých sloupců a r_j lineárně nezávislých řádků, resp. vektorů takových, že jejich součin vytvoří původní matici. Konkrétně lze tedy matice zapsat ve tvaru

$$B_j = U_j V_j^T, \quad \text{kde } U_j \in \mathbb{R}^{m \times r_j}, \quad V_j \in \mathbb{R}^{k \times r_j}, \quad j = 1, 2.$$

Matici B můžeme zapsat

$$B = B_1 + B_2 = U_1 V_1^T + U_2 V_2^T = [U_1, U_2] \begin{bmatrix} V_1^T \\ V_2^T \end{bmatrix}.$$

Hodnost této matice si označíme

$$r = \text{rank}(B), \quad \text{přičemž platí } \max\{r_1, r_2\} \leq r \leq r_1 + r_2.$$

Poznamenejme, že vlevo nastane rovnost např. ve speciálním případě, když bude matice B_1 stejná jako matice B_2 . Tudiž sčítáme dvě stejné matice $B = 2B_1 = 2B_2$ a počet lineárně nezávislých řádků $r_1 = r_2 = r$ se nemění.

V tuto chvíli máme matici B uloženou jako součin sloupců a řádků. Těchto řádků je ale vždy $r_1 + r_2$, nikoliv r , jak bychom chtěli. Po tomto (formálním) sečtení (sečtení se zde zredukovalo jen na sestavení matic $[U_1, U_2]$ a $[V_1, V_2]$ z jejich podmatic U_1, U_2, V_1 a V_2), je potřeba provést tzv. *kompresi*, viz např. [10], nebo [17, sekce 1.3.1]. Tento přístup, který nás bude provázet víceméně i u všech dalších operací, se také často nazývá výrazem *low-rank aritmetika*. Pojem zatím nemá vhodný český ekvivalent; obecně jde o snahu pracovat s objekty – zde maticemi, nebo jejich částmi – relativně nízké hodnoti a tohoto faktu při práci nějak vhodně využít.

5.1.2 Komprese

Kompresse spočívá v tom, že se snažíme zredukovat faktory $[U_1, U_2] \in \mathbb{R}^{m \times (r_1 + r_2)}$ a $[V_1, V_2] \in \mathbb{R}^{k \times (r_1 + r_2)}$ matice B tak, aby měli právě r lineárně nezávislých sloupců. Označme

$$r_U = \text{rank}([U_1, U_2]) \quad \text{a} \quad r_V = \text{rank}([V_1, V_2]),$$

zřejmě opět platí

$$\max\{r_1, r_2\} \leq r_U \leq r_1 + r_2 \quad \text{a} \quad \max\{r_1, r_2\} \leq r_V \leq r_1 + r_2.$$

Pomocí QR rozkladu počítaného např. pomocí Gramovy–Schmidtovy ortogonalizace (viz např. [5, kapitola 3]) dostaneme

$$[U_1, U_2] = Q_U R_U, \quad \text{kde } Q_U \in \mathbb{R}^{m \times r_U}, \quad Q_U^T Q_U = I_{r_U}, \quad R_U \in \mathbb{R}^{r_U \times (r_1 + r_2)},$$

a

$$[V_1, V_2] = Q_V R_V, \quad \text{kde } Q_V \in \mathbb{R}^{k \times r_V}, \quad Q_V^T Q_V = I_{r_V}, \quad R_V \in \mathbb{R}^{r_V \times (r_1 + r_2)},$$

příčemž matice R_U a R_V jsou v tzv. horním schodovitém tvaru, např.

$$\begin{bmatrix} \clubsuit & \heartsuit & \heartsuit & \heartsuit & \heartsuit & \heartsuit & \heartsuit & \heartsuit & \heartsuit & \heartsuit \\ 0 & \clubsuit & \heartsuit & \heartsuit & \heartsuit & \heartsuit & \heartsuit & \heartsuit & \heartsuit & \heartsuit \\ 0 & 0 & \clubsuit & \heartsuit & \heartsuit & \heartsuit & \heartsuit & \heartsuit & \heartsuit & \heartsuit \\ 0 & 0 & 0 & 0 & 0 & \clubsuit & \heartsuit & \heartsuit & \heartsuit & \heartsuit \\ 0 & 0 & 0 & 0 & 0 & 0 & \clubsuit & \heartsuit & \heartsuit & \heartsuit \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \clubsuit & \heartsuit \end{bmatrix}, \quad \text{kde } \clubsuit \neq 0.$$

S využitím těchto dvou QR rozkladů můžeme přepsat součin do tvaru

$$B = [U_1, U_2] \begin{bmatrix} V_1^T \\ V_2^T \end{bmatrix} = Q_U \underbrace{R_U R_V^T}_{W} Q_V^T, \quad \text{kde } W \in \mathbb{R}^{r_U \times r_V}. \quad (5.1)$$

Posledním náročnějším krokem je výpočet singulárního rozkladu matice W , resp. jeho tzv. *ekonomického tvaru* (viz např. [5, kapitola 5]). Nechť tedy

$$W = U_W \Sigma_W V_W^T, \quad \text{kde } U_W \in \mathbb{R}^{r_U \times r}, \Sigma_W \in \mathbb{R}^{r \times r}, V_W \in \mathbb{R}^{r_V \times r}$$

a kde matice Σ_W je navíc diagonální s kladnými čísly (tzv. singulárními čísly matice W) na diagonále, tj. Σ_W je regulární. Zde je vhodné si uvědomit, že singulární čísla matice W jsou zároveň singulárními čísly původní matice B ; hodnota matice W tedy musí být přímo rovna $\text{rank}(W) = \text{rank}(B) = r$. Platí tedy

$$B = [U_1, U_2] \begin{bmatrix} V_1^T \\ V_2^T \end{bmatrix} = Q_U W Q_V^T = Q_U U_W \Sigma_W V_W^T Q_V^T = \underbrace{(Q_U U_W \Sigma_W)}_U \underbrace{(Q_V V_W^T)}_{V^T},$$

kde

$$U = Q_U U_W \Sigma_W \in \mathbb{R}^{m \times r} \quad \text{a} \quad V = Q_V V_W \in \mathbb{R}^{k \times r}.$$

Tím jsme dosáhli vytyčeného cíle, jak zapsat matici B hodnoti r jako součin dvou matic o r řádcích a r sloupcích. Tabulka 5.1 shrnuje výpočetní náročnost součtu dvou mimodiagonálních bloků.

Tabulka 5.1: Výpočetní náklady součtu dvou matic rozměrů $m \times k$ hodnotí r_1 a r_2 . Pro r_U a r_V navíc platí $\max\{r_1, r_2\} \leq r_U \leq r_1 + r_2$ a $\max\{r_1, r_2\} \leq r_V \leq r_1 + r_2$.

Způsob výpočtu	Výpočetní cena
klasicky	mk součtů dvojic reálných čísel
low-rank aritmetika	0 součtů (pouze sestavení matic) + komprese navíc
komprese	2× QR rozklad typu $m \times (r_1 + r_2)$ a $k \times (r_1 + r_2)$ 1× SVD typu $r_U \times r_V$ 4× maticové součiny

Poznamenejme, že alternativně můžeme matici Σ_W zahrnout do matice V případně do obou matic U i V , neboť $(Q_U U_W \Sigma_W)(V_W^T Q_V^T) = (Q_U U_W)(\Sigma_W V_W^T Q_V^T) = (Q_U U_W \Sigma_W^{1/2})(\Sigma_W^{1/2} V_W^T Q_V^T)$, kde diagonální matice $\Sigma_W^{1/2}$ má na diagonále odmocniny z diagonálních prvků Σ_W , tj. $(\Sigma_W^{1/2})^2 = \Sigma_W^{1/2} \Sigma_W^{1/2} = \Sigma_W$.

5.2 Součet hierarchické matice, low-rank matice a husté matice

Abychom mohli následně pracovat s násobením dvou matic, je třeba znát, jakým způsobem lze řešit součty různých typů matic. Přesněji řečeno, potřebujeme vědět, jak jejich součet vypadá.

5.2.1 Součet hierarchické matice a low-rank matice

Začneme součtem hierarchické matice a low-rank matice. Nechť je M_1 hierarchická matice ve tvaru

$$M_1 = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \in \mathbb{R}^{n \times n},$$

a nechť M_2 je low-rank matice zapsaná jako následující součin

$$M_2 = UV^T, \quad \text{kde } U \in \mathbb{R}^{n \times r}, V \in \mathbb{R}^{n \times r}.$$

Rozdělíme-li U a V na bloky vhodných velikostí

$$U = \begin{bmatrix} U_1 \\ U_2 \end{bmatrix}, \quad V = \begin{bmatrix} V_1 \\ V_2 \end{bmatrix},$$

pak lze součet hierarchické matice s low-rank maticí zapsat ve tvaru

$$M_1 + M_2 = \begin{bmatrix} A & B \\ C & D \end{bmatrix} + \begin{bmatrix} U_1 \\ U_2 \end{bmatrix} \begin{bmatrix} V_1^T & V_2^T \end{bmatrix} = \begin{bmatrix} A + U_1V_1^T & B + U_1V_2^T \\ C + U_2V_1^T & D + U_2V_2^T \end{bmatrix} = M.$$

Nyní si rozebereme jednotlivé součty na pravé straně:

1. Součet $A + U_1V_1^T$, resp. $D + U_2V_2^T$ je dvojího typu:
 - (a) Je-li A , resp. D hierarchická matice, jedná se o součet hierarchické matice a low-rank matice. Ten provedeme stejně jako součet $M_1 + M_2$, pokračujeme tedy rekurzivně.
 - (b) Je-li A , resp. D hustá matice (jsme na nejjemnějším dělení), pak součet přímo vyčíslíme.
2. Součet $B + U_1V_2^T$, resp. $C + U_2V_1^T$ reprezentuje součet dvou low-rank matic, který provedeme stejně jako v případě mimodiagonálních bloků v sekci 5.1, resp 5.1.1.

5.2.2 Součet hierarchické a husté, resp. low-rank a husté matice

Sčítáme-li hierarchickou matici s hustou maticí, nebo low-rank maticí s hustou maticí, je třeba si uvědomit, že hustá matice obecně nemá žádnou strukturu, se kterou bychom mohli pracovat. Pokud se tedy nerozhodneme hustou maticí přepočítat a najít její hierarchickou nebo low-rank reprezentaci, bude výsledkem takovéto operace vždy hustá matice. Hierarchický nebo low-rank sčítanec jednoduše vyčíslíme po prvcích a k husté matici přičteme.

5.3 Součin hierarchické matice a vektoru

Než se podíváme na součin dvou matic, bude vhodné se nejprve zaměřit na násobení matice s jedním vektorem. Nechť tedy

$$Mx = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} y \\ z \end{bmatrix} = \begin{bmatrix} Ay + Bz \\ Cy + Dz \end{bmatrix}.$$

Vidíme, že ve výsledku jsou objekty dvou typů:

1. Součiny diagonálních bloků s vektory, Ay a Dz , kde mohou nastat dvě situace:
 - (a) Pokud A , resp. D je hierarchická matice, jedná se opět o součin (menší) hierarchické matice s vektorem. Postupujeme tedy rekurzivně.
 - (b) Pokud jsme na nejjemnější úrovni dělení, tj. pokud A , resp. D je hustá matice, provedeme klasické násobení matice s vektorem.
2. Součiny mimodiagonálních bloků s vektory, Bz a Cy , přičemž tyto bloky jsou nízké hodnoti. Toto jsou jediné součiny, na které se musíme podívat detailněji.

Budeme opět uvažovat mimodiagonální blok $B \in \mathbb{R}^{m \times k}$ hodnoti r , který si můžeme rozepsat do tvaru

$$B = U_B V_B^T, \quad \text{kde } U \in \mathbb{R}^{m \times r}, \quad V \in \mathbb{R}^{k \times r}.$$

Pro blok C bude situace zcela analogická. Součin Bz si zřejmě můžeme rozepsat do tvaru

$$Bz = U_B V_B^T z = U_B z', \quad \text{kde } z' = V_B^T z.$$

Fakticky jsme tím provedli vhodné přezávkování $Bz = (U_B V_B^T)z = U_B (V_B^T z)$. Výpočet se tedy zredukuje na vypočtení dvou klasických součinů matic s vektory, konkrétně matice $r \times k$ s vektorem délky k a pak matice $m \times r$ s vektorem délky r .

Tabulka 5.2: Výpočetní náklady součinu matice rozměrů $m \times k$ hodnoti r s vektorem délky k . Operací rozumíme součin, resp. součet

Způsob výpočtu	Výpočetní cena
klasicky	$2mk - m \sim 2mk$ operací s dvojicí reálných čísel
low-rank aritmetika	$(2rk - r) + (2mr - m) \sim 2r(m + k)$ operací

5.3.1 Součin hierarchické matice a husté matice

Poznamenejme, že analogický postup můžeme použít při součinu hierarchické matice s několika vektory. Respektive, při součinu hierarchické matice s klasickou nehierarchickou (hustou) maticí, neboť tu můžeme interpretovat po sloupcích jako soubor vektorů. V případě součinu klasické matice s hierarchickou (v tomto pořadí) budeme postupovat analogicky, ale po řádcích (jako bychom zleva násobili řádkovými vektory).

5.4 Součin dvou hierarchických matic

Nyní jsme připraveni podívat se na součin dvou matic. Nechť

$$M_1 M_2 = \begin{bmatrix} A_1 & B_1 \\ C_1 & D_1 \end{bmatrix} \begin{bmatrix} A_2 & B_2 \\ C_2 & D_2 \end{bmatrix} = \begin{bmatrix} A_1 A_2 + B_1 C_2 & A_1 B_2 + B_1 D_2 \\ C_1 A_2 + D_1 C_2 & C_1 B_2 + D_1 D_2 \end{bmatrix} = M.$$

Je potřeba si uvědomit, z jakých jednotlivých operací se tento součin matic skládá.

5.4.1 Součiny různých bloků

Nejprve si rozebereme násobení, tj. součiny jednotlivých bloků, pak se podíváme na příslušné součty:

1. Součiny $A_1 A_2$, $D_1 D_2$ mohou představovat několik různých situací. Rozebereme si je pro součin $A_1 A_2$:
 - (a) V případě, že jsou obě matice A_1 i A_2 hierarchické, opakujeme stejný postup jako u M_1 a M_2 , neboli postupujeme opět rekurzivně.
 - (b) V případě, že je A_1 hierarchická a A_2 klasická hustě uložená matice, interpretujeme druhou matici a maticové násobení po sloupcích jakobychom násobili hierarchickou matici souborem vektorů, viz sekce 5.3, resp. 5.3.1 (při opačném pořadí struktur součinitelů postupujeme analogicky).
 - (c) V případě, že jsou obě matice A_1 i A_2 uloženy klasicky (tj. jsme na nejjemnější úrovni v dělení obou matic) provedeme klasické (husté) maticové násobení.
2. Součiny $B_1 C_2$, resp. $C_1 B_2$ dvou mimodiagonálních bloků nízké hodnoty budeme provádět pomocí low-rank operací.
3. Součiny $C_1 A_2$, $D_1 C_2$, $A_1 B_2$, resp. $B_1 D_2$ obsahují buď hierarchické nebo hustě uložené matice v součinu s low-rank maticí.

Protože situace v bodu 1. je jasná, podíváme rovnou na bod 2. Uvažujme součin $B_1 C_2 \in \mathbb{R}^{m \times s}$ (druhý provedeme analogicky) kde,

$$B_1 = U_B V_B^T \in \mathbb{R}^{m \times k}, \quad r_B = \text{rank} B_1, \quad U_B \in \mathbb{R}^{m \times r_B}, \quad V_B \in \mathbb{R}^{k \times r_B},$$

$$C_2 = U_C V_C^T \in \mathbb{R}^{k \times s}, \quad r_C = \text{rank} C_2, \quad U_C \in \mathbb{R}^{k \times r_C}, \quad V_C \in \mathbb{R}^{s \times r_C}.$$

Pak vhodným uzávorkováním dostaneme

$$B_1 C_2 = U_B \underbrace{(V_B^T U_C)}_W V_C^T, \quad \text{kde } W \in \mathbb{R}^{r_B \times r_C}.$$

Dále použijeme *kompresi* na součin $U_B W V_C^T$ analogicky jako byla použita v sekci 5.1.2 na součin $Q_U W Q_V^T$ (5.1). Pro výpočet tedy potřebujeme tři maticové součiny a jeden singulární rozklad. Poznamenejme, že jediný drobný rozdíl je v tom, že matice

Q_U a Q_V (ze sekce 5.1.2) mají ortonormální sloupce, což jsme zde u matic U_B a V_C obecně nevyžadovali. Pokud bychom po kompresi chtěli získat výsledek se stejným vlastnostmi jako při sčítání, museli bychom nejprve provést QR rozklady matic U_B a V_C .

Ve „třetím typu“ násobení v bodu 3. mohou nastat dvě varianty, zaprvé jedna z matic je hustá a druhá nízké hodnoty a zadruhé jedna z matic je hierarchická a druhá nízké hodnoty. Uvažujme např. součin $A_1B_2 \in \mathbb{R}^{m \times k}$ (ostatní provedeme analogicky), kde $A_1 \in \mathbb{R}^{m \times m}$,

$$B_2 = U_B V_B^T \in \mathbb{R}^{m \times k}, \quad r_B = \text{rank} B_2, \quad U_B \in \mathbb{R}^{m \times r_B}, \quad V_B \in \mathbb{R}^{k \times r_B}.$$

Pak zřejmě

$$A_1 B_2 = \underbrace{A_1 U_B}_{U_{AB}} V_B^T, \quad \text{kde } U_{AB} \in \mathbb{R}^{m \times r_B},$$

tedy potřebujeme provést právě r_B součinů matice s vektorem, přičemž matice je buď hustá (a součin provádíme klasicky), nebo hierarchická (a součin provádíme tak jako v sekci 5.3). Výsledek součinu $A_1 B_2$ tak dostaneme přímo v obvyklém low-rank tvaru. Opět poznamenejme, že od matic U_B a V_B jsme obecně nepožadovali, aby měly ortogonální, resp. ortonormální sloupce. Navíc, pokud by měly, matice $U_{AB} = A_1 U_B$ obecně ortogonální sloupce mít nebude. Pokud bychom chtěli její sloupce zortogonalizovat, můžeme opět použít QR rozklad a následnou kompresi tak jako v sekci 5.1.2.

5.4.2 Rekapitulace součinů bloků a součty součinů bloků

Nyní jsme si vyjasnili, jak jednotlivé součiny provádět, důležité ale také bude uvědomit si, co bude výsledkem těchto součinů, neboť s nimi musíme dále pracovat; viz tabulka 5.3. Tyto součiny je potřeba sčítat, přičemž se zde vyskytují součty následujících typů:

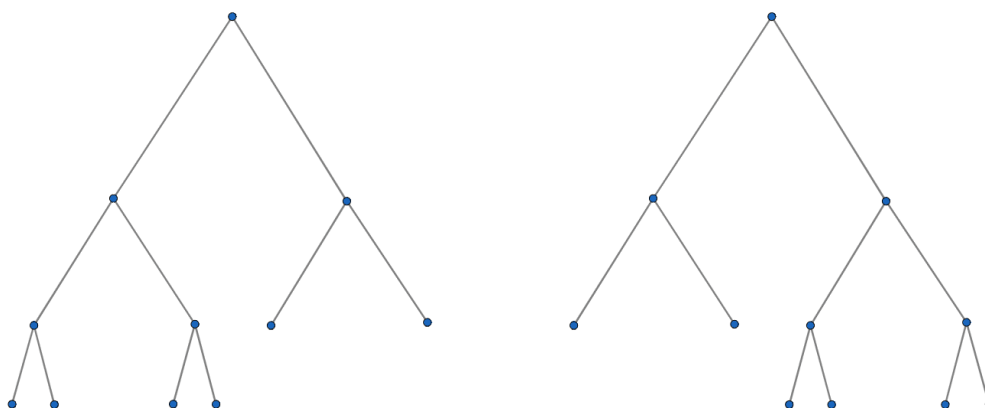
1. Součet $A_1 A_2 + B_1 C_2$, resp. $C_1 B_2 + D_1 D_2$ opět představuje dvě možnosti:
 - (a) První možností je, že sčítáme hierarchickou a low-rank matici, což jsme již podrobně rozebrali v sekci 5.2, resp. 5.2.1. Výsledkem je hierarchická matice.
 - (b) Druhou možností je, že sčítáme hustou matici a low-rank matici, kde nám nezbyvá než výsledek přímo vyčíslit, viz sekci 5.2.2. Výsledkem je hustá matice.
2. Poslední typ součtu, který může vzniknout, připadá na mimodiagonální bloky $C_1 A_2 + D_1 C_2$, resp. $A_1 B_2 + B_1 D_2$. Zde sčítáme dvě low-rank matice, což jsme podrobně rozebrali v sekci 5.1, resp. 5.1.1.

Tabulka 5.3: Tabulka rekapituluje součiny bloků a jejich výsledky. Jednotlivé typy (formáty) matic (obou součinitelů a součinu) jsou označeny zkratkou Hi, Ds a LR, které značí hierarchickou matici, hustou (dense) matici a low-rank matici. (Čísla 1a, 1b, 1c, 2, 3 vedle součinů odkazují na jednotlivé položky v sekci 5.4.1.)

Součin	Formát součinitelů	Formát součinu
A_1A_2, D_1D_2 (viz 1a)	Hi · Hi	Hi
A_1A_2, D_1D_2 (viz 1b)	Hi · Ds	Ds
A_1A_2, D_1D_2 (viz 1c)	Ds · Ds	Ds
B_1C_2, C_1B_2 (viz 2)	LR · LR	LR
$C_1A_2, D_1C_2, A_1B_2, B_1D_2$ (viz 3)	Hi · LR	LR
$C_1A_2, D_1C_2, A_1B_2, B_1D_2$ (viz 3)	Ds · LR	LR

5.5 Poznámka k součtu a součinu dvou hierarchických matic

V sekci 5.2 o součtu dvou hierarchických matic jsme předpokládali, že oba sčítance mají stejnou hierarchickou strukturu. Poznamenejme ale, že ani pro součin ani pro součet dvou hierarchických matic není tento požadavek nutný. Není třeba, aby obě matice měly úplně stejnou stromovou strukturu, musí jen platit, že průnik (v jistém smyslu) jejich stromů je zase strom. Navíc část hierarchického rozkladu odpovídající tomuto společnému podstromu musí odpovídat dělení matic na stejné podmnožiny indexů (bloky musí být stejných rozměrů, aby šlo sčítání, resp. násobení provádět). Větvení jednotlivých stromů však může jít do různě jemných detailů, tj. u jedné matice se může v dané větvi zastavit dříve než u druhé, viz obrázek 5.1.

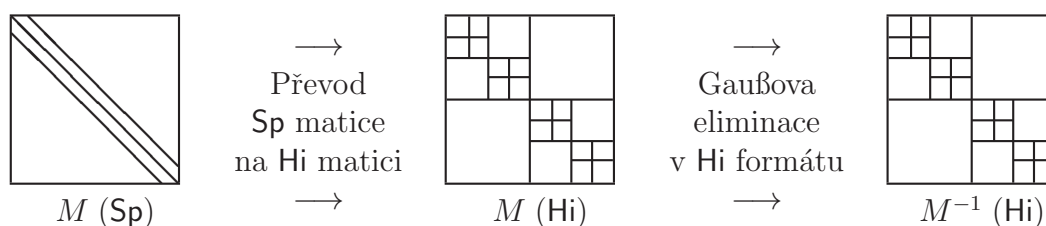


Obrázek 5.1: Dvě hierarchické matice s rozdílnou stromovou strukturou. Pokud společný podstrom odpovídá stejným dělením indexů, lze i takové matice násobit a dokonce i sčítat. Výsledkem operace bude hierarchická matice se strukturou společného podstromu.

Závěr

V této bakalářské práci jsme se snažili případným čtenářům osvětlit úvod do problematiky práce s hierarchickými maticemi. Tyto matice jsou využívány ke strukturovanému ukládání matic pomocí jejího rekurzivního rozdělení do bloků vhodných vlastností. Na příkladu třídiagonální matice T , která je řídká, a s využitím jejích spektrálních vlastností, jsme ukázali, že její inverze T^{-1} má všechny mimodiagonální bloky s omezenou hodnotí, a tedy výborně splňuje předpoklady pro to, abychom ji mohli uložit hierarchicky. Ukázali jsme, že podobný koncept lze aplikovat na velké množství obecných regulárních matic (typicky inverzí řídkých matic, viz [21]).

V praxi, pro obecnou řídkou matici M samozřejmě *nemůžeme* obecně postupovat tak, že bychom nejprve spočítali hustou M^{-1} a pak teprve hledali její hierarchickou reprezentaci. Už ne proto, že bude-li matice velkého řádu, hustá inverze se nemusí vůbec vejít do paměti počítače. V praxi tedy musíme do hierarchického formátu převést už řídkou matici M a její inverzi hledat např. pomocí Gaußovy eliminace, kterou ovšem musíme provozovat také v hierarchickém tvaru, viz obrázek 5.2. Jak takovou Gaußovou eliminaci (LU rozklad) ve skutečnosti provést jsme se v naší práci nezabývali. Je ale zřejmé, že když už umíme dvě hierarchické matice sčítat a násobit, nejsme od popisu hierarchického LU rozkladu daleko. Podrobněji viz např. [13] a [21].



Obrázek 5.2: Řídká (sparse, Sp) matice M a schéma výpočtu její inverze M^{-1} v hierarchickém (Hi) formátu. Díky hierarchickému přístupu je možné spočítat a v počítači uložit i takovou inverzi, která by byla v klasickém formátu neuložitelná z důvodu své velikosti.

Otázkou zůstává, proč bychom vůbec měli hledat inverzi M^{-1} . Důvodem rozhodně není potřeba řešit soustavu rovnic, jak by se možná mohlo zdát (na řešení soustav rovnic existuje řada efektivních metod, viz např. [5], [6], a mnoho dalších). Existují ale úlohy, kde potřebujeme například inverzi (nebo alespoň její aproximaci)

znát explicitně (nebo alespoň v takovém tvaru, abychom s ní mohli dále počítat). Příkladem může být výpočet tzv. sign function iterace, kde počítáme iterace typu

$$M^{(j+1)} = \frac{1}{2} \left(M^{(j)} + (M^{(j)})^{-1} \right),$$

přičemž $M^{(0)} = M$ je daná řídká matice, která je tak velká, že se její obecně hustá inverze nevejde do paměti počítače, viz např. [10].

V této práci jsme podrobněji rozebrali, jak provádět operace sčítání a násobení dvou matic v případě, že jsou tyto matice k dispozici v hierarchickém formátu. Ukázali jsme, že tyto, v hierarchickém formátu poměrně komplikované operace, lze při nejpodrobnějším pohledu rozdělit na méně náročné kroky. U jednotlivých kroků jsme se pokusili, v rámci možností, porovnat jejich náročnost vzhledem ke klasickému přístupu.

Literatura

- [1] M. Bebendorf: *Hierarchical Matrices*, Edice Lecture notes in computational science and engineering (LCNSE) 63, Springer-Verlag, Berlin, Heidelberg, 2008.
doi: [10.1007/978-3-540-77147-0](https://doi.org/10.1007/978-3-540-77147-0)
<http://www.springer.com/la/book/9783540771463>
- [2] S. Börm: *Efficient Numerical Methods for Non-local Operators: \mathcal{H}^2 -Matrix Compression, Algorithms and Analysis*, Edice EMS Tracts in Mathematics 14, European Mathematical Society, Zürich, 2010.
http://www.ems-ph.org/books/book.php?proj_nr=125
- [3] J. Demel: *Grafy a jejich aplikace*, Academia, Praha 2002.
- [4] J. Demel: *Grafy*, SNTL, Státní nakladatelství technické literatury (edice MVŠT, Matematika pro vysoké školy technické, sešit XXXIV), Praha, 1988, resp. 1989 (dotisk).
- [5] E. J. Duintjer Tebbens, I. Hnětynková, M. Plešinger, Z. Strakoš, P. Tichý: *Analýza metod pro maticové výpočty. Základní metody* Matfyzpress, Praha, 2012.
- [6] G. H. Golub, C. F. Van Loan: *Matrix Computations* (Fourth Edition), Johns Hopkins University Press, Baltimore, MD, 2012.
<https://jhupbooks.press.jhu.edu/content/matrix-computations-0>
- [7] W. Hackbusch: *Hierarchische Matrizen: Algorithmen und Analysis*, Springer-Verlag, Berlin, Heidelberg, 2009.
doi: [10.1007/978-3-642-00222-9](https://doi.org/10.1007/978-3-642-00222-9)
<http://www.springer.com/la/book/9783642002212>
- [8] S. Chandrasekaran, P. Dewilde, M. Gu, W. Lyons, T. Pals: *A fast solver for HSS representations via sparse matrices*, SIAM Journal on Matrix Analysis and Applications, Volume 29, Number 1 (2006), pp. 67-81. (15 pages)
doi: [10.1137/050639028](https://doi.org/10.1137/050639028)
<http://epubs.siam.org/doi/abs/10.1137/050639028>
- [9] G. Chartrand: *Introductory graph theory*, Dover Publications, New York, 1985.
- [10] D. Kressner, M. Plešinger, C. Tobler: *A preconditioned low-rank CG method for parameter-dependent Lyapunov matrix equations*, Numerical Linear Algebra with Applications 21(5) (2014), str. 666–684.

- [11] J. Matoušek, J. Nešetřil: *Kapitoly z diskrétní matematiky*, Karolinum, Praha, 2009.
- [12] J. Nešetřil: *Teorie grafů*, SNTL, Státní nakladatelství technické literatury (edice MS, Matematický seminář, č. 13), Praha, 1979.
- [13] S. Pauli: *A numerical solver for Lyapunov equations based on the matrix sign function iteration in HSS arithmetic*, Semester Thesis, SAM, ETH Zürich, 2010.
- [14] J. Sedláček: *Kombinatorika v teorii a praxi. Úvod do teorie grafů*, Academia (edice CV, Cesta k vědění, č. 7), Praha, 1964.
- [15] J. Sedláček: *Úvod do teorie grafů*, Academia (edice CV, Cesta k vědění, č. 25, resp. 29), Praha, 1977 (2. vydání), resp. 1981 (3. vydání).
- [16] R. C. Thompson: *Principal submatrices of normal and Hermitian matrices*, Illinois J. Math., 10 (1966), pp. 296-308
- [17] J. Žáková: *Tenzory a kanonické tenzorové rozklady: Tuckerův rozklad*, bakalářská práce TU v Liberci, 2015.
- [18] HARWELL–BOEING SPARSE MATRIX COLLECTION,
ftp://ftp.numerical.rl.ac.uk/pub/harwell_boeing.
 Sada standardních řídkých matic pro testování algoritmů numerické lineární algebry (původně sestavili Iain Duff, Roger Grimes, John Lewis). Tato sada matic je též k součásti [19] a [20].
- [19] MATRIX MARKET,
<http://math.nist.gov/MatrixMarket>.
 Databáze obsahuje téměř 500 řídkých matic a mnoho nástrojů pro generování matic, jež pocházejí z různých reálných aplikací. Matrix Market je projekt Oddělení matematiky a informatiky (Mathematical and Computational Sciences Division) Laboratoří informatiky (Information Technology Laboratory, ITL) amerického Národního ústavu pro standardy a technologie (National Institute of Standards and Technology, NIST) při ministerstvu obchodu USA.
- [20] THE UNIVERSITY OF FLORIDA SPARSE MATRIX COLLECTION,
<https://sparse.tamu.edu>.
 Sada testovacích matic (sestavili Tim Davis, Yifan Hu (AT&T Research)).
- [21] HIERARCHICAL MATRICES. HLIB PACKAGE,
<http://www.hlib.org>.
 Software pro práci s hierarchickými maticemi.