



TEKNILLINEN KORKEAKOULU
Informaatio- ja luonnontieteiden tiedekunta

Allan Perämäki

Reaalilineaaristen yhtälöryhmien iteratiiviset menetelmät

Diplomityö, joka on jätetty opinnäytteenä tarkastettavaksi diplomi-insinöörin tutkintoa varten teknillisen fysiikan koulutusohjelmassa.

Espoossa, 8.12.2009

Työn valvoja: professori Timo Eirola
Työn ohjaaja: akatemiatutkija Marko Huhtanen

Tekijä:	Allan Perämäki
Koulutusohjelma:	Teknillinen fysiikka
Pääaine:	Matematiikka
Sivuaine:	Tietoliikenneohjelmistot
Työn nimi:	Reaalilineaaristen yhtälöryhmien iteratiiviset menetelmät
Title in English:	Iterative Methods for Real-Linear Systems of Equations
Professuurin koodi ja nimi:	Mat-1 Matematiikka
Työn valvoja:	Professori Timo Eirola
Työn ohjaaja:	Akatemiatutkija Marko Huhtanen
Tiivistelmä:	<p>Työssä esitellään ratkaisumenetelmiä yhtälöille muotoa $Mz + M_{\#}\bar{z} = b$, missä M ja $M_{\#}$ ovat kompleksisia $n \times n$-matriiseja, b annettu vektori ja z ratkaistava vektori. Tällainen yhtälö voitaisiin ratkaista muuttamalla se ensiksi reaaliseksi $2n \times 2n$-kokoiseksi lineaariseksi yhtälöksi $Ax = f$, mutta tässä työssä menetelmät perustuvat alkuperäiseen muotoon.</p> <p>Reaalilineaariset operaattorit ovat kuvauksia $z \mapsto Mz + M_{\#}\bar{z}$. Tässä työssä nämä esitetään taulukkomuodossa matriisien tapaan ja niille annetaan lähteen [1] mukainen reaalilineaarinen LU-hajotelma, Householderin muunnos ja QR-hajotelma. Householderin muunnoksen laskenta-algoritmi annetaan numeerista stabiiliutta parantavalla muutoksella. Lähteessä [1] mainitun LU-hajotelman tukialkion singulaarisuuden ongelman ratkaisuksi ehdotetaan uutta menetelmää.</p> <p>Matriisien Krylovin aliavaruudet ja näihin perustuvia lineaaristen yhtälöryhmien iteratiivisia ratkaisumenetelmiä esitellään lähde [8] noudatellen. Lähteen [1] mukainen reaalilineaarinen GMRES-menetelmä yhtälölle $\kappa z + M_{\#}\bar{z} = b$, missä κ on kompleksiluku, käydään läpi ja lisäksi tapauksille $M_{\#}^T = (-)M_{\#}$ tuodaan esiin Lanczos-tyyppiseen iterointiin perustuvat menetelmät. Ensiksi mainitussa menetelmässä tarvittavan Householderin muunnoksen tilalle tarjotaan uutena mahdollisuutena käyttää reaalilineaarisia Givens-rotatioita. Yleisille reaalilineaarisille yhtälöille ei saada GMRESia, mutta joitakin Galerkinin menetelmään perustuvia iteratiivisia menetelmiä kokeillaan mukaanlukien lähteessä [1] mainittu. Matriisien epätäydelliset LU-hajotelmat (ILU) pohjustusmenetelminä yleistetään reaalilineaarisille operaattoreille. Näistä esimerkkeinä annetaan ILU(0)- ja ILUT-menetelmät.</p> <p>Numeerisia kokeita varten diskretoidaan sähköisessä impedanssitomografiassa esiintyvää $\bar{\partial}$-yhtälöä vastaava integraaliyhtälö päätyen yhtälöön $z + M_{\#}\bar{z} = 1$ lähde [5] noudattaen. MATLABilla suoritettut laskennat osoittavat reaalilineaariset ratkaisumenetelmät tehokkaammiksi kuin GMRES vastaavalle reaaliselle $2n \times 2n$-yhtälöryhmälle. Aikaisemmin ei ole suoraan mainittu edellä olevan yhtälö muuttamista \mathbb{C}-lineaariseen muotoon $(I - M_{\#}\bar{M}_{\#})z = 1 - M_{\#}1$. Matriisin $M_{\#}$ ollessa peräisin edellä mainitun integraaliyhtälön diskretoinnista, on numeeristen kokeiden perusteella ratkaisu tästä muodosta suositeltavaa.</p> <p>Lukijalta edellytetään lineaarialgebran alkeiden hallintaa. Liite A sisältää matriisien osalta käytetyt määritelmät.</p>
Sivumäärä: 84	Avainsanat: reaalilineaarinen, Krylov, iteraatio, Galerkin, GMRES, pohjustus, ILU, impedanssitomografia
Täytetään tiedekunnassa	
Hyväksytty:	Kirjasto:

Author:	Allan Perämäki
Degree Programme:	Engineering Physics
Major subject:	Mathematics
Minor subject:	Data Communications Software
Title:	Iterative Methods for Real-Linear Systems of Equations
Title in Finnish:	Reaalilineaaristen yhtälöryhmien iteratiiviset menetelmät
Chair:	Mat-1 Mathematics
Supervisor:	Professor Timo Eirola
Instructor:	Academy Research Fellow Marko Huhtanen
Abstract:	<p>This thesis presents solution methods for equations of the form $\mathbf{M}\mathbf{z} + \mathbf{M}_{\#}\bar{\mathbf{z}} = \mathbf{b}$, where \mathbf{M} and $\mathbf{M}_{\#}$ are complex $n \times n$-matrices, \mathbf{b} a given vector and \mathbf{z} the vector to be solved. Such an equation could be solved by first rewriting it as a real $2n \times 2n$ system of linear equations $\mathbf{A}\mathbf{x} = \mathbf{f}$, but the methods in this thesis are based on the original formulation.</p> <p>Real-linear operators are mappings $\mathbf{z} \mapsto \mathbf{M}\mathbf{z} + \mathbf{M}_{\#}\bar{\mathbf{z}}$. These can be represented in a table form not unlike matrices and real-linear LU-decomposition, Householder transformation and QR-decomposition are defined following [1]. The Householder transformation is modified to improve numerical stability. A new strategy to avoid the breakdown mentioned in [1] in the real-linear LU-decomposition is suggested.</p> <p>The Krylov subspaces of matrices and some solutions methods based on these are introduced roughly following [8]. The real-linear GMRES method given in [1] for the equation $\kappa\mathbf{z} + \mathbf{M}_{\#}\bar{\mathbf{z}} = \mathbf{b}$, where κ is a complex number, is reconsidered and additionally methods for the cases $\mathbf{M}_{\#}^T = (-)\mathbf{M}_{\#}$ based on Lanczos-type iteration is presented. New real-linear Givens rotations provide an alternative to using Householder transformations during computations of the real-linear GMRES. There's no GMRES for real-linear equations of the general form, but some iterative Galerkin methods are tried including the one mentioned in [1]. The incomplete LU-decompositions (ILU) for matrices as preconditioning methods are generalized to real-linear operators. ILU(0) and ILUT methods are given as examples.</p> <p>Following [5], the integral equation corresponding to the $\bar{\partial}$-equation arising in electrical impedance tomography is discretized resulting in the equation $\mathbf{z} + \mathbf{M}_{\#}\bar{\mathbf{z}} = \mathbf{1}$. Computations with MATLAB provide evidence of the effectiveness for real-linear methods in comparison to the corresponding real $2n \times 2n$ system with GMRES. Previously, there's no direct mention of transforming to the \mathbb{C}-linearized equation $(\mathbf{I} - \mathbf{M}_{\#}\bar{\mathbf{M}}_{\#})\mathbf{z} = \mathbf{1} - \mathbf{M}_{\#}\mathbf{1}$. When the matrix $\mathbf{M}_{\#}$ arises from the discretization of the mentioned integral equation numerical experiments suggest that using this formulation is recommended.</p> <p>The thesis is self-contained requiring only the basics of linear algebra. Appendix A contains the required definitions concerning matrices.</p>
Number of pages: 84	Keywords: real linear, Krylov, iteration, Galerkin, GMRES, preconditioning, ILU, impedance tomography
Faculty fills	
Approved:	Library code:

Sisällys

1	Johdanto	1
1.1	Reaalilineaariset operaattorit	3
1.1.1	Reaalimatriisiesitys	4
1.1.2	Yhdistetty kuvaus ja käänteisoperaattori	6
1.1.3	Liitto-operaattori	7
1.2	QR-hajotelma	8
2	Suora ratkaiseminen ja LU-hajotelma	12
2.1	\mathbb{R} -lineaaristen operaattoreiden LU-hajotelma	12
2.2	Osittaistuenta	14
2.3	Tuenta \mathbb{R} -lineariselle LU-hajotelmalle	16
2.3.1	\mathbb{R} -linearisesti tuetun LU-hajotelman laskeminen	19
3	Iteratiivinen ratkaiseminen	20
3.1	Johdanto	20
3.2	Petrov-Galerkinin menetelmä	21
3.3	Matriisien Krylovin aliavaruudet	23
3.4	GMRES-menetelmä	25
3.4.1	Uudelleenkäynnistys	28
3.5	\mathbb{R} -lineaariset menetelmät	28
3.6	Krylovin aliavaruudet	29
3.7	\mathbb{R} -lineaarinen täyden ortogonalisoinnin menetelmä	33
3.7.1	Menetelmä symmetriselle $M_{\#}$	33
3.7.2	Menetelmä vinosymmetriselle $M_{\#}$	36
3.8	\mathbb{R} -lineaarinen GMRES operaattorille \mathcal{M}_{κ}	38
3.8.1	MINRES symmetriselle $M_{\#}$	42

3.9	\mathbb{C} -linearisoitu yhtälö	42
3.10	Yleisen \mathbb{R} -lineaarisen operaattorin menetelmät	44
4	ILU-pohjustus	47
4.1	Johdanto	47
4.2	Matriisiyhtälöiden ILU-pohjustus	48
4.2.1	ILU-hajotelma ennaltavalitulla nollakuviolla	48
4.2.2	ILU-hajotelma mukautetulla nollakuviolla	49
4.3	\mathbb{R} -lineaaristen yhtälöiden ILU-pohjustus	50
4.3.1	\mathbb{R} -lineaarinen ILU-hajotelma ennaltavalitulla nollakuviolla . .	52
4.3.2	\mathbb{R} -lineaarinen ILU-hajotelma mukautetulla nollakuviolla . . .	53
5	Numeerisia kokeita	56
5.1	GMRES kokeita	56
5.2	GMRES impedanssitomografiassa	57
5.2.1	$\bar{\partial}$ -yhtälön ratkaiseminen	59
5.2.2	Brown-Uhlmann menetelmä	64
5.2.3	Kokeet	65
5.2.4	$\bar{\partial}$ -yhtälö satunnaisella kerroinfunktiolla	65
5.3	Aaltoyhtälö	68
5.4	Epälineaarinen Schrödingerin yhtälö	69
	Yhteenveto	73
	Kirjallisuutta	74
A	Matriisit	76
B	MATLAB-ohjelmalistaukset	78
B.1	Kohta 2.3.1	78
B.2	Kohta 4.2.1	80
B.3	Kohta 4.2.2	81
B.4	Kohta 4.3.1	82
B.5	Algoritmi 4.3.2	83

Luku 1

Johdanto

Matriisien teorian alkuperänä on lineaaristen yhtälöryhmien ratkaisemisessa käytetyt taulukkomuodot. Gaussin eliminaationa tunnettu menetelmä esiintyy jo kolmannelta tai toiselta vuosisadalta eKr peräisin olevassa kiinalaisessa kokoelmassa *Tšiu-tšang suan-su* (Yhdeksän matematiikan taitoja käsittelevää lukua). Japanissa 1700-luvun lopulla Seki Kowa päätyi determinantin käsitteeseen näiden taulukkomuotojen alkeisoperaatioiden pohjalta. Näistä riippumatta samoihin aikoihin Saksassa Gottfried W. Leibniz kehitti myös determinantin, joskin nähtävästi hän käsitteli vain kahden ja kolmen yhtälön ja muuttujan yhtälöryhmiä.

Leibnizin determinantit ja ratkaisutapa kuitenkin vaipuivat unholaan. Hänestä riippumatta sveitsiläinen Gabriel Cramer keksi 1750 hänen nimeään nykyäänkin kantavan yleisen säännön ratkaista lineaarisia yhtälöryhmiä determinanttien avulla. Sääntö teki determinanteista suosittun ratkaisutavan vaikka luultavasti se olikin kehitetty jo tunnetun eliminaatiomenetelmän pohjalta. Tuolloin ei kuitenkaan yleensä edes yritetty ratkaista kovin suuria yhtälöryhmiä.

Carl F. Gaussin taivaanmekaniikkaa käsittelevät julkaisut 1800-luvun alkupuolella sisälsivät pienimmän neliösumman tehtävän ratkaisumenetelmän. Siinä esiintyvän normaaliyhtälön kertoimet muodostavat nykykielellä ilmaisten positiividefiniitin symmetrisen matriisin. Gauss ratkaisi yhtälön eliminointimenetelmällä, joka hyödynsi näitä ominaisuuksia laskutyön nopeuttamiseksi. Esimerkkinä Gauss laskee korjauksen Aurinkokunnan toiseksi suurimman asteroidin Pallaksen rataelementeille, missä ratkaistava normaaliyhtälö sisältää kuusi yhtälöä ja kuusi muuttujaa. Gaussin eliminoinnin nimitys juontaa juurensa näistä julkaisuista.

Matriisi-nimityksen otti käyttöön James J. Sylvester 1850. Hän käytti näitä taulukoi- ta determinanttien teoriassa. Varsinaisen matriisien määritelmän teki Arthur Cayley joitakin vuosia myöhemmin. Hän määritteli matriisien tulon pohtimalla tekijöitä vastaavien lineaarikuvausten yhdistettä. Cayleyn töissä esiintyy kaikki matriisialgebran perusominaisuudet kuten skalaarilla kertoaminen, matriisien yhteenlasku, epäkom- mutoiva tulo, singulaarisuus ja yhteys aikaisempaan determinanttien teoriaan.

1940-luvulle mentäessä olivat erilaiset matriisien hajotelmat (LU, Cholesky) käytös- sä. Mekaanisten pöytälaskukoneiden avulla oli mahdollista ratkaista 10-20 yhtälön ja muuttujan lineaarisia yhtälöryhmiä Gaussin eliminoinnilla. Rajallisen laskentatark- kuuden ja sen pyöristysvirheiden vaikutusta eliminoinnissa ei kuitenkaan tunnettu ja

monet matemaatikot olivatkin pessimistisiä Gaussin eliminoinnin ja suurten yhtälöryhmien suhteen. Aiemmasta pessimistisyydestään huolimatta John von Neumann ja Herman Goldstine sekä Alan Turing julkaisivat 40-luvun lopulla eliminointiprosessin virhearvioita, jotka hälvensivät epäilyksiä pyörästysvirheiden katastrofaalisuudesta. James H. Wilkinson julkaisi 1961 paremmat virhearviot ja otti käyttöön nimitykset osittaistuenta ja täystuenta (näitä menetelmiä oli käytetty jo aiemmin, mm. Wilkinson itse ratkaistessaan pöytälaskukoneella 18 yhtälön ryhmää 1946).

Tietokoneiden kehityksen myötä on Gaussin eliminoinnilla tuhansien yhtälöiden tehtävät nykyään helppo ratkaista. Supertietokoneilla suurimmat tiheiden kerroinmatriisien yhtälöryhmät ovat sisältäneet yli miljoona yhtälöä. Hyvin monissa käytännön ongelmissa kerroinmatriisit kuitenkin koostuvat suurelta osin nolista, minkä johdosta tämän hyödyntämisen mahdollistavat iteratiiviset menetelmät ovat nykyään suosittuja. Iteratiivisia menetelmiä lineaarisen yhtälöryhmän ratkaisemiseksi on esitetty jo 1800-luvulla. Näitä ovat mm. Gauss-Seidelin iteraatio ja Jacobin iteraatio. Tällöin kyseessä ei kuitenkaan ollut matriisin harvuusrakenteen hyödyntäminen, vaan tarkan laskutyön helpottaminen. 1950-luvulla tietokoneilla iteroitiin mm. ylirelaksatiomenetelmällä (SOR), jolla vuonna 1960 kyettiin rutiinomaisesti ratkaisemaan kaksiulotteisia Laplace-tyyppisiä 20000 yhtälön ryhmiä. Yleisten yhtälöryhmien ratkaiseminen Gaussin eliminoinnilla oli tuolloin mahdollista noin 100:lle yhtälölle.

Liittogradienttimenetelmä ja Lanczosin menetelmä syntyivät 1950-luvun alussa ratkaisemaan reaalisien symmetrisen matriisin yhtälöryhmiä. Tarkoilla laskutoimituksilla näillä menetelmillä on ominaisuus, että ratkaisu saavutetaan yhtälöiden lukumäärä vastaavalla iteraatiokierroksella. Tästä syystä näitä pidettiin aluksi mahdollisina suorina menetelminä. Numeeriset kokeet hankalilla yhtälöryhmillä osoittivat ettei edellä mainittu tarkan aritmetiikan ominaisuus pätenyt tietokonearitmetiikassa. Liittogradienttimenetelmä saavutti suuren suosion vasta 70-luvulla, kun sitä alettiin pitäämään aitona iteratiivisena menetelmänä ja se yhdistettiin soveltuvaksi käyttöön harvoille matriiseille. Lisäksi tuolloin tiedettiin, että iteraation suppenemisnopeus riippuu matriisin ominaisarvojakaumasta. Edellä mainitut menetelmät voidaan nähdä Krylovin (venäläisen Aleksei Krylovin vuonna 1931 julkaiseman menetelmän perusteella) aliavaruusmenetelminä, joiden soveltaminen perustuu matriisi-vektoritulojen laskemiseen. Näistä mainittakoon vielä myös epäsymmetrisille matriiseille soveltuva vuonna 1986 Yousef Saadin ja Martin Schultzin esittämä GMRES.

Krylovin aliavaruusmenetelmien kanssa ryhdyttiin 70-luvulla yleisesti käyttämään myös pohjustusmenetelmiä suppenemisen nopeuttamiseksi. Yhtälöryhmän kerroinmatriisin ajatellaan tällöin olevan kerrottu sen likimääräisellä kääntematriisilla, jonka laskemisen on oltava nopeaa. Näistä suosittuja ovat Gaussin eliminaatioon perustuva ILU (incomplete LU) ja sen muunnelmat.

Osittaisdifferentiaaliyhtälöiden ja integraaliyhtälöiden diskretoinnista saadaan usein kerroinmatriiseja, joilla on harvuusrakenteensa tai muun rakenteen johdosta nopea laskea matriisi-vektori tuloja. Täten Krylovin aliavaruusmenetelmät soveltuvat niiden ratkaisuun. Kompleksisten yhtälöiden diskretoinnista saadaan kompleksiluvuisia koostuvia kerroinmatriiseja ja ratkaisuja eikä alun perin reaalille matriiseille kehitettyjen menetelmien suora soveltaminen ole aina järkevää. Liittogradienttimenetelmä voidaan yleistää näille, kunhan matriisi on hermiittinen. Samoin GMRESin voi yleistää suoraviivaisesti kompleksimatriiseille. Kuitenkin on olemassa tapauksia (kompleksinen Helmholtzin yhtälö), joissa kerroinmatriisi on symmetrinen ja

kompleksinen. Tällöin voitaisiin käyttää GMRESia, mutta näille on kehitetty myös symmetristä rakennetta hyödyntävä liittogradienttimenetelmän-tyyppinen menetelmä.

Tämän diplomityön aiheena on matriisien sijaan reaalilineaariset operaattorit ja näiden muodostamat yhtälöryhmät. Tyypiesimerkkinä on sähköisessä impedanssitolmografiassa esiintyvän (kompleksisen) integraaliyhtälön diskretoinnista saatu reaalilineaarinen yhtälö. Seuraavaksi tutustutaan operaattoreihin tarkemmin.

Edellä oleva tiivistelmä lineaaristen yhtälöryhmien ratkaisemisen historiasta on laadittu lähteiden [12], [13], [14] ja [15] pohjalta.

1.1 Reaalilineaariset operaattorit

Lineaarialgebrassa tutkitaan äärellisdimensioisia vektoriavaruuksia ja niiden välisiä lineaarisia kuvauksia. Linearioperaattorilla A kompleksilukukertoimisten vektoriavaruuksien V ja W välillä, $A : V \rightarrow W$, on ominaisuudet $A(v_1 + v_2) = Av_1 + Av_2$ ja $A(\alpha v_1) = \alpha Av_1$ kaikilla vektoreilla $v_1, v_2 \in V$ ja $\alpha \in \mathbb{C}$. Kun avaruuksiin V ja W kiinnitetään kanta, voidaan operaattorille A muodostaa sitä esittävä matriisi.

Tässä työssä käsitellään operaattoreita $\mathcal{M} : V \rightarrow W$, joilla on ominaisuudet

$$\begin{aligned} \mathcal{M}(v_1 + v_2) &= \mathcal{M}(v_1) + \mathcal{M}(v_2), & \text{kaikilla vektoreilla } v_1, v_2 \in V \\ \mathcal{M}(\alpha v) &= \alpha \mathcal{M}(v), & \text{kaikilla vektoreilla } v \in V \text{ ja luvuilla } \alpha \in \mathbb{R}. \end{aligned}$$

Jälkimmäinen ominaisuus vaaditaan nyt vain reaaliluvuille α , mutta avaruudet V ja W ovat edelleen kompleksikertoimisia. Kun avaruuksiin kiinnitetään kannat, voidaan tällainen operaattori esittää kahdella matriisilla. Olkoon v_1, v_2, \dots, v_n avaruuden V kanta ja w_1, w_2, \dots, w_m avaruuden W kanta. Määritellään $m \times n$ -matriisit $A = (a_{jk})$ ja $B = (b_{jk})$ siten, että

$$\sum_{j=1}^m a_{jk} w_j = \frac{1}{2}(\mathcal{M}(v_k) - i\mathcal{M}(iv_k)), \quad \sum_{j=1}^m b_{jk} w_j = \frac{1}{2}(\mathcal{M}(v_k) + i\mathcal{M}(iv_k)).$$

Olkoon sitten $v \in V$ ja $v = \sum_{k=1}^n c_k v_k$. Tällöin

$$\begin{aligned} \mathcal{M}(v) &= \sum_{k=1}^n \mathcal{M}(c_k v_k) = \sum_{k=1}^n \mathcal{M}\left(\frac{1}{2}(c_k + \bar{c}_k)v_k + \frac{1}{2i}(c_k - \bar{c}_k)iv_k\right) \\ &= \sum_{k=1}^n \left(\frac{1}{2}(c_k + \bar{c}_k)\mathcal{M}(v_k) + \frac{1}{2i}(c_k - \bar{c}_k)\mathcal{M}(iv_k)\right) \\ &= \sum_{k=1}^n \left(\frac{1}{2}(\mathcal{M}(v_k) - i\mathcal{M}(iv_k))c_k + \frac{1}{2}(\mathcal{M}(v_k) + i\mathcal{M}(iv_k))\bar{c}_k\right) \\ &= \sum_{k=1}^n \sum_{j=1}^m a_{jk} c_k w_j + \sum_{k=1}^n \sum_{j=1}^m b_{jk} \bar{c}_k w_j = \sum_{j=1}^m \left(\sum_{k=1}^n a_{jk} c_k + \sum_{k=1}^n b_{jk} \bar{c}_k\right) w_j. \end{aligned}$$

Näin ollen $\mathcal{M}(v)$ saadaan laskemalla $Ac + B\bar{c}$, missä c on v :n komponenttien muodostoma pystyvektori. Tämä työ keskittyy numeeriseen laskentaan tällaisilla operaattoreilla, joten abstrakteja vektoriavaruuksia ei jäljempänä enää käsitellä. Asetetaan seuraava määritelmä.

Määritelmä 1.1.1. Olkoot M ja $M_{\#}$ kompleksisia $m \times n$ -matriiseja. Kuvausta

$$\mathcal{M} : \mathbb{C}^n \rightarrow \mathbb{C}^m, \quad \mathcal{M}(z) = Mz + M_{\#}\bar{z}$$

kutsutaan \mathbb{R} -lineaariseksi (reaalilineaariseksi) operaattoriksi. Matriisia M kutsutaan \mathcal{M} :n lineaariseksi osaksi ja matriisia $M_{\#}$ \mathcal{M} :n antilineaariseksi osaksi.

\mathbb{R} -lineaarista operaattoria $\mu : \mathbb{C} \rightarrow \mathbb{C}$ kutsutaan \mathbb{R} -lineaariseksi skalaarioperaattoriksi. Kun $\mathcal{M}(z) = Mz + M_{\#}\bar{z}$ on \mathbb{R} -lineaarinen, niin voidaan määritellä skalaarioperaattorit $\mu_{ij} : \mathbb{C} \rightarrow \mathbb{C}$, $\mu_{ij}(z) = (M)_{ij}z + (M_{\#})_{ij}\bar{z}$. Operaattori \mathcal{M} voidaan silloin ilmaista $m \times n$ -muodossa

$$\mathcal{M} = \begin{bmatrix} \mu_{11} & \mu_{12} & \mu_{13} & \cdots & \mu_{1n} \\ \mu_{21} & \mu_{22} & \mu_{23} & \cdots & \mu_{2n} \\ \mu_{31} & \mu_{32} & \mu_{33} & \cdots & \mu_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mu_{m1} & \mu_{m2} & \mu_{m3} & \cdots & \mu_{mn} \end{bmatrix}, \quad (1.1)$$

$$\mathcal{M}(z) = \begin{bmatrix} \mu_{11}(z_1) + \mu_{12}(z_2) + \cdots + \mu_{1n}(z_n) \\ \mu_{21}(z_1) + \mu_{22}(z_2) + \cdots + \mu_{2n}(z_n) \\ \mu_{31}(z_1) + \mu_{32}(z_2) + \cdots + \mu_{3n}(z_n) \\ \vdots \\ \mu_{m1}(z_1) + \mu_{m2}(z_2) + \cdots + \mu_{mn}(z_n) \end{bmatrix}, \quad (1.2)$$

missä $z = (z_1, z_2, \dots, z_n) \in \mathbb{C}^n$. Tällöin merkitään lyhyesti $\mathcal{M} = (\mu_{ij})$. Näille taulukkumuodoille käytetään vapaasti matriiseista peräisin olevia tuttuja nimityksiä. Esimerkiksi yläkolmio-operaattori tarkoittaa operaattoria, jolle $\mu_{ij} = 0$ kaikilla $i > j$. Käytetään myös nk. MATLAB-notaatiota. Esimerkiksi $\mathcal{M}_{2:4,*}$ on operaattori, joka muodostuu \mathcal{M} :n toisesta, kolmannesta ja neljännestä rivistä. Kun \mathcal{M} on rivi- tai sarakeoperaattori, niin esim. $\mathcal{M}_{1:4}$ on rivi- tai sarakeoperaattori koostuen \mathcal{M} :n neljästä ensimmäisestä skalaarioperaattorista ja \mathcal{M}_2 on \mathcal{M} :n toinen skalaariope-
raattori.

1.1.1 Reaalimatriisiesitys

\mathbb{R} -lineaarinen operaattori $\mathcal{M}(z) = Mz + M_{\#}\bar{z}$ voidaan samaistaa reaalisen $2m \times 2n$ -matriisin kanssa kahdella luonnollisella tavalla. Numeerisessa laskennassa tätä ei välttämättä kannata tehdä, koska matriisien M ja $M_{\#}$ rakenne menetetään. Joitakin reaalilineaaristen operaattoreiden ominaisuuksia tästä kuitenkin nähdään.

Merkitään $z = x + iy \in \mathbb{C}^n$, missä $x = \operatorname{Re}(z)$ on z :n reaaliosa ja $y = \operatorname{Im}(z)$ imaginääriosaa. Samaistetaan sitten $z \in \mathbb{C}^n$ ja vektori $(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) \in \mathbb{R}^{2n}$. Laskemalla saadaan

$$\begin{aligned} \mathcal{M}(z) &= Mz + M_{\#}\bar{z} \\ &= (\operatorname{Re}(M) + i\operatorname{Im}(M))(x + iy) + (\operatorname{Re}(M_{\#}) + i\operatorname{Im}(M_{\#}))(x - iy) \\ &= (\operatorname{Re}(M) + \operatorname{Re}(M_{\#}))x + (-\operatorname{Im}(M) + \operatorname{Im}(M_{\#}))y + \\ &\quad i((\operatorname{Im}(M) + \operatorname{Im}(M_{\#}))x + (\operatorname{Re}(M) - \operatorname{Re}(M_{\#}))y). \end{aligned}$$

Merkitään vektoria $\mathcal{M}(z) = \mathbf{u} + i\mathbf{v}$, missä \mathbf{u} ja \mathbf{v} ovat $\mathcal{M}(z)$:n reaali- ja imaginääriosat. Kun \mathbb{C}^m :n vektorit samaistetaan \mathbb{R}^{2m} :n vektoreiden kanssa vastaavalla tavalla kuin yllä tehtiin \mathbb{C}^n :lle, niin $\mathcal{M}(z) = \mathbf{u} + i\mathbf{v}$ voidaan esittää muodossa

$$\begin{bmatrix} \operatorname{Re}(\mathbf{M}) + \operatorname{Re}(\mathbf{M}_\#) & -\operatorname{Im}(\mathbf{M}) + \operatorname{Im}(\mathbf{M}_\#) \\ \operatorname{Im}(\mathbf{M}) + \operatorname{Im}(\mathbf{M}_\#) & \operatorname{Re}(\mathbf{M}) - \operatorname{Re}(\mathbf{M}_\#) \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix}. \quad (1.3)$$

Tässä esiintyvä $2m \times 2n$ -matriisi on haettu \mathcal{M} :n kanssa samaistettava reaalinen matriisi.

Toinen luonnollinen tapa samaistaa \mathbb{C}^n :n vektorit \mathbb{R}^{2n} :n vektoreiden kanssa on käyttää vektoria

$(x_1, y_1, x_2, y_2, \dots, x_n, y_n) \in \mathbb{R}^{2n}$, kun $z = \mathbf{x} + i\mathbf{y}$. Käytetään vastaavaa järjestystä \mathbb{C}^m :ssä. Kun sitten muodostetaan ensimmäisen tavan mukaisesti reaalinen 2×2 -matriisi \mathbf{m}_{ij} jokaiselle (1.1):n skalaarioperaattorille μ_{ij} , niin operaattoria \mathcal{M} vastaa reaalinen $2m \times 2n$ -matriisi

$$\begin{bmatrix} \mathbf{m}_{11} & \mathbf{m}_{12} & \mathbf{m}_{13} & \dots & \mathbf{m}_{1n} \\ \mathbf{m}_{21} & \mathbf{m}_{22} & \mathbf{m}_{23} & \dots & \mathbf{m}_{2n} \\ \mathbf{m}_{31} & \mathbf{m}_{32} & \mathbf{m}_{33} & \dots & \mathbf{m}_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{m}_{m1} & \mathbf{m}_{m2} & \mathbf{m}_{m3} & \dots & \mathbf{m}_{mn} \end{bmatrix}. \quad (1.4)$$

Esimerkki 1.1.2. Olkoon $\xi_1 : \mathbb{C} \rightarrow \mathbb{C}^3$ 3×1 -operaattori siten, että $\xi_1(z) = (z, 0, 0)$. Kun merkitään nolla- ja yksikköskalaarioperaattoreita $\mathbf{0}(z) = 0$ ja $\mathbf{1}(z) = z$, niin ξ_1 :n 3×1 -muoto ja kohdan (1.4) mukainen matriisi ovat

$$\xi_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \text{ja} \quad \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}.$$

Esimerkki 1.1.3. Olkoot $\mathbf{M} = \begin{bmatrix} 1 + 2i & 3 + 4i \\ 5 + 6i & 7 + 8i \end{bmatrix}$ ja $\mathbf{M}_\# = \begin{bmatrix} 9 + 10i & 11 + 12i \\ 13 + 14i & 15 + 16i \end{bmatrix}$. Kohdan (1.1) mukaiset skalaarioperaattorit ovat

$$\begin{aligned} \mu_{11}(z) &= (1 + 2i)z + (9 + 10i)\bar{z} & \mu_{12}(z) &= (3 + 4i)z + (11 + 12i)\bar{z} \\ \mu_{21}(z) &= (5 + 6i)z + (13 + 14i)\bar{z} & \mu_{22}(z) &= (7 + 8i)z + (15 + 16i)\bar{z}. \end{aligned}$$

Kohdan (1.3) mukainen reaalinen matriisi on

$$\begin{bmatrix} \operatorname{Re}(\mathbf{M}) + \operatorname{Re}(\mathbf{M}_\#) & -\operatorname{Im}(\mathbf{M}) + \operatorname{Im}(\mathbf{M}_\#) \\ \operatorname{Im}(\mathbf{M}) + \operatorname{Im}(\mathbf{M}_\#) & \operatorname{Re}(\mathbf{M}) - \operatorname{Re}(\mathbf{M}_\#) \end{bmatrix} = \begin{bmatrix} 10 & 14 & 8 & 8 \\ 18 & 22 & 8 & 8 \\ 12 & 16 & -8 & -8 \\ 20 & 24 & -8 & -8 \end{bmatrix}$$

ja kohdan (1.4) mukaiset matriisit ovat

$$\begin{aligned} \mathbf{m}_{11} &= \begin{bmatrix} 10 & 8 \\ 12 & -8 \end{bmatrix} & \mathbf{m}_{12} &= \begin{bmatrix} 14 & 8 \\ 16 & -8 \end{bmatrix} & \begin{bmatrix} \mathbf{m}_{11} & \mathbf{m}_{12} \\ \mathbf{m}_{21} & \mathbf{m}_{22} \end{bmatrix} &= \begin{bmatrix} 10 & 8 & 14 & 8 \\ 12 & -8 & 16 & -8 \\ 18 & 8 & 22 & 8 \\ 20 & -8 & 24 & -8 \end{bmatrix}. \end{aligned}$$

1.1.2 Yhdistetty kuvaus ja käänteisoperaattori

\mathbb{R} -linearisista operaattoreista $\mathcal{M} : \mathbb{C}^k \rightarrow \mathbb{C}^m$ ja $\mathcal{N} : \mathbb{C}^n \rightarrow \mathbb{C}^k$ voidaan muodostaa yhdistetty kuvaus

$$\begin{aligned} \mathcal{M} \circ \mathcal{N}(z) &= \mathcal{M}(\mathcal{N}(z)) = M(Nz + N_{\#}\bar{z}) + M_{\#}(\overline{Nz + N_{\#}\bar{z}}) \\ &= (MN + M_{\#}\overline{N_{\#}})z + (MN_{\#} + M_{\#}\overline{N})\bar{z}. \end{aligned}$$

Tämä on \mathbb{R} -lineaarinen ja sitä merkitään myös lyhyemmin $\mathcal{MN}(z) = \mathcal{M} \circ \mathcal{N}(z)$. Kohdan (1.1) mukainen $m \times n$ -skalaarioperaattorimuoto $\mathcal{MN} = (\pi_{ij})$ sille saadaan $\mathcal{M} = (\mu_{ij})$:n ja $\mathcal{N} = (\nu_{ij})$:n muodoista laskemalla

$$\pi_{ij}(z) = \sum_{l=1}^k \mu_{il}(\nu_{lj}(z)) = \sum_{l=1}^k \mu_{il} \circ \nu_{lj}(z) \quad (1 \leq i \leq m, 1 \leq j \leq n).$$

\mathcal{MN} :ää vastaava reaalinen matriisi saadaan \mathcal{M} :ää ja \mathcal{N} :ää vastaavien matriisien tulona.

Määritelmä 1.1.4. *Olkoon $\mathcal{M} : \mathbb{C}^n \rightarrow \mathbb{C}^n$ reaalilineaarinen operaattori. \mathbb{R} -lineaarista operaattoria $\mathcal{N} : \mathbb{C}^n \rightarrow \mathbb{C}^n$ kutsutaan \mathcal{M} :n käänteisoperaattoriksi, jos*

$$\mathcal{M}(\mathcal{N}(z)) = z \quad \text{ja} \quad \mathcal{N}(\mathcal{M}(z)) = z \quad \text{kaikilla } z \in \mathbb{C}^n.$$

Käänteisoperaattoria merkitään $\mathcal{M}^{-1} = \mathcal{N}$. Toisin sanoen, \mathcal{N} on \mathcal{M} :n käänteisoperaattori, jos $\mathcal{MN} = I$ ja $\mathcal{NM} = I$, missä I on yksikköoperaattori (yksikkömatriisi) $Iz = z$ kaikilla $z \in \mathbb{C}^n$.

\mathcal{M} :n käänteisoperaattori \mathcal{M}^{-1} on olemassa täsmälleen silloin kun sitä vastaavalla reaalisella (esim. tavoin (1.3) tai (1.4) muodostetulla) $2n \times 2n$ -matriisilla on käänteismatriisi. Jos tällaiselle reaaliselle matriisille muodostetaan käänteismatriisi, vastaa se operaattoria \mathcal{M}^{-1} . Reaalimatriisisamaistuksen perusteella voidaan myös nähdä, että \mathbb{R} -operaattorilla $\mathcal{M} : \mathbb{C}^n \rightarrow \mathbb{C}^n$ on olemassa käänteisoperaattori mikäli $\mathcal{MN} = I$ tai $\mathcal{NM} = I$ jollakin operaattorilla \mathcal{N} ; toista ehdoista ei siis välttämättä tarvita.

Käänteisoperaattori on yksikäsitteinen, kuten $2n \times 2n$ -reaalimatriisisamaistuksen perusteella voidaan havaita. Voidaan todeta myös suoraan: mikäli \mathcal{N}_1 ja \mathcal{N}_2 ovat kaksi käänteisoperaattoria, niin $\mathcal{N}_1 = \mathcal{N}_1 I = \mathcal{N}_1(\mathcal{MN}_2) = (\mathcal{N}_1 \mathcal{M})\mathcal{N}_2 = I\mathcal{N}_2 = \mathcal{N}_2$.

Esimerkki 1.1.5. *Olkoon $\mu : \mathbb{C} \rightarrow \mathbb{C}$ reaalilineaarinen skalaarioperaattori, $\mu(z) = uz + v\bar{z}$, missä $u, v \in \mathbb{C}$. Tapaa (1.3) vastaava reaalinen matriisi on tällöin*

$$\begin{bmatrix} \operatorname{Re}(u) + \operatorname{Re}(v) & -\operatorname{Im}(u) + \operatorname{Im}(v) \\ \operatorname{Im}(u) + \operatorname{Im}(v) & \operatorname{Re}(u) - \operatorname{Re}(v) \end{bmatrix}.$$

Kun tälle lasketaan käänteismatriisi saadaan

$$\frac{1}{\operatorname{Re}(u)^2 + \operatorname{Im}(u)^2 - \operatorname{Re}(v)^2 - \operatorname{Im}(v)^2} \begin{bmatrix} \operatorname{Re}(u) + \operatorname{Re}(-v) & -(-\operatorname{Im}(u)) + \operatorname{Im}(-v) \\ -\operatorname{Im}(u) + \operatorname{Im}(-v) & \operatorname{Re}(u) - \operatorname{Re}(-v) \end{bmatrix}.$$

Lukemalla tätä ja muotoa (1.3) takaperin, tulee käänteisoperaattori muotoon

$$\mu^{-1}(z) = \frac{1}{|u|^2 - |v|^2}(\bar{u}z - v\bar{z}).$$

Nimittäjässä oleva luku $|u|^2 - |v|^2$ on matriisin determinantti ja siten μ :n käänteisoperaattori on olemassa täsmälleen silloin, kun $|u| \neq |v|$. Voitaisiin myös etsiä suoraan operaattoria $\nu : \mathbb{C} \rightarrow \mathbb{C}$ siten, että $\nu\mu(z) = z$ kaikilla $z \in \mathbb{C}$. Olkoon $\nu(z) = pz + q\bar{z}$ missä $p, q \in \mathbb{C}$. Laskemalla saadaan $\nu\mu(z) = (pu + q\bar{v})z + (pv + q\bar{u})\bar{z}$. Yhtälöparista $pu + q\bar{v} = 1$, $pv + q\bar{u} = 0$ saadaan ratkaistua $p = \bar{u}/(|u|^2 - |v|^2)$, $q = -v/(|u|^2 - |v|^2)$.

1.1.3 Liitto-operaattori

Kun kohdan 1.3 matriisi transponoidaan, on tuloksena

$$\begin{bmatrix} \operatorname{Re}(M^T) + \operatorname{Re}(M_{\#}^T) & -(-\operatorname{Im}(M^T)) + \operatorname{Im}(M_{\#}^T) \\ -\operatorname{Im}(M^T) + \operatorname{Im}(M_{\#}^T) & \operatorname{Re}(M^T) - \operatorname{Re}(M_{\#}^T) \end{bmatrix}.$$

Kulkemalla reaalmatriisiesityksestä takaisin operaattoriksi, saadaan

$$\mathcal{M}^*(z) = M^*z + M_{\#}^T\bar{z}$$

ja sitä kutsutaan operaattorin \mathcal{M} *liitto-operaattoriksi*. Täten erityisesti skalaarioperaattorin $\mu(z) = uz + v\bar{z}$ liitto-operaattori on $\mu^*(z) = \bar{u}z + v\bar{z}$, jolloin skalaarimuodossa (1.1) ilmaisten

$$\mathcal{M}^* = \begin{bmatrix} \mu_{11}^* & \mu_{21}^* & \mu_{31}^* & \cdots & \mu_{m1}^* \\ \mu_{12}^* & \mu_{22}^* & \mu_{32}^* & \cdots & \mu_{m2}^* \\ \mu_{13}^* & \mu_{23}^* & \mu_{33}^* & \cdots & \mu_{m3}^* \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mu_{1n}^* & \mu_{2n}^* & \mu_{3n}^* & \cdots & \mu_{mn}^* \end{bmatrix}.$$

Kahden operaattorin yhdisteelle pätee $(\mathcal{M}\mathcal{N})^* = \mathcal{N}^*\mathcal{M}^*$. Liitto-operaattorista hieman poiketen, operaattorin \mathcal{M} *transpoosi* on

$$\mathcal{M}^T(z) = M^Tz + M_{\#}^T\bar{z}.$$

Tällöin

$$\mathcal{M}^T = \begin{bmatrix} \mu_{11} & \mu_{21} & \mu_{31} & \cdots & \mu_{m1} \\ \mu_{12} & \mu_{22} & \mu_{32} & \cdots & \mu_{m2} \\ \mu_{13} & \mu_{23} & \mu_{33} & \cdots & \mu_{m3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mu_{1n} & \mu_{2n} & \mu_{3n} & \cdots & \mu_{mn} \end{bmatrix}.$$

Yleensä kahden operaattorin yhdisteelle $(\mathcal{M}\mathcal{N})^T \neq \mathcal{N}^T\mathcal{M}^T$.

Määritelmä 1.1.6. *Reaalilineaarinen operaattori $\mathcal{Q} : \mathbb{C}^n \rightarrow \mathbb{C}^m$ on isometria, jos*

$$\|\mathcal{Q}(z)\| = \|z\| \quad \text{kaikilla } z \in \mathbb{C}^n.$$

Sitä kutsutaan unitaariseksi, jos $m = n$.

Lause 1.1.7. *Reaalilineaarinen operaattori $\mathcal{U} : \mathbb{C}^n \rightarrow \mathbb{C}^n$ on unitaarinen jos ja vain jos $\mathcal{U}^*\mathcal{U} = I$.*

Todistus. Olkoon matriisi $\mathbf{A} \in \mathbb{R}^{2n \times 2n}$ unitaarioperaattorin $\mathbf{U} : \mathbb{C}^n \rightarrow \mathbb{C}^n$ reaali-matriisiesitys. Tällöin $\|\mathbf{A}\mathbf{x}\| = \|\mathbf{x}\|$ kaikilla $\mathbf{x} \in \mathbb{R}^{2n}$. Kun merkitään avaruuden \mathbb{R}^{2n} sisätuloa $(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$, niin kaikilla $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{2n}$

$$\begin{aligned} \|\mathbf{A}(\mathbf{x} + \mathbf{y})\|^2 = \|\mathbf{x} + \mathbf{y}\|^2 &\Leftrightarrow (\mathbf{A}(\mathbf{x} + \mathbf{y}), \mathbf{A}(\mathbf{x} + \mathbf{y})) = (\mathbf{x} + \mathbf{y}, \mathbf{x} + \mathbf{y}) \\ &\Leftrightarrow (\mathbf{A}\mathbf{x}, \mathbf{A}\mathbf{y}) = (\mathbf{x}, \mathbf{y}) \Leftrightarrow (\mathbf{A}^T \mathbf{A}\mathbf{x}, \mathbf{y}) = (\mathbf{x}, \mathbf{y}) \\ &\Leftrightarrow (\mathbf{A}^T \mathbf{A}\mathbf{x} - \mathbf{x}, \mathbf{y}) = 0. \end{aligned}$$

Valitsemalla $\mathbf{y} = \mathbf{A}^T \mathbf{A}\mathbf{x} - \mathbf{x}$, todetaan $\mathbf{A}^T \mathbf{A}\mathbf{x} = \mathbf{x}$ kaikilla $\mathbf{x} \in \mathbb{R}^{2n}$. Täten $\mathbf{A}^T \mathbf{A} = \mathbf{I}$.

Toisaalta, jos $\mathbf{A}^T \mathbf{A} = \mathbf{I}$, niin kaikilla \mathbf{x}

$$(\mathbf{A}^T \mathbf{A}\mathbf{x}, \mathbf{x}) = (\mathbf{x}, \mathbf{x}) \Leftrightarrow (\mathbf{A}\mathbf{x}, \mathbf{A}\mathbf{x}) = (\mathbf{x}, \mathbf{x}) \Leftrightarrow \|\mathbf{A}\mathbf{x}\|^2 = \|\mathbf{x}\|^2.$$

□

1.2 QR-hajotelma

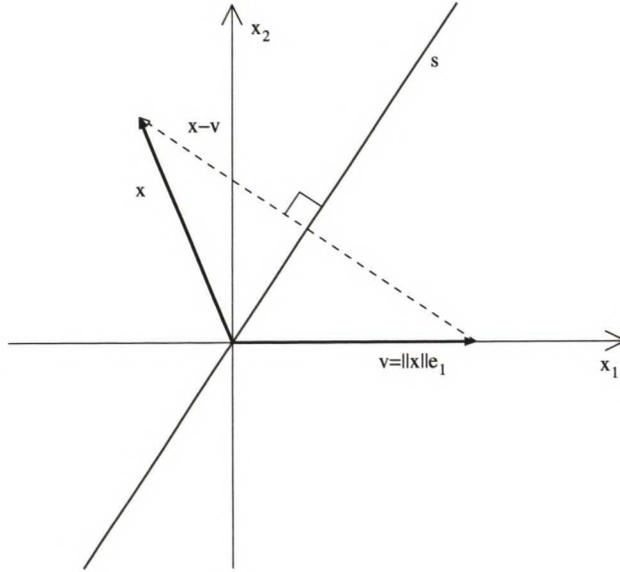
Muun muassa matriisien QR-hajotelman laskemista varten tarvitaan unitaarimatriisia, jolla kertominen kääntää annetun vektorin \mathbf{x} luonnollisen kantavektorin \mathbf{e}_1 suuntaiseksi. Tätä varten tarkastellaan Householderin peilausmatriisia $\mathbf{H} = \mathbf{I} - 2\mathbf{w}\mathbf{w}^*$, missä $\mathbf{w} \in \mathbb{C}^n$ on yksikkövektori. Tämä on selvästi hermiittinen ja unitaarinen $\mathbf{H}^* = \mathbf{H}$, $\mathbf{H}^* \mathbf{H} = \mathbf{H}^2 = \mathbf{I}$. Vektori \mathbf{w} on peilaavan tason normaalivektori ja lausekkeesta $\mathbf{H}\mathbf{x} = \mathbf{x} - 2\mathbf{w}(\mathbf{w}^* \mathbf{x})$ nähdään, että \mathbf{x} :n normaalivektorin suuntainen komponentti vähennetään kaksi kertaa. Siten \mathbf{H} :lla kertominen siirtää \mathbf{x} :n peilaavan tason toiselle puolelle samaan asemaan.

Kuva 1.1 havainnollistaa reaalisen vektorin $\mathbf{x} \in \mathbb{R}^n$ peilaamisen vektoriksi $\mathbf{v} = \|\mathbf{x}\|\mathbf{e}_1$. Peilaavaa tasoa esittää suora s ja nähdään, että normaalivektori on $\mathbf{w} = (\mathbf{x} - \mathbf{v})/\|\mathbf{x} - \mathbf{v}\|$. Kun vektori \mathbf{x} on lähes \mathbf{e}_1 :den suuntainen, on tässä kaavassa nimittäjä lähellä nollaa. Numeerisessa laskennassa kannattaakin tällöin peilata \mathbf{x} vektorin $-\mathbf{e}_1$ suuntaan. Valinta voidaan perustaa \mathbf{x} :n ensimmäisen komponentin etumerkkiin, jolloin saadaan kaava

$$\mathbf{w} = \frac{\mathbf{x} + \alpha\|\mathbf{x}\|\mathbf{e}_1}{\|\mathbf{x} + \alpha\|\mathbf{x}\|\mathbf{e}_1\|}, \quad \alpha = \begin{cases} \frac{\langle \mathbf{x}, \mathbf{e}_1 \rangle}{|\langle \mathbf{x}, \mathbf{e}_1 \rangle|}, & \text{kun } \langle \mathbf{x}, \mathbf{e}_1 \rangle \neq 0, \\ 1, & \text{kun } \langle \mathbf{x}, \mathbf{e}_1 \rangle = 0. \end{cases} \quad (1.5)$$

Kompleksivektorin $\mathbf{x} \in \mathbb{C}^n$ tapauksessa ei voida luottaa geometriseen intuitioon ja joudutaan laskemaan. Yritetään peilata \mathbf{x} vektoriksi \mathbf{v} , jolle $\|\mathbf{v}\| = \|\mathbf{x}\|$ ja $\mathbf{v} \neq \mathbf{x}$. Kokeillaan yllä olevan perusteella $\mathbf{w} = (\mathbf{x} - \mathbf{v})/\|\mathbf{x} - \mathbf{v}\|$, jolloin

$$\begin{aligned} \mathbf{H}\mathbf{x} &= \mathbf{x} - 2 \frac{(\mathbf{x} - \mathbf{v})(\mathbf{x} - \mathbf{v})^*}{\|\mathbf{x} - \mathbf{v}\|^2} \mathbf{x} = \mathbf{v} + (\mathbf{x} - \mathbf{v})(\mathbf{x} - \mathbf{v})^* \frac{(\mathbf{x} - \mathbf{v}) - 2\mathbf{x}}{\|\mathbf{x} - \mathbf{v}\|^2} \\ &= \mathbf{v} - (\mathbf{x} - \mathbf{v}) \frac{(\mathbf{x}^* - \mathbf{v}^*)(\mathbf{x} + \mathbf{v})}{\|\mathbf{x} - \mathbf{v}\|^2} = \mathbf{v} - (\mathbf{x} - \mathbf{v}) \frac{\|\mathbf{x}\|^2 + \mathbf{x}^* \mathbf{v} - \mathbf{v}^* \mathbf{x} - \|\mathbf{v}\|^2}{\|\mathbf{x} - \mathbf{v}\|^2} \\ &= \mathbf{v} + (\mathbf{x} - \mathbf{v}) \frac{2i \operatorname{Im}(\mathbf{v}^* \mathbf{x})}{\|\mathbf{x} - \mathbf{v}\|^2}. \end{aligned}$$



Kuva 1.1: Vektorin \mathbf{x} heijastaminen luonnollisen kantavektorin \mathbf{e}_1 suuntaiseksi. Suora s esittää heijastavaa tasoa.

Täytyy siis lisäksi vaatia $\text{Im} \langle \mathbf{x}, \mathbf{v} \rangle = 0$. Näin onkin, kun $\mathbf{v} = -\alpha \|\mathbf{x}\| \mathbf{e}_1$, missä α on kaavasta (1.5). Täten kaava (1.5) on pätevä myös kompleksivektoreille. Kun vektori $\mathbf{x} \in \mathbb{C}^n$ on annettu, laskemalla \mathbf{w} tästä kaavasta, pätee $\mathbf{H}\mathbf{x} = -\alpha \|\mathbf{x}\| \mathbf{e}_1$.

Siirrytään sitten matriiseista reaalilineaarisiin operaattoreihin. Tällöin tarvitaan unitaarioperaattoria, jolla operointi kääntää *kaksi* annettua vektoria kantavektorin \mathbf{e}_1 suuntaiseksi. Olkoot annetut vektorit $\mathbf{x}, \mathbf{y} \in \mathbb{C}^n$ ja oletetaan ne lineaarisesti riippumattomiksi. Lineaarisesti riippuville \mathbf{x}, \mathbf{y} voidaan käyttää edellä kuvattua tavallista Householderin muunnosta.

Operaattoria $\mathcal{H}(z) = z - \mathbf{U}\mathbf{U}^*z - \mathbf{U}\mathbf{U}^T\bar{z}$, missä $\mathbf{U} \in \mathbb{C}^{n \times 2}$ ja $\text{Re}(\mathbf{U}^*\mathbf{U}) = \mathbf{I}$ kutsutaan reaalilineaariseksi Householderin muunnokseksi. Laskemalla nähdään, että $\mathcal{H}^* = \mathcal{H}$ ja $\mathcal{H}^2 = \mathbf{I}$, joten erityisesti \mathcal{H} on unitaarinen. Lähdetään seuraavaksi etsimään matriisia \mathbf{U} , jolla \mathbf{x} ja \mathbf{y} kääntyvät vektorin \mathbf{e}_1 suuntaan. Merkitään $\mathbf{V} = [\mathbf{x} \ \mathbf{y}]$, $\mathbf{e} = \mathbf{e}_1$ ja $\mathbf{a} = [a_1 \ a_2]^T \in \mathbb{C}^{2 \times 1}$, jolloin ehdot $\mathcal{H}(\mathbf{x}) = a_1 \mathbf{e}$ ja $\mathcal{H}(\mathbf{y}) = a_2 \mathbf{e}$ voidaan kirjoittaa yhtäpitävästi matriisimuotoon

$$\mathbf{V} - \mathbf{U}\mathbf{U}^*\mathbf{V} - \mathbf{U}\mathbf{U}^T\bar{\mathbf{V}} = \mathbf{V} - 2\mathbf{U} \text{Re}(\mathbf{U}^*\mathbf{V}) = \mathbf{e}\mathbf{a}^*. \quad (1.6)$$

Jotta ratkaisu \mathbf{U} voisi olla olemassa, täytyy jollakin matriisilla $\mathbf{R} \in \mathbb{R}^{2 \times 2}$ olla $\mathbf{U}\mathbf{R} = (\mathbf{V} - \mathbf{e}\mathbf{a}^*)$. Tällöin

$$\begin{aligned} \mathbf{V} - 2\mathbf{U} \text{Re}(\mathbf{U}^*\mathbf{V}) &= \mathbf{e}\mathbf{a}^* + \mathbf{U}\mathbf{R} - 2\mathbf{U} \text{Re}(\mathbf{U}^*\mathbf{V}) = \mathbf{e}\mathbf{a}^* + \mathbf{U} \text{Re}(\mathbf{R} - 2\mathbf{U}^*\mathbf{V}) \\ &= \mathbf{e}\mathbf{a}^* + \mathbf{U} \text{Re}(\mathbf{U}^*\mathbf{U}\mathbf{R} - 2\mathbf{U}^*\mathbf{V}) = \mathbf{e}\mathbf{a}^* - \mathbf{U} \text{Re}(\mathbf{U}^*(\mathbf{V} + \mathbf{e}\mathbf{a}^*)). \end{aligned} \quad (1.7)$$

Jos löydetään \mathbf{U} , kääntyvä \mathbf{R} ja vektori \mathbf{a} siten, että

$$\mathbf{U}\mathbf{R} = \mathbf{V} - \mathbf{e}\mathbf{a}^* \quad \text{ja} \quad \text{Re}((\mathbf{V} - \mathbf{e}\mathbf{a}^*)^*(\mathbf{V} + \mathbf{e}\mathbf{a}^*)) = 0, \quad (1.8)$$

niin $\text{Re}(\mathbf{U}^*(\mathbf{V} + \mathbf{e}\mathbf{a}^*)) = (\mathbf{R}^{-1})^* \text{Re}((\mathbf{V} - \mathbf{e}\mathbf{a}^*)^*(\mathbf{V} + \mathbf{e}\mathbf{a}^*)) = 0$, joten kaavan (1.7) mukaan \mathbf{U} on sopiva ratkaisu. Ryhdytään sitten etsimään ehdot (1.8) toteuttavia \mathbf{U} ,

R ja a . Merkitään $w = V^*e$, jolloin jälkimmäinen ehto

$$\operatorname{Re}(V^*V + wa^* - aw^* - aa^*) = 0.$$

Merkitään $c = \operatorname{Re}(V^*V)^{1/2} \begin{bmatrix} 1 \\ i \end{bmatrix}$, jolloin $\operatorname{Re}(V^*V - cc^*) = 0$. Kun $a = \eta c$, missä $\eta \in \mathbb{C}$, $|\eta| = 1$, niin myös $\operatorname{Re}(V^*V - aa^*) = 0$. Etsitään vielä η siten, että $\operatorname{Re}(wa^* - aw^*) = 0$, mikä yhtäpitävästi kirjoitettuna on $\operatorname{Re}(w_1\bar{a}_2 - a_1\bar{w}_2) = 0$ ja

$$\operatorname{Re}(w_1\bar{\eta}c_2 - \eta c_1\bar{w}_2) = 0.$$

Edelleen yhtäpitävästi tämä on $\eta q + \bar{\eta}\bar{q} = 0$, missä $q = \bar{w}_1c_2 - \bar{w}_2c_1$. Kun $q \neq 0$, niin ratkaisut ovat $\eta = \pm i \frac{\bar{q}}{|q|}$. Seuraavan helposti todistettavan lemmän mukaan ainakin toisella näistä matriisin $(V - ea^*)$ sarakkeet ovat lineaarisesti riippumattomat.

Lemma 1.2.1. *Olko x ja y lineaarisesti riippumattomia vektoreita. Kun z on vektori ja α_1, α_2 skalaareita, niin ainakin toinen joukoista $\{x - \alpha_1z, y - \alpha_2z\}$ ja $\{x + \alpha_1z, y + \alpha_2z\}$ on lineaarisesti riippumaton.*

Lopuksi vielä ortonormalisoidaan matriisin $(V - ea^*)$ sarakkeet reaalisen sisätulon $(u, v) = \operatorname{Re}(v^*u)$ suhteen. Yhteenvedona edellä oleva voidaan kirjoittaa seuraaviksi kaavoiksi U :n laskemiseksi.

$$\begin{aligned} c &= \operatorname{Re}(V^*V)^{1/2} \begin{bmatrix} 1 \\ i \end{bmatrix}, \quad q = \langle VJc, e \rangle \quad (J = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}), \\ \alpha &= \begin{cases} \frac{q}{|q|}, & \text{kun } q \neq 0, \\ 1, & \text{kun } q = 0. \end{cases} \\ W_- &= V - i\alpha ec^*, \quad W_+ = V + i\alpha ec^*, \\ W &= \begin{cases} W_+, & \text{kun } \det(\operatorname{Re}(W_+^*W_+)) \geq \det(\operatorname{Re}(W_-^*W_-)), \\ W_-, & \text{kun } \det(\operatorname{Re}(W_+^*W_+)) < \det(\operatorname{Re}(W_-^*W_-)). \end{cases} \\ QR &= \begin{bmatrix} \operatorname{Re}(W) \\ \operatorname{Im}(W) \end{bmatrix} \quad (\text{lasketaan suppea QR-hajotelma}), \\ U &= [I_n \quad 0] Q + i [0 \quad I_n] Q. \end{aligned}$$

Huomautus 1.2.2. Matriisin A sarakkeet ovat lineaarisesti riippumattomat \mathbb{C} :n yli, jos ja vain jos A^*A on positiividefiniitti. Sarakkeet ovat lineaarisesti riippumattomat \mathbb{R} :n yli, jos ja vain jos $\operatorname{Re}(A^*A)$ on positiividefiniitti.

Huomautus 1.2.3. Valinta $W = W_+$ vaikuttaa numeeristen kokeiden perusteella toteutuvan aina ja on analoginen valinta kaavan (1.5) kanssa. Todistus tälle seikalle kuitenkin puuttuu.

Reaalilineaarisen operaattorin \mathcal{M} QR-hajotelman muodostamista varten tarvitaan muunnoksia, joilla muodossa (1.1) ilmaistun operaattorin sarakkeita saadaan ensimmäistä alkioita lukuunottamatta nolllaksi. Toisin sanoen, jos $\mu_i(z) = u_i z + v_i \bar{z}$ ($i = 1, 2, \dots, n$) ovat skalaarioperaattoreita, niin halutaan löytää reaalilineaarinen Householderin muunnos, jolla

$$\mathcal{H} \circ \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} = \begin{bmatrix} \# \\ 0 \\ \vdots \\ 0 \end{bmatrix}. \quad (1.9)$$

1.2. QR-HAJOTELMA

Merkitään $\mathbf{u} = [u_1 \ \cdots \ u_n]^T$ ja $\mathbf{v} = [v_1 \ \cdots \ v_n]^T$, jolloin

$$\begin{aligned}\mathcal{H}(\mathbf{u}z + \mathbf{v}\bar{z}) &= \mathcal{H}(\mathbf{u}z) + \mathcal{H}(\mathbf{v}\bar{z}) \\ &= \operatorname{Re}(z)\mathcal{H}(\mathbf{u}) + \operatorname{Im}(z)\mathcal{H}(i\mathbf{u}) + \operatorname{Re}(z)\mathcal{H}(\mathbf{v}) + \operatorname{Im}(z)\mathcal{H}(-i\mathbf{v}) \\ &= \operatorname{Re}(z)\mathcal{H}(\mathbf{u} + \mathbf{v}) + \operatorname{Im}(z)\mathcal{H}(i(\mathbf{u} - \mathbf{v})).\end{aligned}$$

Ehto (1.9) siis toteutuu, kun \mathcal{H} :ksi valitaan Householderin muunnos, jolla operointi kääntää vektorit $(\mathbf{u} + \mathbf{v})$ ja $i(\mathbf{u} - \mathbf{v})$ kantavektorin \mathbf{e}_1 suuntaisiksi.

Lause 1.2.4 (QR-hajotelma). *Olkkoon $\mathcal{M} : \mathbb{C}^n \rightarrow \mathbb{C}^m$ reaalilineaarinen operaattori ja $m \geq n$. Tällöin on olemassa unitaarinen $\mathcal{Q} : \mathbb{C}^m \rightarrow \mathbb{C}^m$ ja yläkolmio-operaattori $\mathcal{R} : \mathbb{C}^n \rightarrow \mathbb{C}^m$ siten, että*

$$\mathcal{M} = \mathcal{Q}\mathcal{R}.$$

Todistus. Olkkoon operaattori \mathcal{M} ilmaistu muodossa (1.1), missä $m \geq n$. Valitaan Householderin muunnos $\mathcal{H}_1 : \mathbb{C}^m \rightarrow \mathbb{C}^m$ siten, että \mathcal{M} :n ensimmäiselle sarakkeelle pätee (1.9). Tällöin

$$\mathcal{H}_1\mathcal{M} = \begin{bmatrix} \mu_{11}^{(1)} & \mu_{12}^{(1)} & \mu_{13}^{(1)} & \cdots & \mu_{1n}^{(1)} \\ 0 & \mu_{22}^{(1)} & \mu_{23}^{(1)} & \cdots & \mu_{2n}^{(1)} \\ 0 & \mu_{32}^{(1)} & \mu_{33}^{(1)} & \cdots & \mu_{3n}^{(1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \mu_{m2}^{(1)} & \mu_{m3}^{(1)} & \cdots & \mu_{mn}^{(1)} \end{bmatrix}.$$

Valitaan seuraavaksi muunnos $\tilde{\mathcal{H}}_2 : \mathbb{C}^{m-1} \rightarrow \mathbb{C}^{m-1}$ siten, että (1.9) pätee toisen sarakkeen $(m-1)$:lle jälkimmäiselle alkion $[\mu_{22}^{(1)} \ \mu_{32}^{(1)} \ \cdots \ \mu_{m2}^{(1)}]^T$. Asetetaan $\mathcal{H}_2 = \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \tilde{\mathcal{H}}_2 \end{bmatrix}$, missä $\mathbf{1}$ on yksikköskalaarioperaattori $\mathbf{1}(z) = z$. Saadaan

$$\mathcal{H}_2\mathcal{H}_1\mathcal{M} = \begin{bmatrix} \mu_{11}^{(1)} & \mu_{12}^{(1)} & \mu_{13}^{(1)} & \cdots & \mu_{1n}^{(1)} \\ 0 & \mu_{22}^{(2)} & \mu_{23}^{(2)} & \cdots & \mu_{2n}^{(2)} \\ 0 & 0 & \mu_{33}^{(2)} & \cdots & \mu_{3n}^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \mu_{m3}^{(2)} & \cdots & \mu_{mn}^{(2)} \end{bmatrix}.$$

Jatketaan tällä tavoin kunnes päädytään yläkolmiomuotoon

$$\mathcal{H}_n \cdots \mathcal{H}_2\mathcal{H}_1\mathcal{M} = \begin{bmatrix} \mu_{11}^{(1)} & \mu_{12}^{(1)} & \mu_{13}^{(1)} & \cdots & \mu_{1n}^{(1)} \\ 0 & \mu_{22}^{(2)} & \mu_{23}^{(2)} & \cdots & \mu_{2n}^{(2)} \\ 0 & 0 & \mu_{33}^{(3)} & \cdots & \mu_{3n}^{(3)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \mu_{nn}^{(n)} \\ 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix}.$$

Merkitään oikeanpuolimmaista operaattoria \mathcal{R} :llä ja asetetaan $\mathcal{Q} = \mathcal{H}_1^*\mathcal{H}_2^*\cdots\mathcal{H}_n^*$, jolloin $\mathcal{M} = \mathcal{Q}\mathcal{R}$. \square

Luku 2

Suora ratkaiseminen ja LU-hajotelma

Yhtälön suoralla ratkaisumenetelmällä tarkoitetaan sellaista algoritmia, joka ideaalisessa aritmetiikassa tuottaa yhtälön tarkan ratkaisun etukäteen tunnetulla äärellisellä askelmäärällä. Käytännölliset lineaaristen yhtälöiden suorat menetelmät perustuvat Gaussin eliminointiin ja niitä käytetään, kun kerroinmatriisi on tiheä eikä sillä ole erityistä hyödynnettävää rakennetta. Isokokoiset ja suuren määrän nollia sisältävät (ns. harvat) matriisit soveltuvat paremmin seuraavassa luvussa esiteltävin menetelmin ratkaistaviksi. On olemassa myös harvoille matriiseille soveltuvia Gaussin eliminointiin pohjautuvia menetelmiä, mutta niitä ei tässä työssä käsitellä.

Tässä luvussa esitellään suoria ratkaisumenetelmiä reaalilineaarille yhtälöryhmille perustuen Gaussin eliminointiin.

2.1 \mathbb{R} -lineaaristen operaattoreiden LU-hajotelma

Olkoon $\mathcal{M} : \mathbb{C}^n \rightarrow \mathbb{C}^n$ reaalilineaarinen operaattori ja merkitään $\mathcal{M} = (\mu_{ij})_{i,j=1}^n$, missä $\mu_{ij} : \mathbb{C} \rightarrow \mathbb{C}$ ovat reaalilineaarisia skaalarioperaattoreita tai samaistuksen kautta reaalisia 2×2 -matriiseja, jolloin koko operaattorin \mathcal{M} voidaan ajatellaan olevan reaalinen $2n \times 2n$ -matriisi lohkokottuna $n \times n$ määrään 2×2 osia. \mathcal{M} oletetaan kääntyväksi tässä luvussa. Tavoitteena on laskea alakolmio-operaattori \mathcal{L} , jonka lävistäjällä olevat skaalarioperaattorit ovat yksikköoperaattoreita, ja yläkolmio-operaattori \mathcal{U} siten, että $\mathcal{L}\mathcal{U} = \mathcal{M}$.

Kyseistä hajotelmaa kutsutaan, sikäli kun se on olemassa, \mathcal{M} :n LU-hajotelmaksi. Sen avulla voidaan ratkaista lineaarinen yhtälöryhmä

$$\mathcal{M}(z) = b, \tag{2.1}$$

missä $b \in \mathbb{C}^n$ on annettu vektori ja ratkaistava vektori on $z \in \mathbb{C}^n$. Jos \mathcal{M} :llä on LU-hajotelma, niin tällainen yhtälöryhmä voidaan ratkaista seuraavasti. Ensin ratkaistaan w yhtälöryhmästä $\mathcal{L}(w) = b$ eteenpäin sijoittamalla eli ratkaisemalla w :n alkioit järjestyksessä ensimmäisestä aloittaen. Tämän jälkeen ratkaistaan z yhtälöryhmästä $\mathcal{U}(z) = w$ taaksepäin sijoittamalla eli z :n viimeisestä alkioista aloittaen.

Ratkaisussa tarvitaan lävistäjällä olevien skalaarioperaattoreiden käänteisoperaattoreita, ks. esimerkki 1.1.5. Tällöin $\mathcal{M}(z) = \mathcal{L}(\mathbf{U}(z)) = \mathcal{L}(\mathbf{w}) = \mathbf{b}$, joten z ratkaisee yhtälöryhmän (2.1).

LU-hajotelman muodostaminen on kannattavaa esimerkiksi silloin, kun lineaarinen yhtälöryhmä joudutaan ratkaisemaan useaan kertaan eri vektoreilla \mathbf{b} , mutta samalla operaattorilla \mathcal{M} . Tällöin tarvitaan vähemmän laskutoimituksia, kun ensin muodostetaan \mathcal{M} :n LU-hajotelma, kuin ratkaistaessa jokainen yhtälöryhmä erikseen esim. Gaussin eliminoinneilla.

Esitetään seuraavaksi tapa muodostaa LU-hajotelma. Suoritetaan reaali-linearisia Gaussin eliminaatioita kertomalla operaattori $\mathcal{M} = (\mu_{ij})$ sopivalla operaattorilla \mathcal{L}_1 seuraavasti. Olettaen, että \mathcal{M} :n ensimmäinen tukialkio μ_{11} on kääntyvä saadaan

$$\begin{aligned} \mathcal{L}_1 \mathcal{M} &= \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ -\mu_{21}\mu_{11}^{-1} & 1 & 0 & \cdots & 0 \\ -\mu_{31}\mu_{11}^{-1} & 0 & 1 & & \\ \vdots & \vdots & & \ddots & \\ -\mu_{n1}\mu_{11}^{-1} & 0 & & & 1 \end{bmatrix} \begin{bmatrix} \mu_{11} & \mu_{12} & \mu_{13} & \cdots & \mu_{1n} \\ \mu_{21} & \mu_{22} & \mu_{23} & \cdots & \mu_{2n} \\ \mu_{31} & \mu_{32} & \mu_{33} & \cdots & \mu_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mu_{n1} & \mu_{n2} & \mu_{n3} & \cdots & \mu_{nn} \end{bmatrix} = \\ &= \begin{bmatrix} \mu_{11} & \mu_{12} & \mu_{13} & \cdots & \mu_{1n} \\ 0 & \mu_{22} - \mu_{21}\mu_{11}^{-1}\mu_{12} & \mu_{23} - \mu_{21}\mu_{11}^{-1}\mu_{13} & \cdots & \mu_{2n} - \mu_{21}\mu_{11}^{-1}\mu_{1n} \\ 0 & \mu_{32} - \mu_{31}\mu_{11}^{-1}\mu_{12} & \mu_{33} - \mu_{31}\mu_{11}^{-1}\mu_{13} & \cdots & \mu_{3n} - \mu_{31}\mu_{11}^{-1}\mu_{1n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \mu_{n2} - \mu_{n1}\mu_{11}^{-1}\mu_{12} & \mu_{n3} - \mu_{n1}\mu_{11}^{-1}\mu_{13} & \cdots & \mu_{nn} - \mu_{n1}\mu_{11}^{-1}\mu_{1n} \end{bmatrix} = \\ &= \begin{bmatrix} \mu_{11} & \mu_{12} & \mu_{13} & \cdots & \mu_{1n} \\ 0 & \nu_{22} & \nu_{23} & \cdots & \nu_{2n} \\ 0 & \nu_{32} & \nu_{33} & \cdots & \nu_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \nu_{n2} & \nu_{n3} & \cdots & \nu_{nn} \end{bmatrix}, \end{aligned} \quad (2.2)$$

missä ν_{ij} :llä on merkitty eliminoinnin tuloksena syntyneitä uusia operaattoreita. Jos nyt seuraava \mathcal{M} :n tukialkio ν_{22} on kääntyvä, niin eliminointia voidaan jatkaa. Saadaan

$$\begin{aligned} \mathcal{L}_2 \mathcal{L}_1 \mathcal{M} &= \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ 0 & -\nu_{32}\nu_{22}^{-1} & 1 & & \\ \vdots & \vdots & & \ddots & \\ 0 & -\nu_{n2}\nu_{22}^{-1} & & & 1 \end{bmatrix} \begin{bmatrix} \mu_{11} & \mu_{12} & \mu_{13} & \cdots & \mu_{1n} \\ 0 & \nu_{22} & \nu_{23} & \cdots & \nu_{2n} \\ 0 & \nu_{32} & \nu_{33} & \cdots & \nu_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \nu_{n2} & \nu_{n3} & \cdots & \nu_{nn} \end{bmatrix} \\ &= \begin{bmatrix} \mu_{11} & \mu_{12} & \mu_{13} & \cdots & \mu_{1n} \\ 0 & \nu_{22} & \nu_{23} & \cdots & \nu_{2n} \\ 0 & 0 & \pi_{33} & \cdots & \pi_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \pi_{n3} & \cdots & \pi_{nn} \end{bmatrix}. \end{aligned} \quad (2.3)$$

Jos kaikki syntyvät \mathcal{M} :n tukialkiot ovat kääntyviä, niin tätä toistamalla oikealle puolelle muodostuu lopulta yläkolmio-operaattori; merkitään tätä \mathbf{U} :lla. Yllä olevia alakolmio-operaattoreita on merkitty \mathcal{L}_1 :llä (ensimmäinen näistä) ja \mathcal{L}_2 :lla (toinen).

Yllä olevan kaltaista prosessia toistettaessa käytettyjä loppuja alakolmio-operaattoreita merkitään \mathcal{L}_k :lla. Nämä operaattorit ovat muotoa $\mathcal{L}_k = \mathcal{I} - \eta_k \xi_k^T$, missä η_k on $n \times 1$ -operaattori, ξ_k on $n \times 1$ -operaattori, jonka k :s komponentti on yksikköoperaattori muiden ollessa nollaoperaattoreita, ja lisäksi pätee $\xi_j^T \eta_k = 0$, kun $j \leq k$. Esimerkiksi $\eta_1 = [0 \ \mu_{21}\mu_{11}^{-1} \ \mu_{31}\mu_{11}^{-1} \ \cdots \ \mu_{n1}\mu_{11}^{-1}]^T$.

Edellä olevin merkinnöin

$$\mathcal{L}_{n-1}\mathcal{L}_{n-2}\cdots\mathcal{L}_1\mathcal{M} = \mathcal{U}.$$

Operaattorin $\mathcal{L}_k = \mathcal{I} - \eta_k \xi_k^T$ käänteisoperaattori on $\mathcal{L}_k^{-1} = \mathcal{I} + \eta_k \xi_k^T$. Tämä nähdään laskemalla

$$(\mathcal{I} - \eta_k \xi_k^T)(\mathcal{I} + \eta_k \xi_k^T) = \mathcal{I} + \eta_k \xi_k^T - \eta_k \xi_k^T - \eta_k \xi_k^T \eta_k \xi_k^T = \mathcal{I} - \eta_k (\xi_k^T \eta_k) \xi_k^T = \mathcal{I}.$$

Olkoon $\mathcal{L} = \mathcal{L}_1^{-1} \cdots \mathcal{L}_{n-2}^{-1} \mathcal{L}_{n-1}^{-1}$, joka on alakolmio-operaattori, koska sellaisten tulot säilyvät alakolmio-operaattoreina. Lisäksi \mathcal{L} :n lävistjäalkiot ovat yksikköoperaattoreita. Myös $\mathcal{M} = \mathcal{L}\mathcal{U}$, joten tämä on \mathcal{M} :n LU-hajotelma.

Operaattori \mathcal{L} voidaan laskea suoraan Gaussin eliminaatiossa käytetyistä kertoimisista. Operaattori $\mathcal{R}_k = \mathcal{L}_k^{-1} \cdots \mathcal{L}_{n-2}^{-1} \mathcal{L}_{n-1}^{-1}$ voidaan kirjoittaa seuraavasti:

$$\mathcal{R}_k = \mathcal{I} + \eta_k \xi_k^T + \cdots + \eta_{n-2} \xi_{n-2}^T + \eta_{n-1} \xi_{n-1}^T.$$

Tämä nähdään induktiolla. Selvästi \mathcal{R}_{n-1} on tätä muotoa. Jos \mathcal{R}_{k+1} on edellä mainittua muotoa, niin

$$\begin{aligned} \mathcal{R}_k &= \mathcal{L}_k^{-1} \mathcal{R}_{k+1} = (\mathcal{I} + \eta_k \xi_k^T)(\mathcal{I} + \eta_{k+1} \xi_{k+1}^T + \cdots + \eta_{n-1} \xi_{n-1}^T) \\ &= \mathcal{I} + \eta_{k+1} \xi_{k+1}^T + \cdots + \eta_{n-1} \xi_{n-1}^T + \eta_k \xi_k^T + \eta_k (\xi_k^T \eta_{k+1}) \xi_{k+1}^T + \cdots + \\ &\quad \eta_k (\xi_k^T \eta_{n-1}) \xi_{n-1}^T = \mathcal{I} + \eta_k \xi_k^T + \eta_{k+1} \xi_{k+1}^T + \cdots + \eta_{n-1} \xi_{n-1}^T, \end{aligned}$$

missä käytettiin ominaisuutta $\xi_j^T \eta_l = 0$, kun $j \leq l$. Siten saadaan

$$\mathcal{L} = \mathcal{R}_1 = \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ \mu_{21}\mu_{11}^{-1} & 1 & 0 & \cdots & 0 \\ \mu_{31}\mu_{11}^{-1} & \nu_{32}\nu_{22}^{-1} & 1 & & \\ \vdots & \vdots & & \ddots & \\ \mu_{n1}\mu_{11}^{-1} & \nu_{n2}\nu_{22}^{-1} & & & 1 \end{bmatrix}.$$

Operaattoreiden \mathcal{L} ja \mathcal{U} laskemiseksi edellä kuvatulla tavalla täytyy tehdä noin $4n^3/3$ kompleksilukukertolaskua ja saman verran yhteenlaskuja. Tietokoneella laskettaessa kompleksiliukulaskutoimituksia (complex flops) tehdään noin $4n^3/3$, koska määritelmän mukaan yksi complex flop tarkoittaa yhden kompleksilukukertolaskun ja -yhteenlaskun yhdistelmää.

Reaalimatriisiesityksestä (1.4) havaitaan, että \mathbb{R} -lineaarinen LU-hajotelma voidaan nähdään myös matriisien LU-lohkohajotelmana, missä lohkot ovat 2×2 -kokoisia.

2.2 Osittaistuenta

Edellä oletettiin kaikki \mathcal{M} :n tukialkiot μ_{11} , ν_{22} jne. kääntyviksi. Tämä ei välttämättä päde, vaikka \mathcal{M} onkin kääntyvä. Lisäksi muihin skalaarioperaattorialkioihin

nähdessä lähes singulaarinen, mutta kuitenkin kääntyvä, tukialkio voi aiheuttaa äärellistarkkuisessa tietokonearitmetiikassa merkitsevien numeroiden menetyksiä huonontamalla lopputuloksen tarkkuutta tai tehdä sen jopa hyödyttömäksi. Matriisien LU-hajotelman laskennassa molemmat vastaavat ongelmat voidaan ratkaista käyttämällä tuentamenetelmiä, kuten osittais-, torni- ja täystuenta. Reaalilineaarisille operaattoreille näiden tuentamenetelmien analogioilla ei välttämättä päästä eroon kääntymättömistä tukialkioista, kuten seuraava esimerkki osoittaa.

Esimerkki 2.2.1. Tarkastellaan seuraavaa operaattoria.

$$\mu_{11}(z) = z + \bar{z}, \quad \mu_{12}(z) = iz + i\bar{z}, \quad \mu_{21}(z) = z - \bar{z}, \quad \mu_{22}(z) = -iz + i\bar{z},$$

$$\mathcal{M} = \begin{bmatrix} \mu_{11} & \mu_{12} \\ \mu_{21} & \mu_{22} \end{bmatrix}.$$

Skalaarioperaattori $\mu : \mathbb{C} \rightarrow \mathbb{C}$, $\mu(z) = uz + v\bar{z}$ on kääntyvä täsmälleen silloin, kun $|u| \neq |v|$. Mikään yllä olevista skalaarioperaattoreista ei siis ole kääntyvä. Skalaarioperaattoria μ vastaava reaalinen matriisi on

$$\begin{bmatrix} \operatorname{Re}(u+v) & -\operatorname{Im}(u-v) \\ \operatorname{Im}(u+v) & \operatorname{Re}(u-v) \end{bmatrix}.$$

\mathcal{M} :ää vastaa matriisi

$$\begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 2 \\ 0 & 2 & 0 & 0 \end{bmatrix}.$$

Tämä on kääntyvä, joten \mathcal{M} on kääntyvä operaattori. Tarvittavaa kääntävää tukialkiota ei siis välttämättä ole mahdollista löytää operaattorille \mathcal{M} , vaikka \mathcal{M} olisikin kääntyvä. Esimerkin skalaarioperaattoreista μ_{ij} mikään ei kelpaa ensimmäiseksi tukialkioksi, joten edes matriisien täydellisen tuennan analogia ei toimi kaikilla reaalilineaarisilla operaattoreilla.

Esitellään kuitenkin seuraavaksi matriisien osittaistuennan vastine. Menetelmä laajennetaan myöhemmin ottamaan huomioon myös edellä kuvatun esimerkin kaltaisia tilanteita. Tässä tukialkion valinta suoritetaan rivien järjestyksestä vaihtamalla. Sarakkeiden järjestyksen vaihtaminen olisi myös mahdollista (saraketuenta) kuten matriisien täydellisen tuennan vastinekin (etsitään jäljellä olevasta alamatriisista itseisarvoltaan suurin alkio ja vaihdetaan sekä rivejä että sarakkeita).

Skalaarioperaattori $\mu(z) = uz + v\bar{z}$ on kääntyvä täsmälleen silloin, kun $|u| \neq |v|$ ja tällöin käänteisoperaattori on

$$\mu^{-1}(z) = \frac{1}{|u|^2 - |v|^2}(\bar{u}z - v\bar{z}).$$

Pyritään välttämään pienellä luvulla jakamista, jottei μ^{-1} :n lineaarisesta eikä antilinearisesta osasta tulisi itseisarvoltaan suuri. Yritetään siis saada hajotelman laskemisessa käännettävien skalaarioperaattoreiden $|u|^2 - |v|^2$ itseisarvoltaan mahdollisimman suureksi. Merkitään $N(\mu) = |u|^2 - |v|^2$.

Osittaistuennassa etsitään aluksi operaattorin $\mathcal{M} = (\mu_{ij})$ ensimmäiseltä sarakeelta alkio $\mu_{k,1}$, jolle $|N(\mu_{k,1})|$ on suurin. Vaihdetaan operaattorin \mathcal{M} ensimmäisen ja k :nnen rivin paikkaa. Tämä voidaan tehdä laskemalla $P_1\mathcal{M}$, missä matriisi $P_1 \in \mathbb{C}^{n \times n}$ on permutaatiomatriisi, joka on saatu vaihtamalla yksikkömatriisiin

ensimmäisen ja k :nnen sarakkeen paikat. Oletetaan, että löydetylle alkioille pätee $|N(\boldsymbol{\mu}_{k,1})| > 0$, jolloin operaattorin $\mathbf{P}_1\mathcal{M}$ vasemman yläkulman alkio on kääntyvä.

Tämän jälkeen voidaan tehdä Gaussin eliminaatio kuten kaavassa (2.2), mutta nyt \mathcal{M} :n paikalla on $\mathbf{P}_1\mathcal{M}$. Tämän tuloksena saadaan operaattori $\mathcal{L}_1\mathbf{P}_1\mathcal{M}$.

Seuraavaksi etsitään matriisin $\mathcal{L}_1\mathbf{P}_1\mathcal{A}$ toisen sarakkeen toiselta riviltä alkaen suurinta $|N(\boldsymbol{\mu}_{k,2})|$ eli olkoon $2 \leq k \leq n$ sellainen, että

$$|(\mathcal{L}_1\mathbf{P}_1\mathcal{M})_{k,2}| = \max_{2 \leq i \leq n} |(\mathcal{L}_1\mathbf{P}_1\mathcal{M})_{i,2}|.$$

Oletetaan tämä jälleen positiiviseksi. Olkoon $\mathbf{P}_2 \in \mathbb{C}^{n \times n}$ permutaatiomatriisi, joka on saatu vaihtamalla yksikkömatriisin sarakkeet 2 ja k keskenään. Gaussin eliminaation tuloksena saadaan $\mathcal{L}_2\mathbf{P}_2\mathcal{L}_1\mathbf{P}_1\mathcal{M}$, kuten kaavassa (2.3), mutta nyt $\mathcal{L}_1\mathcal{M}$:n paikalla on $\mathbf{P}_2\mathcal{L}_1\mathbf{P}_1\mathcal{M}$.

Kun edellä kuvattua menettelyä jatketaan, saadaan yläkolmio-operaattori

$$\mathcal{L}_{n-1}\mathbf{P}_{n-1}\mathcal{L}_{n-2}\mathbf{P}_{n-2} \cdots \mathcal{L}_1\mathbf{P}_1\mathcal{M} = \mathbf{U}.$$

Koska tässä olevat matriisit \mathbf{P}_k ovat yksikkömatriiseja mahdollisesti k :s ja jokin sen oikealla puolella oleva sarake vaihdettuna, niin $\mathbf{P}_k\boldsymbol{\xi}_j = \boldsymbol{\xi}_j$, kun $j < k$. Myös $\mathbf{P}_k = \mathbf{P}_k^T$.

Edellisessä kohdassa nähtiin, että $\mathcal{L}_k = \mathcal{I} - \boldsymbol{\eta}_k\boldsymbol{\xi}_k^T$. Matriisit \mathbf{P}_j voidaan kuljettaa kaikkien \mathcal{L}_k operaattoreiden perään seuraavalla tavalla. Merkitään $\boldsymbol{\eta}'_{n-2} = \mathbf{P}_{n-1}\boldsymbol{\eta}_{n-2}$ ja $\mathcal{L}'_{n-2} = \mathcal{I} - \boldsymbol{\eta}'_{n-2}\boldsymbol{\xi}_{n-2}^T$.

$$\begin{aligned} \mathbf{P}_{n-1}\mathcal{L}_{n-2} &= \mathbf{P}_{n-1}(\mathcal{I} - \boldsymbol{\eta}_{n-2}\boldsymbol{\xi}_{n-2}^T) = \mathbf{P}_{n-1} - \boldsymbol{\eta}'_{n-2}\boldsymbol{\xi}_{n-2}^T \\ &= \mathbf{P}_{n-1} - \boldsymbol{\eta}'_{n-2}(\mathbf{P}_{n-1}\boldsymbol{\xi}_{n-2})^T = \mathbf{P}_{n-1} - \boldsymbol{\eta}'_{n-2}\boldsymbol{\xi}_{n-2}^T\mathbf{P}_{n-1}^T \\ &= \mathbf{P}_{n-1} - \boldsymbol{\eta}'_{n-2}\boldsymbol{\xi}_{n-2}^T\mathbf{P}_{n-1} = \mathcal{L}'_{n-2}\mathbf{P}_{n-1}. \end{aligned}$$

Merkitään lisäksi $\boldsymbol{\eta}'_k = \mathbf{P}_{n-1}\mathbf{P}_{n-2} \cdots \mathbf{P}_{k+1}\boldsymbol{\eta}_k$, $\mathcal{L}'_k = \mathcal{I} - \boldsymbol{\eta}'_k\boldsymbol{\xi}_k^T$ ja

$$\mathbf{P} = \mathbf{P}_{n-1}\mathbf{P}_{n-2} \cdots \mathbf{P}_1.$$

Edellä olevan kaltaisten laskujen jälkeen tuloksena on $\mathcal{L}_{n-1}\mathcal{L}'_{n-2} \cdots \mathcal{L}'_1\mathbf{P}\mathcal{M} = \mathbf{U}$. Merkitään vielä $\mathcal{L} = \mathcal{L}'_1 \cdots \mathcal{L}'_{n-2}\mathcal{L}_{n-1}$, jolloin $\mathcal{L}\mathbf{U} = \mathbf{P}\mathcal{M}$. Operaattori \mathcal{L} koostuu Gaussin eliminaatiossa käytetyistä kertoimista

$$\mathcal{L} = \mathcal{I} + \boldsymbol{\eta}'_1\boldsymbol{\xi}_1^T + \boldsymbol{\eta}'_2\boldsymbol{\xi}_2^T + \cdots + \boldsymbol{\eta}'_{n-2}\boldsymbol{\xi}_{n-2}^T + \boldsymbol{\eta}_{n-1}\boldsymbol{\xi}_{n-1}^T.$$

2.3 Tuenta \mathbb{R} -lineaariseksi LU-hajotelmalle

Seuraavaksi esitellään strategia, jonka avulla reaali-lineaariseksi operaattorille voidaan aina muodostaa LU-hajotelman kaltainen hajotelma. Ensiksi osoitetaan hajotelman olemassaolo, jonka jälkeen tarkastellaan kuinka sen voi laskea numeerisesti mielekkäällä tavalla. Tavoitteena on päästä eroon edellisessä kohdassa mainittujen kääntymättömien tukialkioiden ongelmasta.

Lemma 2.3.1. *Olkooot $\mu_1 \neq \mathbf{0}$ ja μ_2 singulaarisia reaalioperaattoreita. Jos $\mu_1 + q\mu_2$ ei ole kääntyvä millään $q \in \mathbb{C}$, niin on olemassa $r \in \mathbb{C}$ siten, että $\mu_2 - r\mu_1 = \mathbf{0}$.*

Todistus. Jos $\mu_2 = \mathbf{0}$, niin valitaan $r = 0$. Olkoon sitten $\mu_2 \neq \mathbf{0}$. Merkitsemällä $\mu_1(z) = u_1z + v_1\bar{z}$, $\mu_2(z) = u_2z + v_2\bar{z}$ voidaan laskea

$$0 = |u_1 + qu_2|^2 - |v_1 + qv_2|^2 = (|u_1|^2 - |v_1|^2) + |q|^2(|u_2|^2 - |v_2|^2) + 2\operatorname{Re}(\bar{q}(u_1\bar{u}_2 - v_1\bar{v}_2)) = 2\operatorname{Re}(\bar{q}(u_1\bar{u}_2 - v_1\bar{v}_2)).$$

Valitsemalla $q = (u_1\bar{u}_2 - v_1\bar{v}_2)$ seuraa tästä $u_1\bar{u}_2 = v_1\bar{v}_2$. Olkoon $r = u_2/u_1$. Tällöin $u_2 - ru_1 = 0$ ja

$$v_2 - rv_1 = v_2 - \frac{u_2}{u_1} \frac{u_1\bar{u}_2}{\bar{v}_2} = v_2 - \frac{u_2\bar{u}_2v_2}{\bar{v}_2v_2} = 0.$$

Eli $\mu_2 - r\mu_1 = \mathbf{0}$. □

Lause 2.3.2. *Olkoon $\mathcal{M} : \mathbb{C}^n \rightarrow \mathbb{C}^n$ reaalioperaattori. Tällöin on olemassa alakolmio- $\mathcal{L} : \mathbb{C}^n \rightarrow \mathbb{C}^n$, jonka lävistäjä koostuu skalaarisista yksikköoperaattoreista, yläkolmio- $\mathcal{U} : \mathbb{C}^n \rightarrow \mathbb{C}^n$, yläkolmio- $\mathcal{V} : \mathbb{C}^n \rightarrow \mathbb{C}^n$, jonka jokaisella rivillä on yksikkölävistäjäalkion lisäksi korkeintaan yksi toinen nollasta poikkeava kompleksiluku, ja permutaatio- $\mathcal{P} : \mathbb{C}^n \rightarrow \mathbb{C}^n$ siten, että*

$$\mathcal{L}\mathcal{U} = \mathcal{V}\mathcal{P}\mathcal{M}.$$

Todistus. Olkoon $\mathcal{M} : \mathbb{C}^n \rightarrow \mathbb{C}^n$, $\mathcal{M} = (\mu_{ij})$ reaalioperaattori. Jos tämän ensimmäisessä sarakkeessa rivillä k on kääntyvä operaattori, niin kerrotaan \mathcal{M} vasemmalta rivinvaihto- \mathcal{P}_1 , joka vaihtaa \mathcal{M} :n ensimmäisen ja k :nnen rivin paikat keskenään. Tämän jälkeen voidaan suorittaa Gaussin eliminointi kuten edellä ja saadaan $\mathcal{L}_1\mathcal{P}_1\mathcal{M}$.

Jos ensimmäisessä sarakkeessa ei ole kääntyvää operaattoria, niin etsitään sieltä jokin nollasta poikkeava operaattori μ_{k1} . Jos tällaista ei ole, niin ensimmäinen sarake on jo eliminoitu ja voidaan siirtyä eliminoimaan seuraavaa saraketta, jolloin $\mathcal{L}_1 = \mathcal{P}_1 = \mathcal{I}$. Jos on ja löytyy jokin toinen μ_{j1} ja $q \in \mathbb{C}$ siten, että $\mu_{k1} + q\mu_{j1}$ on kääntyvä operaattori, niin kerrotaan \mathcal{M} rivinvaihto- \mathcal{P}_1 , joka vaihtaa ensimmäisen ja k :nnen rivin paikat, ja sen jälkeen kerrotaan operaattorilla $\mathcal{V}_1 = \mathcal{I} + \xi_1\kappa_1^T$, missä κ_j on $n \times 1$ -operaattori, jonka j :s komponentti on q . Operaattorin $\mathcal{V}_1\mathcal{P}_1\mathcal{M}$ vasemmassa yläkulmassa on nyt $\mu_{k1} + q\mu_{j1}$, joka on kääntyvä. Voidaan siis eliminoida Gaussin algoritmilla, jolloin saadaan $\mathcal{L}_1\mathcal{V}_1\mathcal{P}_1\mathcal{M}$.

Jos $\mu_{k1} + q\mu_{j1}$ ei ole kääntyvä millään $1 \leq j \leq n$, $j \neq k$, $q \in \mathbb{C}$, niin lemmän 2.3.1 mukaan on olemassa luvut $r_j \in \mathbb{C}$ siten, että $\mu_{j1} - r_j\mu_{k1} = \mathbf{0}$. Olkoon \mathcal{P}_1 jälleen rivinvaihto- \mathcal{P}_1 , joka vasemmalta kertomalla vaihtaa ensimmäisen ja k :nnen rivin paikat, ja olkoon $\mathcal{L}_1 = \mathcal{I} - \eta_1\xi_1^T$, jossa η_1 sisältää komponentteina r_j :t paitsi ensimmäisessä komponentissa 0 ja k :nnessä r_1 . Operaattorin $\mathcal{P}_1\mathcal{M}$ ensimmäisen sarakkeen, ensimmäisen rivin alapuoliset, komponentit voidaan nyt eliminoida kertomalla \mathcal{L}_1 :llä vasemmalta. Saadaan $\mathcal{L}_1\mathcal{P}_1\mathcal{M}$.

Voidaan siirtyä eliminoimaan seuraava sarake samaan tapaan. Lopulta saadaan

$$\mathcal{L}_{n-1}\mathcal{V}_{n-1}\mathcal{P}_{n-1}\mathcal{L}_{n-2}\mathcal{V}_{n-2}\mathcal{P}_{n-2}\cdots\mathcal{L}_1\mathcal{V}_1\mathcal{P}_1\mathcal{M}=\mathcal{U}.$$

Operaattorit \mathcal{L}_k ja \mathcal{V}_k ovat muotoa

$$\mathcal{L}_k=\mathcal{I}-\eta_k\xi_k^T, \quad \mathcal{V}_k=\mathcal{I}+\xi_k\kappa_k^T \quad k=1,\dots,n-1,$$

missä η_k, κ_k ovat $n \times 1$ -operaattoreita, joille $\xi_j^T\eta_k=\xi_j^T\kappa_k=0$, kun $j \leq k$.

Operaattorit \mathcal{V}_k ja \mathcal{P}_k voidaan kuljettaa operaattorien \mathcal{L}_k perään seuraavasti. Pätee

$$\mathcal{V}_{n-1}\mathcal{P}_{n-1}\mathcal{L}_{n-2}=\mathcal{V}_{n-1}\mathcal{L}'_{n-2}\mathcal{P}_{n-1},$$

missä $\eta'_{n-2}=\mathcal{P}_{n-1}\eta_{n-2}$ ja $\mathcal{L}'_{n-2}=\mathcal{I}-\eta'_{n-2}\xi_{n-2}^T$. Olkoon $\eta''_{n-2}=\mathcal{V}_{n-1}\eta'_{n-2}$ ja $\mathcal{L}''_{n-2}=\mathcal{I}-\eta''_{n-2}\xi_{n-2}^T$. Tällöin

$$\begin{aligned} \mathcal{V}_{n-1}\mathcal{L}'_{n-2} &= \mathcal{V}_{n-1}-\eta''_{n-2}\xi_{n-2}^T = \mathcal{I}+\xi_{n-1}^T\kappa_{n-1}-\eta''_{n-2}\xi_{n-2}^T \\ &= (\mathcal{I}-\eta''_{n-2}\xi_{n-2}^T)(\mathcal{I}+\xi_{n-1}^T\kappa_{n-1}) = \mathcal{L}''_{n-2}\mathcal{V}_{n-1}. \end{aligned}$$

Merkitään vielä $\eta''_k=\mathcal{V}_{n-1}\mathcal{P}_{n-1}\mathcal{V}_{n-2}\mathcal{P}_{n-2}\cdots\mathcal{V}_{k+1}\mathcal{P}_{k+1}\eta_k$ ja $\mathcal{L}''_k=\mathcal{I}-\eta''_k\xi_k^T$. Tällöin edellä kuvatusen kaltaisesti saadaan

$$\mathcal{L}_{n-1}\mathcal{L}''_{n-2}\cdots\mathcal{L}''_1\mathcal{V}_{n-1}\mathcal{P}_{n-1}\mathcal{V}_{n-2}\mathcal{P}_{n-2}\cdots\mathcal{V}_1\mathcal{P}_1\mathcal{M}=\mathcal{U}.$$

Operaattorit \mathcal{P}_k voidaan vielä kuljettaa operaattoreiden \mathcal{V}_k eteen. Kun merkitään $\kappa'_{n-2}=\mathcal{P}_{n-1}\kappa_{n-2}$ ja $\mathcal{V}'_{n-2}=\mathcal{I}+\xi_{n-2}\kappa'_{n-2}$, niin

$$\begin{aligned} \mathcal{P}_{n-1}\mathcal{V}_{n-2} &= \mathcal{P}_{n-1}+\mathcal{P}_{n-1}\xi_{n-2}\kappa_{n-2}^T = \mathcal{P}_{n-1}+\xi_{n-2}\kappa_{n-2}^T \\ &= \mathcal{P}_{n-1}+\xi_{n-2}(\mathcal{P}_{n-1}^2\kappa_{n-2})^T = \mathcal{P}_{n-1}+\xi_{n-2}(\mathcal{P}_{n-1}\kappa_{n-2})^T\mathcal{P}_{n-1}^T \\ &= \mathcal{P}_{n-1}+\xi_{n-2}\kappa'_{n-2}\mathcal{P}_{n-1} = \mathcal{V}'_{n-2}\mathcal{P}_{n-1}. \end{aligned}$$

Jatketaan tätä kunnes kaikki \mathcal{P}_k :t ovat perässä. Kun $\kappa'_k=\mathcal{P}_{n-1}\mathcal{P}_{n-2}\cdots\mathcal{P}_{k+1}\kappa_k$ ja $\mathcal{V}'_k=\mathcal{I}+\xi_k\kappa_k'^T$, niin

$$\mathcal{L}_{n-1}\mathcal{L}''_{n-2}\cdots\mathcal{L}''_1\mathcal{V}_{n-1}\mathcal{V}'_{n-2}\cdots\mathcal{V}'_1\mathcal{P}_{n-1}\mathcal{P}_{n-2}\cdots\mathcal{P}_1\mathcal{M}=\mathcal{U}.$$

Merkitsemällä

$$\mathcal{L}=\mathcal{L}_1''^{-1}\cdots\mathcal{L}_{n-2}''^{-1}\mathcal{L}_{n-1}^{-1}, \quad \mathcal{V}=\mathcal{V}_{n-1}\mathcal{V}'_{n-2}\cdots\mathcal{V}'_1 \quad \text{ja} \quad \mathcal{P}=\mathcal{P}_{n-1}\mathcal{P}_{n-2}\cdots\mathcal{P}_1$$

voidaan edellä oleva lausua $\mathcal{L}\mathcal{U}=\mathcal{V}\mathcal{P}\mathcal{M}$. Tässä operaattori \mathcal{L} voidaan muodostaa suoraan Gaussin eliminoinnissa käytetyistä kertoimista kuten aiemmin:

$$\mathcal{L}=\mathcal{I}+\eta_1''\xi_1^T+\cdots+\eta_{n-2}''\xi_{n-2}^T+\eta_{n-1}''\xi_{n-1}^T.$$

Myös operaattori \mathcal{V} voidaan muodostaa suoraan seuraavasti. Merkitään $\mathcal{S}_k=\mathcal{V}'_k\cdots\mathcal{V}'_2\mathcal{V}'_1$, jolloin induktiolla voidaan osoittaa

$$\mathcal{S}_k=\mathcal{I}+\xi_1\kappa_1'^T+\xi_2\kappa_2'^T+\cdots+\xi_k\kappa_k'^T.$$

Nimittäin selvästi $\mathcal{S}_1=\mathcal{V}'_1$ on tätä muotoa. Jos \mathcal{S}_k on yllä olevaa muotoa, niin

$$\begin{aligned} \mathcal{S}_{k+1} &= \mathcal{V}'_{k+1}\mathcal{S}_k = (\mathcal{I}+\xi_{k+1}\kappa_{k+1}'^T)\mathcal{S}_k = \mathcal{I}+\xi_1\kappa_1'^T+\xi_2\kappa_2'^T+\cdots+\xi_k\kappa_k'^T+ \\ &\quad \xi_{k+1}\kappa_{k+1}'^T+\xi_{k+1}(\kappa_{k+1}'^T\xi_1)\kappa_1'^T+\xi_{k+1}(\kappa_{k+1}'^T\xi_2)\kappa_2'^T+\cdots+\xi_{k+1}(\kappa_{k+1}'^T\xi_k)\kappa_k'^T \\ &= \mathcal{I}+\xi_1\kappa_1'^T+\xi_2\kappa_2'^T+\cdots+\xi_k\kappa_k'^T+\xi_{k+1}\kappa_{k+1}'^T, \end{aligned}$$

joka on haluttua muotoa. Täten

$$\mathcal{V}=\mathcal{I}+\xi_1\kappa_1'^T+\xi_2\kappa_2'^T+\cdots+\xi_{n-2}\kappa_{n-2}'^T+\xi_{n-1}\kappa_{n-1}'^T.$$

Lause on nyt osoitettu. \square

2.3.1 \mathbb{R} -lineaarisesti tuetun LU-hajotelman laskeminen

Esitellään sitten eräs tapa laskea edellä kuvattu hajotelma reaalin lineaarisille operaattoreille tavalla, joka pyrkii ottamaan huomioon numeeriset näkökohdat. Oletetaan seuraavassa, että operaattori \mathcal{M} on kääntyvä.

Lauseen 2.3.2 mukaiset \mathcal{L} , \mathcal{V} ja \mathcal{P} ovat kääntyviä, joten myös \mathcal{U} on kääntyvä. Erityisesti tästä seuraa, että kaikki \mathcal{U} :n skalaariset lävistäjäoperaattorit ovat kääntyviä. Täten lauseen 2.3.2 todistuksen alussa sarakkeelta löytyy joko heti kääntyvä operaattori tai sitten nolasta poikkeava μ_{k1} ja $q \in \mathbb{C}$, $j \neq k$ siten, että $\mu_{k1} + q\mu_{j1}$ on kääntyvä.

Skalaarioperaattorille $\mu(z) = uz + v\bar{z}$ määritellään kohdan 2.2 tavoin suure $N(\mu) = |u|^2 - |v|^2$ ja käytetään seuraavaa strategiaa. Etsitään ensin operaattorin \mathcal{M} ensimmäiseltä sarakkeelta alkio μ_{k1} siten, että $|N(\mu_{k1})|$ on suurin mahdollinen. Seuraavaksi lasketaan jokaiselle $i, j \in \mathbb{N}$, $1 \leq i, j \leq n$, $i \neq j$ luku $N(\mu_{i1} + q_{ij}\mu_{j1})$, jossa $q_{ij} = s_{ij}r_{ij}$, $r_{ij} = (u_{i1}\bar{u}_{j1} - v_{i1}\bar{v}_{j1})$ ja $s_{ij} \in \mathbb{R}$ vielä määräämätön reaaliluku. Tällöin

$$\begin{aligned} N(\mu_{i1} + q_{ij}\mu_{j1}) &= |u_{i1} + q_{ij}u_{j1}|^2 - |v_{i1} + q_{ij}v_{j1}|^2 \\ &= (|u_{i1}|^2 - |v_{i1}|^2) + |q_{ij}|^2(|u_{j1}|^2 - |v_{j1}|^2) + \\ &\quad 2\operatorname{Re}(\bar{q}_{ij}(u_{i1}\bar{u}_{j1} - v_{i1}\bar{v}_{j1})) = N(\mu_{i1}) + s_{ij}^2|r_{ij}|^2N(\mu_{j1}) + 2s_{ij}|r_{ij}|^2. \end{aligned}$$

Huomaa, että pidettäessä $|q_{ij}|$ kiinnitettynä, tehty valinta q_{ij} :lle maksimoi arvon $N(\mu_{i1} + q_{ij}\mu_{j1})$.

Valittaessa lukuja s_{ij} pyritään siihen ettei operaattorin $\mathcal{V}_1\mathcal{P}_1\mathcal{M}$ ensimmäisen rivin alkioden suuruusluokka kasvaisi kohtuuttomasti operaattorin $\mathcal{P}_1\mathcal{M}$ ensimmäiseen riviin nähden. Valitaan tässä yksinkertaisesti joko $s_{ij} = 1$ tai $s_{ij} = -1$ sen perusteella kummalla valinnalla lausekkeen $|N(\mu_{i1} + q_{ij}\mu_{j1})|$ arvo on suurempi.

Lopuksi poimitaan suurempi arvoista $|N(\mu_{k1})|$ ja $|N(\mu_{i1} + q_{ij}\mu_{j1})|$. Liitteessä B oleva MATLAB-funktio `r1_luvp` laskee yllä kuvatulla tavalla hajotelman $\mathcal{LU} = \mathcal{VPM}$.

Luku 3

Iteratiivinen ratkaiseminen

3.1 Johdanto

Olkoon $\mathbf{A} \in \mathbb{C}^{n \times n}$, $\mathbf{b} \in \mathbb{C}^n$ ja tehtävänä ratkaista lineaarinen yhtälöryhmä

$$\mathbf{Ax} = \mathbf{b}. \quad (3.1)$$

Tätä tarkoitusta varten edellisessä luvussa esiteltiin Gaussin menetelmään perustuva \mathbf{A} :n LU-hajotelma. Tämä algoritmi vaatii $O(n^3)$ laskuoperaatiota ja on esimerkiksi ns. suorasta menetelmästä lineaarisen yhtälöryhmän ratkaisemiseksi. Muistia se vaatii $O(n^2)$ alkioita. Suurilla n toinen tai molemmat näistä vaatimuksista saattavat ylittää käytettävissä olevat tietokoneresurssit. Erityisesti osittaisdifferentiaaliyhtälöiden ja integraaliyhtälöiden diskretoinnista seuraavissa yhtälöryhmissä n kasvaa haettaessa lisää tarkkuutta lopulta resurssit ylittäen. Toisaalta näiden diskretointien matriiseilla \mathbf{A} on usein rakenne, jota voidaan käyttää hyödyksi. Matriisia \mathbf{A} kutsutaan harvaksi, jos niin suuri määrä sen alkioista on nollia, että tämän hyödyntäminen on mahdollista.

Iteratiivisissa menetelmissä lähdetään liikkeelle ratkaisun approksimaatiosta, jota pyritään tarkentamaan algoritmin jokaisella iteraatiokierroksella. Iterointi pysäytetään, kun riittävä tarkkuus on saavutettu. Erilaisia menetelmiä on 1950-luvulta lähtien esitetty lukuisia ja niiden suppenemisnopeudet kohti ratkaisua riippuvat paitsi menetelmästä myös matriisin \mathbf{A} ominaisuuksista. Suppenemisnopeuteen voidaan huomattavasti vaikuttaa myös ns. pohjustuksella, mistä enemmän seuraavassa luvussa.

Usein käytännön tehtävissä esiintyvien yhtälöryhmien matriiseilla \mathbf{A} on ominaisuuksia, jotka tekevät iteratiiviset menetelmät suoraa menetelmiä laskennallisesti edullisemmiksi. Tyypillisiä iteratiivisten menetelmien perusoperaatioita jokaisella iteraatiokierroksella ovat matriisilla \mathbf{A} kertominen ja sisätulon laskeminen. Harvoista matriiseista nauhamatriisien tulot vektoreiden kanssa vievät vain $O(n)$ laskuoperaatiota. Suorat menetelmät eivät myöskään tule kyseeseen mikäli itse matriisia \mathbf{A} ei ole mielekästä muodostaa ja tallentaa tietokoneen muistiin, mutta operaatio $\mathbf{x} \mapsto \mathbf{Ax}$ voidaan suorittaa mille tahansa vektorille \mathbf{x} . Toisaalta joskus yhtälön ratkaisun ei tarvitse olla tarkka, jolloin iteroinnin voi pysäyttää ratkaisun saavutettua saman virherajan kuin mittauservoista peräisin olevilla \mathbf{A} :lla ja \mathbf{b} :llä. Suoralla menetelmäl-

lä ratkaiseminen täytyy viedä loppuun asti, koska laskennan välitulokset eivät ole käyttökelpoisia likimääräisratkaisuja.

Myös tässä luvussa esitellään aluksi joitakin iteratiivisia menetelmiä matriiseille ja sen jälkeen siirrytään käsittelemään reaalilineaaristen operaattoreiden menetelmiä.

3.2 Petrov-Galerkinin menetelmä

Abstrakti Galerkinin menetelmä on pohjana useille eri menetelmille. Yhtälö (3.1) voidaan kirjoittaa yhtäpitävästi muotoon

$$\langle \mathbf{Ax}, \mathbf{y} \rangle = \langle \mathbf{b}, \mathbf{y} \rangle \quad \text{kaikilla } \mathbf{y} \in \mathbb{C}^n.$$

Notaatio $\langle \cdot, \cdot \rangle$ tarkoittaa \mathbb{C}^n :n sisätuloa $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{y}^* \mathbf{x}$. Kirjoitetaan edellinen yhtälö muotoon

$$\langle \mathbf{b} - \mathbf{Ax}, \mathbf{y} \rangle = 0 \quad \text{kaikilla } \mathbf{y} \in \mathbb{C}^n. \quad (3.2)$$

Korvaamalla tässä avaruus \mathbb{C}^n sopivasti aliavaruuksilla, saadaan approksimoiva yhtälö. Olkoot $X_k, Y_k \subset \mathbb{C}^n$ joitakin \mathbb{C}^n :n aliavaruuksia ja \mathbf{x}_0 ratkaisun alkuarvaus. Avaruuksien X_k, Y_k dimensiot ovat tyypillisesti huomattavasti pienemmät kuin n . Petrov-Galerkinin ehto ratkaisun approksimaatioksi $\mathbf{x}_k \in \mathbf{x}_0 + X_k$ on

$$\langle \mathbf{b} - \mathbf{Ax}_k, \mathbf{y} \rangle = 0 \quad \text{kaikilla } \mathbf{y} \in Y_k. \quad (3.3)$$

Ratkaisua \mathbf{x}_k etsitään siis affiinista aliavaruudesta $\mathbf{x}_0 + X_k$ siten, että jäännös $\mathbf{r}_k = \mathbf{b} - \mathbf{Ax}_k$ on kohtisuorassa Y_k :ta vastaan.

Ratkaisun \mathbf{x}_k laskemiseksi, olkoot matriisin $\mathbf{V}_k \in \mathbb{C}^{n \times k}$ sarakkeet avaruuden X_k jotkin kantavektorit ja matriisin $\mathbf{W}_k \in \mathbb{C}^{n \times k}$ sarakkeet Y_k :n kanta. Olkoon $\mathbf{z}_k \in \mathbb{C}^k$ siten, että

$$\mathbf{x}_k = \mathbf{x}_0 + \mathbf{V}_k \mathbf{z}_k. \quad (3.4)$$

Merkitään vielä $\mathbf{v} = \mathbf{b} - \mathbf{Ax}_0$, jolloin yhtälö (3.3) saadaan muotoon

$$\mathbf{W}_k^* \mathbf{AV}_k \mathbf{z}_k = \mathbf{W}_k^* \mathbf{v}. \quad (3.5)$$

Ratkaisemalla tästä \mathbf{z}_k ja sijoittamalla kaavaan (3.4) saadaan yhtälön (3.1) ratkaisun approksimaatio \mathbf{x}_k .

Matriisi $\mathbf{W}_k^* \mathbf{AV}_k$ ei välttämättä ole kääntövä vaikka \mathbf{A} olisikin. Seuraavan lauseen kahdessa tapauksessa näin kuitenkin on.

Lause 3.2.1. *Olkoon $\mathbf{V} \in \mathbb{C}^{n \times k}$ aliavaruuden X_k kanta ja $\mathbf{W} \in \mathbb{C}^{n \times k}$ aliavaruuden Y_k kanta. Tällöin $\mathbf{W}^* \mathbf{AV}$ on kääntövä, jos*

- (1) $\langle \mathbf{Ax}, \mathbf{x} \rangle \neq 0$ kaikilla $\mathbf{x} \neq 0$ ja $X_k = Y_k$, tai
- (2) \mathbf{A} on kääntövä ja $Y_k = \mathbf{AX}_k$.

Todistus. Ensimmäisessä tapauksessa on olemassa kannanvaihtomatriisi $\mathbf{K} \in \mathbb{C}^{k \times k}$ siten, että $\mathbf{W} = \mathbf{VK}$. Tällöin $\mathbf{W}^* \mathbf{AV} = \mathbf{K}^* \mathbf{V}^* \mathbf{AV}$. Koska

$$\langle \mathbf{V}^* \mathbf{AV} \mathbf{z}, \mathbf{z} \rangle = \langle \mathbf{AV} \mathbf{z}, \mathbf{V} \mathbf{z} \rangle \neq 0 \quad \text{kaikilla } \mathbf{z} \in \mathbb{C}^k, \mathbf{z} \neq \mathbf{0},$$

niin matriisi V^*AV on kääntyvä, koska sen nolla-avaruus sisältää vain nollavektorin. Koska myös kannanvaihtomatriisi K on kääntyvä, väite on saatu osoitettua.

Toista tapausta varten huomataan, että AV on Y_k :n kanta. Täten on olemassa kannanvaihtomatriisi $K \in \mathbb{C}^{k \times k}$ siten, että $AV = WK$. Tällöin $W^*AV = W^*WK$. Matriisi K on kääntyvä ja myös W^*W on kääntyvä, koska

$$\langle W^*Wz, z \rangle = \langle Wz, Wz \rangle = \|Wz\|^2 > 0 \quad \text{kaikilla } z \in \mathbb{C}^k, z \neq \mathbf{0}.$$

□

Petrov-Galerkinin yhtälön ratkaisulle pätee seuraavat yleiset optimaalisuustulokset.

Lause 3.2.2. *Olkoon $A \in \mathbb{C}^{n \times n}$ hermiittinen ja positiividefiniitti matriisi ja $X_k = Y_k$. Tällöin Petrov-Galerkinin yhtälön ratkaisun x_k virhe on kohtisuorassa aliavaruutta Y_k vastaan A -normissa eli*

$$\langle x_* - x_k, y \rangle_A = 0 \quad \text{kaikilla } y \in Y_k,$$

missä $\langle u, v \rangle_A = \langle Au, v \rangle$ ja x_* on yhtälön (3.1) tarkka ratkaisu. Approksimaatio x_k minimoi virheen A -normissa

$$\|x_* - x_k\|_A = \min_{y \in X_k} \|x_* - y\|_A.$$

Todistus. Asetetaan yhtälössä (3.2) $x = x_*$ ja vähennetään se puolittain (3.3):sta, jolloin

$$\langle A(x_* - x_k), y \rangle = 0 \quad \text{kaikilla } y \in Y_k.$$

Virheen minimoituminen seuraa Pythagoraan lauseesta.

$$\|x_* - y\|_A^2 = \|x_* - x_k + x_k - y\|_A^2 = \|x_* - x_k\|_A^2 + \|x_k - y\|_A^2 \quad (y \in X_k),$$

missä on käytetty kohtisuoruutta $(x_* - x_k) \perp (x_k - y)$. Virhe minimoituu, kun $y = x_k$. □

Toisessa tapauksessa ($Y_k = AX_k$) nähdään yhtälöstä (3.3), että approksimaation x_k antama jäännös $r_k = b - Ax_k$ on kohtisuorassa avaruutta AX_k vastaan. Pythagoraan lauseesta seuraa (kuten edellisen lauseen todistuksessa), että jäännöksen 2-normi $\|r_k\|_2$ on minimissään täsmälleen silloin, kun $r_k \perp AX_k$. Saadaan seuraava lause.

Lause 3.2.3. *Olkoon $Y_k = AX_k$. Tällöin Petrov-Galerkinin yhtälön ratkaisun x_k jäännös $r_k = b - Ax_k$ on kohtisuorassa avaruutta AX_k vastaan. Tämä tapahtuu täsmälleen silloin, kun jäännöksen 2-normi on minimissään*

$$\|b - Ax_k\|_2 = \min_{y \in AX_k} \|b - Ay\|_2.$$

Seuraavaksi esitellään Krylovin aliavaruudet, joita käyttäen yhdessä Petrov-Galerkinin yhtälön kanssa saadaan muodostettua iteratiivisia ratkaisumenetelmiä yhtälölle (3.1). Liittogradienttimenetelmä on esimerkki tällaisesta menetelmästä ($Y_k = X_k$ ja A hermiittinen ja positiividefiniitti) ja sille pätee lause 3.2.2. GMRES on esimerkki menetelmästä, jolle pätee lause 3.2.3. Näistä esitellään vain GMRES, minkä jälkeen siirrytään käsittelemään reaalilineaarisia menetelmiä.

3.3 Matriisien Krylovin aliavaruudet

Valitaan X_k :ksi Krylovin aliavaruudet

$$\mathcal{K}_k(\mathbf{A}, \mathbf{v}) = \text{span} \left\{ \mathbf{v}, \mathbf{A}\mathbf{v}, \mathbf{A}^2\mathbf{v}, \dots, \mathbf{A}^{k-1}\mathbf{v} \right\}, \quad k = 1, 2, \dots$$

Näille pätee $\mathcal{K}_k(\mathbf{A}, \mathbf{v}) \subset \mathcal{K}_{k+1}(\mathbf{A}, \mathbf{v})$ kaikilla k . Vektorin \mathbf{v} minimipolynomi \mathbf{A} :n suhteen on pienintä astetta oleva mooninen polynomi p (korkeimman asteen termin kerroin on yksi) siten, että $p(\mathbf{A})\mathbf{v} = \mathbf{0}$. Tällainen minimipolynomi on aina olemassa, koska erityisesti \mathbf{A} :n karakteristisella polynomilla p on $p(\mathbf{A}) = 0$ Cayley-Hamiltonin lauseen mukaan. Seuraava lause kertoo, että Krylovin aliavaruudet kasvavat aidosti minimipolynomin astelukuun asti eivätkä sen jälkeen muutu.

Lause 3.3.1. *Olkoot $\mathbf{A} \in \mathbb{C}^{n \times n}$, $\mathbf{v} \in \mathbb{C}^n$ ja p on \mathbf{v} :n minimipolynomi. Kun merkitään $\mu = \deg p$, niin*

- (1) $\dim \mathcal{K}_k(\mathbf{A}, \mathbf{v}) = k$, kun $k \leq \mu$,
- (2) $\mathcal{K}_k(\mathbf{A}, \mathbf{v}) = \mathcal{K}_\mu(\mathbf{A}, \mathbf{v})$, kun $k \geq \mu$.

Todistus. Olkoon $k \leq \mu$. Selvästi $\dim \mathcal{K}_k(\mathbf{A}, \mathbf{v}) \leq k$. Kun jokin kertoimista c_0, c_1, \dots, c_{k-1} poikkeaa nolasta, täytyy olla $c_0\mathbf{v} + c_1\mathbf{A}\mathbf{v} + \dots + c_{k-1}\mathbf{A}^{k-1}\mathbf{v} \neq \mathbf{0}$. Muutoin tästä voitaisiin muodostaa minimipolynomi, jonka asteluku olisi pienempi kuin μ vastoin p :n valintaa. Vektorit $\mathbf{v}, \mathbf{A}\mathbf{v}, \dots, \mathbf{A}^{k-1}\mathbf{v}$ ovat siten lineaarisesti riippumattomat, joten kohta (1) on osoitettu.

Olkoon sitten $k \geq \mu$. Minimipolynomin mukaan $\mathbf{A}^\mu\mathbf{v}$ on lineaarikombinaatio vektoreista $\mathbf{v}, \mathbf{A}\mathbf{v}, \dots, \mathbf{A}^{\mu-1}\mathbf{v}$. Huomaamalla, että $\mathbf{A}^j\mathbf{v} = \mathbf{A}\mathbf{A}^{j-1}\mathbf{v}$, saadaan induktiolla, että myös $\mathbf{A}^k\mathbf{v}$ on lineaarikombinaatio samoista vektoreista. Kohta (2) on täten osoitettu. \square

Krylovin aliavaruuksiin $\mathcal{K}_k(\mathbf{A}, \mathbf{v})$ voidaan muodostaa ortonormaali kanta Arnoldin menetelmällä. Olkoon ensimmäinen kantavektori $\mathbf{q}_1 = \mathbf{v}/\|\mathbf{v}\|$. Oletetaan sitten, että kantavektorit $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_j$ on laskettu ja näytetään kuinka \mathbf{q}_{j+1} saadaan. Lasketaan vektori $\mathbf{A}\mathbf{q}_j$ ja ortogonalisoidaan se käyttäen Gram-Schmidt menetelmää

$$h_{ij} = \langle \mathbf{A}\mathbf{q}_j, \mathbf{q}_i \rangle \quad (i = 1, 2, \dots, j) \quad (3.6)$$

$$\mathbf{w}_j = \mathbf{A}\mathbf{q}_j - \sum_{i=1}^j h_{ij}\mathbf{q}_i \quad (3.7)$$

$$h_{j+1,j} = \|\mathbf{w}_j\| \quad (3.8)$$

Menetelmä pysäytetään, jos $\mathbf{w}_j = \mathbf{0}$. Tällöin sanotaan menetelmän murtuneen ja tätä j :nnettä askelta murtumisaskelleeksi. Muutoin asetetaan $\mathbf{q}_{j+1} = \mathbf{w}_j/\|\mathbf{w}_j\|$.

Algoritmi 3.3.1 laskee saman kannan käyttäen muunnettua Gram-Schmidt menetelmää. Tämä on tarkassa aritmetiikassa yhtäpitävä edellä kuvatun kanssa, mutta on numeerisessa laskennassa tapahtuvien pyöristysvirheiden läsnäollessa luotettavampi. Vektorin $\mathbf{A}\mathbf{q}_j$ ortogonalisointi tapahtuu siinä yksitellen, jolloin vektoreissa $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_j$ jo oleva pyöristysvirheistä aiheutunut epäortogonaalisuus tulee huomioiduksi.

Algoritmi 3.3.1 Arnoldin menetelmä ortonormaalien kannan muodostamiseksi Krylovin aliaruudelle $\mathcal{K}_k(\mathbf{A}, \mathbf{v})$ käyttäen muunnettua Gram-Schmidt menetelmää. Laskee ortonormaalit kantavektorit $\mathbf{q}_1, \mathbf{q}_2, \dots$ ja tallentaa ortogonalisoinnissa käytetyt sisätulot matriisiin $\mathbf{H} = (h_{ij})$. Algoritmi pysähtyy muodostettuaan $\min\{k + 1, \dim \mathcal{K}_k(\mathbf{A}, \mathbf{v})\}$ kantavektoria.

```

1:  $\mathbf{q}_1 \leftarrow \mathbf{v}/\|\mathbf{v}\|$ 
2: for  $j = 1$  to  $k$  do
3:    $\mathbf{w} \leftarrow \mathbf{A}\mathbf{q}_j$ 
4:   for  $i = 1$  to  $j$  do
5:      $h_{ij} \leftarrow \langle \mathbf{w}, \mathbf{q}_i \rangle$ 
6:      $\mathbf{w} \leftarrow \mathbf{w} - h_{ij}\mathbf{q}_i$ 
7:   end for
8:    $h_{j+1,j} \leftarrow \|\mathbf{w}\|$ 
9:   if  $h_{j+1,j} = 0$  then
10:    lopeta algoritmi
11:  end if
12:   $\mathbf{q}_{j+1} \leftarrow \mathbf{w}/h_{j+1,j}$ 
13: end for

```

Lause 3.3.2. Oletetaan, että Arnoldin menetelmällä on saatu muodostettua vektorit $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_k$. Nämä vektorit muodostavat ortonormaalien kannan avaruuteen $\mathcal{K}_k(\mathbf{A}, \mathbf{v})$.

Todistus. Vektorit ovat ortonormaalit Gram-Schmidt prosessin ja normeerauksen myötä. Osoitetaan, että jokaiselle \mathbf{q}_j on olemassa astetta $j - 1$ oleva polynomi p_{j-1} siten, että $\mathbf{q}_j = p_{j-1}(\mathbf{A})\mathbf{v}$. Tämä on selvästi totta \mathbf{q}_1 :lle. Oletetaan sitten, että näin on $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_j$:lle ja osoitetaan todeksi myös \mathbf{q}_{j+1} :lle. Kaavasta (3.7) saadaan

$$\mathbf{q}_{j+1} = \frac{1}{\|\mathbf{w}_j\|} \mathbf{A}\mathbf{q}_j - \sum_{i=1}^j \frac{1}{\|\mathbf{w}_j\|} h_{ij} \mathbf{q}_i = \frac{1}{\|\mathbf{w}_j\|} \mathbf{A}p_{j-1}(\mathbf{A})\mathbf{v} - \sum_{i=1}^j \frac{1}{\|\mathbf{w}_j\|} h_{ij} p_{i-1}(\mathbf{A})\mathbf{v},$$

mistä nähdään, että oikealla puolella on j -asteinen polynomi.

Ollaan saatu, että $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_k$ kuuluvat avaruuteen $\mathcal{K}_k(\mathbf{A}, \mathbf{v})$. Krylovin aliaruuden määritelmästä saadaan $\dim \mathcal{K}_k(\mathbf{A}, \mathbf{v}) \leq k$, joten vektorit virittävät avaruuden ja $\dim \mathcal{K}_k(\mathbf{A}, \mathbf{v}) = k$. Lisäksi nähdään, että \mathbf{v} :n minimipolynomi on vähintään astetta k . \square

Lause 3.3.3. Arnoldin menetelmässä $\mathbf{w}_j = 0$ jos ja vain jos j on \mathbf{v} :n minimipolynomin asteluku.

Todistus. Oletetaan aluksi $\mathbf{w}_j = 0$. Kuten lauseen 3.3.2 todistuksessa, olkoon p_{j-1} astetta j oleva polynomi siten, että $\mathbf{q}_j = p_{j-1}(\mathbf{A})\mathbf{v}$. Kaavasta (3.7) saadaan

$$0 = \mathbf{w}_j = \mathbf{A}\mathbf{q}_j - \sum_{i=1}^j h_{ij} \mathbf{q}_i = \mathbf{A}p_{j-1}(\mathbf{A})\mathbf{v} - \sum_{i=1}^j h_{ij} p_{i-1}(\mathbf{A})\mathbf{v},$$

jonka oikealla puolella on j -asteinen polynomi. Siispä \mathbf{v} :n minimipolynomin asteluku on korkeintaan j . Koska $\mathbf{w}_{j-1} \neq 0$, on toisaalta saatu muodostettua $\mathbf{q}_1, \dots, \mathbf{q}_j$, joten

lauseen 3.3.2 perusteella $\dim \mathcal{K}_j(\mathbf{A}, \mathbf{v}) = j$ ja lauseen 3.3.1 mukaan minimipolynomin asteluku ei voi olla pienempi kuin j .

Oletetaan sitten, että j on \mathbf{v} :n minimipolynomin asteluku. Lauseiden 3.3.1 ja 3.3.2 perusteella Arnoldin menetelmällä on mahdotonta muodostaa enemmän kuin j ortonormaalia vektoria. Täytyy siis olla $\mathbf{w}_i = 0$ jollakin $i \leq j$. Toisaalta, jos $i < j$, niin tämän todistuksen ensimmäisen osan mukaan asteluvunkin täytyisi olla i . Täten $\mathbf{w}_j = 0$. \square

Edeltävien lauseiden sisältönä on, että Arnoldin menetelmä toimii niin pitkälle kuin mahdollista. Se pysähtyy täsmälleen silloin, kun Krylovin aliavaruuksien jono ($k = 1, 2, \dots$) lakkaa kasvamasta.

Olkoot $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_k$ muodostettu Arnoldin menetelmällä ja asetetaan ne sarakkeiksi matriisiin $\mathbf{Q}_k = [\mathbf{q}_1 \ \mathbf{q}_2 \ \dots \ \mathbf{q}_k]$. Olkoon $\widehat{\mathbf{H}}_k$ Arnoldin menetelmän alkioista h_{ij} muodostettu $(k+1) \times k$ -matriisi, $(\widehat{\mathbf{H}}_k)_{ij} = h_{ij}$ ($i \leq j+1$) ja muut alkiot nolliina. Olkoon vielä \mathbf{H}_k se $k \times k$ -matriisi, kun matriisista $\widehat{\mathbf{H}}_k$ poistetaan viimeinen rivi. Kaavasta (3.7) saadaan

$$h_{j+1,j} \mathbf{q}_{j+1} = \mathbf{A} \mathbf{q}_j - \sum_{i=1}^j h_{ij} \mathbf{q}_i \quad \Leftrightarrow \quad \mathbf{A} \mathbf{q}_j = \sum_{i=1}^{j+1} h_{ij} \mathbf{q}_i,$$

joka voidaan kirjoittaa matriisimuotoihin

$$\mathbf{A} \mathbf{Q}_j = \mathbf{Q}_{j+1} \widehat{\mathbf{H}}_j, \quad (3.9)$$

$$\mathbf{A} \mathbf{Q}_j = \mathbf{Q}_j \mathbf{H}_j + \mathbf{w}_j \mathbf{e}_j^T, \quad (3.10)$$

missä \mathbf{e}_j on \mathbb{C}^j :n j :s luonnollinen kantavektori. Kertomalla jälkimmäinen puolittain \mathbf{Q}_j^* :llä ja hyödyntämällä \mathbf{Q}_j :n sarakkeiden ortonormalisuutta, saadaan lisäksi

$$\mathbf{Q}_j^* \mathbf{A} \mathbf{Q}_j = \mathbf{H}_j. \quad (3.11)$$

3.4 GMRES-menetelmä

GMRES (Generalized Minimal Residual) menetelmässä likimääräisratkaisua \mathbf{x}_k etsitään avaruudesta $\mathbf{x}_0 + \mathcal{K}_k(\mathbf{A}, \mathbf{v})$ ja vaaditaan, että jäännösvektorin $\mathbf{r}_k = \mathbf{b} - \mathbf{A} \mathbf{x}_k$ 2-normi $\|\mathbf{r}_k\|$ on pienin mahdollinen. Tässä $\mathbf{v} = \mathbf{b} - \mathbf{A} \mathbf{x}_0$ ja \mathbf{x}_0 on alkuarvaus. Lauseen 3.2.3 mukaan normin $\|\mathbf{r}_k\|$ minimointi on yhtäpitävä kohtisuoruuden ehdon (3.3) kanssa, kun merkitään $\mathbf{X}_k = \mathcal{K}_k(\mathbf{A}, \mathbf{v})$, $\mathbf{Y}_k = \mathbf{A} \mathbf{X}_k$.

Valitaan k ja lasketaan Arnoldin menetelmällä - käyttäen esimerkiksi algoritmia 3.3.1 - ortonormaalit kantavektorit Krylovin aliavaruudelle $\mathcal{K}_k(\mathbf{A}, \mathbf{v})$ ja asetetaan ne matriisiin \mathbf{Q}_k sarakkeiksi. Tällöin yhtälöstä (3.9)

$$\begin{aligned} \|\mathbf{r}_k\| &= \|\mathbf{b} - \mathbf{A} \mathbf{x}_k\| = \|\mathbf{v} - \mathbf{A} \mathbf{Q}_k \mathbf{z}_k\| = \|\mathbf{v} - \mathbf{Q}_{k+1} \widehat{\mathbf{H}}_k \mathbf{z}_k\| \\ &= \left\| \mathbf{Q}_{k+1} (\|\mathbf{v}\| \mathbf{e}_1 - \widehat{\mathbf{H}}_k \mathbf{z}_k) \right\| = \left\| \|\mathbf{v}\| \mathbf{e}_1 - \widehat{\mathbf{H}}_k \mathbf{z}_k \right\|, \end{aligned} \quad (3.12)$$

missä jälkimmäisin yhtäsuuruus seuraa \mathbf{Q}_{k+1} :n sarakkaiden ortonormalisuudesta. Jäljelle on jäänyt pienimmän neliösumman tehtävä. On löydettävä $\mathbf{z}_k \in \mathbb{C}^k$ siten,

että $\|\|\mathbf{v}\|e_1 - \widehat{\mathbf{H}}_k \mathbf{z}_k\|$ on pienin. Tässä $\widehat{\mathbf{H}}_k = [h_{ij}]$ on $(k+1) \times k$ -matriisi, jonka alkioit ovat nollia $h_{ij} = 0$, kun $i > j + 1$.

Tehtävä voidaan ratkaista laskemalla QR-hajotelma Givensin rotaatioilla. Olkoon $\mathbf{G} \in \mathbb{C}^{(k+1) \times (k+1)}$ laskettu ortonormaali matriisi siten, että $\mathbf{G}\widehat{\mathbf{H}}_k = \widehat{\mathbf{R}}_k$ on yläkolmiomatriisi. Merkitään $\widehat{\boldsymbol{\eta}}_k = \mathbf{G}e_1$ ja olkoon vektori $\boldsymbol{\eta}_k$ tämän k ensimmäistä komponenttia ja ζ_k viimeinen komponentti, jolloin $\widehat{\boldsymbol{\eta}}_k = \begin{bmatrix} \boldsymbol{\eta}_k \\ \zeta_k \end{bmatrix}$. Merkitään vielä \mathbf{R}_k :lla sitä $k \times k$ -matriisia, joka saadaan poistamalla $\widehat{\mathbf{R}}_k$:sta viimeinen (nolla)rivi.

Käyttäen \mathbf{G} :n ortonormalisuutta, tulee (3.12) muotoon

$$\begin{aligned} \|\mathbf{r}_k\|^2 &= \|\|\mathbf{v}\|e_1 - \widehat{\mathbf{H}}_k \mathbf{z}_k\|^2 = \|\mathbf{G}(\|\mathbf{v}\|e_1 - \widehat{\mathbf{H}}_k \mathbf{z}_k)\|^2 \\ &= \|\|\mathbf{v}\|\widehat{\boldsymbol{\eta}}_k - \widehat{\mathbf{R}}_k \mathbf{z}_k\|^2 = \|\|\mathbf{v}\|\boldsymbol{\eta}_k - \mathbf{R}_k \mathbf{z}_k\|^2 + \|\mathbf{v}\|^2 |\zeta_k|^2. \end{aligned} \quad (3.13)$$

Kun Arnoldin menetelmä ei ole vielä murtunut, on $\widehat{\mathbf{H}}_k$:n rangi määrittäen (3.8) perusteella k . Kun \mathbf{A} on kääntövä, on $\widehat{\mathbf{H}}_k$:n rangi k myös murtumisaskeleella yhtälöstä (3.10) johtuen. Näissä tapauksissa $\widehat{\mathbf{R}}_k$:n ja \mathbf{R}_k :n rangi on myös k , joten \mathbf{R}_k on kääntövä. Ratkaistaan yhtälö

$$\mathbf{R}_k \mathbf{z}_k = \|\mathbf{v}\|\boldsymbol{\eta}_k$$

taaksepäin sijoituksin, jonka jälkeen GMRES-menetelmän likimääräisratkaisu saadaan kaavasta $\mathbf{x}_k = \mathbf{x}_0 + \mathbf{Q}_k \mathbf{z}_k$. Tällöin $\|\mathbf{r}_k\| = \|\mathbf{v}\||\zeta_k|$.

Huomataan lisäksi, että Arnoldin menetelmän murtuessa $h_{k+1,k} = 0$ ($\mathbf{w}_k = 0$), jolloin GMRES antaa tarkan ratkaisun yhtälölle. Tällöin \mathbf{G}_k on yksikkömatriisi, joten $\zeta_k = 0$ ja kaavasta (3.13) nähdään $\mathbf{r}_k = 0$.

Seuraava lause antaa mahdollisuuden arvioida GMRESin suppenemisnopeutta.

Lause 3.4.1. *Olkoon \mathbf{A} diagonalisoituva, $\mathbf{A} = \mathbf{X}\boldsymbol{\Lambda}\mathbf{X}^{-1}$. Tällöin jäännökselle pätee arvio*

$$\|\mathbf{r}_k\| \leq \|\mathbf{r}_0\| \kappa_2(\mathbf{X}) \min_{p \in \mathbb{P}_k, p(0)=1} \max_{\lambda \in \sigma(\mathbf{A})} |p(\lambda)|,$$

missä \mathbb{P}_k on korkeintaan k -asteisten kompleksikertoimisten polynomien joukko ja $\kappa_2(\mathbf{X}) = \|\mathbf{X}\|_2 \|\mathbf{X}^{-1}\|_2$ on häiriöalttius 2-normissa.

Todistus. Olkoon p jokin korkeintaan k -asteinen polynomi siten, että $p(0) = 1$. Asetetaan $q(\lambda) = (1 - p(\lambda))/\lambda$ ja $\mathbf{x} = \mathbf{x}_0 + q(\mathbf{A})\mathbf{r}_0$. Tällöin

$$\mathbf{b} - \mathbf{A}\mathbf{x} = \mathbf{b} - \mathbf{A}\mathbf{x}_0 - (\mathbf{I} - p(\mathbf{A}))\mathbf{r}_0 = p(\mathbf{A})\mathbf{r}_0.$$

Koska q on $(k-1)$ -asteinen polynomi, kuuluu $\mathbf{x} \in \mathbf{x}_0 + \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$. GMRES minimoi jäännöksen normin $\|\mathbf{r}_k\|$ avaruudessa $\mathbf{x}_0 + \mathcal{K}_k(\mathbf{A}, \mathbf{r}_0)$, joten

$$\begin{aligned} \|\mathbf{r}_k\| &\leq \|\mathbf{b} - \mathbf{A}\mathbf{x}\| = \|p(\mathbf{A})\mathbf{r}_0\| = \|\mathbf{X}p(\boldsymbol{\Lambda})\mathbf{X}^{-1}\mathbf{r}_0\| \\ &\leq \|\mathbf{X}\| \|p(\boldsymbol{\Lambda})\| \|\mathbf{X}^{-1}\| \|\mathbf{r}_0\| = \|\mathbf{r}_0\| \kappa_2(\mathbf{X}) \max_{\lambda \in \sigma(\mathbf{A})} |p(\lambda)|, \end{aligned}$$

koska diagonaalimatriisin 2-normi on itseisarvoltaan suurin diagonaalialkio. On lisäksi helppo nähdä, että oikeanpuolimmaisena lausekkeen minimi saavutetaan jollakin polynomilla p . \square

Algoritmi 3.4.1 GMRES-menetelmä yhtälön $\mathbf{Ax} = \mathbf{b}$ likimääräiseksi ratkaisemiseksi. Suorittaa k iteraatiota tai pysähtyy, jos tarkka ratkaisu saavutetaan.

```

1:  $\mathbf{r}_0 \leftarrow \mathbf{b} - \mathbf{Ax}_0$ ,  $\mathbf{q}_1 \leftarrow \mathbf{r}_0 / \|\mathbf{r}_0\|$ ,  $\hat{\boldsymbol{\eta}} \leftarrow [1] \in \mathbb{C}^{1 \times 1}$ 
2: for  $j = 1$  to  $k$  do
3:    $\mathbf{w} \leftarrow \mathbf{Aq}_j$ 
4:   for  $i = 1$  to  $j$  do
5:      $h_{ij} \leftarrow \langle \mathbf{w}, \mathbf{q}_i \rangle$ 
6:      $\mathbf{w} \leftarrow \mathbf{w} - h_{ij}\mathbf{q}_i$ 
7:   end for
8:    $h_{j+1,j} \leftarrow \|\mathbf{w}\|$ 
9:    $\mathbf{s} \leftarrow [h_{1j} \ h_{2j} \ \dots \ h_{j+1,j}]^T$ 
10:  for  $i = 1$  to  $j - 1$  do
11:     $\mathbf{s}_{i:(i+1)} \leftarrow \mathbf{G}_i \mathbf{s}_{i:(i+1)}$ 
12:  end for
13:  Olkoon  $\mathbf{G}_j \in \mathbb{C}^{2 \times 2}$  Givensin rotaatio siten, että  $\mathbf{G}_j \mathbf{s}_{j:(j+1)} = [\# \ 0]^T$ 
14:   $\mathbf{s}_{j:(j+1)} \leftarrow \mathbf{G}_j \mathbf{s}_{j:(j+1)}$ 
15:   $\hat{\mathbf{R}} \leftarrow \begin{bmatrix} \hat{\mathbf{R}} \\ 0 \end{bmatrix} \in \mathbb{C}^{(j+1) \times (j-1)}$ ,  $\hat{\mathbf{R}} \leftarrow \begin{bmatrix} \hat{\mathbf{R}} & \mathbf{s} \end{bmatrix}$ ,  $\hat{\boldsymbol{\eta}} \leftarrow \begin{bmatrix} \hat{\boldsymbol{\eta}} \\ 0 \end{bmatrix}$ ,  $\hat{\boldsymbol{\eta}}_{j:(j+1)} \leftarrow \mathbf{G}_j \hat{\boldsymbol{\eta}}_{j:(j+1)}$ 
16:  if  $h_{j+1,j} = 0$  then
17:    break
18:  end if
19:   $\mathbf{q}_{j+1} \leftarrow \mathbf{w} / h_{j+1,j}$ 
20: end for
21:  $\mathbf{Q} \leftarrow [\mathbf{q}_1 \ \dots \ \mathbf{q}_j]$ ,  $\mathbf{R} \leftarrow \hat{\mathbf{R}}_{1:j,*}$ ,  $\boldsymbol{\eta} \leftarrow \hat{\boldsymbol{\eta}}_{1:j}$ 
22: Ratkaistaan  $\mathbf{z}$  yhtälöstä  $\mathbf{Rz} = \|\mathbf{r}_0\| \boldsymbol{\eta}$ 
23:  $\mathbf{x} \leftarrow \mathbf{x}_0 + \mathbf{Qz}$ , jäännöksen normi on  $\|\mathbf{r}_0\| \|\hat{\boldsymbol{\eta}}_{j+1}\|$ 

```

Seuraus 3.4.2. *Olkoon \mathbf{A} diagonalisoituva ja kuulukoon sen ominaisarvot ζ -keskiseen R -säteiseen suljettuun kiekkoon. Oletetaan lisäksi, että origo ei kuulu tähän kiekkoon. Tällöin*

$$\|\mathbf{r}_k\| \leq \|\mathbf{r}_0\| \kappa_2(\mathbf{X}) \left(\frac{R}{|\zeta|} \right)^k.$$

Todistus. Olkoon polynomi $p(\lambda) = (1 - \lambda/\zeta)^k$. Saadaan

$$|p(\lambda)| = \left(\frac{1}{|\zeta|} |\zeta - \lambda| \right)^k \leq \left(\frac{R}{|\zeta|} \right)^k.$$

□

Edellisen lauseen ja sen seurauksen sisältämä häiriöalttius $\kappa_2(\mathbf{X})$ on yleensä vaikea laskea, mikä tekee arvioiden käyttämisestä hankalaa. Kuitenkin normaalimatriisille (tai lähes normaaleille) pätee $\kappa_2(\mathbf{X}) = 1$, koska matriisi on normaali jos ja vain jos se on diagonalisoituva ortonormeeratulla muunnosmatriisilla \mathbf{X} .

3.4.1 Uudelleenkäynnistys

GMRESin suorittama täysi ortogonalisointi vaatii kaikkien vektorien $\mathbf{q}_j \in \mathbb{C}^n$ tallentamisen tietokoneen muistiin. Vaadittu muistin määrä kierrosten k kasvaessa on siten $O(nk)$, jolloin tietokoneen muistikapasiteetin rajat saattavat tulla vastaan laskennan edistyessä. Jokaisen kierroksen laskenta myös hidastuu, koska uusi vektori täytyy ortogonalisoida kaikkia aikaisempia vastaan. Näistä syistä voi olla edullista käynnistää algoritmi uudelleen käyttäen alkuarvauksena \mathbf{x}_0 edellisen algoritmin ajon likimääräisratkaisua \mathbf{x}_k . Merkitään k :n iteraatiokierroksen välein uudelleenkäynnistettyä menetelmää GMRES(k):lla.

Lause 3.4.3. Jos $\langle \mathbf{Ax}, \mathbf{x} \rangle \neq 0$ kaikilla $\mathbf{x} \neq 0$, niin GMRES(k) suppenee.

Todistus. Olkoon \mathbf{x}_0 alkuarvaus ja $\mathbf{r}_0 = \mathbf{b} - \mathbf{Ax}_0$. Olkoon \mathbf{x}_1 ensimmäinen GMRESin antama likimääräisratkaisu eli $\|\mathbf{b} - \mathbf{Ax}\|$ ($\mathbf{x} \in \mathbf{x}_0 + \mathcal{K}_1(\mathbf{A}, \mathbf{r}_0)$) on minimissään, kun $\mathbf{x} = \mathbf{x}_1$. Tällöin $\mathbf{x}_1 = \mathbf{x}_0 + \alpha \mathbf{r}_0$ ja jäännös $\mathbf{r}_1 = \mathbf{b} - \mathbf{Ax}_1 = \mathbf{r}_0 - \alpha \mathbf{Ar}_0$. Normille

$$\|\mathbf{r}_1\|^2 = \langle \mathbf{r}_0 - \alpha \mathbf{Ar}_0, \mathbf{r}_0 - \alpha \mathbf{Ar}_0 \rangle = \langle \mathbf{r}_0 - \alpha \mathbf{Ar}_0, \mathbf{r}_0 \rangle - \bar{\alpha} \langle \mathbf{r}_0 - \alpha \mathbf{Ar}_0, \mathbf{Ar}_0 \rangle, \quad (3.14)$$

missä $\langle \mathbf{r}_0 - \alpha \mathbf{Ar}_0, \mathbf{Ar}_0 \rangle = 0$, koska $\mathbf{r}_1 \perp \mathcal{AK}_1(\mathbf{A}, \mathbf{r}_0)$. Tästä saadaan myös

$$\alpha = \frac{\langle \mathbf{r}_0, \mathbf{Ar}_0 \rangle}{\langle \mathbf{Ar}_0, \mathbf{Ar}_0 \rangle}.$$

Jatkaen kohdasta (3.14)

$$\|\mathbf{r}_1\|^2 = \|\mathbf{r}_0\|^2 - \alpha \langle \mathbf{Ar}_0, \mathbf{r}_0 \rangle = \|\mathbf{r}_0\|^2 \left(1 - \frac{|\langle \mathbf{Ar}_0, \mathbf{r}_0 \rangle|^2 \langle \mathbf{r}_0, \mathbf{r}_0 \rangle}{\langle \mathbf{r}_0, \mathbf{r}_0 \rangle^2 \langle \mathbf{Ar}_0, \mathbf{Ar}_0 \rangle} \right).$$

Asetetaan $\mu = \min_{\|\mathbf{x}\|=1} |\langle \mathbf{Ax}, \mathbf{x} \rangle|$ ja käytetään epäyhtälöä $\|\mathbf{Ar}_0\| \leq \|\mathbf{A}\| \|\mathbf{r}_0\|$, jolloin

$$\|\mathbf{r}_1\|^2 \leq \|\mathbf{r}_0\|^2 \left(1 - \frac{\mu^2}{\|\mathbf{A}\|^2} \right).$$

Täten GMRESin ensimmäinen iteraatio pienentää jäännöksen normia $\|\mathbf{r}_0\|$ ainakin tekijällä $(1 - \mu^2/\|\mathbf{A}\|^2)^{1/2} < 1$. Muissa iteraatioissa jäännöksen normi ei kasva. Seuraavissakin uudelleenkäynnistyksissä ensimmäinen iteraatio aina pienentää normia samalla tekijällä. Näin ollen uudelleenkäynnistetty GMRES suppenee. \square

3.5 \mathbb{R} -lineaariset menetelmät

Olkoon $\mathcal{M} : \mathbb{C}^n \rightarrow \mathbb{C}^n$, $\mathcal{M}(\mathbf{z}) = \mathbf{M}\mathbf{z} + \mathbf{M}_{\#}\bar{\mathbf{z}}$ reaali-lineaarinen operaattori, $\mathbf{b} \in \mathbb{C}^n$ ja tuntematon vektori \mathbf{z} on ratkaistava yhtälöstä

$$\mathcal{M}(\mathbf{z}) = \mathbf{b}. \quad (3.15)$$

Tällainen yhtälö voitaisiin ratkaista muodostamalla reaali-matriisiesitys ja käyttämällä edellä kuvattuja menettelyjä matriisiyhtälöiden ratkaisemiseksi. Pitämällä \mathbb{R} -lineaarinen yhtälö annetussa muodossaan tarjoaa kuitenkin etuja. Etsittäessä reaali-matriisiesityksen likimääräisratkaisua jostakin k :n vektorin virittämästä \mathbb{R}^{2n} :n aliavaruudesta tilanne voidaan siirtää alkuperäisen yhtälön (3.15) ratkaisemiseen.

Viritettyä aliavaruutta vastaa tällöin k :n vektorin virittämä \mathbb{R} -kertoiminen \mathbb{C}^n :n aliavaruus. Tässä kappaleessa esitellään menetelmiä, joissa ratkaisua etsitään näiden k :n vektorin virittämästä \mathbb{C} -kertoimisesta aliavaruudesta. Likimääräisratkaisun haakuvaruus kullakin iteraatiolla on siten laajempi ja pyrkimyksenä on saavuttaa reaalimatriisimenetelmiä nopeampi suppeneminen.

Yleisen reaalin lineaarisen operaattori sijaan keskitytään pääasiassa muotoa $\mathcal{M}_\kappa(\mathbf{z}) = \kappa\mathbf{z} + \mathbf{M}_\# \bar{\mathbf{z}}$ ($\kappa \in \mathbb{C}$) oleviin operaattoreihin, koska näiden ominaisuudet soveltuvat Krylovin aliavaruuksien käyttöön. Tämän muotoisilla operaattoreilla on sovelluksia mm. lääketieteellisessä kuvantamisessa. Sähköisessä impedanssitomografiassa muodostetaan poikkileikkauskuva esimerkiksi potilaan rintakehästä syöttämällä ihon pinnalle sähkövirtoja ja mittaamalla vasteena saatu sähköpotentiaali. Kyseessä on inversio-ongelma, jonka ratkaisuna on sähköjohtavuus poikkileikkauksen jokaisessa pisteessä. Kuva keuhkoista ja sydäimestä voidaan muodostaa, koska eri kudoksilla on erilaiset johtavuudet. Ongelmaa ratkaistaessa joudutaan heikosti singulaarisiin Fredholmian toisen lajin integraaliyhtälöihin, jolloin diskretoinnin myötä päädytään yhtälöihin muotoa $\mathbf{z} + \mathbf{M}_\# \bar{\mathbf{z}} = \mathbf{1}$, missä $\mathbf{1}$ on luvuista 1 koostuva vektori. Tässä matriiseja $\mathbf{M}_\#$ ei tallenneta tietokoneen muistiin vaan operaatio $\mathbf{z} \rightarrow \mathbf{M}_\# \bar{\mathbf{z}}$ on peräisin funktioiden konvoluutiosta ja diskretoitu operaatio voidaan laskea käyttäen nopeaa Fourier-muunnosta (FFT). Yhtälöt ovat yleensä vähintään kokoa $n = 2^{16}$ ja yhden kuvan laskemiseen sellainen joudutaan ratkaisemaan jokaisessa diskretointipisteessä. Suorat ratkaisumenetelmät eivät tule kyseeseen ja iteratiivisissa menetelmissä suppenemisnopeuden parantaminen antaa kuvan nopeammin lääkärin käyttöön.

Muotoillaan seuraavaksi Petrov-Galerkinin yhtälö matriiseja vastaavalla tavalla. Olkoon X_k jokin \mathbb{C}^n :n \mathbb{C} -kertoiminen aliavaruus ja Y_k jokin \mathbb{R} -kertoiminen aliavaruus. Kun $\mathbf{z}_0 \in \mathbb{C}^n$ on ratkaisun alkuarvaus, Petrov-Galerkinin ehto likimääräisratkaisuksi $\mathbf{z}_k \in \mathbf{z}_0 + X_k$ on

$$\langle \mathbf{b} - \mathcal{M}(\mathbf{z}_k), \mathbf{y} \rangle = 0 \quad \text{kaikilla } \mathbf{y} \in Y_k. \quad (3.16)$$

Merkitään $\mathbf{r}_0 = \mathbf{b} - \mathcal{M}(\mathbf{z}_0)$. Kun matriisien \mathbf{V}_k ja \mathbf{W}_k sarakkeet muodostavat kannat avaruuksille X_k ja Y_k (vastaavasti), niin edellä oleva ehto voidaan kirjoittaa muotoon

$$\mathbf{W}_k^* \mathbf{M} \mathbf{V}_k \mathbf{u}_k + \mathbf{W}_k^* \mathbf{M}_\# \bar{\mathbf{V}}_k \bar{\mathbf{u}}_k = \mathbf{W}_k^* \mathbf{r}_0, \quad (3.17)$$

mistä ratkaisemalla \mathbf{u}_k saadaan likimääräisratkaisu alkuperäiselle yhtälölle (3.15) laskemalla

$$\mathbf{z}_k = \mathbf{z}_0 + \mathbf{V}_k \mathbf{u}_k. \quad (3.18)$$

3.6 Krylovin aliavaruudet

Krylovin aliavaruudet voidaan määritellä matriiseja vastaavalla tavalla. Olkoon

$$\mathcal{K}_k(\mathcal{M}, \mathbf{r}_0) = \text{span} \left\{ \mathbf{r}_0, \mathcal{M}(\mathbf{r}_0), \mathcal{M}^2(\mathbf{r}_0), \dots, \mathcal{M}^{k-1}(\mathbf{r}_0) \right\} \quad (k = 1, 2, \dots).$$

Vektorin \mathbf{r}_0 minimipolynomi \mathcal{M} :n suhteen on pienintä astetta oleva kompleksiker-toiminen mooninen polynomi p siten, että $p(\mathcal{M})(\mathbf{r}_0) = 0$. Kun operaattorille \mathcal{M} muodostetaan sitä vastaava reaalin $2n \times 2n$ -matriisi \mathbf{A} , niin Cayley-Hamiltonin lauseen mukaan on olemassa (reaaliker-toiminen) polynomi p siten, että $p(\mathbf{A}) = 0$. Tällä samalla polynomilla myös $p(\mathcal{M}) = 0$, joten minimipolynomi on hyvin määritelty.

Seuraavassa keskitytään siirrettyihin antilineaarisiin operaattoreihin $\mathcal{M}_\kappa(z) = \kappa z + M_\# \bar{z}$, koska yleisemmillä reaalilineaarisilla operaattoreilla ei ole tarvittavia ominaisuuksia. Kuitenkin mikäli operaattorin $\mathcal{M}(z) = Mz + M_\# \bar{z}$ matriisi M on helposti kääntyvä, niin yhtälö $\mathcal{M}(z) = b$ voidaan kertoa puolittain M^{-1} :llä ja siten saadaan operaattori muotoon $\widetilde{\mathcal{M}}(z) = z + M^{-1}M_\# \bar{z}$.

Lemma 3.6.1. *Olkoon $\mathcal{M}_\kappa(z) = \kappa z + M_\# \bar{z}$ siirretty antilineaarinen operaattori ja k ei-negatiivinen kokonaisluku. Tällöin*

$$\mathcal{M}_\kappa^k = p(\mathcal{M}_0), \quad (3.19)$$

$$\mathcal{M}_0^k = q(\mathcal{M}_\kappa), \quad (3.20)$$

joillakin k -asteisilla moonisilla polynomeilla p ja q . Lisäksi, kun \tilde{p} on k -asteinen polynomi, niin

$$\mathcal{M}_\kappa \tilde{p}(\mathcal{M}_\kappa) = \tilde{q}(\mathcal{M}_\kappa), \quad (3.21)$$

jollakin $(k+1)$ -asteisella polynomilla \tilde{q} .

Todistus. Kirjoitetaan ensimmäistä kaavaa varten

$$\mathcal{M}_\kappa^k = (\kappa I + \mathcal{M}_0)^k = (\kappa I + \mathcal{M}_0) \cdots (\kappa I + \mathcal{M}_0).$$

Kerrotaan tämä auki ja käytetään kaavoja $\mathcal{M}_0 \circ \kappa I = \bar{\kappa} \mathcal{M}_0$ ja $\kappa I \circ \mathcal{M}_0 = \kappa \mathcal{M}_0$, niin saadaan (3.19).

Jo osoitetun kaavan perusteella on olemassa korkeintaan $(k-1)$ -asteinen polynomi r siten, että $\mathcal{M}_0^k = \mathcal{M}_\kappa^k - r(\mathcal{M}_0)$. Kaava (3.20) seuraa siten induktiosta, koska se on selvästi tosi $k=0$:lla.

Olkoon vielä \tilde{p} jokin k -asteinen polynomi. Tällöin (3.19) perusteella on olemassa k -asteinen polynomi \tilde{r} siten, että

$$\mathcal{M}_\kappa \tilde{p}(\mathcal{M}_\kappa) = \mathcal{M}_\kappa \tilde{r}(\mathcal{M}_0) = \kappa \tilde{r}(\mathcal{M}_0) + \mathcal{M}_0 \tilde{r}(\mathcal{M}_0) = \kappa \tilde{r}(\mathcal{M}_0) + \tilde{s}(\mathcal{M}_0),$$

jollakin $(k+1)$ -asteisella polynomilla \tilde{s} , koska $\mathcal{M}_0 \circ \alpha_j \mathcal{M}_0^j = \bar{\alpha}_j \mathcal{M}_0^{j+1}$. Käyttämällä vielä kaavaa (3.20), saadaan viimeinen väite. \square

Huomautus 3.6.2. Lemman perusteella erityisesti $\mathcal{K}_k(\mathcal{M}_\kappa, r_0) = \mathcal{K}_k(\mathcal{M}_0, r_0)$. Antilineaarille operaattorille $\mathcal{M}_0(z) = M_\# \bar{z}$ ja vektorille $v \in \mathbb{C}^n$ pätee

$$\begin{aligned} \mathcal{M}_0^0(v) &= v, \\ \mathcal{M}_0^{k+1}(v) &= M_\# \overline{\mathcal{M}_0^k(v)}. \end{aligned}$$

Tällöin esimerkiksi $\mathcal{M}_0(v) = v$, $\mathcal{M}_0^1(v) = M_\# \bar{v}$, $\mathcal{M}_0^2(v) = M_\# \overline{M_\# \bar{v}}$, ja $\mathcal{M}_0^3(v) = M_\# \overline{M_\# \overline{M_\# \bar{v}}}$. Täten

$$\mathcal{K}_k(\mathcal{M}_\kappa, r_0) = \text{span} \left\{ r_0, M_\# \bar{r}_0, M_\# \overline{M_\# r_0}, \dots, \mathcal{M}_0^{k-1}(r_0) \right\}.$$

Seuraava lause siirretyille antilineaarille operaattoreille vastaa matriisien lausetta

3.3.1. Lause ei yleisty kaikille reaalilineaarisille operaattoreille.

Lause 3.6.3. *Olkoon \mathcal{M}_κ siirretty antilineaarinen operaattori, $r_0 \in \mathbb{C}^n$ ja p on r_0 :n minimipolynomi. Kun merkitään $\mu = \deg p$, niin*

- (1) $\dim \mathcal{K}_k(\mathcal{M}_\kappa, \mathbf{r}_0) = k$, kun $k \leq \mu$,
 (2) $\mathcal{K}_k(\mathcal{M}_\kappa, \mathbf{r}_0) = \mathcal{K}_\mu(\mathcal{M}_\kappa, \mathbf{r}_0)$, kun $k \geq \mu$.

Lisäksi $\mathcal{K} = \mathcal{K}_\mu(\mathcal{M}_\kappa, \mathbf{r}_0)$ on \mathcal{M}_κ :n invariantti aliavaruus, $\mathcal{M}_\kappa(\mathcal{K}) \subset \mathcal{K}$.

Todistus. Ensimmäisen kohdan todistus on sama kuin lauseessa 3.3.1. Vektorin \mathbf{r}_0 minimipolynomin perusteella $\mathcal{M}_\kappa^\mu(\mathbf{r}_0) \in \mathcal{K}_\mu(\mathcal{M}_\kappa, \mathbf{r}_0)$, mistä saadaan $(\mu-1)$ -asteinen polynomi \tilde{p} siten, että

$$\mathcal{M}_\kappa^\mu(\mathbf{r}_0) = \tilde{p}(\mathcal{M}_\kappa)(\mathbf{r}_0). \quad (3.22)$$

Nyt $\mathcal{M}_\kappa^{\mu+1}(\mathbf{r}_0) = \mathcal{M}_\kappa(\mathcal{M}_\kappa^\mu(\mathbf{r}_0)) = \mathcal{M}_\kappa(\tilde{p}(\mathcal{M}_\kappa)(\mathbf{r}_0))$. Kohdan (3.21) perusteella on olemassa μ -asteinen polynomi \tilde{q} siten, että $\mathcal{M}_\kappa^{\mu+1}(\mathbf{r}_0) = \tilde{q}(\mathcal{M}_\kappa)(\mathbf{r}_0)$. Käyttämällä vielä kaavaa (3.22), saadaan $\mathcal{M}_\kappa^{\mu+1}(\mathbf{r}_0) \in \mathcal{K}_\mu(\mathcal{M}_\kappa, \mathbf{r}_0)$.

Koska $\mathcal{M}_\kappa^{\mu+2}(\mathbf{r}_0) = \mathcal{M}_\kappa(\mathcal{M}_\kappa^{\mu+1}(\mathbf{r}_0))$, voidaan edellistä jatkaa ja kohdan 2 väite seuraa induktiosta. Avaruuden $\mathcal{K}_\mu(\mathcal{M}_\kappa, \mathbf{r}_0)$ invarianssi on selvä edeltävän perusteella. \square

Ryhdytään sitten laskemaan ortonormaaleja vektoreita noudattamalla matriisien Arnoldin menetelmää. Valitaan $\mathbf{q}_1 = \mathbf{r}_0 / \|\mathbf{r}_0\|$. Kun kantavektorit $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_j$ on laskettu, niin \mathbf{q}_{j+1} saadaan laskemalla $\mathcal{M}(\mathbf{q}_j)$ ja ortogonalisoidamalla se käyttäen Gram-Schmidt menetelmää

$$h_{ij} = \langle \mathcal{M}(\mathbf{q}_j), \mathbf{q}_i \rangle \quad (i = 1, 2, \dots, j) \quad (3.23)$$

$$\mathbf{w}_j = \mathcal{M}(\mathbf{q}_j) - \sum_{i=1}^j h_{ij} \mathbf{q}_i \quad (3.24)$$

$$h_{j+1,j} = \|\mathbf{w}_j\| \quad (3.25)$$

Menetelmä pysäytetään, jos $\mathbf{w}_j = 0$. Muutoin asetetaan $\mathbf{q}_{j+1} = \mathbf{w}_j / \|\mathbf{w}_j\|$.

On huomattava, että sisätulon arvo kaavassa (3.23) on yleensä kompleksiluku. Reaalilineaarisille operaattoreille pätee vain $\mathcal{M}(\alpha z) = \alpha \mathcal{M}(z)$, missä α on reaalinen. Tästä syystä vektori \mathbf{w}_j ei välttämättä enää kuulu Krylovin aliavaruuteen eikä algoritmi siten generoi sille kantaa. Operaattoreille \mathcal{M}_κ pätee kuitenkin matriiseja vastaavat lauseet.

Lause 3.6.4. Oletetaan, että operaattorilla \mathcal{M}_κ on Arnoldin menetelmällä laskettu vektorit $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_k$. Nämä vektorit muodostavat ortonormaalin kannan avaruuteen $\mathcal{K}_k(\mathcal{M}_\kappa, \mathbf{r}_0)$.

Todistus. Vektorit ovat ortonormaalit Gram-Schmidt prosessin ja normeerauksen myötä. Osoitetaan, että jokaiselle \mathbf{q}_j on olemassa astetta $j-1$ oleva polynomi p_{j-1} siten, että $\mathbf{q}_j = p_{j-1}(\mathcal{M}_\kappa)(\mathbf{r}_0)$. Tämä on selvästi totta \mathbf{q}_1 :lle. Oletetaan sitten, että näin on $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_j$:lle ja osoitetaan todeksi myös \mathbf{q}_{j+1} :lle. Kaavasta (3.24) saadaan

$$\begin{aligned} \mathbf{q}_{j+1} &= \frac{1}{\|\mathbf{w}_j\|} \mathcal{M}_\kappa(\mathbf{q}_j) - \sum_{i=1}^j \frac{1}{\|\mathbf{w}_j\|} h_{ij} \mathbf{q}_i \\ &= \frac{1}{\|\mathbf{w}_j\|} \mathcal{M}_\kappa p_{j-1}(\mathcal{M}_\kappa)(\mathbf{r}_0) - \sum_{i=1}^j \frac{1}{\|\mathbf{w}_j\|} h_{ij} p_{i-1}(\mathcal{M}_\kappa)(\mathbf{r}_0), \end{aligned}$$

missä $\mathcal{M}_{\kappa p_{j-1}}(\mathcal{M}_{\kappa}) = q(\mathcal{M}_{\kappa})$ jollakin j -asteisella polynomilla lemmän 3.6.1 mukaan. Täten oikealla puolella on j -asteinen polynomi, joten induktioaskel on osoitettu.

Ollaan saatu, että $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_k$ kuuluvat avaruuteen $\mathcal{K}_k(\mathcal{M}_{\kappa}, \mathbf{r}_0)$. Selvästi aina $\dim \mathcal{K}_k(\mathcal{M}_{\kappa}, \mathbf{r}_0) \leq k$, joten vektorit virittävät avaruuden ja $\dim \mathcal{K}_k(\mathcal{M}_{\kappa}, \mathbf{r}_0) = k$. \square

Lause 3.6.5. *Operaattorilla \mathcal{M}_{κ} Arnoldin menetelmässä $\mathbf{w}_j = 0$ jos ja vain jos j on \mathbf{v} :n minimipolynomin asteluku.*

Todistus. Oletetaan aluksi $\mathbf{w}_j = 0$. Kuten lauseen 3.6.4 todistuksessa, olkoon p_{j-1} astetta j oleva polynomi siten, että $\mathbf{q}_j = p_{j-1}(\mathcal{M}_{\kappa})(\mathbf{r}_0)$. Kaavasta (3.24) saadaan

$$0 = \mathbf{w}_j = \mathcal{M}_{\kappa}(\mathbf{q}_j) - \sum_{i=1}^j h_{ij} \mathbf{q}_i = \mathcal{M}_{\kappa p_{j-1}}(\mathcal{M}_{\kappa})(\mathbf{r}_0) - \sum_{i=1}^j h_{ij} p_{i-1}(\mathcal{M}_{\kappa})(\mathbf{r}_0),$$

jonka oikealla puolella on j -asteinen polynomi lemmän 3.6.1 perusteella. Siispä \mathbf{r}_0 :n minimipolynomin asteluku on korkeintaan j . Koska $\mathbf{w}_{j-1} \neq 0$, on toisaalta saatu muodostettua $\mathbf{q}_1, \dots, \mathbf{q}_j$, joten lauseen 3.6.4 perusteella $\dim \mathcal{K}_j(\mathcal{M}_{\kappa}, \mathbf{r}_0) = j$ ja lauseen 3.6.3 mukaan minimipolynomin asteluku ei voi olla pienempi kuin j .

Oletetaan sitten, että j on \mathbf{r}_0 :n minimipolynomin asteluku. Lauseiden 3.6.3 ja 3.6.4 perusteella Arnoldin menetelmällä on mahdotonta muodostaa enemmän kuin j ortonormaalia vektoria. Täytyy siis olla $\mathbf{w}_i = 0$ jollakin $i \leq j$. Toisaalta, jos $i < j$, niin tämän todistuksen ensimmäisen osan mukaan asteluvunkin täytyisi olla i . Täten $\mathbf{w}_j = 0$. \square

Huomautuksen 3.6.2 mukaisesti operaattorin \mathcal{M}_{κ} Krylovin aliavaruudet ovat samoja kuin \mathcal{M}_0 :n, joten yksinkertaisuuden vuoksi muodostetaan kantavektorit käyttäen \mathcal{M}_0 :aa. Olkoot $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_k$ saatu operaattorilla \mathcal{M}_0 Arnoldin menetelmästä ja asetetaan $\mathbf{Q}_k = [\mathbf{q}_1 \ \mathbf{q}_2 \ \dots \ \mathbf{q}_k]$. Olkoon $\widehat{\mathbf{H}}_k$ Arnoldin menetelmän alkioista h_{ij} muodostettu $(k+1) \times k$ -matriisi, $(\widehat{\mathbf{H}}_k)_{ij} = h_{ij}$ ($i \leq j+1$) ja muut alkioit nolliia. Olkoon vielä \mathbf{H}_k se $k \times k$ -matriisi, kun matriisista $\widehat{\mathbf{H}}_k$ poistetaan viimeinen rivi. Kaavasta (3.24) saadaan

$$h_{j+1,j} \mathbf{q}_{j+1} = \mathbf{M}_{\#} \bar{\mathbf{q}}_j - \sum_{i=1}^j h_{ij} \mathbf{q}_i \quad \Leftrightarrow \quad \mathbf{M}_{\#} \bar{\mathbf{q}}_j = \sum_{i=1}^{j+1} h_{ij} \mathbf{q}_i,$$

joka voidaan kirjoittaa matriisimuotoihin

$$\mathbf{M}_{\#} \bar{\mathbf{Q}}_j = \mathbf{Q}_{j+1} \widehat{\mathbf{H}}_j, \quad (3.26)$$

$$\mathbf{M}_{\#} \bar{\mathbf{Q}}_j = \mathbf{Q}_j \mathbf{H}_j + \mathbf{w}_j \mathbf{e}_j^T, \quad (3.27)$$

missä \mathbf{e}_j on \mathbb{C}^j :n j :s luonnollinen kantavektori. Kertomalla jälkimmäinen puolittain \mathbf{Q}_j^* :llä, saadaan

$$\mathbf{Q}_j^* \mathbf{M}_{\#} \bar{\mathbf{Q}}_j = \mathbf{H}_j. \quad (3.28)$$

3.7 \mathbb{R} -lineaarinen täyden ortogonalisoinnin menetelmä

Olkoon \mathbf{z}_0 alkuarvaus yhtälön $\mathcal{M}_\kappa(\mathbf{z}) = \mathbf{b}$ ratkaisulle ja asetetaan $\mathbf{r}_0 = \mathbf{b} - \mathcal{M}_\kappa(\mathbf{z}_0)$. Täyden ortogonalisaation menetelmässä (Full Orthogonalization Method, FOM) lasketaan Arnoldin menetelmällä kantavektoreita Krylovin aliavaruuksiin $\mathcal{K}_k(\mathcal{M}_\kappa, \mathbf{r}_0)$, jolloin likimääräisratkaisu $\mathbf{z}_k \in \mathbf{z}_0 + \mathcal{K}_k(\mathcal{M}_\kappa, \mathbf{r}_0)$ saadaan kohtisuoruusvaatimukselta

$$\langle \mathbf{b} - \mathcal{M}_\kappa(\mathbf{z}_k), \mathbf{y} \rangle = 0 \quad \text{kaikilla } \mathbf{y} \in \mathcal{K}_k(\mathcal{M}_\kappa, \mathbf{r}_0).$$

Kun matriisin \mathbf{Q}_k sarakkeet ovat Arnoldin menetelmän ortonormaalit kantavektorit, yhtälöksi tulee

$$\kappa \mathbf{u}_k + \mathbf{Q}_k^* \mathbf{M}_\# \overline{\mathbf{Q}_k} \overline{\mathbf{u}_k} = \mathbf{Q}_k^* \mathbf{r}_0,$$

mistä ratkaisemalla \mathbf{u}_k saadaan likimääräisratkaisuksi $\mathbf{z}_k = \mathbf{z}_0 + \mathbf{Q}_k \mathbf{u}_k$. Käyttäen yhtälöä (3.28) yksinkertaistuu tämä vielä muotoon

$$\kappa \mathbf{u}_k + \mathbf{H}_k \overline{\mathbf{u}_k} = \|\mathbf{r}_0\| \mathbf{e}_1, \quad (3.29)$$

$$\mathbf{z}_k = \mathbf{z}_0 + \mathbf{Q}_k \mathbf{u}_k. \quad (3.30)$$

Likimääräisratkaisun virheen arvioimiseksi käytetään edellä olevia yhtälöitä (3.27), (3.29) ja (3.30), jolloin

$$\begin{aligned} \mathbf{b} - \mathcal{M}_\kappa(\mathbf{z}_k) &= \mathbf{r}_0 - \mathcal{M}_\kappa(\mathbf{Q}_k \mathbf{u}_k) = \mathbf{r}_0 - \kappa \mathbf{Q}_k \mathbf{u}_k - \mathbf{M}_\# \overline{\mathbf{Q}_k} \overline{\mathbf{u}_k} \\ &= \mathbf{r}_0 - \kappa \mathbf{Q}_k \mathbf{u}_k - \mathbf{Q}_k \mathbf{H}_k \overline{\mathbf{u}_k} - (\mathbf{w}_k \mathbf{e}_k^T) \overline{\mathbf{u}_k} \\ &= \mathbf{Q}_k (\|\mathbf{r}_0\| \mathbf{e}_1 - \kappa \mathbf{u}_k - \mathbf{H}_k \overline{\mathbf{u}_k}) - \mathbf{w}_k (\mathbf{e}_k^T \overline{\mathbf{u}_k}) = -(\mathbf{e}_k^T \overline{\mathbf{u}_k}) \mathbf{w}_k, \end{aligned} \quad (3.31)$$

mistä edelleen

$$\mathbf{b} - \mathcal{M}_\kappa \mathbf{z}_k = -h_{k+1,k} (\mathbf{e}_k^T \overline{\mathbf{u}_k}) \mathbf{q}_{k+1}, \quad (3.32)$$

$$\|\mathbf{b} - \mathcal{M}_\kappa \mathbf{z}_k\| = h_{k+1,k} |\mathbf{e}_k^T \overline{\mathbf{u}_k}|. \quad (3.33)$$

Tätä on mahdollista käyttää algoritmin pysäyttämiseksi, kun riittävä tarkkuus on saavutettu. Kaavasta 3.31 myös nähdään, että ratkaisu \mathbf{z}_k on tarkka Arnoldin menetelmän murtuessa (kun $\mathbf{w}_k = 0$).

3.7.1 Menetelmä symmetriselle $\mathbf{M}_\#$

Symmetrisellä matriisilla $\mathbf{M}_\#$ saadaan edellä kuvattua FOMia yksinkertaistettua siten, että tietokoneen muistissa tarvitsee pitää ainoastaan kaksi viimeisintä kantavektori \mathbf{q}_j ja likimääräisratkaisu saadaan eräänlaisen päivitysprosessin kautta.

Kun operaattorille \mathcal{M}_κ ajetaan Arnoldin menetelmää symmetrisellä matriisilla $\mathbf{M}_\#$, niin yhtälöstä (3.28) nähdään

$$\mathbf{H}_k^T = (\mathbf{Q}_k^* \mathbf{M}_\# \overline{\mathbf{Q}_k})^T = \overline{\mathbf{Q}_k}^T \mathbf{M}_\#^T (\mathbf{Q}_k^*)^T = \mathbf{Q}_k^* \mathbf{M}_\# \overline{\mathbf{Q}_k} = \mathbf{H}_k,$$

joten Hessenbergin matriisi \mathbf{H}_k on myös symmetrinen. Täten se on kolmilävistäjämuotoa

$$\mathbf{T}_k = \begin{bmatrix} \alpha_1 & \beta_2 & & & \\ \beta_2 & \alpha_2 & \beta_3 & & \\ & \ddots & \ddots & \ddots & \\ & & \beta_{k-1} & \alpha_{k-1} & \beta_k \\ & & & \beta_k & \alpha_k \end{bmatrix}.$$

3.7. \mathbb{R} -LINEAARINEN TÄYDEN ORTOGONALISOINNIN MENETELMÄ

Tällöin $\alpha_j = h_{jj}$ ja $\beta_j = h_{j,j-1}$ ja Arnoldin menetelmä (3.23)-(3.25) yksinkertaistuu seuraavaksi Lanczosin menetelmäksi. Asetetaan $\beta_1 = 0$, $\mathbf{q}_1 = \mathbf{r}_0/\|\mathbf{r}_0\|$ ja lasketaan vektorit $\mathbf{q}_2, \mathbf{q}_3, \dots$ iteroimalla seuraavia kaavoja aloittaen arvosta $j = 1$.

$$\alpha_j = \langle \mathbf{M}_\# \bar{\mathbf{q}}_j, \mathbf{q}_j \rangle, \quad (3.34)$$

$$\mathbf{w}_j = \mathbf{M}_\# \bar{\mathbf{q}}_j - \alpha_j \mathbf{q}_j - \beta_j \mathbf{q}_{j-1}, \quad (3.35)$$

$$\beta_{j+1} = \|\mathbf{w}_j\|. \quad (3.36)$$

Menetelmä pysäytetään, jos $\beta_{j+1} = 0$. Muutoin asetetaan $\mathbf{q}_{j+1} = \mathbf{w}_j/\beta_{j+1}$. Muunnettua Gram-Schmidt menetelmää käyttäen tästä tulee algoritmi 3.7.1.

Algoritmi 3.7.1 Lanczosin menetelmä ortonormaalin kannan muodostamiseksi Krylovin aliavaruudelle $\mathcal{K}_k(\mathcal{M}_0, \mathbf{r}_0)$ symmetrisellä $\mathbf{M}_\#$ käyttäen muunnettua Gram-Schmidt menetelmää. Laskee ortonormaalit kantavektorit $\mathbf{q}_1, \mathbf{q}_2, \dots$. Algoritmi pysähtyy muodostettuaan $\min\{k+1, \dim \mathcal{K}_k(\mathbf{A}, \mathbf{v})\}$ kantavektoria.

```

1:  $\mathbf{q}_1 \leftarrow \mathbf{r}_0/\|\mathbf{r}_0\|$ ,  $\beta_1 \leftarrow 0$ ,  $\mathbf{q}_0 \leftarrow 0$ 
2: for  $j = 1$  to  $k$  do
3:    $\mathbf{w} \leftarrow \mathbf{M}_\# \bar{\mathbf{q}}_j - \beta_j \mathbf{q}_{j-1}$ 
4:    $\alpha_j \leftarrow \langle \mathbf{w}, \mathbf{q}_j \rangle$ 
5:    $\mathbf{w} \leftarrow \mathbf{w} - \alpha_j \mathbf{q}_j$ 
6:    $\beta_{j+1} \leftarrow \|\mathbf{w}\|$ 
7:   if  $\beta_{j+1} = 0$  then
8:     lopeta algoritmi
9:   end if
10:   $\mathbf{q}_{j+1} \leftarrow \mathbf{w}/\beta_{j+1}$ 
11: end for

```

Lanczosin menetelmässä ortonormalisointiin tarvitaan vain kahta edellistä vektoria. FOMin likimääräisratkaisun laskemiseen kaavan (3.30) mukaan tarvitaan kuitenkin kaikkia vektoreita \mathbf{q}_j . Näytetään seuraavaksi miten tästä vaatimuksesta päästään eroon.

Merkitään $\mathcal{T}_k(z) = \kappa z + \mathbf{T}_k \bar{z}$, joka on tavalla (1.1) ilmaistuna kolmilävistäjämuotoa

$$\mathcal{T}_k = \begin{bmatrix} \hat{\alpha}_1 & \hat{\beta}_2 & & & & \\ \hat{\beta}_2 & \hat{\alpha}_2 & \hat{\beta}_3 & & & \\ & \ddots & \ddots & \ddots & & \\ & & \hat{\beta}_{k-1} & \hat{\alpha}_{k-1} & \hat{\beta}_k & \\ & & & \hat{\beta}_k & \hat{\alpha}_k & \end{bmatrix},$$

missä skalaarioperaattorit ovat $\hat{\alpha}_j(z) = \kappa z + \alpha_j \bar{z}$, $\hat{\beta}_j(z) = \beta_j \bar{z}$. Oletetaan, että sillä on reaaliineaarinen LU-hajotelma, jonka on tällöin oltava muotoa

$$\mathcal{T}_k = \mathcal{L}_k \mathcal{U}_k = \begin{bmatrix} 1 & & & & & \\ \lambda_2 & 1 & & & & \\ & \lambda_3 & \ddots & & & \\ & & \ddots & 1 & & \\ & & & \lambda_k & 1 & \end{bmatrix} \begin{bmatrix} \nu_1 & \hat{\beta}_2 & & & & \\ & \nu_2 & \hat{\beta}_3 & & & \\ & & \ddots & \ddots & & \\ & & & \nu_{k-1} & \hat{\beta}_k & \\ & & & & \nu_k & \end{bmatrix},$$

missä λ_k ja ν_k syntyvät Gaussin eliminoinnissa ja niille saadaan kaavat

$$\begin{aligned}\nu_1 &= \hat{\alpha}_1, \\ \lambda_k &= \hat{\beta}_k \nu_{k-1}^{-1}, \\ \nu_k &= \hat{\alpha}_k - \lambda_k \hat{\beta}_k.\end{aligned}$$

Oletetaan \mathcal{T}_k vielä kääntyväksi, jolloin kaavasta (3.30) tulee

$$z_k = z_0 + Q_k U_k^{-1} \mathcal{L}_k^{-1}(\|r_0\|e_1).$$

Merkitään

$$\mathcal{P}_k = Q_k U_k^{-1}, \quad y_k = \mathcal{L}_k^{-1}(\|r_0\|e_1), \quad (3.37)$$

jonka jälkeen likimääräisratkaisu saadaan kaavasta $z_k = z_0 + \mathcal{P}_k(y_k)$. Olkoot π_1, \dots, π_k operaattorin \mathcal{P}_k sarakkeet. Tällöin ensimmäisestä kaavasta (3.37) nähdään, että

$$\pi_k = (q_k - \pi_{k-1} \hat{\beta}_k) \nu_k^{-1}. \quad (3.38)$$

Merkitään η_k :lla vektorin y_k viimeistä komponenttia, jolloin toisesta kaavasta (3.37) saadaan

$$\eta_k = -\lambda_k(\eta_{k-1}), \quad (3.39)$$

koska alakolmio-operaattoreille pätee

$$\mathcal{L}_k = \begin{bmatrix} 1 & & & & & & \\ \lambda_2 & 1 & & & & & \\ & \lambda_3 & \ddots & & & & \\ & & \ddots & \ddots & & & \\ & & & \ddots & 1 & & \\ & & & & \lambda_{k-1} & 1 & \\ & & & & & & 1 \end{bmatrix} \begin{bmatrix} 1 & & & & & & \\ & 1 & & & & & \\ & & \ddots & & & & \\ & & & \ddots & & & \\ & & & & 1 & & \\ & & & & & 1 & \\ & & & & & & \lambda_k & 1 \end{bmatrix}.$$

Nyt likimääräisratkaisulle saadaan kaava

$$z_k = z_0 + \mathcal{P}_k(y_k) = z_0 + [\mathcal{P}_{k-1} \quad \pi_k] \begin{bmatrix} y_{k-1} \\ \eta_k \end{bmatrix} = z_{k-1} + \pi_k(\eta_k). \quad (3.40)$$

Kun likimääräisratkaisu z_{k-1} on laskettu, se voidaan päivittää ratkaisuksi z_k laske-
malla uusi kantavektori q_k yhdellä kierroksella algoritmia 3.7.1. Muistissa tarvitsee
tällöin olla vain q_{k-1} ja q_{k-2} . Sen jälkeen voidaan käyttää kaavaa (3.38) π_k :n päi-
vittämiseksi. Muistissa tarvitsee olla edellinen π_{k-1} . Vielä kaavasta (3.39) saadaan
päivitettyä η_k edellisen η_{k-1} :n perusteella. Lopuksi z_k lasketaan kaavasta (3.40).
Tämä on kirjoitettu algoritmiksi 3.7.2.

Ratkaisun tarkkuutta on mahdollista arvioida kaavalla (3.33). Jäännökselle $r_k = b - \mathcal{M}_k(z_k)$ pätee

$$\|r_k\| = \beta_{k+1} |e_k^T u_k|,$$

missä $u_k = U_k^{-1} \mathcal{L}_k^{-1}(\|r_0\|e_1) = U_k^{-1}(y_k)$. Täten

$$\|r_k\| = \beta_{k+1} |\nu_k^{-1}(\eta_k)| = |\lambda_{k+1}(\eta_k)| = |\eta_{k+1}|. \quad (3.41)$$

Algoritmi 3.7.2 Symmetrisen $M_{\#}$ ratkaisumenetelmä. K on suurin sallittu iteraatioiden lukumäärä.

```

1:  $\mathbf{q}_1 \leftarrow \mathbf{r}_0 / \|\mathbf{r}_0\|$ ,  $\beta_1 \leftarrow 0$ ,  $\mathbf{q}_0 \leftarrow 0$ ,  $\mathbf{p}_0 \leftarrow 0$ ,  $\mathbf{p}_{a,0} \leftarrow 0$ 
2:  $\eta_1 \leftarrow \|\mathbf{r}_0\|$ ,  $\lambda_1 \leftarrow 0$ ,  $\lambda_{a,1} \leftarrow 0$ 
3: for  $k = 1$  to  $K$  do
4:    $\mathbf{w} \leftarrow M_{\#}\bar{\mathbf{q}}_k - \beta_k \mathbf{q}_{k-1}$ 
5:    $\alpha_k \leftarrow \langle \mathbf{w}, \mathbf{q}_k \rangle$ 
6:    $\begin{bmatrix} u_k & u_{a,k} \end{bmatrix} \leftarrow \begin{bmatrix} \kappa & \alpha_k \end{bmatrix} - \beta_k \begin{bmatrix} \lambda_{a,k} & \lambda_k \end{bmatrix}$ 
7:    $\begin{bmatrix} \mathbf{p}_k & \mathbf{p}_{a,k} \end{bmatrix} \leftarrow \begin{bmatrix} \mathbf{q}_k - \beta_k \mathbf{p}_{a,k-1} & -\beta_k \mathbf{p}_{k-1} \end{bmatrix} \begin{bmatrix} u_k & u_{a,k} \\ \bar{u}_{a,k} & \bar{u}_k \end{bmatrix}^{-1}$ 
8:    $\mathbf{z}_k \leftarrow \mathbf{z}_{k-1} + \eta_k \mathbf{p}_k + \bar{\eta}_k \mathbf{p}_{a,k}$ 
9:   Jos  $\mathbf{z}_k$  on riittävän tarkka, lopetetaan algoritmi.
10:   $\mathbf{w} \leftarrow \mathbf{w} - \alpha_k \mathbf{q}_k$ 
11:   $\beta_{k+1} \leftarrow \|\mathbf{w}\|$ 
12:   $\mathbf{q}_{k+1} \leftarrow \mathbf{w} / \beta_{k+1}$ 
13:   $\begin{bmatrix} \lambda_{k+1} & \lambda_{a,k+1} \end{bmatrix} \leftarrow \begin{bmatrix} 0 & \beta_{k+1} \end{bmatrix} \begin{bmatrix} u_k & u_{a,k} \\ \bar{u}_{a,k} & \bar{u}_k \end{bmatrix}^{-1}$ 
14:   $\eta_{k+1} \leftarrow -\lambda_{k+1} \eta_k - \lambda_{a,k+1} \bar{\eta}_k$ 
15: end for
    
```

3.7.2 Menetelmä vinosymmetriselle $M_{\#}$

Edellä oleva symmetrisen $M_{\#}$ matriisin menetelmä saattaa rikkoutua, koska oletettua LU-hajotelmaa ei ole. Vinosymmetrisillä matriiseilla ei tätä ongelmaa esiinny, kunhan $\kappa \neq 0$. Seuraavassa asetetaan $\kappa = 1$ ja ratkaistavana on siten yhtälö $\mathbf{z} + M_{\#}\bar{\mathbf{z}} = \mathbf{b}$, missä $M_{\#}$ on vinosymmetrinen. Tuloksena on matriisien liittogradienttimenetelmää muistuttava menetelmä.

Yhtälöstä (3.28) saadaan nyt

$$\mathbf{H}_k^T = (\mathbf{Q}_k^* M_{\#} \bar{\mathbf{Q}}_k)^T = \bar{\mathbf{Q}}_k^T M_{\#}^T (\mathbf{Q}_k^*)^T = \mathbf{Q}_k^* (-M_{\#}) \bar{\mathbf{Q}}_k = -\mathbf{H}_k,$$

joten Hessenbergin matriisi \mathbf{H}_k on myös vinosymmetrinen. Täten se on muotoa

$$\mathbf{T}_k = \begin{bmatrix} 0 & -\beta_2 & & & & \\ \beta_2 & 0 & -\beta_3 & & & \\ & \ddots & \ddots & \ddots & & \\ & & & \beta_{k-1} & 0 & -\beta_k \\ & & & & \beta_k & 0 \end{bmatrix}.$$

Tällöin $\beta_j = h_{j,j-1}$ ja Arnoldin menetelmästä (3.23)-(3.25) tulee vielä edellistä kohtaakin yksinkertaisempi, koska $\langle M_{\#}\bar{\mathbf{q}}, \mathbf{q} \rangle = 0$ kaikilla vektoreilla \mathbf{q} . Asetetaan $\beta_1 = 0$, $\mathbf{q}_1 = \mathbf{r}_0 / \|\mathbf{r}_0\|$ ja lasketaan vektorit $\mathbf{q}_2, \mathbf{q}_3, \dots$ iteroimalla seuraavia kaavoja aloittaen arvosta $j = 1$.

$$\mathbf{w}_j = M_{\#}\bar{\mathbf{q}}_j + \beta_j \mathbf{q}_{j-1}, \quad (3.42)$$

$$\beta_{j+1} = \|\mathbf{w}_j\|. \quad (3.43)$$

Menetelmä pysäytetään, jos $\beta_{j+1} = 0$. Muutoin asetetaan $\mathbf{q}_{j+1} = \mathbf{w}_j / \beta_{j+1}$. Vaikka tarkassa aritmetiikassa ortonormalisointia ei enää tarvitakaan, niin numeerisessa

laskennassa tästä saattaa olla hyötyä. Tätä varten voidaan vaihtaa (3.42) tilalle kaava

$$\tilde{\mathbf{w}}_j = \mathbf{M}_{\#} \bar{\mathbf{q}}_j + \beta_j \mathbf{q}_{j-1}, \quad \mathbf{w}_j = \tilde{\mathbf{w}}_j - \langle \tilde{\mathbf{w}}_j, \mathbf{q}_j \rangle \mathbf{q}_j. \quad (3.44)$$

Merkitään $\mathcal{T}_k(\mathbf{z}) = \mathbf{z} + \mathbf{T}_k \bar{\mathbf{z}}$, joka on kolmilävistäjämuotoa

$$\mathcal{T}_k = \begin{bmatrix} \mathbf{1} & -\widehat{\beta}_2 & & & & \\ \widehat{\beta}_2 & \mathbf{1} & -\widehat{\beta}_3 & & & \\ & \ddots & \ddots & \ddots & & \\ & & \widehat{\beta}_{k-1} & \mathbf{1} & -\widehat{\beta}_k & \\ & & & \widehat{\beta}_k & \mathbf{1} & \\ & & & & & \mathbf{1} \end{bmatrix},$$

missä skalaarioperaattorit ovat $\widehat{\beta}_j(z) = \beta_j \bar{z}$. Sen LU-hajotelman (mikäli olemassa) on oltava muotoa

$$\mathcal{T}_k = \mathcal{L}_k \mathcal{U}_k = \begin{bmatrix} \mathbf{1} & & & & & \\ \lambda_2 & \mathbf{1} & & & & \\ & \lambda_3 & \ddots & & & \\ & & \ddots & \mathbf{1} & & \\ & & & \lambda_k & \mathbf{1} & \\ & & & & & \mathbf{1} \end{bmatrix} \begin{bmatrix} \nu_1 & -\widehat{\beta}_2 & & & & \\ & \nu_2 & -\widehat{\beta}_3 & & & \\ & & \ddots & \ddots & & \\ & & & \nu_{k-1} & -\widehat{\beta}_k & \\ & & & & \nu_k & \end{bmatrix},$$

missä λ_k ja ν_k syntyvät Gaussin eliminoinnissa, mistä niille saadaan kaavat

$$\begin{aligned} \nu_1 &= \mathbf{1}, \\ \lambda_k &= \widehat{\beta}_k \nu_{k-1}^{-1}, \\ \nu_k &= \mathbf{1} + \lambda_k \widehat{\beta}_k. \end{aligned}$$

Nähdään helposti, että nämä saadaan muotoon

$$\begin{aligned} \gamma_1 &= \mathbf{1}, \\ l_k &= \beta_k / \gamma_{k-1}, \\ \gamma_k &= \mathbf{1} + l_k \beta_k, \end{aligned}$$

missä $\nu_k(z) = \gamma_k z$ ja $\lambda_k(z) = l_k \bar{z}$. Kaikilla j on $\gamma_j \geq 1$, joten LU-hajotelma on aina olemassa ja \mathcal{T}_k on kääntyvä. Kaavasta (3.30) saadaan

$$\mathbf{z}_k = \mathbf{z}_0 + \mathbf{Q}_k \mathcal{U}_k^{-1} \mathcal{L}_k^{-1} (\|\mathbf{r}_0\| \mathbf{e}_1)$$

ja kaavoista (3.38)-(3.40) tulee

$$\begin{aligned} \pi_k &= (\mathbf{q}_k - \pi_{k-1} \widehat{\beta}_k) \nu_k^{-1} = (1/\gamma_k) (\mathbf{q}_k - \pi_{k-1} \widehat{\beta}_k), \\ \eta_k &= -\lambda_k (\eta_{k-1}) = -l_k \eta_{k-1}, \\ \mathbf{z}_k &= \mathbf{z}_{k-1} + \pi_k (\eta_k), \end{aligned}$$

missä on käytetty huomiota, että nyt $\eta_j \in \mathbb{R}$ kaikilla j . Kaavoja saadaan vielä yksinkertaistettua merkitsemällä $\pi_k(z) = \tilde{\mathbf{p}}_k z + \tilde{\mathbf{p}}_{\#,k} \bar{z}$ ja $\mathbf{p}_k = \tilde{\mathbf{p}}_k + \tilde{\mathbf{p}}_{\#,k}$. Koska $\eta_k \in \mathbb{R}$ ja $\widehat{\beta}_k$ on anti-lineaarinen, niin kaavat saadaan muotoon

$$\mathbf{p}_k = \frac{1}{\gamma_k} (\mathbf{q}_k - \beta_k \mathbf{p}_{k-1}), \quad (3.45)$$

$$\eta_k = -l_k \eta_{k-1}, \quad (3.46)$$

$$\mathbf{z}_k = \mathbf{z}_{k-1} + \eta_k \mathbf{p}_k. \quad (3.47)$$

Näin saadaan algoritmi 3.7.3. Ratkaisun jäännöksen normille (3.41) pätee edelleen

$$\|\mathbf{r}_k\| = |\eta_{k+1}|.$$

Kunkin kierroksen laskuja varten täytyy muistissa pitää vain kaksi edellistä \mathbf{q} -vektoria ja yksi edellinen \mathbf{p} -vektori.

Algoritmi 3.7.3 Vinosymmetrisen $M_\#$ ratkaisumenetelmä ($\kappa = 1$). K on suurin sallittu iteraatioiden lukumäärä.

- 1: $\mathbf{q}_1 \leftarrow \mathbf{r}_0 / \|\mathbf{r}_0\|$, $\beta_1 \leftarrow 0$, $\mathbf{q}_0 \leftarrow 0$, $\mathbf{p}_0 \leftarrow 0$, $\eta_1 \leftarrow \|\mathbf{r}_0\|$, $\gamma_1 \leftarrow 1$
 - 2: **for** $k = 1$ to K **do**
 - 3: $\mathbf{p}_k \leftarrow \frac{1}{\gamma_k} (\mathbf{q}_k - \beta_k \mathbf{p}_{k-1})$
 - 4: $\mathbf{z}_k \leftarrow \mathbf{z}_{k-1} + \eta_k \mathbf{p}_k$
 - 5: Jos \mathbf{z}_k on riittävän tarkka, lopetetaan algoritmi.
 - 6: $\mathbf{w} \leftarrow M_\# \bar{\mathbf{q}}_k + \beta_k \mathbf{q}_{k-1}$
 - 7: $\mathbf{w} \leftarrow \mathbf{w} - \langle \mathbf{w}, \mathbf{q}_k \rangle \mathbf{q}_k$
 - 8: $\beta_{k+1} \leftarrow \|\mathbf{w}\|$
 - 9: $\mathbf{q}_{k+1} \leftarrow \mathbf{w} / \beta_{k+1}$
 - 10: $\lambda_{k+1} \leftarrow \beta_{k+1} / \gamma_k$, $\gamma_{k+1} \leftarrow 1 + \lambda_{k+1} \beta_{k+1}$
 - 11: $\eta_{k+1} \leftarrow -\lambda_{k+1} \eta_k$
 - 12: **end for**
-

3.8 \mathbb{R} -lineaarinen GMRES operaattorille \mathcal{M}_κ

Tämän kohdan menetelmä toimii yleisille matriiseille $M_\#$. Likimääräisratkaisua yhtälölle $\mathcal{M}_\kappa(\mathbf{z}) = \mathbf{b}$ etsitään jälleen muodossa $\mathbf{z}_k \in \mathbf{z}_0 + \mathcal{K}_k(\mathcal{M}_\kappa, \mathbf{r}_0)$, mutta nyt vaaditaan jäännöksen $\mathbf{r}_k = \mathbf{b} - \mathcal{M}_\kappa(\mathbf{z}_k)$ normi pienimmäksi mahdolliseksi. Tämä on yhtäpitävää ehdon $\mathbf{r}_k \perp \mathcal{M}_\kappa(\mathcal{K}_k(\mathcal{M}_\kappa, \mathbf{r}_0))$ kanssa. Kun matriisin \mathbf{Q}_k sarakkeet on laskettu Arnoldin menetelmällä, normille saadaan yhtälön (3.26) mukaan

$$\begin{aligned} \|\mathbf{r}_k\| &= \|\mathbf{r}_0 - \kappa \mathbf{Q}_k \mathbf{u}_k - M_\# \bar{\mathbf{Q}}_k \bar{\mathbf{u}}_k\| = \|\mathbf{r}_0 - \kappa \mathbf{Q}_k \mathbf{u}_k - \mathbf{Q}_{k+1} \widehat{\mathbf{H}}_k \bar{\mathbf{u}}_k\| \\ &= \left\| \mathbf{Q}_{k+1} (\|\mathbf{r}_0\| \widehat{\mathbf{e}}_1 - \kappa \widehat{\mathbf{I}}_k \mathbf{u}_k - \widehat{\mathbf{H}}_k \bar{\mathbf{u}}_k) \right\| = \left\| \|\mathbf{r}_0\| \widehat{\mathbf{e}}_1 - \kappa \widehat{\mathbf{I}}_k \mathbf{u}_k - \widehat{\mathbf{H}}_k \bar{\mathbf{u}}_k \right\|, \end{aligned} \quad (3.48)$$

missä $\widehat{\mathbf{e}}_1$ on \mathbb{C}^{k+1} :n ensimmäinen luonnollinen kantavektori ja $\widehat{\mathbf{I}}_k$ on yksikkömatriisi \mathbf{I}_k lisättyä yhdellä nollarivillä. Tämä pienimmän neliösumman tehtävä voidaan ratkaista laskemalla QR-hajotelma operaattorille $\widehat{\mathcal{H}}_k : \mathbb{C}^k \rightarrow \mathbb{C}^{k+1}$, $\widehat{\mathcal{H}}_k(\mathbf{u}) = \kappa \widehat{\mathbf{I}}_k \mathbf{u} + \widehat{\mathbf{H}}_k \bar{\mathbf{u}}$. Käyttäen skalaarioperaattorimuotoa (1.1), voidaan kirjoittaa

$$\widehat{\mathcal{H}}_k = \begin{bmatrix} \mu_{11} & \mu_{12} & \mu_{13} & \cdots & \mu_{1k} \\ \mu_{21} & \mu_{22} & \mu_{23} & & \mu_{2k} \\ & \mu_{32} & \mu_{33} & & \mu_{3k} \\ & & \ddots & & \vdots \\ & & & \mu_{k,k-1} & \mu_{k,k} \\ & & & & \mu_{k+1,k} \end{bmatrix}, \quad (3.49)$$

Unitaarisina operaattoreina (3.50) voidaan käyttää kohdassa 1.2 esiteltyjä Householderin muunnoksia. Toinen mahdollisuus on käyttää seuraavan lauseen tarjoamia Givensin rotaatioiden tapaisia reaali lineaarisia operaattoreita. Huomaa, että kaavan (3.25) mukaan $h_{j+1,j}$ on reaaliluku.

Lause 3.8.1. *Olkooot skalaarioperaattorit*

$$\mu(z) = \kappa z + a\bar{z}, \quad \nu(z) = b\bar{z},$$

missä $\kappa, a \in \mathbb{C}$ ja $b \in \mathbb{R}$, $b \neq 0$. Asetetaan

$$\widehat{U} = \begin{bmatrix} \bar{\kappa} & 0 \\ 0 & \kappa \end{bmatrix}, \quad \widehat{U}_\# = \begin{bmatrix} a & b \\ -b & \bar{a} \end{bmatrix}, \quad \widehat{U}(z) = \widehat{U}z + \widehat{U}_\#\bar{z},$$

$$x = |\kappa|^2 + |a|^2 + b^2, \quad y = 2a\kappa,$$

$$r_1 = r_2 = \sqrt{\frac{x + \sqrt{x^2 - |y|^2}}{2}}, \quad s_1 = \frac{a\kappa}{r_1}, \quad s_2 = \frac{\bar{a}\kappa}{r_2},$$

$$\beta_i(z) = r_i z + s_i \bar{z}.$$

Tällöin operaattori $\mathcal{U} : \mathbb{C}^2 \rightarrow \mathbb{C}^2$,

$$\mathcal{U} = \widehat{U} \circ \begin{bmatrix} \beta_1^{-1} & \\ & \beta_2^{-1} \end{bmatrix}$$

on unitaarinen ja

$$\mathcal{U} \circ \begin{bmatrix} \mu \\ \nu \end{bmatrix} = \begin{bmatrix} \# \\ 0 \end{bmatrix}.$$

Todistus. Suoraan laskemalla nähdään

$$\widehat{U}^* \widehat{U} = \begin{bmatrix} \beta_1^2 & \\ & \beta_2^2 \end{bmatrix},$$

joten

$$\mathcal{U}^* \mathcal{U} = \begin{bmatrix} \beta_1^{-1} & \\ & \beta_2^{-1} \end{bmatrix} \widehat{U}^* \widehat{U} \begin{bmatrix} \beta_1^{-1} & \\ & \beta_2^{-1} \end{bmatrix} = I$$

ja täten \mathcal{U} on unitaarinen. Niin ikään suoralla laskulla todetaan

$$\widehat{U} \circ \begin{bmatrix} \beta_1^{-1} \circ \mu \\ \beta_2^{-1} \circ \nu \end{bmatrix} = \begin{bmatrix} \# \\ 0 \end{bmatrix}.$$

□

Esimerkki 1.1.5 sisältää kaavan edellisessä lauseessa olevien β_i^{-1} laskemiseksi.

Yhtälö $\mathcal{M}_\kappa(z) = \mathbf{b}$ voitaisiin myös ratkaista kirjoittamalla se kaksi kertaa suuremmaksi reaalisiksi yhtälöryhmäksi. Käyttäen kohdan (1.3) esitystä saadaan ($\kappa = c + id$)

$$\mathbf{A} = \begin{bmatrix} c\mathbf{I} + \operatorname{Re}(\mathcal{M}_\#) & -d\mathbf{I} + \operatorname{Im}(\mathcal{M}_\#) \\ d\mathbf{I} + \operatorname{Im}(\mathcal{M}_\#) & c\mathbf{I} - \operatorname{Re}(\mathcal{M}_\#) \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} \operatorname{Re}(z) \\ \operatorname{Im}(z) \end{bmatrix}, \quad \mathbf{f} = \begin{bmatrix} \operatorname{Re}(\mathbf{b}) \\ \operatorname{Im}(\mathbf{b}) \end{bmatrix},$$

$$\mathbf{A}\mathbf{x} = \mathbf{f}.$$

Algoritmi 3.8.1 GMRES-menetelmä yhtälön $\mathcal{M}_\kappa(\mathbf{z}) = \mathbf{b}$ likimääräiseksi ratkaisemiseksi. Suorittaa k iteraatiota tai pysähtyy, jos tarkka ratkaisu saavutetaan. Alkuarvaus on \mathbf{z}_0 .

```

1:  $\mathbf{r}_0 \leftarrow \mathbf{b} - \mathcal{M}_\kappa(\mathbf{z}_0)$ ,  $\beta \leftarrow \|\mathbf{r}_0\|$ ,  $\mathbf{q}_1 \leftarrow \mathbf{r}_0/\beta$ ,  $\mathbf{v} \leftarrow 1$ ,  $\widehat{\mathcal{R}} \leftarrow []$ 
2: for  $j = 1$  to  $k$  do
3:    $\mathbf{w} \leftarrow M_{\#} \bar{\mathbf{q}}_j$ 
4:   for  $i = 1$  to  $j$  do
5:      $h_{ij} \leftarrow \langle \mathbf{w}, \mathbf{q}_i \rangle$ 
6:      $\mathbf{w} \leftarrow \mathbf{w} - h_{ij} \mathbf{q}_i$ 
7:   end for
8:    $h_{j+1,j} \leftarrow \|\mathbf{w}\|$ 
9:    $\mathbf{s} \leftarrow [h_{1j} \ h_{2j} \ \dots \ h_{j+1,j}]^T$ 
10:  Olkoon  $\mathcal{S} : \mathbb{C} \rightarrow \mathbb{C}^{j+1}$  siten, että  $\mathcal{S}(z) = \kappa e_j z + s \bar{z}$ 
11:  for  $i = 1$  to  $j - 1$  do
12:     $\mathcal{S}_{i:(i+1)} \leftarrow \mathcal{U}_i \circ \mathcal{S}_{i:(i+1)}$ 
13:  end for
14:  Valitaan unitaari  $\mathcal{U}_j : \mathbb{C}^2 \rightarrow \mathbb{C}^2$  siten, että  $\mathcal{U}_j \circ \mathcal{S}_{j:(j+1)} = [\# \ 0]^T$ 
15:   $\mathcal{S}_{j:(j+1)} \leftarrow \mathcal{U}_j \circ \mathcal{S}_{j:(j+1)}$ 
16:   $\widehat{\mathcal{R}} \leftarrow \begin{bmatrix} \widehat{\mathcal{R}} \\ \mathbf{0} \end{bmatrix}$ ,  $\widehat{\mathcal{R}} \leftarrow \begin{bmatrix} \widehat{\mathcal{R}} & \mathbf{s} \end{bmatrix}$ , missä  $\mathbf{0}$  on  $1 \times (j - 1)$ -nollaoperaattori
17:   $\mathbf{v} \leftarrow \begin{bmatrix} \mathbf{v} \\ 0 \end{bmatrix} \in \mathbb{C}^{j+1}$ ,  $\mathbf{v}_{j:(j+1)} \leftarrow \mathcal{U}_j(\mathbf{v}_{j:(j+1)})$ 
18:  if  $h_{j+1,j} = 0$  then
19:    break
20:  end if
21:   $\mathbf{q}_{j+1} \leftarrow \mathbf{w}/h_{j+1,j}$ 
22: end for
23:  $\mathbf{Q} \leftarrow [\mathbf{q}_1 \ \dots \ \mathbf{q}_j]$ ,  $\mathcal{R} \leftarrow \widehat{\mathcal{R}}_{1:j,*}$ 
24: Ratkaistaan  $\mathbf{u}$  yhtälöstä  $\mathcal{R}(\mathbf{u}) = \beta \mathbf{v}_{1:j}$ 
25:  $\mathbf{z} \leftarrow \mathbf{z}_0 + \mathbf{Q}\mathbf{u}$ , jäännöksen normi on  $\beta|\mathbf{v}_{j+1}|$ 

```

Kun matriisien GMRES etsii jäännöksen minimoivan ratkaisun Krylovin aliavaruudesta ($\mathbf{v} = \mathbf{f} - \mathbf{A}\mathbf{x}_0$)

$$\mathcal{K}_k(\mathbf{A}, \mathbf{v}) = \text{span} \left\{ \mathbf{v}, \mathbf{A}\mathbf{v}, \mathbf{A}^2\mathbf{v}, \dots, \mathbf{A}^{k-1}\mathbf{v} \right\},$$

niin reaalityyppinen GMRES etsii ratkaisun avaruudesta ($\mathbf{r}_0 = \mathbf{b} - \mathcal{M}_\kappa(\mathbf{z}_0)$)

$$\mathcal{K}_k(\mathcal{M}_\kappa, \mathbf{r}_0) = \text{span} \left\{ \mathbf{r}_0, \mathcal{M}_0(\mathbf{r}_0), \mathcal{M}_0^2(\mathbf{r}_0), \dots, \mathcal{M}_0^{k-1}(\mathbf{r}_0) \right\}.$$

Edellinen on \mathbb{R} -kertoiminen ja jälkimmäinen on \mathbb{C} -kertoiminen, joten samaistuksen $\mathbf{z} \leftrightarrow \begin{bmatrix} \text{Re}(\mathbf{z}) \\ \text{Im}(\mathbf{z}) \end{bmatrix}$ mielessä on aidosti $\mathcal{K}_k(\mathbf{A}, \mathbf{v}) \subset \mathcal{K}_k(\mathcal{M}_\kappa, \mathbf{r}_0)$ ja siten reaalityyppinen GMRES etsii kullakin iteraatiokierroksella jäännöksen minimoivan ratkaisun laajemmasta joukosta. Täten se suppenee aina vähintään yhtä nopeasti kuin matriisien GMRES. Numeeristen kokeiden mukaan se on yleensä nopeampi (ks. kohta 5.1).

3.8.1 MINRES symmetriselle $M_{\#}$

Kohdan 3.7.1 menetelmän ongelmana on sen mahdollinen rikkoutuminen LU-hajotelman puuttumisen vuoksi. Edellisen kohdan GMRES voidaan erikoistaa symmetrisille $M_{\#}$ ilman tällaista ongelmaa ja silti saavuttaa kohtaa 3.7.1 vastaava säästö kantavektoreiden suhteen.

Kohdassa 3.7.1 kuvatulla Lanczosin menetelmällä edellisen kohdan GMRESin Hessenbergin operaattori (3.49) tulee kolmilävistäjämuotoon

$$\widehat{\mathcal{T}}_k = \begin{bmatrix} \widehat{\alpha}_1 & \widehat{\beta}_2 & & & & \\ \widehat{\beta}_2 & \widehat{\alpha}_2 & \widehat{\beta}_3 & & & \\ & \ddots & \ddots & \ddots & & \\ & & \widehat{\beta}_{k-1} & \widehat{\alpha}_{k-1} & \widehat{\beta}_k & \\ & & & \widehat{\beta}_k & \widehat{\alpha}_k & \\ & & & & \widehat{\beta}_{k+1} & \end{bmatrix},$$

missä skalaarioperaattorit ovat $\widehat{\alpha}_j(z) = \kappa z + \alpha_j \bar{z}$, $\widehat{\beta}_j(z) = \beta_j \bar{z}$. Kun tämä muunnetaan yläkolmiomuotoon edellisen kohdan muotoisilla unitaarioperaattoreilla (3.51), saadaan kolmilävistäjinen yläkolmio-operaattori $\widehat{\mathcal{R}}_k$. Merkitään edelleen $\mathcal{R}_k = (\rho_{ij})$ sitä operaattoria, joka saadaan poistamalla $\widehat{\mathcal{R}}_k$:n viimeinen rivi. Määritellään sitten operaattori $\mathcal{P}_k = \mathcal{Q}_k \mathcal{R}_k^{-1}$ ja olkoot π_1, \dots, π_k sen sarakkeet. Koska $\mathcal{P}_k \mathcal{R}_k = \mathcal{Q}_k$, niin

$$\pi_k = (\mathbf{q}_k - \pi_{k-1} \rho_{k-1,k} - \pi_{k-2} \rho_{k-2,k}) \rho_{k,k}^{-1}.$$

Näin ollen, muistamalla vain kaksi edellistä π_{k-2} ja π_{k-1} , saadaan laskettua seuraava π_k . Likimääräisratkaisua voidaan päivittää kullakin kierroksella kaavalla

$$\begin{aligned} \mathbf{z}_k &= \mathbf{z}_0 + \mathcal{Q}_k \mathcal{R}_k^{-1} (\|\mathbf{r}_0\| \boldsymbol{\eta}_k) = \mathbf{z}_0 + [\mathcal{P}_{k-1} \quad \pi_k] (\|\mathbf{r}_0\| \boldsymbol{\eta}_k) \\ &= \mathbf{z}_{k-1} + \|\mathbf{r}_0\| \pi_k ((\boldsymbol{\eta}_k)_k). \end{aligned}$$

Edellä kuvatusta menetelmästä saadaan algoritmi 3.8.2. Jäännöksen normi $\|\mathbf{r}_k\|$ olisi nytkin helposti käytettävissä jokaisella kierroksella laskemalla $\beta |(\mathbf{v}_j)_{j+1}|$.

3.9 C-linearisoitu yhtälö

Kerrotaan yhtälö $\mathcal{M}_0(\mathbf{z}) = \mathbf{b}$ vasemmalta puolittain operaattorilla \mathcal{M}_0^{k-1} , jolloin parillisilla k saadaan C-lineaarinen (matriisi-)yhtälö $(M_{\#} \overline{M_{\#}})^{k/2} \mathbf{z} = \widehat{\mathbf{b}}$, missä $\widehat{\mathbf{b}} = (M_{\#} \overline{M_{\#}})^{\frac{k}{2}-1} M_{\#} \mathbf{b}$. Kun $\kappa \neq 0$ ja tarkastellaan yhtälöä $\mathcal{M}_{\kappa}(\mathbf{z}) = \mathbf{b}$, voidaan kertoa vasemmalta polynomilla

$$p(\mathcal{M}_0) = \sum_{j=0}^{k-1} \left(-\frac{1}{\kappa} \mathcal{M}_0\right)^j \frac{1}{\kappa},$$

mikä saadaan katkaisemalla käänteisoperaattorin $(\kappa \mathbf{I} + \mathcal{M}_0)^{-1} = \left(\kappa \left(\mathbf{I} + \frac{1}{\kappa} \mathcal{M}_0\right)\right)^{-1} = \left(\mathbf{I} + \frac{1}{\kappa} \mathcal{M}_0\right)^{-1} \frac{1}{\kappa}$ geometrinen sarja. Saadaan matriisiyhtälö

$$\left(\mathbf{I} - \frac{1}{\kappa} (M_{\#} \overline{M_{\#}})^{k/2}\right) \mathbf{z} = \widehat{\mathbf{b}}, \quad \widehat{\mathbf{b}} = p(\mathcal{M}_0)(\mathbf{b}). \quad (3.53)$$

Algoritmi 3.8.2 MINRES menetelmä symmetrisille $M_{\#}$. Suorittaa k iteraatiota tai pysähtyy, jos tarkka ratkaisu saavutetaan. Alkuarvaus on z_0 .

- 1: $r_0 \leftarrow b - \mathcal{M}_{\kappa}(z_0)$, $\beta \leftarrow \|r_0\|$, $q_1 \leftarrow r_0/\beta$, $v_0 \leftarrow 1$
- 2: $\beta_1 \leftarrow 0$, $q_0 \leftarrow 0$, $\pi_0 \leftarrow 0$, $\pi_{-1} \leftarrow 0$, $z \leftarrow z_0$
- 3: **for** $j = 1$ to k **do**
- 4: $w \leftarrow M_{\#}\bar{q}_j - \beta_j q_{j-1}$
- 5: $\alpha_j \leftarrow \langle w, q_j \rangle$
- 6: $w \leftarrow w - \alpha_j q_j$
- 7: $\beta_{j+1} \leftarrow \|w\|$
- 8: $s \leftarrow [0 \ \beta_j \ \alpha_j \ \beta_{j+1}]^T$
- 9: Olkoon $\mathcal{S} : \mathbb{C} \rightarrow \mathbb{C}^4$ siten, että $\mathcal{S}(z) = \kappa e_3 z + s\bar{z}$
- 10: Jos $j > 2$, $\mathcal{S}_{1:2} \leftarrow \mathcal{U}_{j-2} \circ \mathcal{S}_{1:2}$
- 11: Jos $j > 1$, $\mathcal{S}_{2:3} \leftarrow \mathcal{U}_{j-1} \circ \mathcal{S}_{2:3}$
- 12: Valitaan unitaari $\mathcal{U}_j : \mathbb{C}^2 \rightarrow \mathbb{C}^2$ siten, että $\mathcal{U}_j \circ \mathcal{S}_{3:4} = [\# \ 0]^T$
- 13: $\mathcal{S}_{3:4} \leftarrow \mathcal{U}_j \circ \mathcal{S}_{3:4}$
- 14: $\pi_j \leftarrow (q_j - \pi_{j-1}\mathcal{S}_2 - \pi_{j-2}\mathcal{S}_1)\mathcal{S}_3^{-1}$
- 15: $v_j \leftarrow \begin{bmatrix} v_{j-1} \\ 0 \end{bmatrix} \in \mathbb{C}^{j+1}$, $(v_j)_{j:(j+1)} \leftarrow \mathcal{U}_j((v_j)_{j:(j+1)})$
- 16: $z \leftarrow z + \beta\pi_j((v_j)_j)$
- 17: **if** $\beta_{j+1} = 0$ **then**
- 18: lopeta algoritmi
- 19: **end if**
- 20: $q_{j+1} \leftarrow w/\beta_{j+1}$
- 21: **end for**
- 22: jäännöksen normi on $\beta|(v_j)_{j+1}|$

Erityisesti kertomalla yhtälö $\mathcal{M}_{\kappa}(z) = b$ vasemmalta operaattorilla $(\bar{\kappa}I - \mathcal{M}_0)$, saadaan $(\bar{\kappa}I - \mathcal{M}_0)\mathcal{M}_{\kappa} = (\bar{\kappa}I - \mathcal{M}_0)(\kappa I + \mathcal{M}_0) = |\kappa|^2 I - \mathcal{M}_0^2$, joten ratkaistavaksi yhtälöksi tulee

$$(|\kappa|^2 I - M_{\#}\overline{M_{\#}})z = \hat{b}, \quad \hat{b} = \bar{\kappa}b - M_{\#}\bar{b}. \quad (3.54)$$

Koska $(\bar{\kappa}I - \mathcal{M}_0)$ on kääntövä jos ja vain jos $(\kappa I + \mathcal{M}_0)$ on kääntövä, niin (3.54) on yhtäpitävä alkuperäisen yhtälön $\mathcal{M}_{\kappa}(z) = b$ kanssa, jos \mathcal{M}_{κ} on kääntövä. Samoin (3.53), jos $p(\mathcal{M}_0)$ on kääntövä.

Matriisien GMRESia voidaan nyt käyttää näiden ratkaisuun. Suppenemisnopeus riippuu matriisista $M_{\#}$ ja voi olla nopeampi tai hitaampi kuin edellä kuvattu reaali-lineaarinen GMRES (ks. kohta 5.1). Tässä täytyy myös laskea k kertaa matriisi-vektoritulo $M_{\#}$:llä jokaista uutta Krylovin aliavaruuden kantavektoria kohden.

On mielenkiintoista vertailla matriisien iteratiivisia menetelmiä tässä luvussa esiteltyihin reaali-lineaarisiin menetelmiin. Kun ratkaistavana on lineaarinen yhtälöryhmä $Ax = f$, käytetään matriisien iteratiivisia menetelmiä usein seuraavasti. Hermiittiselle ja positiividefiniitille A voidaan käyttää liittogradienttimenetelmää CG (myös MINRES tulisi kyseeseen). Hermiittiselle ja indefiniitille A käytetään MINRESiä (CG:n toimivuus ei tällöin ole taattu tarkassakaan aritmetiikassa). Kun A ei ole hermiittinen, on vaihtoehtojen useita, mutta tässä luvussa on esitelty vain GMRES.

Olkoon nyt $\kappa = 1$. Vinosymmetrisellä $M_{\#}$ on $(I - M_{\#}\overline{M_{\#}}) = (I + M_{\#}M_{\#}^*)$, joka on hermiittinen ja positiividefiniitti. Siten CG soveltuu tähän. Aiemmin jo todet-

tiin kohdan 3.7.2 vinosymmetrisen $M_{\#}$ reaalin lineaarisen menetelmän muistuttavan matriisien CG:tä. Seuraavassa on tälle menetelmälle lausetta 3.2.2 vastaava optimaalisuustulos. Tätä tapausta varten olisi myös mahdollista tehdä kohtaa 3.8.1 vastaava reaalin lineaarinen MINRES vinosymmetrisille $M_{\#}$.

Lause 3.9.1. *Olkoon $M_{\#}$ vinosymmetrinen ja ratkaistavana on yhtälö $z + M_{\#}\bar{z} = b$. Tällöin kohdan 3.7.2 menetelmän likimääräisratkaisulle $z_k \in z_0 + \mathcal{K}_k(\mathcal{M}_1, r_0)$ pätee*

$$\|z_* - z_k\|_A \leq \|z_* - y\|_A \quad \text{kaikilla } y \in z_0 + \mathcal{K}_k(\mathcal{M}_1, r_0),$$

missä $A = I - M_{\#}\overline{M_{\#}}$ on hermiittinen ja positiividefiniitti, z_* on tarkka ratkaisu ja $\|z\|_A = \langle Az, z \rangle^{1/2}$.

Todistus. Kohdassa 3.7.2 nähtiin, että likimääräisratkaisu $z_k = z_0 + Q_k u_k$ saadaan ratkaisemalla yhtälö

$$u_k + T_k \bar{u}_k = \|r_0\| e_1, \quad (3.55)$$

jolla on yksikäsitteinen ratkaisu u_k . Tässä $\|r_0\| e_1 = Q_k^* r_0$ ja $T_k = Q_k^* M_{\#} \overline{Q_k}$. Tällöin operaattori $u \mapsto u - T_k \bar{u}$ on kääntyvä, joten (3.55) on yhtäpitävä yhtälön

$$(I - T_k \overline{T_k}) u_k = Q_k^* r_0 - T_k \overline{Q_k^* r_0}$$

kanssa. Sijoittamalla $T_k = Q_k^* M_{\#} \overline{Q_k}$ tästä saadaan

$$Q_k^* (I - M_{\#} \overline{M_{\#}}) Q_k u_k = Q_k^* (r_0 - M_{\#} \overline{r_0}).$$

Väite seuraa nyt lauseesta 3.2.2. □

Symmetrisellä $M_{\#}$ on $(I - M_{\#} \overline{M_{\#}}) = (I - M_{\#} M_{\#}^*)$ ja tämä on mahdollisesti indefiniitti. Täten tulee kyseeseen matriisien MINRES eikä CG välttämättä toimi. Kohdan 3.7.1 symmetrisen $M_{\#}$ menetelmä muistuttaa kohdan 3.7.2 menetelmää ja siten matriisien CG:tä. Takeita toimivuudesta ei kuitenkaan ole mahdollisesti puuttuvan LU-hajotelman vuoksi. Sen sijaan voidaan käyttää kohdan 3.8.1 reaalin lineaarista MINRESiä.

Yleisen $M_{\#}$ tapauksessa matriisille $(I - M_{\#} \overline{M_{\#}})$ voidaan käyttää GMRESiä ja operaattorille $z \mapsto z + M_{\#} \bar{z}$ kohdan 3.8 reaalin lineaarista GMRESiä.

3.10 Yleisen \mathbb{R} -lineaarisen operaattorin menetelmät

Yleisen operaattorin $\mathcal{M}(z) = Mz + M_{\#} \bar{z}$ tapauksessa palataan Petrov-Galerkinin ehdosta saatuun yhtälöön (3.17). Kun z_0 on alkuarvaus ja $r_0 = b - \mathcal{M}(z_0)$ sen jäännös, voidaan ryhtyä laskemaan Arnoldin menetelmällä ortonormaaleja vektoreita kaavoilla (3.23)-(3.25) kuten algoritmossa 3.10.1. Olkoon Q_k :n sarakkeet lasketut vektorit. Etsitään ratkaisua Q_k :n sarakkeiden virittämästä aliavaruudesta ja vaaditaan jäännös kohtisuoraksi samaa aliavaruutta vastaan, jolloin saadaan yhtälö

$$Q_k^* M Q_k u_k + Q_k^* M_{\#} \overline{Q_k} \bar{u}_k = Q_k^* r_0, \quad (3.56)$$

mistä ratkaisemalla u_k saadaan likimääräisratkaisuksi $z_k = z_0 + Q_k u_k$.

Edellä kuvattu vastaa kohdan 3.7 mukaista operaattorin \mathcal{M}_κ täyden ortogonalisointimenetelmää, mutta matriisit $\mathbf{Q}_k^* \mathbf{M} \mathbf{Q}_k$ ja $\mathbf{Q}_k^* \mathbf{M}_\# \overline{\mathbf{Q}}_k$ eivät nyt yksinkertaistu Hessenbergin muotoon. Kaavan (3.28) sijaan saadaan $\mathbf{Q}_k^* \mathbf{M} \mathbf{Q}_k + \mathbf{Q}_k^* \mathbf{M}_\# \overline{\mathbf{Q}}_k = \mathbf{H}_k$, mitä ei voida hyödyntää. Lisäksi kantavektorit \mathbf{Q}_k eivät välttämättä viritä Krylovin aliavaruutta $\mathcal{K}_k(\mathcal{M}, \mathbf{r}_0)$. Tästä huolimatta tätä kantaa käyttäen yhtälö (3.56) voidaan ratkaista luvun 2 menettelyin. Näin saadaan viitteessä [1] esitetty Galerkinin approksimaatiomenetelmä `r1_Ga1` yleisille \mathbb{R} -lineaarille operaattoreille.

Algoritmi 3.10.1 Yksinkertainen menetelmä ortonormaalin kannan muodostamiseksi \mathbb{R} -operaattorille $\mathcal{M} : \mathbb{C}^n \rightarrow \mathbb{C}^n$. Laskee matriisin $\mathbf{Q} \in \mathbb{C}^{n \times k}$, jonka sarakkeet ovat kantavektorit.

```

1:  $\mathbf{q}_1 \leftarrow \mathbf{r}_0 / \|\mathbf{r}_0\|$ ,  $\mathbf{Q} \leftarrow [\mathbf{q}_1]$ 
2: for  $j = 2$  to  $k$  do
3:    $\mathbf{v} \leftarrow \mathcal{M}(\mathbf{q}_{j-1})$ 
4:    $\mathbf{w} \leftarrow \mathbf{v} - \sum_{i=1}^{j-1} \langle \mathbf{v}, \mathbf{q}_i \rangle \mathbf{q}_i$ 
5:    $\mathbf{q}_j \leftarrow \mathbf{w} / \|\mathbf{w}\|$ 
6:    $\mathbf{Q} \leftarrow [\mathbf{Q} \ \mathbf{q}_j]$ 
7: end for
    
```

Kompleksisen sisätulon sijaan voidaan myös käyttää reaalista sisätuloa, jolloin saadaan algoritmi 3.10.2. Tällöin saadut vektorit sisältyvät (reaalikertoimisiin) Krylovin aliavaruuksiin, mutta nyt ortogonalisointi ei generoi ortonormaalia kantaa kompleksisen sisätulon mielessä. Reaalisen sisätulon mielessä kohtisuorat vektorit saattavat myös olla kompleksisesti lineaarisesti riippuvat, kuten vektorit $[1 \ 1 \ \dots \ 1]^T$ ja $[i \ i \ \dots \ i]^T$ osoittavat. Petrov-Galerkinin ehdosta (3.16) saatua yhtälöä (3.56) voidaan kuitenkin ryhtyä ratkaisemaan tälläkin \mathbf{Q}_k vaikka menetelmän rikkoutuminen onkin mahdollista.

On huomattava, että algoritmi 3.10.2 muodostaa saman kannan kuin algoritmi 3.3.1, kunhan \mathbf{A} on operaattoria \mathcal{M} vastaava $2n \times 2n$ -matriisi muodostettuna esimerkiksi käyttäen tapaa (1.3) tai (1.4). Tämän voi nähdä seuraavasti. Olkoot $\mathbf{z}, \mathbf{w} \in \mathbb{C}^n$, $\mathbf{z} = \mathbf{x} + i\mathbf{y}$, $\mathbf{w} = \mathbf{u} + i\mathbf{v}$. Tällöin

$$\operatorname{Re} \langle \mathbf{z}, \mathbf{w} \rangle = \operatorname{Re} ((\mathbf{u}^T - i\mathbf{v}^T)(\mathbf{x} + i\mathbf{y})) = \mathbf{u}^T \mathbf{x} + \mathbf{v}^T \mathbf{y} = \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix}^T \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix},$$

joka on vastaava reaalinen sisätulo.

Algoritmi 3.10.2 \mathbb{R} -lineaarinen menetelmä kannan muodostamiseksi Krylovin aliavaruudelle $\mathcal{K}_k(\mathcal{M}, \mathbf{r}_0)$. Laskee matriisin $\mathbf{Q} \in \mathbb{C}^{n \times k}$, jonka sarakkeet ovat kantavektorit.

```

1:  $\mathbf{q}_1 \leftarrow \mathbf{r}_0 / \|\mathbf{r}_0\|$ ,  $\mathbf{Q} \leftarrow [\mathbf{q}_1]$ 
2: for  $j = 2$  to  $k$  do
3:    $\mathbf{v} \leftarrow \mathcal{M}(\mathbf{q}_{j-1})$ 
4:    $\mathbf{w} \leftarrow \mathbf{v} - \sum_{i=1}^{j-1} \operatorname{Re}(\langle \mathbf{v}, \mathbf{q}_i \rangle) \mathbf{q}_i$ 
5:    $\mathbf{q}_j \leftarrow \mathbf{w} / \|\mathbf{w}\|$ 
6:    $\mathbf{Q} \leftarrow [\mathbf{Q} \ \mathbf{q}_j]$ 
7: end for
    
```

Sekä algoritmin 3.10.1 että 3.10.2 ortogonalisoinnissa voidaan (ja tietokonetoteutuksessa kannattaa) käyttää muunnettua Gram-Schmidt prosessia.

Mikäli algoritmin 3.10.2 tarjoamat kantavektorit ortonormeerataan lopuksi myös kompleksisen sisätulon suhteen, saadaan ortonormaali kanta Krylovin aliavaruuteen. Tämän tavoitteena on lisätä numeerista stabiiliutta. Näin saadaan algoritmi 3.10.3, joka on työn valvojan ehdottama [9].

Algoritmi 3.10.3 \mathbb{R} -lineaarinen menetelmä ortonormaalin kannan muodostamiseksi Krylovin aliavaruudelle $\mathcal{K}_k(\mathcal{M}, \mathbf{r}_0)$. Laskee matriisin $\mathbf{Q} \in \mathbb{C}^{n \times k}$, jonka sarakkeet ovat kantavektorit.

- 1: $\mathbf{U} \leftarrow$ Algoritmi 3.10.2
 - 2: Lasketaan (suppea) QR-hajotelma, $\mathbf{QR} = \mathbf{U}$. Hajotelman \mathbf{Q} on tämän algoritmin tulos.
-

3.3.1 Esitetyistä kolmesta vaihtoehdoisesta algoritmista kantavektorimatriisin \mathbf{Q}_k muodostamiseksi, saadaan kolme mahdollista menetelmää yhdistettäessä yhtälön (3.56) ratkaisemisen kanssa. Näiden kaikkien hankaluutena on, ettei ratkaistavaksi tulevaan yhtälöön muodostu Hessenbergin matriiseja. Sikäli nämä menetelmät jäävät epätydyttäväiksi verrattuna aikaisemmin tässä luvussa esitettyihin operaattorin \mathcal{M}_k menetelmiin.

Luku 4

ILU-pohjustus

4.1 Johdanto

Olkoon $A \in \mathbb{C}^{n \times n}$, $b \in \mathbb{C}^n$ ja tarkastellaan jälleen lineaarista yhtälöryhmää

$$Ax = b. \quad (4.1)$$

Edellisessä luvussa esiteltiin joitakin iteratiivisia ratkaisumenetelmiä, jotka pyrkivät tarkentamaan tämän yhtälön likimääräisratkaisua jokaisella kierroksellaan. Pohjustuksen tavoitteena on suppenemisen nopeuttaminen. Suppenemisnopeus riippuu matriisin A ominaisuuksista, kuten esimerkiksi GMRESille saatiin ominaisarvoista riippuva seuraus 3.4.2. Pohjustuksessa yhtälö muunnetaan sopivasti yhtäpitävään muotoon siten, että kerroinmatriisin ominaisuudet soveltuvat käytettävälle iteratiiviselle menetelmälle paremmin. GMRESin tapauksessa voidaan odottaa nopeutumista, mikäli pohjustetun kerroinmatriisin ominaisarvot kasautuvat selkeästi origosta irti olevan pisteen ympärille.

Vasemmanpuoleisessa pohjustuksessa korvataan yhtälö (4.1) yhtälöllä

$$M^{-1}Ax = M^{-1}b.$$

Tähän matriisi M^{-1} saadaan varsinaisesta pohjustusmenetelmästä, joita on kehitetty lukuisia vaihdellen yleispätevyyteen pyrkivistä menetelmistä tarkoin tiettyä tehtävää varten säädettyihin menetelmiin. Tavallisesti matriisin M tulee olla jossain mielessä lähellä matriisia A .

Oikeanpuoleisessa pohjustuksessa toisaalta korvataan yhtälö (4.1) yhtälöllä

$$AM^{-1}y = b.$$

Alkuperäisen yhtälön ratkaisu saadaan laskemalla lopuksi $x = M^{-1}y$. Tässäkin tulisi olla $M \approx A$.

Yllä mainittua matriisia M^{-1} ei yleensä lasketa eksplisiittisesti ja kerrota matriisin A kanssa, vaan näihin suhtaudutaan operaatioina $v \mapsto M^{-1}v$ ja $v \mapsto Av$. Iteratiivisissa menetelmissä tyypillisesti tarvitseekin vain laskea jokaisella kierroksella yhdistetyn operaattorin $v \mapsto M^{-1}Av$ tulos jollakin edellisestä iteraatiokierroksesta riippuvalla vektorilla v . Mahdollisesti monimutkaista käänteismatriisin laskentaa ja

matriisiin A harvuusrakenteen sotkevaa kertolaskua ei tarvitse tehdä. Lisäksi vaatimuksena on, että operaatio $v \mapsto M^{-1}v$ on nopea suorittaa. Kokonaislaskenta-aikaa kasvattava pohjustusmenetelmä ei palvele mitään tarkoitusta, vaikka se iteraatioiden lukumäärää vähentäisikin.

Tässä luvussa käsitellään ILU-pohjustusta, joka pyrkii olemaan yleispätevä pohjustin. Edellisten lukujen tapaan esitellään aluksi matriisitapaus ja tämän jälkeen yleistetään reaalilineaarisille operaattoreille.

4.2 Matriisiyhtälöiden ILU-pohjustus

ILU-pohjustuksessa ryhdytään muodostamaan ala- ja yläkolmiomatriiseja L ja U tavallisen LU-hajotelman tapaan. Kuitenkin laskennan edetessä näistä pudotetaan pois joitakin alkioita asettamalla ne nolliksi. Tästä tuleekin menetelmän nimitys epä-täydellinen LU-hajotelma (Incomplete LU). Lopputulokselle pätee vain suunnilleen $A \approx LU$.

Edellä olevassa johdannossa oleva matriisi M on nyt $M = LU$, jota ei kuitenkaan tarvitse muodostaa eksplisiittisesti tietokoneen muistiin. Käänteisoperaatio $v \mapsto M^{-1}v$ voidaan laskea yhtälöistä

$$Ly = v$$

$$Ux = y$$

eteen- ja taaksepäin sijoituksin. Pohjustusoperaation tulos on $x = M^{-1}v$. Kun matriiseja L ja U muodostettaessa alkioita pudotetaan pois riittävästi, tulee operaatiosta $v \mapsto U^{-1}L^{-1}v$ nopea. Toisaalta, jos alkioita pudotetaan pois liikaa, tulee approksimaatiosta $A \approx LU$ huono eikä menetelmä nopeuta iteraation suppenemistä.

ILU-hajotelman poispudotettavat alkiot voidaan valita eri strategioin ja näistä saadaan kokonainen perhe pohjustusmenetelmiä. Seuraavassa esitellään kaksi mahdollisuutta. Huomattakoon, ettei kummassakaan käytetä tuentaa luvun 2 tapaan. Näin hajotelman laskeminen nopeutuu, mutta ei ole numeerisesti stabiili. Lisäksi tukialkio saattaa olla nolla. Jälkimmäisestä ongelmasta pääsisi eroon korvaamalla tukialkion jollakin nollostä poikkeavalla luvulla, mutta tämä huonontaa $LU \approx A$ tarkkuutta. Tuenta on mahdollista myös ILU-hajotelmille, mutta sen käsittely on tämän työn ulkopuolella.

4.2.1 ILU-hajotelma ennaltavalitulla nollakuviolla

Olkoon $A \in \mathbb{C}^{n \times n}$ ja $Z \subset \{(i, j) \in \mathbb{Z} \times \mathbb{Z} \mid 1 \leq i, j \leq n, i \neq j\}$. Tulkitaan seuraavassa, että joukko Z on ILU-hajotelman ns. nollakuvio. Jos indeksipari (i, j) kuuluu joukkoon Z , niin tätä indeksiparia vastaava alkio pudotetaan pois. Mikäli (i, j) on aidossa alakolmiossa ($i > j$), niin asetetaan matriisiin L alkio $l_{ij} = 0$. Jos taas (i, j) on aidossa yläkolmiossa ($i < j$), niin asetetaan $u_{ij} = 0$ matriisissa U . Matriisien lävistäjälkioita ei voida asettaa nolliksi, koska L :n ja U :n täytyy olla kääntyviä.

Aluksi asetetaan kaikki A :n alkiot $a_{ij} = 0$, jos (i, j) kuuluu nollakuvioon. Sitten

ryhdytään laskemaan tämän A :n LU-hajotelmaa Gaussin eliminoinnein. Tämä suoritetaan muuten tuttuun tapaan, mutta rivien yhteenlaskussa nollakuvioon kuuluvia (aluksi nolliksi asetettuja) alkioita ei ylikirjoiteta. Näin saadaan algoritmi 4.2.1. Tässä matriisit L ja U on pidetty erillään. Yleensä LU-toteutuksissa hajotelma muodostetaan suoraan matriisin A päälle, jolloin sen yläkolmioon muodostuu U :n alkiot ja aitoon alakolmioon L :n aito alakolmio eikä L :n yksiköksi tiedettyä lävistäjää tarvitse tallentaa.

Algoritmi 4.2.1 Laskee ILU-hajotelman matriisit $L = (l_{ij}), U = (u_{ij})$ matriisille $A \in \mathbb{C}^{n \times n}$ nollakuviolla Z .

```

1:  $L \leftarrow I$ 
2:  $U \leftarrow A$ 
3:  $u_{ij} \leftarrow 0$  kaikille  $(i, j) \in Z$ .
4: for  $i = 1$  to  $n$  do
5:   for  $k = i + 1$  to  $n$  do
6:     if  $(k, i) \notin Z$  then
7:        $l_{ki} \leftarrow u_{ki}/u_{ii}$ 
8:        $u_{kr} \leftarrow u_{kr} - l_{ki}u_{ir}$  kaikilla  $r > i$  s.e.  $(k, r) \notin Z$ 
9:     end if
10:     $u_{ki} \leftarrow 0$ 
11:  end for
12: end for

```

Algoritmi 4.2.2 laskee saman ILU-hajotelman, mutta kahden uloimman silmukan järjestyksestä on vaihdettu. Vastaava MATLAB-funktio `ilu_v2` löytyy liitteestä B. Silmukoiden järjestyksen vaihto vastaa Gaussin eliminoinnissa sitä, että se suoritetaan valmiiksi järjestyksessä rivi kerrallaan ensimmäisestä aloittaen. Järjestyksen vaihdolla saadaan paremmin harvoille matriiseille soveltuva algoritmi, koska ainoastaan työn alla olevaa riviä muokataan kerrallaan. Rivin laskennan valmistuttua se voidaan tallentaa tietokoneen muistiin harvan matriisin tietorakenteen muotoon ja siirtyä käsittelemään seuraavaa riviä. Ennaltavalitulla nollakuviolla ei tällä ole laskemisen kannalta suurta merkitystä. Seuraavassa kohdassa esitettävälle ILUT-hajotelmalle silmukoiden vaihto on kuitenkin tärkeä muunnos ja sitä käytetään tästä lähtien yksinomaan (ILUT:lle laskennan lopputulos myös riippuu silmukoiden järjestyksestä).

Tähän mennessä nollakuvion valinnasta ei olla vielä sanottu mitään. Seuraava määritelmä tekee yksinkertaisen valinnan.

Määritelmä 4.2.1. Matriisin $A \in \mathbb{C}^{n \times n}$ ILU(0)-hajotelmaksi kutsutaan ILU-hajotelmaa, missä nollakuvioksi on valittu A :n nollakuvio. Tarkemmin $(i, j) \in Z \Leftrightarrow a_{ij} = 0$.

Liitteessä olevien MATLAB-koodien avulla ILU(0)-hajotelman voi laskea suorittamalla `ilu_v2(A, A)`.

4.2.2 ILU-hajotelma mukautetulla nollakuviolla

Tarkastellaan edellisen kohdan algoritmia 4.2.2. Siinä alkioiden pudotuksia tapahtuu kahdessa paikassa; sekä **for** k-silmukassa että **for** j-silmukassa. Alkioiden pu-

Algoritmi 4.2.2 Laskee ILU-hajotelman matriisit $L = (l_{ij}), U = (u_{ij})$ matriisille $A \in \mathbb{C}^{n \times n}$ nollakuviolla Z .

```

1:  $L \leftarrow I$ 
2:  $U \leftarrow A$ 
3: for  $i = 1$  to  $n$  do
4:   for  $k = 1$  to  $i - 1$  do
5:     if  $(i, k) \in Z$  then
6:        $l_{ik} \leftarrow 0$ 
7:     else
8:        $l_{ik} \leftarrow u_{ik}/u_{kk}$ 
9:        $u_{i*} \leftarrow u_{i*} - l_{ik}u_{k*}$ 
10:    end if
11:     $u_{ik} \leftarrow 0$ 
12:  end for
13:  for  $j = i + 1$  to  $n$  do
14:    if  $(i, j) \in Z$  then
15:       $u_{ij} \leftarrow 0$ 
16:    end if
17:  end for
18: end for

```

dotus perustuu ennaltavalittuun nollakuviioon Z . Mukautetulla nollakuviolla tarkoitetaan tässä sellaista alkioiden pudotusmenettelyä, joka ottaa huomioon niiden (itseis)arvojen suuruuden. Hajotelmassa säilytetään itseisarvoltaan suuret alkiot pois pudotettavien ollessa pieniä, jolloin approksimaatiosta $A \approx LU$ saadaan parempi kuin suuruudet huomiotta jättävissä ennaltavalitun nollakuvion menetelmissä.

Matriisien L ja U alkiosta pois pudotetaan sellaiset, joiden suhteellinen koko (verrattuna A :n vastaavan koko rivin normiin) on riittävän pieni. Olkoon $\epsilon > 0$ ja a_{i*} matriisin A rivi i . Tällöin hajotelman laskettu alkiot l_{ij} pudotetaan pois mikäli $|l_{ij}| < \epsilon \|a_{i*}\|$. Samoin laskettu u_{ij} pudotetaan pois, jos $|u_{ij}| < \epsilon \|a_{i*}\|$. Tässä normi on tavallinen 2-normi, mutta myös muiden normien käyttö on mahdollista.

Lisäksi valitaan argumentti p , joka määrää sekä L :n aidossa alakolmiossa että U :n aidossa yläkolmiossa kullakin rivillä olevien nolasta poikkeavien alkioiden suurimman määrän. Alkiosta säilytetään p itseisarvoltaan suurinta ja loput pudotetaan pois. Tällä saadaan lisää hallittavuutta L :n ja U :n harvuusrakenteeseen.

Tällaisella menettelyllä muodostettua hajotelmaa kutsutaan ILUT:ksi ja se on esitetty algoritmina 4.2.3. Vastaava MATLAB-funktio `ilut` löytyy liitteestä B.

4.3 \mathbb{R} -lineaaristen yhtälöiden ILU-pohjustus

Olkoon $\mathcal{M} : \mathbb{C}^n \rightarrow \mathbb{C}^n$ reaalilineaarinen operaattori ja $b \in \mathbb{C}^n$. Yhtälölle

$$\mathcal{M}(z) = b \tag{4.2}$$

on edellisessä luvussa esitetty iteratiivisia menetelmiä, joita matriisiyhtälöiden tapaan lähdetään nyt pohjustamaan. Tavoitteena on iteraatiokierrosten ja kokonaislaskenta-ajan pienentäminen käyttäen nopeasti laskettavaa pohjustinta.

Algoritmi 4.2.3 ILUT(ϵ, p). Laskee ILU-hajotelman matriisit $L = (l_{ij}), U = (u_{ij})$ matriisille $A \in \mathbb{C}^{n \times n}$ käyttäen mukautettua nollakuviota. Argumentti ϵ määrää alkioiden suuruustoleranssin suhteessa rivikokoon. Tuloksena olevissa L :ssä ja U :ssa on lävistäjäalkioiden lisäksi korkeintaan p nollasta poikkeavaa alkia kullakin rivillä.

```
1:  $L \leftarrow I$ 
2:  $U \leftarrow A$ 
3: for  $i = 1$  to  $n$  do
4:    $m \leftarrow \|u_{i*}\|$ 
5:   for  $k = 1$  to  $i - 1$  do
6:      $c \leftarrow u_{ik}/u_{kk}$ 
7:     if  $|c| < \epsilon m$  then
8:        $l_{ik} \leftarrow 0$ 
9:     else
10:       $l_{ik} \leftarrow c$ 
11:       $u_{i*} \leftarrow u_{i*} - cu_{k*}$ 
12:    end if
13:     $u_{ik} \leftarrow 0$ 
14:  end for
15:  for  $j = i + 1$  to  $n$  do
16:    if  $|u_{ij}| < \epsilon m$  then
17:       $u_{ij} \leftarrow 0$ 
18:    end if
19:  end for
20:  if  $i > p + 1$  then
21:    Säilytetään alkiosta  $l_{i,1}, \dots, l_{i,i-1}$   $p$  itseisarvoltaan suurinta; muut asetetaan nolliksi.
22:  end if
23:  if  $i < n - p$  then
24:    Säilytetään alkiosta  $u_{i,i+1}, \dots, u_{i,n}$   $p$  itseisarvoltaan suurinta; muut asetetaan nolliksi.
25:  end if
26: end for
```

Olkoon $\mathcal{P} : \mathbb{C}^n \rightarrow \mathbb{C}^n$ jonkin pohjustusmenetelmän antama operaattori siten, että $\mathcal{P} \approx \mathcal{M}$. Vasemmanpuoleisessa pohjustuksessa ratkotaan yhtälön (4.2) sijaan yhtälöä

$$\mathcal{P}^{-1}\mathcal{M}(z) = \mathcal{P}^{-1}(b).$$

Matriisitapausta vastaten, oikeanpuoleinen pohjustus on

$$\mathcal{M}\mathcal{P}^{-1}(w) = b,$$

mistä alkuperäisen yhtälön (4.2) ratkaisu saadaan laskemalla $z = \mathcal{P}^{-1}(w)$.

On huomattava, ettei siirretty antilineaarinen operaattori $\mathcal{M}_\kappa(z) = \kappa z + M_\# \bar{z}$ edellä kuvatulla pohjustuksella säilytä muotoaan. Olkoon $\mathcal{P}^{-1}(z) = Sz + S_\# \bar{z}$, jolloin laskemalla saadaan

$$\mathcal{P}^{-1}\mathcal{M}_\kappa(z) = (\kappa S + S_\# \overline{M_\#})z + (SM_\# + \bar{\kappa} S_\#) \bar{z},$$

mistä nähdään ettei lineaarinen osa yleensä ole muotoa $\widehat{\kappa}I$. Edellisessä luvussa tälle tärkeälle erikoistapaukselle esiteltiin reaalilineaarinen GMRES, mutta sen pohjustaminen on tällä hetkellä avoin ongelma. Pohjustus koskeekin tässä työssä ainoastaan yleiselle operaattorille kohdassa 3.10 esitetyjä menetelmiä. Näidenkin menetelmien hankaluutena on, että lineaarisella ja antilinearisella osalla täytyy kyetä operoimaan erikseen. Kun operaattorin \mathcal{M} matriisit M ja $M_\#$ ovat tunnettuja, täytyisi tämän jälkeen pohjustetusta operaattorista $\mathcal{P}^{-1}\mathcal{M}$ vielä erotella vastaavat matriisit, mikä ei yleensä ole mielekäästä. Tämän vuoksi joudutaan käyttämään seuraavia kaavoja, jotka pätevät mille tahansa operaattorille $\mathcal{N}(z) = Nz + N_\# \bar{z}$,

$$\begin{aligned} Nz &= \frac{1}{2}(\mathcal{N}(z) - i\mathcal{N}(iz)), \\ N_\# \bar{z} &= \frac{1}{2}(\mathcal{N}(z) + i\mathcal{N}(iz)). \end{aligned}$$

Valitsemalla $\mathcal{N} = \mathcal{P}^{-1}\mathcal{M}$ nähdään, että pohjustetun operaattorin osilla voidaan iteratiivista menetelmää varten aina laskea, mutta tähän tarvitaan kaksi täyttä operaattori-vektori tuloa. Algoritmia 3.10.3 vastaavan menetelmän pohjustaminen on erityisen työlästä, koska pohjustusoperaatiot täytyy laskea yllä kuvatulla tavalla ensin algoritmia 3.10.2 ajettaessa ja sen jälkeen uudestaan QR-hajotelmasta saatujen vektoreiden kanssa.

\mathbb{R} -linearisessa ILU-pohjustuksessa operaattoriksi \mathcal{P} valitaan $\mathcal{P} = \mathcal{L}\mathcal{U}$, missä operaattorit $\mathcal{L} : \mathbb{C}^n \rightarrow \mathbb{C}^n$ ja $\mathcal{U} : \mathbb{C}^n \rightarrow \mathbb{C}^n$ muodostavat \mathbb{R} -lineaarisen epätäydellisen LU-hajotelman. Perustana on luvussa 2 esitetty reaalilineaarinen LU-hajotelma.

4.3.1 \mathbb{R} -lineaarinen ILU-hajotelma ennaltavalitulla nollakuviolla

Olkoon $\mathcal{M} : \mathbb{C}^n \rightarrow \mathbb{C}^n$, $\mathcal{M}(z) = Mz + M_\# \bar{z}$, reaalilineaarinen operaattori ja

$$Z \subset \{(i, j) \in \mathbb{Z} \times \mathbb{Z} \mid 1 \leq i, j \leq n, i \neq j\}$$

nollakuvio. Ryhdytään muodostamaan \mathcal{M} :n LU-hajotelmaa kohdan 2.1 tavalla, mutta asettamalla nollassi nollakuviossa olevat alkioit. Menetelmä on täysin analoginen kohdan 4.2.1 matriisitapauksen kanssa. Matriisien tilalla voidaan ajatella manipuloitavan muodossa (1.1) esitetyjä reaalilineaarisia operaattoreita.

Aluksi asetetaan nolliksi ne matriisien M ja $M_{\#}$ alkiot, joiden indeksipari (i, j) kuuluu nollakuvioon. Sitten suoritetaan Gaussin eliminointeja luvun 2 tapaan, mutta rivejä yhteenlaskiessa nollakuvioon kuuluvaa indeksiparia (i, j) vastaavaa (alussa nollassi asetettua) skalaarioperaattoria ei koskaan ylikirjoiteta. Algoritmi 4.3.1 laskee ILU-hajotelman vastaten matriisitapauksen algoritmia 4.2.2, missä eliminointi suoritetaan valmiiksi rivi kerrallaan. Vastaava MATLAB-funktio `rl_ilu` löytyy liitteestä B.

Algoritmi 4.3.1 Laskee \mathbb{R} -ILU-hajotelman operaattorit $\mathcal{L}(z) = Lz + L_{\#}\bar{z}$, $\mathcal{U}(z) = Uz + U_{\#}\bar{z}$ operaattorille $\mathcal{M}(z) = Mz + M_{\#}\bar{z}$ nollakuviolla Z .

```

1:  $L \leftarrow I, L_{\#} \leftarrow 0$ 
2:  $U \leftarrow M, U_{\#} \leftarrow M_{\#}$ 
3: for  $i = 1$  to  $n$  do
4:   for  $k = 1$  to  $i - 1$  do
5:     if  $(i, k) \in Z$  then
6:        $l_{ik} \leftarrow 0, l_{\#,ik} \leftarrow 0$ 
7:     else
8:        $[l_{ik} \ l_{\#,ik}] \leftarrow [u_{ik} \ u_{\#,ik}] \begin{bmatrix} u_{kk} & u_{\#,kk} \\ \bar{u}_{\#kk} & \bar{u}_{kk} \end{bmatrix}^{-1}$ 
9:        $u_{i*} \leftarrow u_{i*} - [l_{ik} \ l_{\#,ik}] \begin{bmatrix} u_{k*} \\ \bar{u}_{\#k*} \end{bmatrix}$ 
10:       $u_{\#,i*} \leftarrow u_{\#,i*} - [l_{ik} \ l_{\#,ik}] \begin{bmatrix} u_{\#,k*} \\ \bar{u}_{k*} \end{bmatrix}$ 
11:     end if
12:    $u_{ik} \leftarrow 0, u_{\#,ik} \leftarrow 0$ 
13:   end for
14:   for  $j = i + 1$  to  $n$  do
15:     if  $(i, j) \in Z$  then
16:        $u_{ij} \leftarrow 0, u_{\#,ij} \leftarrow 0$ 
17:     end if
18:   end for
19: end for

```

Yksinkertainen nollakuviio saadaan seuraavasta määritelmästä.

Määritelmä 4.3.1. Operaattorin \mathcal{M} \mathbb{R} -ILU(0):ksi kutsutaan hajotelmaa, jossa nollakuvioksi on valittu matriisien M ja $M_{\#}$ yhteinen nollakuviio. Tarkemmin $(i, j) \in Z \Leftrightarrow ((M)_{ij} = 0 \text{ ja } (M_{\#})_{ij} = 0)$.

Kun M ja $M_{\#}$ ovat MATLABin harvoja matriiseja, funktion `rl_ilu` avulla \mathbb{R} -ILU(0):n voi laskea suorittamalla `rl_ilu(M, Ma, max(spones(M), spones(Ma)))`.

4.3.2 \mathbb{R} -lineaarinen ILU-hajotelma mukautetulla nollakuviolla

Reaalilineaarinen ILUT, jota tämän jälkeen kutsutaan lyhyesti \mathbb{R} -ILUT:ksi, saadaan edellä olevan kohdan koodista vastaavin muutoksin kuin matriisien tapauksessa kohdassa 4.2.2.

Olkoon $\epsilon > 0$ ja $m_{i*}, m_{\#,i*}$ matriisien M ja $M_{\#}$ rivit i . Merkitään reaalilineaarisen LU-hajotelman operaattoreita $\mathcal{L}(z) = Lz + L_{\#}\bar{z}$ ja $\mathcal{U}(z) = Uz + U_{\#}\bar{z}$. Alkiot

l_{ij} ja $l_{\#,ij}$ pudotetaan pois, jos $|l_{ij}| + |l_{\#,ij}| < \epsilon(\|l_{i*}\|_1 + \|l_{\#,i*}\|_1)$. Samoin u_{ij} ja $u_{\#,ij}$ pudotetaan pois, jos $|u_{ij}| + |u_{\#,ij}| < \epsilon(\|u_{i*}\|_1 + \|u_{\#,i*}\|_1)$. Näissä on käytetty 1-normia, mutta toisenlainenkin normin valinta on mahdollista.

Jälleen, argumentti p määrää määrää sekä \mathcal{L} :n aidossa alakolmiossa että \mathcal{U} :n aidossa yläkolmiossa kullakin rivillä olevien nollassa poikkeavien skalaarioperaattoreiden suurimman määrän. Skalaareista säilytetään p normiltaan suurinta ja loput pudotetaan pois.

Tuloksena on algoritmi 4.3.2, joka vastaa matriisien algoritmia 4.2.3. Liitteessä B on vastaava MATLAB-funktio nimellä `r1_ilut`.

Yleisestä reaalmatriisiesityksestä (1.4) nähdään, että reaalilineaariset ILU-hajotelmat voidaan nähdä matriisien lohko-ILU-hajotelmina, missä lohkot ovat 2×2 -kokoisia. Täten näiden tehokkuuden voidaan olettaa olevan samaa luokkaa tunnettujen matriisien lohkohajotelmiin perustuvien pohjustusmenetelmien kanssa. Tehtyjä numeerisia kokeita esitetyille ILU-hajotelmille käydään läpi seuraavassa luvussa.

Algoritmi 4.3.2 \mathbb{R} -ILUT(ϵ, p). Laskee \mathbb{R} -ILUT-hajotelman operaattorit $\mathcal{L}(z) = Lz + L_{\#}\bar{z}$, $\mathcal{U}(z) = Uz + U_{\#}\bar{z}$ operaattorille $\mathcal{M}(z) = Mz + M_{\#}\bar{z}$ käyttäen muokattua nollakuviota. Argumentti ϵ määrää alkioiden suuruustoleranssin suhteessa rivikokoon. Tuloksena olevissa $L, L_{\#}$:ssä ja $U, U_{\#}$:ssa on lävistäjäalkioiden lisäksi korkeintaan p nollasta poikkeavaa alkioita kullakin rivillä.

```

1:  $L \leftarrow I, L_{\#} \leftarrow 0$ 
2:  $U \leftarrow M, U_{\#} \leftarrow M_{\#}$ 
3: for  $i = 1$  to  $n$  do
4:    $m \leftarrow \|u_{i*}\|_1 + \|u_{\#,i*}\|_1$ 
5:   for  $k = 1$  to  $i - 1$  do
6:      $\begin{bmatrix} c & c_{\#} \end{bmatrix} \leftarrow \begin{bmatrix} u_{ik} & u_{\#,ik} \end{bmatrix} \begin{bmatrix} u_{kk} & u_{\#,kk} \\ \bar{u}_{\#kk} & \bar{u}_{kk} \end{bmatrix}^{-1}$ 
7:     if  $|c| + |c_{\#}| < \epsilon m$  then
8:        $l_{ik} \leftarrow 0, l_{\#,ik} \leftarrow 0$ 
9:     else
10:       $l_{ik} \leftarrow c, l_{\#,ik} \leftarrow c_{\#}$ 
11:       $u_{i*} \leftarrow u_{i*} - \begin{bmatrix} c & c_{\#} \end{bmatrix} \begin{bmatrix} u_{k*} \\ \bar{u}_{\#k*} \end{bmatrix}$ 
12:       $u_{\#,i*} \leftarrow u_{\#,i*} - \begin{bmatrix} c & c_{\#} \end{bmatrix} \begin{bmatrix} u_{\#,k*} \\ \bar{u}_{k*} \end{bmatrix}$ 
13:     end if
14:      $u_{ik} \leftarrow 0, u_{\#,ik} \leftarrow 0$ 
15:   end for
16:   for  $j = i + 1$  to  $n$  do
17:     if  $|u_{ij}| + |u_{\#,ij}| < \epsilon m$  then
18:        $u_{ij} \leftarrow 0, u_{\#,ij} \leftarrow 0$ 
19:     end if
20:   end for
21:   if  $i > p + 1$  then
22:     Säilytetään alkiosta  $l_{i,1}, l_{\#,i,1}, \dots, l_{i,i-1}, l_{\#,i,i-1}$  ne, joiden indeksi ovat samat kuin  $p$ :n suurimman listassa  $|l_{i,1}| + |l_{\#,i,1}|, \dots, |l_{i,i-1}| + |l_{\#,i,i-1}|$ ; muut asetetaan nolliksi.
23:   end if
24:   if  $i < n - p$  then
25:     Säilytetään alkiosta  $u_{i,i+1}, u_{\#,i,i+1}, \dots, u_{i,n}, u_{\#,i,n}$  ne, joiden indeksi ovat samat kuin  $p$ :n suurimman listassa  $|u_{i,i+1}| + |u_{\#,i,i+1}|, \dots, |u_{i,n}| + |u_{\#,i,n}|$ ; muut asetetaan nolliksi.
26:   end if
27: end for

```

Luku 5

Numeerisia kokeita

Tässä luvussa esitetään suoritettujen numeeristen kokeiden koejärjestelyt ja tulokset. Kokeet jakaantuvat kahteen osaan. Aluksi tarkastellaan ja vertaillaan luvun 3 GMRES-menetelmiä, jonka jälkeen siirrytään luvun 4 pohjustusmenetelmien kokeisiin. Kaikki numeerinen laskenta tehtiin MATLABin versiolla 7.6.0.324 (R2008a).

5.1 GMRES kokeita

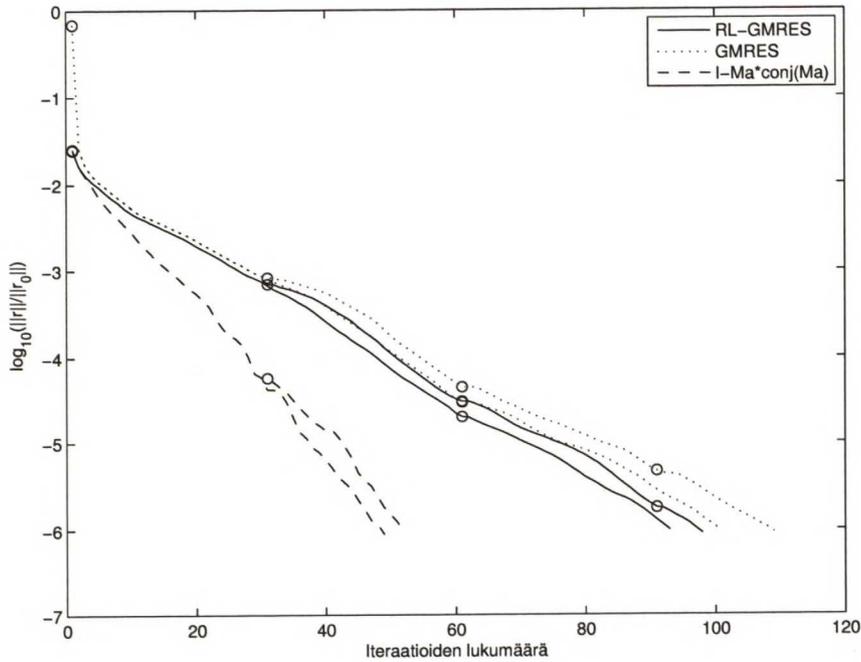
Kokeissa ratkaistiin yhtälöryhmä $\mathbf{z} + \mathbf{M}_\# \bar{\mathbf{z}} = \mathbf{b}$, missä vektori \mathbf{b} koostuu luvuista 1. Menetelminä olivat reaalilineaarinen GMRES (RL-GMRES, kohta 3.8) ja matriisien GMRES (kohta 3.4) sekä kaksinkertaiseksi kirjoitetulle reaalilaiselle yhtälöryhmälle (kohdan 3.8 loppu) että yhtälölle (3.54).

Kokeet suoritettiin MATLABilla ja ohjelmakoodit perustuivat viitteen [1] `r1_GMRES`-funktioon. Matriisien GMRES toteutuksena käytettiin `r1_GMRES`ista muokattua versiota, joten molempien toteutustapa oli erittäin samankaltainen. Kokeissa matriiseja $\mathbf{M}_\#$ oli kolmea tyyppiä. Merkitään $\mathbf{R}_n^{(j)} \in \mathbb{R}^{n \times n}$ matriisia, jonka alkiot ovat valittu tasaisesti satunnaisesti väliltä $(0, 1)$. Nämä matriisit muodostettiin MATLABilla käyttäen $\mathbf{R}_n^{(j)} = \text{rand}(n)$ ja ne vaihtuivat satunnaisesti jokaisessa testiajossa.

Ensimmäisessä kokeessa oli $\mathbf{M}_\# = (\mathbf{I}_{500} + \mathbf{R}_{500}^{(1)} + i\mathbf{R}_{500}^{(2)})/10 \in \mathbb{C}^{500 \times 500}$. Tyypillisen ajon jäännösvektoreiden normit $\|\mathbf{r}_k\|/\|\mathbf{r}_0\|$ on piirretty kuvaan 5.1 \log_{10} -skaalassa. Toisessa kokeessa $\mathbf{M}_\# = (5\mathbf{I}_{100} + \mathbf{R}_{100}^{(1)} + i\mathbf{R}_{100}^{(2)}) \in \mathbb{C}^{100 \times 100}$. Jäännösvektoreiden normit on piirretty kuvaan 5.2. Viimeisessä kokeessa ei käytetty satunnaismatriiseja vaan valittiin $\mathbf{M}_\# = \text{diag}(2, 3, \dots, 101) \in \mathbb{C}^{100 \times 100}$, mistä saatiin kuva 5.3. Alkuarvauksena oli $\mathbf{z}_0 = 0$ kaikissa kokeissa.

Tasapuolisen vertailun vuoksi yhtälön (3.54) jäännösvektoreiden normit laskettiin kaavalla $\|\mathbf{r}_k\| = \|\mathbf{b} - \mathcal{M}_1(\mathbf{z}_k)\|$ jokaisella iteraatiokierroksella. Tulokset eivät kuitenkaan juurikaan poikenneet käytännöllisen iteroinnin kaavasta $\|\hat{\mathbf{r}}_k\| = \|(\mathbf{I} - \mathcal{M}_0)(\mathbf{b} - \mathcal{M}_1(\mathbf{z}_k))\|$.

Kuvista nähdään, että RL-GMRES suppenee kaksinkertaista reaalista systeemiä hieman nopeammin lukuunottamatta viimeistä koetta, jossa suppeneminen on yhtä nopeaa. Yhtälö (3.54) vaatii vain suunnilleen puolet edellisten iteraatioista lukuunotta-



Kuva 5.1: Iteraatioiden jäännöksen suhteellinen virhe, kun $\kappa = 1$ ja $M_{\#} = (\text{eye}(500) + \text{rand}(500) + i * \text{rand}(500))/10$ ja vektori \mathbf{b} koostuu luvuista 1. Kuvaan on piirretty sekä 30 että 60 iteraation välein uudelleenkäynnistetyt GMRESit. Käynnistyskohdat on merkitty symbolilla o.

matta viimeistä koetta, jossa se ei suppene. Toisaalta tämän yhtälön ratkaisu vaatii kahden matriisi-vektoritulon laskemista $M_{\#}$:lla.

5.2 GMRES impedanssitomografiassa

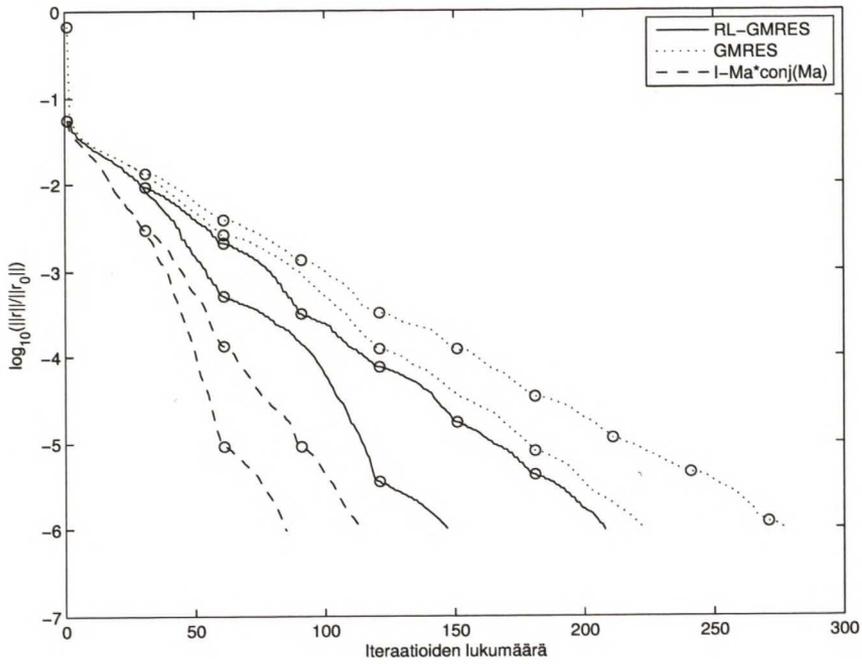
Edellä olevissa GMRES-kokeissa nähtiin, että suppenemisnopeus riippuu hyvin paljon teennäisesti valitusta matriisista $M_{\#}$. Käytännön suorituskyvyn selvittämiseksi kokeita suoritettiin myös impedanssitomografiassa esiintyvillä yhtälöillä ja simuloidulla datalla.

Olkoon Ω tason \mathbb{R}^2 avoin yksikkökierros ja tarkastellaan siellä johtavuusyhtälön reuna-arvotettavaa

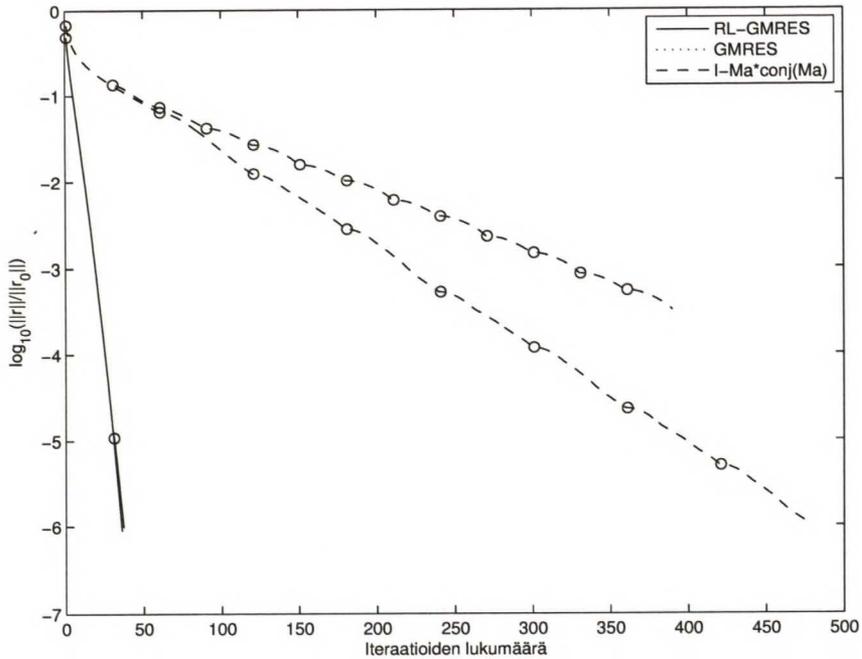
$$\begin{cases} \nabla \cdot \gamma(\mathbf{x}) \nabla u(\mathbf{x}) = 0, & \mathbf{x} \in \Omega, \\ u(\mathbf{x}) = f(\mathbf{x}), & \mathbf{x} \in \partial\Omega, \end{cases}$$

missä sähkönjohtavuus γ on vakioiden rajoittama $0 < c \leq \gamma \leq C < \infty$. Kun (sähkö)potentiaali f kiekon reunalla on annettu ja johtavuus γ tunnetaan, ratkaisu u on potentiaali koko kiekossa. Tällöin voidaan laskea kiekon reunalla virrantiheys $\gamma \frac{\partial u}{\partial \nu} \Big|_{\partial\Omega}$ funktio, missä ν on yksikkönormaalivektori ulospäin. Tunnetun johtavuuden γ perusteella voidaan täten määritellä kuvaus Λ_{γ} reunalla annetuista potentiaaleista f reunan virrantiheyksiin, $\Lambda_{\gamma}(f) = \gamma \frac{\partial u}{\partial \nu} \Big|_{\partial\Omega}$.

Sähköisessä impedanssitomografiassa ratkaistavana on inversio-ongelma, missä ku-



Kuva 5.2: Iteraatioiden jäännöksen suhteellinen virhe, kun $\kappa = 1$ ja $M_{\#} = 5 \cdot \text{eye}(100) + \text{rand}(100) + i \cdot \text{rand}(100)$ ja vektori \mathbf{b} koostuu luvuista 1. Kuvaan on piirretty sekä 30 että 60 iteraation välein uudelleenkäynnistetyt GMRESit. Käynnistyskohdat on merkitty symbolilla \circ .



Kuva 5.3: Iteraatioiden jäännöksen suhteellinen virhe, kun $\kappa = 1$ ja $M_{\#} = \text{diag}(2, 3, \dots, 101)$. Kuvaan on piirretty sekä 30 että 60 iteraation välein uudelleenkäynnistetyt GMRESit. Käynnistyskohdat on merkitty symbolilla \circ . RL-GMRES ja GMRES viivat ovat tässä päällekkäin.

vaus Λ_γ on annettu ja tätä vastaava johtavuus γ on ratkaistava. Varsin yleisin oletuksien voidaan osoittaa, että γ on yksikäsitteinen [7]. Tätä ennen vaativammin oletuksien yksikäsitteisyyden on todistanut mm. Brown ja Uhlmann (mm. γ :lla täytyy olla ensimmäiset osittaisderivaatat). Kuten [5], tässä esityksessä käytetään Brown-Uhlmannin todistuksessa käytettyä menetelmää γ :n rekonstruointiin. Johtavuuden γ uudelleenrakentaminen jakautuu kahteen vaiheeseen. Ensimmäisessä vaiheessa kuvauksen Λ_γ perusteella lasketaan nk. sirontamuunnos ja toisessa vaiheessa sirontamuunnoksesta rakennetaan johtavuus γ . Jälkimmäinen vaihe tapahtuu ratkaisemalla $\bar{\partial}$ -yhtälöitä (5.1). Tässä tarkoituksena on kokeilla luvun 3 menetelmiä, joten kokonaisen inversiotehtävän ratkaisemisen sijaan keskitytään tähän jälkimmäiseen vaiheeseen. Tällöin on sallittua hieman oikaista ja laskea sirontamuunnos suoraan annetusta johtavuudesta γ . Brown-Uhlmannin menetelmässä tämäkin voidaan tehdä ratkaisemalla $\bar{\partial}$ -yhtälöitä.

5.2.1 $\bar{\partial}$ -yhtälön ratkaiseminen

Tarkastellaan seuraavaa nk. $\bar{\partial}$ -yhtälöä

$$\bar{\partial}_z w(z) = -T(z)\overline{w(z)}, \quad \lim_{|z| \rightarrow \infty} w(z) = 1 \quad (5.1)$$

missä kerroinfunktiolla $T : \mathbb{R}^2 \rightarrow \mathbb{C}$ on kompakti kantaja. Tästä on ratkaistava $w : \mathbb{R}^2 \rightarrow \mathbb{C}$ asymptoottisella ehdolla $\lim_{|z| \rightarrow \infty} w(z) = 1$. Kompleksiluku $z = z_1 + iz_2$ samaistetaan tason \mathbb{R}^2 pisteen (z_1, z_2) kanssa ja käytetään derivaatta-operaattoreita

$$\partial_z = \frac{\partial}{\partial z} = \frac{1}{2} \left(\frac{\partial}{\partial z_1} - i \frac{\partial}{\partial z_2} \right), \quad \bar{\partial}_z = \frac{\partial}{\partial \bar{z}} = \frac{1}{2} \left(\frac{\partial}{\partial z_1} + i \frac{\partial}{\partial z_2} \right),$$

missä $z = (z_1, z_2)$. Yhtälö (5.1) ja w :n asymptoottinen ehto ovat yhtäpitäviä integraaliyhtälön

$$w(z) = 1 - \frac{1}{\pi} \int_{\mathbb{R}^2} \frac{T(z') \overline{w(z')}}{z - z'} dz'_1 dz'_2 \quad (5.2)$$

kanssa, missä $g(z) = 1/(\pi z)$ on operaattorin $\bar{\partial}_z$ Greenin funktio. Voidaan kirjoittaa myös lyhyemmin

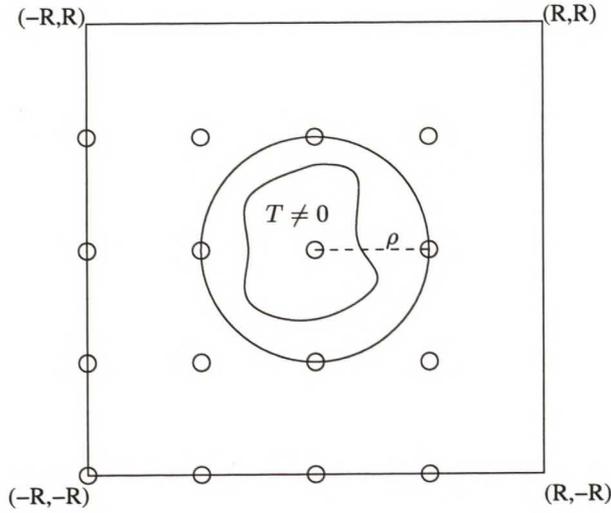
$$w + g * (T\bar{w}) = 1, \quad (5.3)$$

missä $*$ tarkoittaa funktioiden konvoluutiota.

Numeerista laskentaa varten diskretoidaan integraaliyhtälö (5.2) viitteen [5] tavalla. Olkoon $\bar{B}(0, \rho)$ origokeskinen ρ -säteinen suljettu kiekko siten, että T :n kantaja sisältyy siihen. Olkoon $R \geq 2\rho$ ja asetetaan $S_R = [-R, R] \times [-R, R]$ (ks. kuva 5.4). Yhtälön (5.2) ratkaisu w kiekkossa $\bar{B}(0, \rho)$ riippuu Greenin funktion g arvoista vain kiekkossa $\bar{B}(0, 2\rho)$, joten g voidaan haluttaessa määritellä uudelleen tämän kiekon ulkopuolella. Olkoon

$$K(z) = \begin{cases} \frac{1}{\pi z}, & \text{kun } |z| \leq R, \\ 0, & \text{kun } |z| > R. \end{cases} \quad (5.4)$$

Katkaisu voitaisiin tehdä myös sileästi valitsemalla $R = 2\rho + \epsilon$ jollakin pienellä $\epsilon > 0$ ja asettamalla $K(z) = g(z)\chi(z)$, missä χ on sileä, $0 \leq \chi \leq 1$ ja saa arvon 1 kiekkossa $\bar{B}(0, 2\rho)$ ja 0 kiekon $\bar{B}(0, R)$ ulkopuolella. Tässä esityksessä käytetään vain katkaisua (5.4).



Kuva 5.4: Keskellä avoin kiekko $B(0, \rho)$ ja $R = 2\rho$. Pienet ympyrät esittävät 4×4 hilaa ($m = 2$).

Tarkastellaan periodista integraaliyhtälöä

$$w(z) = 1 - \int_{S_R} K(z - z') T(z') \overline{w(z')} dz'_1 dz'_2 \quad (z \in S_R), \quad (5.5)$$

missä K ja T on ensin rajoitettu neliöön S_R ja sitten jatkettu $2R$ -periodisiksi funktioiksi koko tasoon. Tämän ratkaisu w on niin ikään $2R$ -periodinen funktio.

Periodisen yhtälön (5.5) ratkaisulle \tilde{w} ja alkuperäisen yhtälön (5.2) ratkaisulle w pätee $\tilde{w}(z) = w(z)$ kaikilla $z \in B(0, \rho)$. Kun periodisesta yhtälöstä (5.5) on saatu ratkaistua \tilde{w} , voitaisiin se laajentaa kiekosta $B(0, \rho)$ yhtälön (5.2) ratkaisuksi koko \mathbb{R}^2 :een laskemalla

$$w(z) = 1 - \frac{1}{\pi} \int_{B(0, \rho)} \frac{T(z')}{z - z'} \overline{\tilde{w}(z')} dz'_1 dz'_2 \quad (z \in \mathbb{R}^2).$$

Tässä työssä ratkaisun kiinnostava osuus kuitenkin sijaitsee kiekossa $B(0, \rho)$ eikä tätä laajennusta tarvitse laskea.

Ryhdytään sitten ratkaisemaan yhtälöä (5.5). Olkoon m positiivinen kokonaisluku ja asetetaan $N = 2^m$, $h = 2R/N$. Muodostetaan hila pisteistä jh , missä $j \in \mathbb{Z}_N^2$

$$\mathbb{Z}_N^2 = \left\{ (j_1, j_2) \in \mathbb{Z}^2 \mid -\frac{N}{2} \leq j_1, j_2 < \frac{N}{2} \right\}. \quad (5.6)$$

Tämä on $N \times N$ -hila, jonka tapaus $m = 2$ on piirretty kuvaan 5.4. Tällainen N :n valinta sopii erityisen hyvin jäljempänä esiteltävän nopean Fourier-muunnoksen käyttöön.

Kun $j \in \mathbb{Z}^2$, merkitään

$$T_j = T(jh), \quad (5.7)$$

$$w_j = w(jh), \quad (5.8)$$

$$K_j = \begin{cases} 0, & \text{kun } j_1 \equiv j_2 \equiv 0 \pmod{N}, \\ K(jh), & \text{muutoin.} \end{cases} \quad (5.9)$$

Funktion K diskretoinnissa asetettiin nolla singulariteettipisteissä. Singulariteettipisteet ovat integroituvia, joten tästä aiheutuva diskreetointivirhe pienenee hilaa tihentäessä. Koska funktiot T , w ja K ovat $2R$ -periodisia, ovat diskretoidut kaksiindeksiset jonot T_j , w_j ja K_j N -periodisia molempien indeksien suhteen.

Diskretoidaan seuraavaksi yhtälössä (5.5) esiintyvä integraali muotoon

$$h^2 \sum_{k \in \mathbb{Z}_N^2} K_{j-k} T_k \overline{w_k}. \quad (5.10)$$

Tässä $K_{j-k} = K((j-k)h) = K(jh - kh)$, kun $j \neq k$. Summa (5.10) on periodisten jonojen konvoluutio, joka voidaan merkitä lyhyemmin $h^2(K_j) * (T_j \overline{w_j})$. Yhtälö (5.5) on diskretoitu nyt muotoon

$$w_j + h^2(K_j) * (T_j \overline{w_j}) = 1 \quad (j \in \mathbb{Z}_N^2).$$

Konvoluutiota ei kannata laskea suoraan kaavasta (5.10) vaan käyttäen nopeaa Fourier-muunnosta (FFT:stä voi lukea esim. kirjasta [10]). Merkitään \mathcal{F} :llä kaksiulotteista diskreettiä Fourier-muunnosta, jolloin edellä olevasta yhtälöstä tulee

$$w_j + h^2 \mathcal{F}^{-1}(\mathcal{F}(K_j) \cdot \mathcal{F}(T_j \overline{w_j})) = 1, \quad (j \in \mathbb{Z}_N^2), \quad (5.11)$$

missä \cdot tarkoittaa Fourier-kertoimien kertomista keskenään alkioittain.

Matriisiyhtälöksi tämä voidaan kirjoittaa pakkaamalla ratkaistavat w_j vektoriksi $\mathbf{u} \in \mathbb{C}^{N^2}$ asettamalla $u_{(j_1+N/2)+N(j_2+N/2)+1} = w_{j_1, j_2}$, missä $j = (j_1, j_2) \in \mathbb{Z}_N^2$. Tämän pakkauksen myötä summa (5.10) voidaan esittää matriisien avulla. Yhtälöksi saadaan $\mathbf{u} + \mathbf{M}_\# \overline{\mathbf{u}} = \mathbf{1}$, missä $\mathbf{1} \in \mathbb{C}^{N^2}$ koostuu luvuista 1 ja $\mathbf{M}_\# = \mathbf{C}\mathbf{D}$, missä $\mathbf{C}, \mathbf{D} \in \mathbb{C}^{N^2 \times N^2}$. Matriisi \mathbf{D} on diagonaalimatriisi koostuen alkioista T_j ja alkioista K_j koostuvalla matriisilla \mathbf{C} kertominen suorittaa konvoluution. Tämä yhtälö ratkeaa operaattorille \mathcal{M}_κ luvussa 3 esitellyin menettelyin.

Trigonometrinen kollokaatio

Kerrotaan yhtälö (5.5) puolittain $\overline{T(z)}$:lla, jolloin

$$\overline{T(z)} w(z) = \overline{T(z)} - \overline{T(z)} \int_{S_R} K(z-z') T(z') \overline{w(z')} dz'_1 dz'_2 \quad (z \in S_R).$$

Merkitsemällä $v(z) = \overline{T(z)} w(z)$ saadaan periodinen integraaliyhtälö

$$v(z) = \overline{T(z)} - \overline{T(z)} \int_{S_R} K(z-z') \overline{v(z')} dz'_1 dz'_2. \quad (5.12)$$

Kun tämä ratkaistaan ja lasketaan funktio w kaavasta

$$w(z) = 1 - \int_{S_R} K(z-z') \overline{v(z')} dz'_1 dz'_2,$$

niin w on yhtälön (5.5) ratkaisu. Diskretoidaan (5.12) käyttäen trigonometristä kollokaatiota viitteen [6] tapaan.

Määritellään ortonormaali kanta avaruuteen $L^2(S_R)$ funktioilla

$$\varphi_j(z) = \frac{1}{2R} \exp\left(i\pi j \cdot \frac{z}{R}\right), \quad j \in \mathbb{Z}^2, \quad (5.13)$$

ja koostukoon avaruus \mathcal{T}_N trigonometrisistä polynomeista $v_N = \sum_{j \in \mathbb{Z}_N^2} c_j \varphi_j$, missä $c_j \in \mathbb{C}$. Kun v on $2R$ -periodinen funktio, määritellään sen trigonometrinen interpolatio $Q_N v$ vaatimalla

$$Q_N v \in \mathcal{T}_N, \quad (Q_N v)(jh) = v(jh), \quad j \in \mathbb{Z}_N^2.$$

Yhtälö (5.12) diskretoidaan nyt muotoon

$$v_N = Q_N(\bar{T}) - Q_N(\bar{T}K\bar{v}_N), \quad (5.14)$$

missä $(Kv)(z) = \int_{S_R} K(z-z')v(z') dz'_1 dz'_2$ ja on ratkaistava $v_N \in \mathcal{T}_N$. Nähdään, että

$$K\varphi_j = \widehat{K}(j)\varphi_j \quad (j \in \mathbb{Z}^2),$$

missä $\widehat{K}(j) = \int_{S_R} K(z)\varphi_{-j}(z) dz_1 dz_2$ ovat funktion K Fourier-kertoimet.

Trigonometrisellä polynomilla $v_N \in \mathcal{T}_N$ on seuraavat kaksi mahdollista esitystapaa

$$\begin{aligned} v_N(z) &= \sum_{k \in \mathbb{Z}_N^2} \widehat{v}_N(k)\varphi_k(z), \\ v_N(z) &= \sum_{j \in \mathbb{Z}_N^2} v_N(jh)\varphi_{N,j}(z), \end{aligned}$$

missä $\varphi_{N,j} \in \mathcal{T}_N$ toteuttaen $\varphi_{N,j}(kh) = \delta_{j,k}$ (Kroneckerin δ) ja $j, k \in \mathbb{Z}_N^2$. Ensimmäinen näistä on Fourier-kertoimien avulla esitetty ja jälkimmäinen hilapisteissä olevien arvojen avulla. Funktiot $\varphi_{N,j}$ saadaan kaavasta

$$\varphi_{N,j}(z) = \frac{1}{N^2} \sum_{k \in \mathbb{Z}_N^2} \exp\left(i\pi k \cdot \left(\frac{x}{R} - \frac{2j}{N}\right)\right).$$

Esityksestä toiseen voidaan siirtyä diskreetin Fourier-muunnoksen avulla. Laskemalla saadaan

$$\widehat{v}_N(k) = \int_{S_R} v_N(z)\varphi_{-j}(z) dz_1 dz_2 = \frac{2R}{N^2} \sum_{j \in \mathbb{Z}_N^2} v_N(jh) \exp\left(-i\pi k \cdot \frac{2j}{N}\right). \quad (5.15)$$

Arvot hilapisteissä lasketaan diskreetillä käänteismuunnoksella

$$v_N(jh) = \frac{1}{2R} \sum_{k \in \mathbb{Z}_N^2} \widehat{v}_N(k) \exp\left(i\pi k \cdot \frac{2j}{N}\right). \quad (5.16)$$

Muodostetaan sitten yhtälön (5.14) matriisimuoto. Kun ratkaisun $v_N \in \mathcal{T}_N$ arvot $v_N(jh)$, $j \in \mathbb{Z}_N^2$, tunnetaan, niin $v_N = \sum_{k \in \mathbb{Z}_N^2} \widehat{v}_N(k)\varphi_k$, missä kertoimet $\widehat{v}_N(k)$ saadaan kaavasta (5.15). Joukon (5.6) lisäksi määritellään

$$\widetilde{\mathbb{Z}}_N^2 = \left\{ (j_1, j_2) \in \mathbb{Z}^2 \mid -\frac{N}{2} < j_1, j_2 \leq \frac{N}{2} \right\}.$$

Sijoitetaan ja lasketaan

$$\begin{aligned}\overline{TKv_N} &= \overline{T} \sum_{k \in \mathbb{Z}_N^2} \widehat{v}_N(k) \widehat{K}(-k) \varphi_{-k} = \overline{T} \sum_{k \in \mathbb{Z}_N^2} \widehat{v}_N(-k) \widehat{K}(-k) \varphi_{-k} \\ &= \overline{T} \sum_{k \in \mathbb{Z}_N^2} \widehat{K}(k) \widehat{v}_N(k) \varphi_k,\end{aligned}$$

joten matriisimuodoksi saadaan

$$v_j + \overline{T}_j \mathcal{F}_h^{-1}(\widehat{K}(j) \cdot \mathcal{F}_h(\overline{v}_j)) = \overline{T}_j, \quad (j \in \mathbb{Z}_N^2), \quad (5.17)$$

missä $v_j = v_N(jh)$ ja $\mathcal{F}_h, \mathcal{F}_h^{-1}$ ovat diskreetit Fourier-muunnokset (5.15) ja (5.16). Kun v_j pakataan vektoriksi \mathbf{y} edellisen kohdan lopussa olevaa tapaa käyttäen, niin tämä yhtälö on matriisimuotoa

$$\mathbf{y} + \overline{DC}\overline{\mathbf{y}} = \overline{D}\mathbf{1}. \quad (5.18)$$

Kun tästä ratkaistaan \mathbf{y} ja lasketaan $\mathbf{u} = \mathbf{1} - \overline{C}\overline{\mathbf{y}}$, niin nähdään, että \mathbf{u} toteuttaa yhtälön

$$\mathbf{u} + \overline{CD}\overline{\mathbf{u}} = \mathbf{1}. \quad (5.19)$$

Toisaalta kertomalla yhtälö (5.19) puolittain vasemmalta \overline{D} :llä ja ottamalla uudeksi muuttujaksi $\mathbf{y} = \overline{D}\mathbf{u}$ päädytään yhtälöön (5.18). Nämä ovat siten yhtäpitävät. Yhtälöön (5.19) päädyttiin myös edellisessä kohdassa. Erona on, että yhtälön (5.11) Fourier-kertoimet $\mathcal{F}(K_j)$ on laskettu diskreetillä Fourier-muunnoksella diskreetoinista (5.9) ja tämän paikalla on nyt tarkasti laskettavissa olevat kertoimet $\widehat{K}(j)$. Lasketaan nämä seuraavaksi. Huomautettakoon, ettei yhtälössä (5.17) ole tekijää h^2 yhtälön (5.11) tapaan. Tämä tekijä sisältyy nyt kertoimiin $\widehat{K}(j)$.

Katkaistun Greenin funktion Fourier-kertoimet

Olkoon $j = (j_1, j_2) \neq (0, 0)$. Derivoimalla nähdään, että $\overline{\partial}_z \varphi_j(z) = \lambda_j^{-1} \varphi_j(z)$, missä $\lambda_j = \left(\frac{\pi}{2R}(-j_2 + ij_1)\right)^{-1}$, joten

$$\begin{aligned}\widehat{K}(j) &= \int_{S_R} K(z) \varphi_{-j}(z) dz_1 dz_2 \\ &= \int_{B(0,R)} \frac{1}{\pi z} \varphi_{-j}(z) dz_1 dz_2 = -\lambda_j \int_{B(0,R)} \frac{1}{\pi z} \overline{\partial}_z \varphi_{-j}(z) dz_1 dz_2 \\ &= -\lambda_j \lim_{\delta \rightarrow 0} \int_{B(0,R) \setminus B(0,\delta)} \frac{1}{\pi z} \overline{\partial}_z \varphi_{-j}(z) dz_1 dz_2 \\ &= -\lambda_j \lim_{\delta \rightarrow 0} \left\{ \left(\int_{S(0,R)} - \int_{S(0,\delta)} \right) \left(\frac{1}{\pi z} \varphi_{-j}(z) \frac{z}{2|z|} \right) ds \right. \\ &\quad \left. - \int_{B(0,R) \setminus B(0,\delta)} \left(\overline{\partial}_z \frac{1}{\pi z} \right) \varphi_{-j}(z) dz_1 dz_2 \right\},\end{aligned}$$

missä on käytetty Gauss-Green lausetta. Jälkimmäisin integraali on nolla, koska $1/(\pi z)$ on operaattorin $\overline{\partial}$ Greenin funktio. Saadaan

$$\widehat{K}(j) = \lambda_j \left(\varphi_{-j}(0) - \frac{1}{2\pi R} \int_{S(0,R)} \varphi_{-j}(z) ds \right).$$

Lasketaan

$$\begin{aligned} \int_{S(0,R)} e^{i\pi j \cdot z/R} ds &= R \int_{S(0,1)} e^{i\pi j \cdot z} ds = R \int_{S(0,1)} e^{i\pi |j| z_1} ds \\ &= 2R \int_{-1}^1 e^{i\pi |j| z_1} (1 - z_1^2)^{-1/2} dz_1 \\ &= 2R \int_{-1}^1 \cos(\pi |j| z_1) (1 - z_1^2)^{-1/2} dz_1 = 2\pi R J_0(\pi |j|), \end{aligned}$$

missä Besselin funktiolle J_0 käytetty kaava löytyy esim. kirjasta [16]. Täten

$$\widehat{K}(j) = \frac{\lambda_j}{2R} (1 - J_0(\pi |j|)).$$

Kun $j = (0, 0)$, nähdään helposti $\widehat{K}(0) = \int_{B(0,R)} \frac{1}{\pi z} dz_1 dz_2 = 0$. Fourier-kertoimiksi on näin saatu

$$\widehat{K}(j) = \begin{cases} \frac{1 - J_0(\pi |j|)}{\pi(-j_2 + ij_1)}, & \text{kun } j \neq (0, 0), \\ 0, & \text{kun } j = (0, 0). \end{cases}$$

5.2.2 Brown-Uhlmann menetelmä

Seuraavassa esitetään vain ratkaistavat yhtälöt ja laskentakaavat, katso [5] ja sen viitteet tarkempaan matemaattiseen kuvaukseen. Olkoon sitten johtavuusfunktio γ annettu kiekossa Ω ja oletetaan, että $\gamma = 1$ lähellä kiekon reunaa. Lasketaan (eräänlainen potentiaali)

$$q = -\gamma^{-1/2} \partial_z \gamma^{1/2}. \quad (5.20)$$

Koska oletetaan, että $\gamma = 1$ lähellä kiekon reunaa $\partial\Omega$, voidaan q laajentaa sileästi koko tasoon asettamalla $q(z) = 0$ kaikilla $z \in \mathbb{C} \setminus \Omega$. Ratkaistaan kaikki seuraavat $\bar{\partial}$ -yhtälöt parametrina $k \in \mathbb{C}$

$$\bar{\partial}_z m_{\pm}(z, k) = \pm q(z) e(z, -k) \overline{m_{\pm}(z, k)}, \quad \lim_{|z| \rightarrow \infty} m_{\pm}(z, k) = 1, \quad (5.21)$$

missä $e(z, k) = \exp(i(zk + \bar{z}k)) = \exp(2i \operatorname{Re}(zk))$. Lasketaan $m_1(z, k) = \frac{1}{2}(m_+(z, k) + m_-(z, k))$ ja potentiaalin q sirontamuunnos kaavasta

$$S(k) = \frac{-i}{\pi} \int_{\Omega} e(z, k) \overline{q(z)} m_1(z, k) dz_1 dz_2. \quad (5.22)$$

Johtavuus γ uudelleenrakennetaan sirontamuunnoksesta S ratkaisemalla $\bar{\partial}$ -yhtälöt parametrina $z \in \Omega$

$$\bar{\partial}_k \tilde{m}^+(z, k) = \overline{S(-k)} e(z, -k) \overline{\tilde{m}^+(z, k)}, \quad \lim_{|k| \rightarrow \infty} \tilde{m}^+(z, k) = 1. \quad (5.23)$$

Johtavuus saadaan sitten laskemalla

$$\gamma(z) = (\operatorname{Re}(\tilde{m}^+(z, 0)))^2. \quad (5.24)$$

5.2.3 Kokeet

Kokeissa toistettiin viitteen [5, kohta 5.1] järjestely. Ensiksi valmisteltiin johtavuusfunktio γ yksikkökiekkoon. Taustajohtavuus on 1, sydämen johtavuus on 2, keuhkojen johtavuus on 0.33 ja toisessa keuhkossa olevan poikkeaman johtavuus 1.33. Kuvassa 5.5 on tämä alkuperäinen johtavuus jo uudelleenrakennettuna. Potentiaali q laskettiin käyttäen kaavaa (5.20) ja numeerista derivointia.

Yhtälöt (5.21) ratkaistiin käyttäen kohdan 5.2.1 menetelmää. Neliöön $[-2, 2] \times [-2, 2]$ muodostettiin 128×128 -hila muuttujan z arvoille ja neliöön $[-40, 40] \times [-40, 40]$ niin ikään 128×128 -hila muuttujan k arvoille. Sirontamuunnos laskettiin kaavasta (5.22) numeerisella integroinnilla. Koska sirontamuunnos S ei poikkea nolasta vain rajoitetussa joukossa, se täytyy katkaista seuraavaa vaihetta varten. Tästä aiheutuu systemaattinen virhe, mutta tämä jätetään huomiotta, koska S on $|k|$:n kasvaessa nopeasti vaimeneva. Edellä neliö $[-40, 40] \times [-40, 40]$ on jo valittu tätä katkaisua varten eikä $S(k)$ arvoja laskettu, kun $|k| > 20$, vaan asetettiin nolliksi.

Laskettu sirontamuunnos tallennettiin ja seuraavan vaiheen johtavuuden rekonstruointiin käytetyt iteraatiot ja laskenta-ajat mitattiin eri menetelmin. Jokaisessa kokeessa ratkaistiin $\bar{\partial}$ -yhtälöt (5.23) kaikille z :n arvoille neliössä $[-1, 1] \times [-1, 1]$ 32×32 -hilassa lukuunottamatta hilapisteitä, joille $|z| > 1$. Jokaisessa hilapisteessä ratkaistiin reaaliineaarinen yhtälö $\mathbf{u} + \mathbf{M}_{\#} \bar{\mathbf{u}} = \mathbf{1}$ iteratiivisin menetelmin ja iterointi pysäytettiin, kun jäännös $\|\mathbf{r}_k\|/\|\mathbf{r}_0\| = \|\mathbf{1} - \mathbf{u}_k - \mathbf{M}_{\#} \bar{\mathbf{u}}_k\|/\|\mathbf{r}_0\| < 10^{-12}$.

Jokaisen kokeen tuloksena oli kuva 5.5 ja mittaukset eri menetelmin on kirjattu taulukkoon 5.1. Laskenta-ajat mitattiin MATLABin TIC-TOC komennoilla kahdessa eri ajossa. Taulukossa $\mathbf{I} - \mathbf{M}_{\#} \overline{\mathbf{M}_{\#}}$ tarkoittaa yhtälön (3.54) ratkaisemista matriisien GMRESillä ja $\mathbf{I} - \mathbf{M}_{\#} \overline{\mathbf{M}_{\#} \mathbf{M}_{\#} \overline{\mathbf{M}_{\#}}}$ yhtälön (3.53) ratkaisua arvolla $k = 4$. Näille kahdelle suoritettiin kokeet myös käyttäen epäkäytännöllistä jäännöksen normia $\|\mathbf{r}_k\| = \|\mathbf{b} - \mathcal{M}_1(\mathbf{u}_k)\|$. Nämä on merkitty J:llä. Tämä tehtiin, jotta iterointi tulisi pysäytettyä samassa tarkkuudessa muiden menetelmien kanssa. Kuten taulukosta nähdään, ei tällä ole iteraatioiden määrään juurikaan vaikutusta. Laskenta-aika J:llä merkityissä kohdissa johtuu kohtuuttomasta ylimääräisestä laskennasta. Kuvaan 5.6 on piirretty yhden $\bar{\partial}$ -yhtälön ratkaisun suhteelliset jäännökset.

Vertailun vuoksi laskenta suoritettiin myös seuraavalla yksinkertaisella kiintopisteiteraatiolla

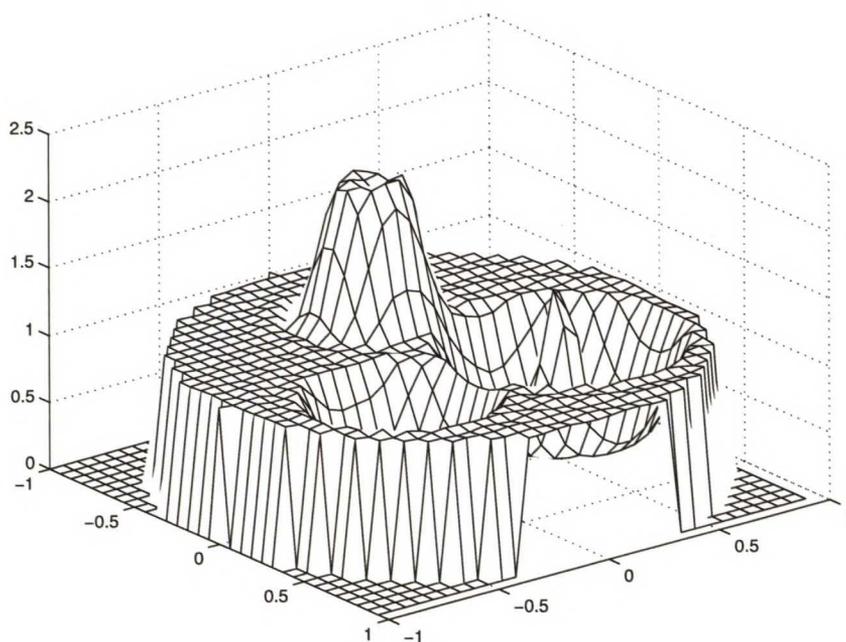
$$\mathbf{u}_{k+1} = -\mathbf{M}_{\#} \bar{\mathbf{u}}_k + \mathbf{1}.$$

Tämän iterointi pysäytettiin, kun $\|\mathbf{1} - \mathbf{u}_{k+1} - \mathbf{M}_{\#} \bar{\mathbf{u}}_k\|/\|\mathbf{r}_0\| < 10^{-12}$. Kuten taulukosta nähdään, iteraatioiden lukumäärä on suurempi kuin GMRES-menetelmissä, mutta laskenta-aika on tästä huolimatta hyvä.

5.2.4 $\bar{\partial}$ -yhtälö satunnaisella kerroinfunktiolla

Edellisen kohdan kiintopisteiteraation ja GMRES-menetelmän menestys matriisille $\mathbf{I} - \mathbf{M}_{\#} \overline{\mathbf{M}_{\#} \mathbf{M}_{\#} \overline{\mathbf{M}_{\#}}}$ antaa aiheutta epäillä, että tämä johtuu matriisin $\mathbf{M}_{\#}$ normin pienyydestä.

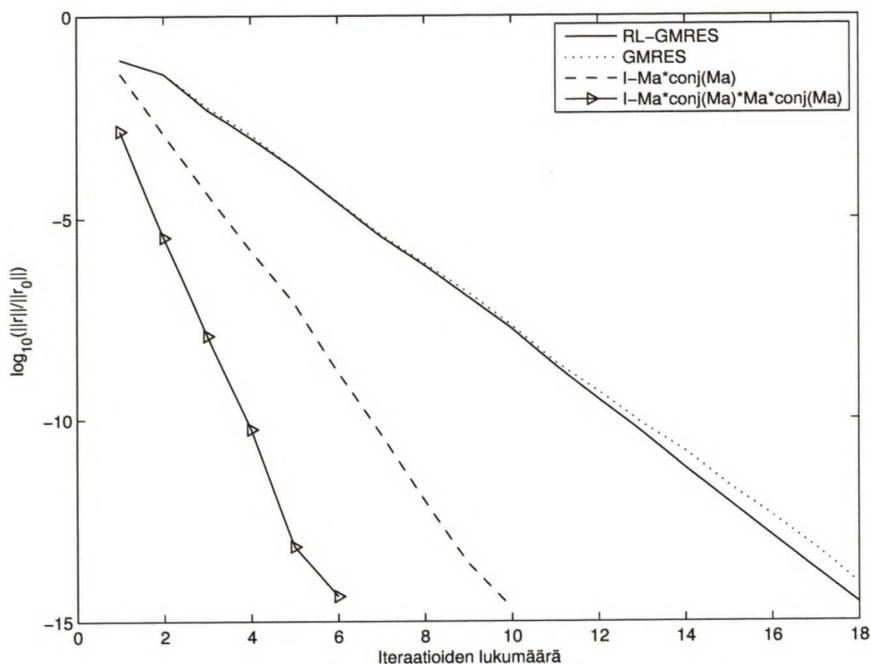
Seuraavissa kokeissa yhtälö (5.1) ratkaistiin neliössä $[-40, 40] \times [-40, 40]$ 256×256 -hilassa ja kerroinfunktion T arvot hilapisteissä $|z| < 20$ asetettiin tasaisesti satunnai-



Kuva 5.5: Johtavuus yksikköympyrässä.

Menetelmä	Iteraatiot	Aika (s)
RL-GMRES	12192	105.1
	12192	104.9
GMRES	12500	120.9
	12500	119.0
$I - M_{\#} \overline{M_{\#}}$	6742	69.0
	6742	68.3
$I - M_{\#} \overline{M_{\#}} J$	6744	97.7
	6744	98.5
$I - M_{\#} \overline{M_{\#} M_{\#} M_{\#}}$	3756	61.3
	3756	61.1
$I - M_{\#} \overline{M_{\#} M_{\#} M_{\#}} J$	3764	76.1
	3764	75.3
Kiintopiste	15991	55.1
	15991	55.2

Taulukko 5.1: Sirontamuunnoksesta lasketun rekonstruoinnin aika.



Kuva 5.6: Iteraatioiden jäännöksen suhteellinen virhe, kun $\bar{\delta}$ -yhtälö on ratkaistu hilan keskipisteessä ($z_1 = z_2 = 0$).

Menetelmä	0.01	0.1	0.5	1.0	5.0	10.0	20.0
RL-GMRES	5	10	25	39	160	313	625
GMRES	5	10	25	41	179	355	713
$I - M_{\#} \overline{M_{\#}} J$	3	5	13	22	80	155	352
$I - M_{\#} \overline{M_{\#}} M_{\#} \overline{M_{\#}} J$	2	4	36	> 960			
Kiintopiste	7	75	∞	∞			

Taulukko 5.2: $\bar{\delta}$ -yhtälö ratkaistuna T :llä, jonka arvot ovat tasaisesti satunnaisia väliltä 0 ylimmällä rivillä olevaan arvoon. Arvot ovat iteraatioiden lukumääriä ja merkinnällä ∞ menetelmä ei supennut.

sesti 0:n ja annetun ylärajan välillä MATLABin `rand`-funktiolla. Muissa hilapisteissä T asetettiin nolllaksi. Kaikki GMRESit uudelleenkäynnistettiin 60 iteraation välein ja pysäytettiin, kun jäännöksen normi oli $< 10^{-6}$. Tulokset ovat taulukossa 5.2, mistä nähdään kiintopisteiteraation ja matriisin $I - M_{\#} \overline{M_{\#}} M_{\#} \overline{M_{\#}}$ GMRESin suppenemisen hidastuvan rajusti T :n kasvaessa.

On huomattava, että $I - M_{\#} \overline{M_{\#}}$ säilytti nopeussuhteensa RL-GMRESiin ja kaksinkertaiseen reaaliseen GMRESiin nähden.

5.3 Aaltoyhtälö

Tarkastellaan seuraavaa kaksiulotteisen aaltoyhtälön alkuarvo- ja reuna-arvotettavaa alueessa $\Omega = (0, 1) \times (0, 1) \subset \mathbb{R}^2$

$$\begin{cases} u_{tt}(\mathbf{x}, t) = \Delta u(\mathbf{x}, t), & \mathbf{x} \in \Omega, t > 0, \\ u(\mathbf{x}, t) = 0, & \mathbf{x} \in \partial\Omega, t \geq 0, \\ u(\mathbf{x}, 0) = f(\mathbf{x}), u_t(\mathbf{x}, 0) = g(\mathbf{x}), & \mathbf{x} \in \Omega. \end{cases}$$

Ratkaistaan tämä numeerisesti käyttämällä aluksi differenssiapproksimaatioita paikamuuttujan suhteen. Olkoon $n_1, n_2 \in \mathbb{Z}^+$ ja asetetaan $h_1 = \frac{1}{n_1}$, $h_2 = \frac{1}{n_2}$. Merkitään $u_{(i,j)}(t) = u(ih_1, jh_2, t)$, $0 \leq i \leq n_1 + 1$, $0 \leq j \leq n_2 + 1$. Reunaehdoista johtuen $u_{(0,j)}(t) = u_{(n_1+1,j)}(t) = u_{(i,0)}(t) = u_{(i,n_2+1)}(t) = 0$. Määritellään $n \times n$ -matriisi

$$\Delta_h = \frac{1}{h^2} \begin{bmatrix} -2 & 1 & & & & & \\ & 1 & -2 & 1 & & & \\ & & 1 & -2 & 1 & & \\ & & & \ddots & \ddots & \ddots & \\ & & & & & 1 & -2 & 1 \\ & & & & & & 1 & -2 \end{bmatrix},$$

missä $h = \frac{1}{n}$. Tämä on itse asiassa yksidimensioisen Laplacen operaattorin (toinen derivaatta) differenssiapproksimaatio, kun molemmissa päissä on Dirichlet'n reunaehto 0.

Kaksidimensioista tapausta varten järjestetään tuntemattomat $u_{(i,j)}(t)$, $1 \leq i \leq n_1$, $1 \leq j \leq n_2$, pystyvektoriksi. Merkitään $v_{i+(j-1)n_1}(t) = u_{(i,j)}(t)$. Kroneckerin tulon avulla muodostettua matriisiä

$$\Delta_{h_1, h_2} = \mathbf{I}_{n_2} \otimes \Delta_{h_1} + \Delta_{h_2} \otimes \mathbf{I}_{n_1}$$

käyttämällä tulee tehtäväksi ratkaista differentiaaliyhtälösystemi

$$\mathbf{v}''(t) = \Delta_{h_1, h_2} \mathbf{v}(t), \quad \mathbf{v}(0) = \mathbf{f}_{h_1, h_2}, \quad \mathbf{v}'(0) = \mathbf{g}_{h_1, h_2},$$

missä $(\mathbf{f}_{h_1, h_2})_{i+(j-1)n_1} = f(ih_1, jh_2)$ ja $(\mathbf{g}_{h_1, h_2})_{i+(j-1)n_1} = g(ih_1, jh_2)$. Kirjoitetaan tämä ensimmäisen kertaluvun systeemiksi. Merkitään $\mathbf{w}(t) = \mathbf{v}'(t)$, jolloin systeemiksi saadaan

$$\mathbf{y}'(t) = \mathbf{A}\mathbf{y}(t), \quad \mathbf{y}(0) = \mathbf{y}_0,$$

missä

$$\mathbf{A} = \begin{bmatrix} 0 & \mathbf{I}_{n_1 n_2} \\ \Delta_{h_1, h_2} & 0 \end{bmatrix} \quad \text{ja} \quad \mathbf{y}(t) = \begin{bmatrix} \mathbf{v}(t) \\ \mathbf{w}(t) \end{bmatrix}, \quad \mathbf{y}_0 = \begin{bmatrix} \mathbf{f}_{h_1, h_2} \\ \mathbf{g}_{h_1, h_2} \end{bmatrix}.$$

Suoritetaan seuraavaksi aikadiskretointi. Olkoon diskretointiväli $\delta > 0$ ja merkitään $\mathbf{y}_j = \mathbf{y}(j\delta)$. Tällöin implisiittinen keskipistesääntö antaa yhtälöt

$$\mathbf{y}_{j+1} = \mathbf{y}_j + \delta \mathbf{A} \left(\frac{1}{2} (\mathbf{y}_j + \mathbf{y}_{j+1}) \right), \quad j = 0, 1, 2, \dots,$$

jotka voidaan kirjoittaa muotoon

$$(\mathbf{I} - \frac{\delta}{2} \mathbf{A}) \mathbf{y}_{j+1} = (\mathbf{I} + \frac{\delta}{2} \mathbf{A}) \mathbf{y}_j, \quad j = 0, 1, 2, \dots \quad (5.25)$$

Tässä \mathbf{y}_{j+1} saadaan siis edellisen \mathbf{y}_j avulla ratkaisemalla lineaarinen yhtälöryhmä.

Kokeessa laskettiin \mathbf{y}_1 GMRESillä käyttäen sekä tavallista matriisien ILU(0)-pohjustusta että \mathbb{R} -lineaarista ILU(0)-pohjustusta. Lisäksi laskettiin kohdan 3.10 mukaisilla reaali-lineaarilla menetelmillä, jolloin edellä olevan yhtälön matriisi muunnettiin \mathbb{R} -operaattoriksi käyttäen tapaa (1.3). Diskretointiparametrit olivat $n_1 = n_2 = 32$ ja $\delta = 0.1$. Alkuehtona oli

$$f(\mathbf{x}) = \begin{cases} \frac{1}{2}(1 + \cos(\pi\|\mathbf{x} - [0.5 \ 0.5]^T\|/0.1)), & \text{kun } \|\mathbf{x} - [0.5 \ 0.5]^T\| < 0.1, \\ 0, & \text{muutoin,} \end{cases}$$

$$g(\mathbf{x}) = \begin{cases} 2 - 2\|\mathbf{x} - [1 \ 1]^T\|/0.3, & \text{kun } \|\mathbf{x} - [1 \ 1]^T\| < 0.3, \\ 0, & \text{muutoin.} \end{cases}$$

Matriisien ILU(0)-pohjustus laskettiin käyttäen liitteen B funktiota `ilu_v2`. \mathbb{R} -lineaarisen pohjustuksen tapauksessa kerroinmatriisista $(\mathbf{I} - \frac{\delta}{2}\mathbf{A})$ tehtiin ensin \mathbb{R} -lineaarinen operaattori $\mathcal{M} : \mathbb{C}^{n_1 n_2} \rightarrow \mathbb{C}^{n_1 n_2}$ kohdan (1.3) mukaan. Tämän jälkeen laskettiin \mathcal{M} :lle \mathbb{R} -lineaarinen ILU(0)-hajotelma käyttäen funktiota `r1_ilu` ja tätä hajotelmaa käytettiin pohjustimena MATLABin `gmres`-funktiolle. Kuvassa 5.7 näkyy ratkaisun jäännös piirrettynä iteraatioiden lukumäärää vasten. Nähdään, että \mathbb{R} -ILU(0) parantaa selkeästi suppenemisnopeutta verrattuna sekä pohjustamattomaan iterointiin että matriisien ILU(0):aan nähden. Algoritmi 3.10.1 ei kuitenkaan suppene \mathbb{R} -ILU(0):n kanssa. Nähdään myös, että tässä tehtävässä kohdan 3.10 reaali-lineaariset iteratiiviset menetelmät eivät nopeuta suppenemistä.

Diskretoinnista $n_1 = n_2 = 4$, saadaan kuva 5.8. Algoritmin 3.10.1 tapauksessa iteraatio ei suppene. Sen rivillä 3 laskettava vektori \mathbf{v} on (numeerisesti lähes) nolla ensimmäisten iteraatioiden jälkeen, joten tässä tapauksessa algoritmi ei toimi.

5.4 Epälineaarinen Schrödingerin yhtälö

Tarkastellaan seuraavaa epälineaarista Schrödingerin yhtälöä

$$i\psi_t(\mathbf{x}, t) = \Delta\psi(\mathbf{x}, t) + |\psi(\mathbf{x}, t)|^2\psi(\mathbf{x}, t).$$

Rajoitutaan tässä myös neliöön $[0, 1] \times [0, 1] \subset \mathbb{R}^2$. Merkitään

$$\psi_{(i,j)}(t) = \psi(ih_1, jh_2), \quad 0 \leq i \leq n_1 + 1, \quad 0 \leq j \leq n_2 + 1.$$

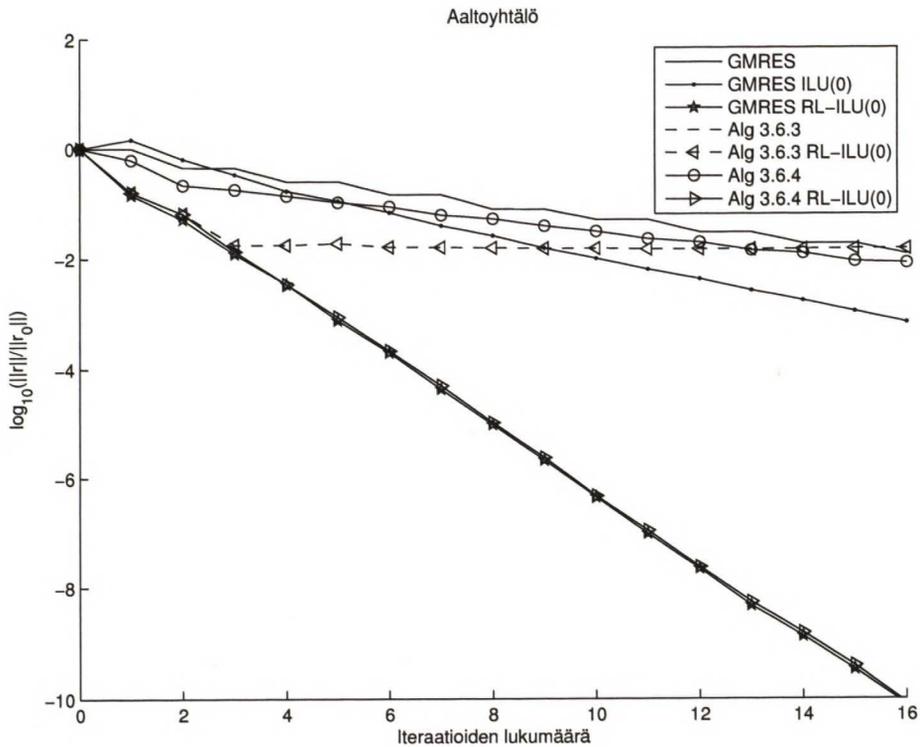
Olkoon tässäkin $\psi_{(0,j)}(t) = \psi_{(n_1+1,j)}(t) = \psi_{(i,0)}(t) = \psi_{(i,n_2+1)}(t) = 0$. Merkitään edelleen $w_{i+(j-1)n_1}(t) = \psi_{(i,j)}(t)$, $1 \leq i \leq n_1$, $1 \leq j \leq n_2$. Diskretoitu yhtälö on nyt

$$i\mathbf{w}'(t) = \mathbf{\Delta}_{h_1, h_2}\mathbf{w}(t) + \mathbf{G}(\mathbf{w}(t)), \quad (5.26)$$

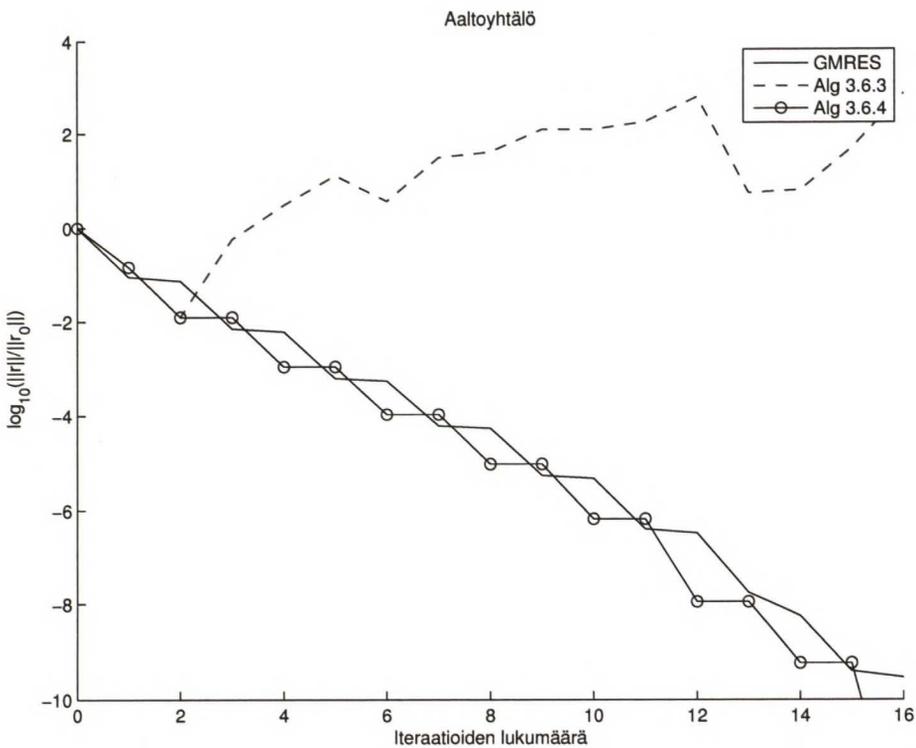
missä $\mathbf{G} : \mathbb{C}^{n_1 n_2} \rightarrow \mathbb{C}^{n_1 n_2}$, $\mathbf{G}(\mathbf{w})_i = |w_i|^2 w_i$.

Tehdään seuraavaksi aikadiskretointi askelpituudella $\delta > 0$. Käytetään differentiaaliyhtälölle $\mathbf{w}'(t) = \mathcal{M}(\mathbf{w}(t))$, missä $\mathcal{M} : \mathbb{C}^n \rightarrow \mathbb{C}^n$ on annettu ja \mathbf{w} tuntematon, menetelmää

$$\mathbf{w}^{k+1} = \mathbf{w}^k + \delta\mathcal{L}(\mathbf{w}^k) \left(\frac{1}{2}(\mathbf{w}^k + \mathbf{w}^{k+1}) \right), \quad k = 0, 1, 2, \dots, \quad (5.27)$$



Kuva 5.7: Aaltoyhtälön ratkaisun suhteellinen jäännös (tapaus $n_1 = n_2 = 32$). RL-ILU(0) on reaalin lineaarinen ILU(0). Alg 3.10.1 ja Alg 3.10.2 viivat ovat päällekkäin, samoin GMRES RL-ILU(0) ja Alg 3.10.2 RL-ILU(0).



Kuva 5.8: Aaltoyhtälön ratkaisun suhteellinen jäännös (tapaus $n_1 = n_2 = 4$).

missä $\mathcal{L}(\mathbf{w}^k)$ on \mathcal{M} :n \mathbb{R} -linearisointi pisteessä \mathbf{w}^k . Tällöin siis $\mathcal{L}(\mathbf{w}^k)$ on vakion ja \mathbb{R} -lineaarisen operaattorin summa ja

$$\mathcal{M}(\mathbf{w}) = \mathcal{L}(\mathbf{w}^k)(\mathbf{w}) + o(|\mathbf{w} - \mathbf{w}^k|).$$

Differentiaaliyhtälön (5.26) tapauksessa on $\mathcal{M}(\mathbf{w}) = -i\Delta_{h_1, h_2}\mathbf{w} - i\mathbf{G}(\mathbf{w})$ ja, koska

$$|v + u|^2(v + u) = (v + u)^2(\bar{v} + \bar{u}) = |v|^2v + 2|v|^2u + v^2\bar{u} + o(|u|),$$

sen \mathbb{R} -linearisointi pisteessä \mathbf{w}^k on

$$\mathcal{L}(\mathbf{w}^k)(\mathbf{w}) = -i\left(\Delta_{h_1, h_2}\mathbf{w}^k + \mathbf{G}(\mathbf{w}^k) + (\Delta_{h_1, h_2} + \mathbf{B}(\mathbf{w}^k))(\mathbf{w} - \mathbf{w}^k) + \mathbf{C}(\mathbf{w}^k)(\overline{\mathbf{w} - \mathbf{w}^k})\right),$$

missä $\mathbf{B}(\mathbf{w})_{i,i} = 2|w_i|^2$, $\mathbf{C}(\mathbf{w})_{i,i} = w_i^2$ ja $\mathbf{B}(\mathbf{w}) \in \mathbb{C}^{n_1 n_2 \times n_1 n_2}$:n ja $\mathbf{C}(\mathbf{w}) \in \mathbb{C}^{n_1 n_2 \times n_1 n_2}$:n muut alkioit nolli. Sijoittamalla tämä menetelmään (5.27) saadaan

$$\begin{aligned} & \left(i\mathbf{I} - \frac{\delta}{2}(\Delta_{h_1, h_2} + \mathbf{B}(\mathbf{w}^k))\right)(\mathbf{w}^{k+1} - \mathbf{w}^k) - \frac{\delta}{2}\mathbf{C}(\mathbf{w}^k)(\overline{\mathbf{w}^{k+1} - \mathbf{w}^k}) \\ & = \delta\Delta_{h_1, h_2}\mathbf{w}^k + \delta\mathbf{G}(\mathbf{w}^k). \end{aligned}$$

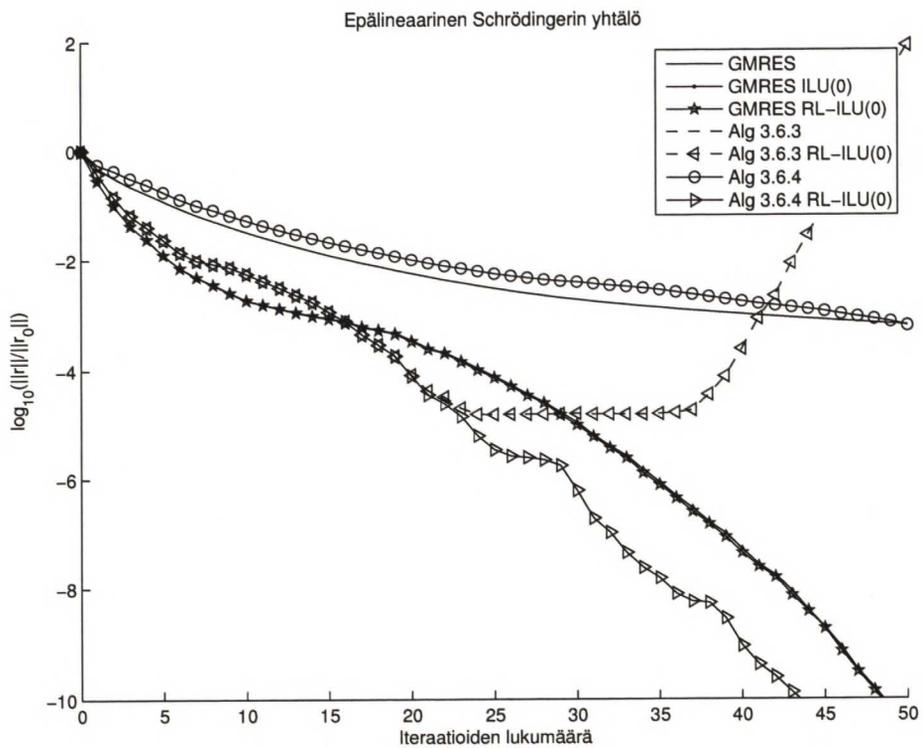
Reaalilineaarisiksi yhtälöksi saadaan näin $\mathbf{M}\mathbf{z} + \mathbf{M}_\# \bar{\mathbf{z}} = \mathbf{b}$, missä

$$\begin{aligned} \mathbf{M} &= i\mathbf{I} - \frac{\delta}{2}(\Delta_{h_1, h_2} + \mathbf{B}(\mathbf{w}^k)), & \mathbf{M}_\# &= -\frac{\delta}{2}\mathbf{C}(\mathbf{w}^k), \\ \mathbf{b} &= \delta\Delta_{h_1, h_2}\mathbf{w}^k + \delta\mathbf{G}(\mathbf{w}^k), \end{aligned}$$

ja merkitään $\mathbf{z} = \mathbf{w}^{k+1} - \mathbf{w}^k$.

Kun alkuehto $\psi(\mathbf{x}, 0)$ on annettu, on se diskretoituna $(\mathbf{w}^0)_{i+(j-1)n_1} = \psi(ih_1, jh_2, 0)$. Numeerisessa kokeessa suoritettiin vastaavat laskut kuin edellisen kohdan aaltoyhtälön tapauksessa. Edellä olevasta reaalilineaarisesta yhtälöstä ratkaistiin \mathbf{w}^1 annetun alkuehtovektorin \mathbf{w}^0 perusteella. Parametrit olivat $n_1 = n_2 = 32$ ja $\delta = 0.1$ ja alkuehdoksi valittiin edellisen kohdan 5.3 alkuehtovektoreita käyttäen $\mathbf{w}^0 = \mathbf{f}_{h_1, h_2} + i\mathbf{g}_{h_1, h_2}$ (huom. i on tässä imaginaariyksikkö eikä indeksiluku). Matriisien GMRESia varten muunnettiin reaalilineaarinen yhtälö käyttäen tapaa (1.3).

Kuvassa 5.9 näkyy ratkaisun suhteellinen jäännös erilaisin menetelmin. Edelleen \mathbb{R} -ILU(0) tuo suppenemiseen parannusta, mutta nyt ILU(0) suppenee yhtä hyvin matriisien GMRESin kanssa. Lisäksi havaitaan, että algoritmia 3.10.1 vastaava menetelmä käyttäytyy edelleen huonosti.



Kuva 5.9: Epälineaarisen Schrödingerin yhtälön ratkaisun suhteellinen jäännös. Alg 3.10.1 ja Alg 3.10.2 viivat ovat päällekkäin, samoin GMRES ILU(0) ja RL-ILU(0) ovat päällekkäin.

Yhteenveto

Tässä työssä on jatkettu viitteen [1] aloittamaa reaalilineaaristen menetelmien tutkimusta. QR-hajotelman laskennassa tarvittava Householderin reaalilineaarinen muunnos tarkentui kohdassa 1.2 numeerisesti stabiilimpaan muotoon. Kohdissa 2.3 ja 2.3.1 ehdotettiin uutta tapaa tukea reaalilineaarista LU-hajotelmaa. Viitteen [1] GMRES-menetelmä esiteltiin kohdassa 3.8 ja Householderin muunnoksien tilalle vaihtoehdoksi tarjottiin uudet lauseen 3.8.1 esittämät reaalilineaariset Givens-rotaatiot. Luvun 5 numeeriset kokeet vahvistivat reaalilineaarisen GMRES-menetelmän nopeusedun reaalisten matriisien menetelmään verrattuna. Uutena kokeellisena havaintona kohdan 3.9 \mathbb{C} -linearisoitu yhtälö (3.54) on luvun 5 numeeristen kokeiden perusteella mielenkiintoinen mahdollisuus sähköisen impedanssitomografian tyyppisen $\bar{\delta}$ -yhtälön ratkaisussa.

Toisaalta viitteessä [1] mainittu mahdollisuus reaalilineaaristen ILU-pohjustimien käytöstä reaalisten $2n \times 2n$ -matriisien pohjustimina vaikuttaa hedelmättömältä. Reaalilineaarinen LU-hajotelma ja ILU-hajotelmat ovat tällöin erikoistapauksia yleisemmistä matriisien lohko hajotelmista. Tällaiset hajotelmat ovat yleisesti, riippuen sovelluskohteesta, toimiviksi tunnettuja. Tästä osoituksena myös luvun 5 kokeet vahvistavat \mathbb{R} -ILU(0) hajotelman pienentävän iteraatiokierroksia.

Kirjallisuutta

- [1] T. Eirola, M. Huhtanen, J. von Pfafer. *Solution methods for \mathbb{R} -linear problems in \mathbb{C}^n* . SIAM J. Matrix Anal. Appl. 25 (2004), pp. 804-828.
- [2] Marko Huhtanen, Olavi Nevanlinna. *Real linear matrix analysis*. Banach Center Publ. 75 (2007), pp. 171-189.
- [3] J.L. Mueller, S. Siltanen. *Direct reconstruction of conductivities from boundary measurements*. SIAM J. Sci. Comput. 24 (2003), pp. 1232-1266.
- [4] J.L. Mueller, S. Siltanen, D. Isaacson. *An implementation of the reconstruction algorithm of A. Nachman for the 2D inverse conductivity problem*. Inverse Problems 16 (2000), pp. 681-699.
- [5] K. Knudsen, J.L. Mueller, S. Siltanen. *Numerical solution method for the d -bar equation in the plane*. Journal of Computational Physics 198 (2004), pp. 500-517.
- [6] G. Vainikko. *Fast solvers of the Lippmann-Schwinger equation*. Direct and Inverse Problems of Mathematical Physics, Int. Soc. Anal. Appl. Comput., 5, Kluwer Acad. Publ., Dordrecht, 2000, pp. 423-440.
- [7] Kari Astala, Lassi Päiväranta. *Calderón's inverse conductivity problem in the plane*. Annals of Mathematics 163 (2006), pp. 265-299.
- [8] Yousef Saad. *Iterative Methods for Sparse Linear Systems*. 2nd Edition, Society for Industrial and Applied Mathematics, Philadelphia, USA, 2003. ISBN 0-89871-534-2.
- [9] Timo Eirola. *Krylov integrators for Hamiltonian systems*. Workshop on Exponential Integrators, 20-23.10.2004, Innsbruck, Itävalta. <http://techmath.uibk.ac.at/numbau/alex/events/files04/slides/timo.pdf>.
- [10] William H. Press, Saul A. Teukolsky, William T. Vetterling, Brian P. Flannery. *Numerical Recipes in C: the art of scientific computing*. 2nd Edition, Cambridge University Press, New York, USA, 1992. ISBN 0-521-43108-5.
- [11] Roger A. Horn, Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, 1987. ISBN 0-521-30586-1.
- [12] Carl B. Boyer. *Tieteiden kuningatar*. Toinen painos, 2 osaa, Art House, 1995. ISBN 951-884-150-0. Englanninkielinen alkuperäisteos Carl B. Boyer, *A history of Mathematics*, 2nd Edition, Revised by Uta C. Merzbach, suomentanut Kimmo Pietiläinen.

- [13] Carl D. Meyer. *History of Gaussian Elimination*. The 1982 Mathematical Calendar, Rome Press. http://meyer.math.ncsu.edu/Meyer/PS_Files/GaussianEliminationHistory.pdf.
- [14] Henk A. van der Vorst. *Krylov subspace iteration*. Computing in Science and Engineering, vol. 2, no. 1 (2000), pp. 32-37.
- [15] James H. Wilkinson. *Some Comments from a Numerical Analyst*. Journal of the ACM, vol. 18, no. 2 (1971), pp. 137-147.
- [16] N. N. Lebedev. *Special Functions and Their Applications*. Dover, New York, 1972. Revised English Edition Translated and Edited by Richard A. Silverman.

Liite A

Matriisit

Matriisi \mathbf{A} on $m \times n$ -kokoinen kompleksilukualkioista a_{ij} koostuva taulukko. Riviindeksi i saa arvot $1, 2, \dots, m$ ja sarakeindeksi j arvot $1, 2, \dots, n$. Kaikkien $m \times n$ -matriisien joukkoa merkitään $\mathbb{C}^{m \times n}$. Reaalisten matriisien joukkoa merkitään $\mathbb{R}^{m \times n}$.

Indeksien $i = j$ paikkoja sanotaan (pää)diagonaaliksi (l. lävistäjäksi). Matriisin transpoosi \mathbf{A}^T saadaan kirjoittamalla rivit sarakkeiksi, tarkemmin $(\mathbf{A}^T)_{ij} = a_{ji}$. Matriisin hermitoinnissa otetaan lisäksi kompleksikonjugaatti $\mathbf{A}^* = \overline{\mathbf{A}}^T$.

Matriisi $\mathbf{A} \in \mathbb{C}^{m \times n}$ on

- (1) diagonaalimatriisi (l. lävistäjämatriisi), jos se on neliömatriisi ja $a_{ij} = 0$ kaikilla $i \neq j$,
- (2) yläkolmiomatriisi, jos $a_{ij} = 0$ kaikilla $i > j$,
- (3) alakolmiomatriisi, jos $a_{ij} = 0$ kaikilla $i < j$,
- (4) Hessenbergin matriisi, jos $a_{ij} = 0$ kaikilla $i > j + 1$.

Matriisia $\mathbf{A} \in \mathbb{C}^{n \times n}$ sanotaan neliömatriisiksi. Se on

- (1) yksikkömatriisi, jos $a_{ii} = 1$ kaikilla i ja $a_{ij} = 0$ kaikilla $i \neq j$. Sitä merkitään \mathbf{I} ja joskus tarkemmin \mathbf{I}_n .
- (2) kääntyvä, jos on olemassa $\mathbf{B} \in \mathbb{C}^{n \times n}$ siten, että $\mathbf{AB} = \mathbf{BA} = \mathbf{I}$. Tällöin merkitään $\mathbf{A}^{-1} = \mathbf{B}$.
- (3) (vino)symmetrinen, jos $\mathbf{A}^T = (-)\mathbf{A}$.
- (4) hermiittinen, jos $\mathbf{A}^* = \mathbf{A}$.
- (5) normaali, jos $\mathbf{A}^* \mathbf{A} = \mathbf{AA}^*$.
- (6) unitaarinen, jos $\mathbf{A}^* \mathbf{A} = \mathbf{I}$.
- (7) ortogonaalinen, jos \mathbf{A} on reaalinen ja $\mathbf{A}^T \mathbf{A} = \mathbf{I}$.

Avaruuden \mathbb{C}^n vektorit esitetään yleensä $n \times 1$ -(pysty)vektoreina. Luonnollinen kantavektori \mathbf{e}_i on $n \times 1$ -pystyvektori, jonka i :s alkio on 1 ja muut nollia. Vektoreiden sisätuloa merkitään $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{y}^* \mathbf{x}$ ja vektorin normia $\|\mathbf{x}\| = (\mathbf{x}^* \mathbf{x})^{1/2}$. Vektorin $\mathbf{x} = [x_1 \ \dots \ x_n]^T$ p -normi on $\|\mathbf{x}\|_p = (\sum_{k=1}^n |x_k|^p)^{1/p}$, kun $1 \leq p < \infty$, ja ∞ -normi on $\|\mathbf{x}\|_\infty = \max_k |x_k|$.

Hermiittinen matriisi \mathbf{A} on positiividefiniitti, jos $\langle \mathbf{Ax}, \mathbf{x} \rangle > 0$ kaikilla $\mathbf{x} \neq 0$. Luku $\lambda \in \mathbb{C}$ on neliömatriisin \mathbf{A} ominaisarvo, jos on olemassa (ominaisvektori) $\mathbf{x} \neq 0$ siten,

että $A\mathbf{x} = \lambda\mathbf{x}$. Kaikkien ominaisarvojen joukkoa kutsutaan A :n spektriiksi ja sitä merkitään $\sigma(A)$.

Matriisille $A \in \mathbb{C}^{m \times n}$ voidaan käyttää lohkomerkintää

$$A = \begin{bmatrix} B & C \\ D & E \end{bmatrix},$$

missä $B \in \mathbb{C}^{p \times q}$, $C \in \mathbb{C}^{p \times s}$, $D \in \mathbb{C}^{r \times q}$ ja $E \in \mathbb{C}^{r \times s}$, $m = p + r$, $n = q + s$. Tällöin alkiot ovat

$$a_{ij} = \begin{cases} b_{ij}, & \text{kun } 1 \leq i \leq p \text{ ja } 1 \leq j \leq q, \\ c_{i,j-q}, & \text{kun } 1 \leq i \leq p \text{ ja } q < j \leq n, \\ d_{i-p,j}, & \text{kun } p < i \leq m \text{ ja } 1 \leq j \leq q, \\ e_{i-p,j-q}, & \text{kun } p < i \leq m \text{ ja } q < j \leq n. \end{cases}$$

Lohkomerkinnän erikoistapauksena on $Q = [q_1 \ q_2 \ \cdots \ q_k]$, missä $q_1, q_2, \dots, q_k \in \mathbb{C}^n$ ovat pystyvektoreita. Tällöin Q on $n \times k$ -matriisi, jonka sarakkeet koostuvat vektoreista q_j . Matriisi $[Q \ q_{k+1}]$ saadaan lisäämällä matriisiin Q viimeiseksi sarakkeeksi uusi pystyvektori q_{k+1} .

Matriisin $A = (a_{ij})$ saraketta j vastaavaa vektoria merkitään \mathbf{a}_{*j} ja riviä i vastaavaa \mathbf{a}_{i*} . Matriisi $A_{r:s,t:u}$ koostuu A :n alkiosta a_{ij} , joille $r \leq i \leq s$ ja $t \leq j \leq u$. Vektori $\mathbf{v}_{r:s}$ koostuu vektorin \mathbf{v} komponenteista i , joille $r \leq i \leq s$.

Matriisien $A \in \mathbb{C}^{m \times n}$ ja $B \in \mathbb{C}^{p \times q}$ Kroneckerin tulo on $(mp) \times (nq)$ -matriisi

$$A \otimes B = \begin{bmatrix} a_{11}B & a_{12}B & \cdots & a_{1n}B \\ a_{21}B & a_{22}B & \cdots & a_{2n}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}B & a_{m2}B & \cdots & a_{mn}B \end{bmatrix}.$$

Jos matriisien A_1 ja A_2 tulo voidaan laskea ja samoin matriisien B_1 ja B_2 tulo, niin $(A_1 \otimes B_1)(A_2 \otimes B_2) = A_1 A_2 \otimes B_1 B_2$. Myös $(A \otimes B)^* = A^* \otimes B^*$.

Liite B

MATLAB-ohjelmalistaukset

B.1 Kohta 2.3.1

```
function [L,La,U,Ua,V,P] = rl_luvp(M,Ma)

% Laskee LU-hajotelman kääntyvälle reaalityyppiselle operaattorille
% käyttäen Gaussin eliminointia osittaistuennalla.
%
% Palauttaa:
% L,La - Alakolmiomatriiseja. Matriisin L lävistäjä koostuu
%        ykkösistä ja matriisin La lävistäjä nolista.
% U,Ua - Yläkolmiomatriiseja.
% V     - Yläkolmiomatriisi, jonka lävistäjä koostuu ykkösistä ja
%        jokainen rivi sisältää korkeintaan yhden muun nolista
%        poikkeavan alkion.
% P     - Permutaatiomatriisi.
%
% Laskettu hajotelma toteuttaa: op(L,La)*op(U,Ua)=V*P*op(M,Ma), missä
% esim. op(M,Ma) viittaa reaalityyppiseen operaattoriin, jonka
% lineaarinen osa on M ja anti-lineaarinen osa on Ma.
%

n=size(M,1); L=eye(n); La=zeros(n); U=M; Ua=Ma; V=eye(n); rperm=1:n;

for k=2:n

    % Etsitään nykyisestä sarakkeesta suurimman normin operaattori.

    ls=real(U(k-1:n,k-1)).^2+imag(U(k-1:n,k-1)).^2;
    ls=ls-real(Ua(k-1:n,k-1)).^2-imag(Ua(k-1:n,k-1)).^2;

    [sk,si]=max(abs(ls));
    si=si+k-2; sj=0;

    % Etsitään nykyisestä sarakkeesta pari skalaarioperaattoreita,
```

% joista voidaan muodostaa kääntyvä operaattori. Yritetään myös
 % valita tämä pari siten, että niistä muodostetulla operaattorilla
 % on iso normi parantaen yllä etsittyä.

```

for i=k-1:n-1
    for j=i+1:n
        qt=U(i,k-1)*conj(U(j,k-1))-Ua(i,k-1)*conj(Ua(j,k-1));
        qt2=conj(qt);
        qst=abs(qt)^2;
        skt=ls(i-k+2)+qst*ls(j-k+2);

        if (skt > 0)
            skt=skt+2*qst;
        else
            skt=2*qst-skt;
            qt=-qt;
        end

        if (skt > sk)
            sk=skt; si=i; sj=j; q=qt;
        end

        skt=ls(j-k+2)+qst*ls(i-k+2);

        if (skt > 0)
            skt=skt+2*qst;
        else
            skt=2*qst-skt;
            qt2=-qt2;
        end

        if (skt > sk)
            sk=skt; si=j; sj=i; q=qt2;
        end
    end
end

% Permutoidaan rivejä tarvittaessa. V:lle permutoidaan sarakkeita.

if (si ~= k-1)
    r=U(k-1,1:n); U(k-1,1:n)=U(si,1:n); U(si,1:n)=r;
    r=Ua(k-1,1:n); Ua(k-1,1:n)=Ua(si,1:n); Ua(si,1:n)=r;

    if (k > 2)
        r=L(k-1,1:k-2); L(k-1,1:k-2)=L(si,1:k-2); L(si,1:k-2)=r;
        r=La(k-1,1:k-2); La(k-1,1:k-2)=La(si,1:k-2); La(si,1:k-2)=r;
        r=V(1:k-2,k-1); V(1:k-2,k-1)=V(1:k-2,si); V(1:k-2,si)=r;
    end

    rs=rperm(k-1); rperm(k-1)=rperm(si); rperm(si)=rs;

```

```

    if (sj == k-1)
        sj = si;
    end
end

% Lisätään alapuolinen rivi (kerrottuna sopivalla kompleksiluvulla)
% nykyiseen riviin, jolloin saadaan kääntyvä tukialkio.

if (sj ~= 0)
    V(k-1,sj)=q;
    U(k-1,k-1:n)=U(k-1,k-1:n)+q*U(sj,k-1:n);
    Ua(k-1,k-1:n)=Ua(k-1,k-1:n)+q*Ua(sj,k-1:n);
    if (k > 2)
        L(k-1,1:k-2)=L(k-1,1:k-2)+q*L(sj,1:k-2);
        La(k-1,1:k-2)=La(k-1,1:k-2)+q*La(sj,1:k-2);
    end
end

a=U(k-1,k-1); b=Ua(k-1,k-1);
w=[U(k:n,k-1),Ua(k:n,k-1)]/[a,b;b',a'];
L(k:n,k-1)=w(:,1); La(k:n,k-1)=w(:,2);
U(k:n,k:n)=U(k:n,k:n)-w*[U(k-1,k:n);conj(Ua(k-1,k:n))];
Ua(k:n,k:n)=Ua(k:n,k:n)-w*[Ua(k-1,k:n);conj(U(k-1,k:n))];
U(k:n,k-1)=zeros(n-k+1,1); Ua(k:n,k-1)=zeros(n-k+1,1);
end

% Luodaan permutaatiomatriisi permutaatioindeksivektorista.

P=zeros(n);

for k=1:n
    P(k,rperm(k))=1;
end

```

B.2 Kohta 4.2.1

Alla on MATLAB-esimerkkikoodi kohtaan 4.2.1. Tässä NZ on nollakuviomatriisi. Alkion arvo nolla tarkoittaa, että sitä vastaava indeksipari kuuluu nollakuvioon.

```

function [L, U] = ilu_v2(A, NZ)

    n = size(A, 1); L = eye(n); U = A;

    for i = 1:n
        for k = 1:i-1
            if NZ(i, k) == 0
                L(i, k) = 0;
            end
        end
    end

```

```

        else
            L(i, k) = U(i, k) / U(k, k);
            U(i, k+1:n) = U(i, k+1:n) - L(i, k)*U(k, k+1:n);
        end
    end
end

U(i, 1:i-1) = zeros(1, i-1);
for j = i+1:n
    if NZ(i, j) == 0
        U(i, j) = 0;
    end
end
end
end
end

```

B.3 Kohta 4.2.2

Alla kohdan 4.2.2 funktio. Argumentti `droptol` on sama kuin em. kohdan ϵ .

```

function [L, U] = ilut(A, p, droptol)

    n = size(A, 1); L = eye(n); U = A;

    for i = 1:n
        rownorm = norm(A(i, :));

        for k = 1:i-1
            c = U(i, k) / U(k, k);
            if abs(c) < droptol*rownorm
                L(i, k) = 0;
            else
                L(i, k) = c;
                U(i, k+1:n) = U(i, k+1:n) - c*U(k, k+1:n);
            end
        end
    end

    if i > 1
        U(i, 1:i-1) = zeros(1, i-1);
    end

    for j = i+1:n
        if abs(U(i, j)) < droptol*rownorm
            U(i, j) = 0;
        end
    end
end

if i > p+1
    [s, si] = sort(abs(L(i, 1:i-1)));
    L(i, si(1:i-1-p)) = 0;
end

```

```

end

if i < n-p
    [s, si] = sort(abs(U(i, i+1:n)));
    U(i, i + si(1:n-i-p)) = 0;
end
end

```

B.4 Kohta 4.3.1

```

%
% Lasketaan reaalityyppisen operaattorin  $z \rightarrow Mz + Ma \cdot \text{conj}(z)$ 
% ILU-hajotelma käyttäen nollakuviota NZ.
%
function [L, La, U, Ua] = rl_ilu(M, Ma, NZ)

n = size(M, 1);
L = speye(n); La = spalloc(n, n, 3*n);
U = spalloc(n, n, 3*n); Ua = spalloc(n, n, 3*n);

for i = 1:n
    w = full(M(i,:)); wa = full(Ma(i,:));

    for k = 1:i-1
        if NZ(i, k) == 0
            w(k) = 0; wa(k) = 0;
        else
            a = U(k, k); b = Ua(k, k);
            c = [w(k), wa(k)] / [a, b; b', a'];
            w(k) = c(1); wa(k) = c(2);
            w(k+1:n) = w(k+1:n) - ...
                c*[U(k, k+1:n); conj(Ua(k, k+1:n))];
            wa(k+1:n) = wa(k+1:n) - ...
                c*[Ua(k, k+1:n); conj(U(k, k+1:n))];
        end
    end

    for k = i+1:n
        if NZ(i, k) == 0
            w(k) = 0; wa(k) = 0;
        end
    end

    w = sparse(w); wa = sparse(wa);
    if i > 1
        L(i, 1:i-1) = w(1:i-1); La(i, 1:i-1) = wa(1:i-1);
    end
end

```

```

    U(i, i:n) = w(i:n); Ua(i, i:n) = wa(i:n);
end

```

B.5 Algoritmi 4.3.2

Alla algoritmia 4.3.2 vastaava funktio. Tämä hyödyntää MATLABin harvojen matriisien käsittelyfunktioita.

```

function [L, La, U, Ua] = rl_ilut(M, Ma, p, droptol)

    n = size(M, 1);
    L = speye(n); La = spalloc(n, n, 3*n);
    U = spalloc(n, n, 3*n); Ua = spalloc(n, n, 3*n);

    for i = 1:n
        w = full(M(i, :)); wa = full(Ma(i, :));
        rownorm = sum(abs(w) + abs(wa));

        for k = 1:i-1
            a = U(k, k); b = Ua(k, k);
            c = [w(k), wa(k)] / [a, b; b', a'];
            if (abs(c(1)) + abs(c(2))) < droptol*rownorm
                w(k) = 0; wa(k) = 0;
            else
                w(k) = c(1); wa(k) = c(2);
                w(k+1:n) = w(k+1:n) - ...
                    c*[U(k, k+1:n); conj(Ua(k, k+1:n))];
                wa(k+1:n) = wa(k+1:n) - ...
                    c*[Ua(k, k+1:n); conj(U(k, k+1:n))];
            end
        end

        for k = i:n
            if (abs(w(k)) + abs(wa(k))) < droptol*rownorm
                w(k) = 0; wa(k) = 0;
            end
        end

        if i > p+1
            [s, si] = sort(abs(w(1:i-1)) + abs(wa(1:i-1)));
            L(i, si(i-p:i-1)) = w(si(i-p:i-1));
            La(i, si(i-p:i-1)) = wa(si(i-p:i-1));
        else
            if (i > 1)
                L(i, 1:i-1) = w(1:i-1); La(i, 1:i-1) = wa(1:i-1);
            end
        end
    end
end

```

```
if i < n-p
    [s, si] = sort(abs(w(i+1:n)) + abs(wa(i+1:n)));
    U(i, i) = w(i); Ua(i, i) = wa(i);
    U(i, i + si(n-i-p+1:n-i)) = w(i + si(n-i-p+1:n-i));
    Ua(i, i + si(n-i-p+1:n-i)) = wa(i + si(n-i-p+1:n-i));
else
    U(i, i:n) = w(i:n); Ua(i, i:n) = wa(i:n);
end
end
```